

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1 Same data, different analysts: variation in effect sizes due to analytical
2 decisions in ecology and evolutionary biology.

3 Elliot Gould, School of Agriculture Food and Ecosystem Sciences, University of Melbourne, Australia

4 Hannah S. Fraser, School of Historical and Philosophical Studies, University of Melbourne, Australia

5 Timothy H. Parker, Department of Biology, Whitman College, USA. Author for Correspondence:
6 parkerth@whitman.edu

7 Shinichi Nakagawa, School of Biological, Earth & Environmental Sciences, University of New South
8 Wales, Australia

9 Simon C. Griffith, School of Natural Sciences, Macquarie University, Australia

10 Peter A. Vesk, School of Agriculture Food and Ecosystem Sciences, University of Melbourne, Australia

11 Fiona Fidler, School of Historical and Philosophical Studies, University of Melbourne, Australia

12 Daniel G. Hamilton, School of Public Health and Preventive Medicine, Monash University, Australia

13 Robin N Abbey-Lee, Länsstyrelsen Östergötland, Sweden

14 Jessica K. Abbott, Biology Department, Lund University, Sweden

15 Luis A. Aguirre, Department of Biology, University of Massachusetts, USA

16 Carles Alcaraz, Marine and Continental Waters, IRTA, Spain

17 Irith Aloni, Department of Life Sciences, Ben Gurion University of the Negev, Israel

18 Drew Altschul, Department of Psychology, The University of Edinburgh, UK

19 Kunal Arekar, Centre for Ecological Sciences, Indian Institute of Science, India

20 Jeff W. Atkins, Southern Research Station, USDA Forest Service, USA

21 Joe Atkinson, Center for Ecological Dynamics in a Novel Biosphere (ECONOVO), Department of
22 Biology, Aarhus University, Denmark

23 Christopher M. Baker, School of Mathematics and Statistics, University of Melbourne, Australia

24 Meghan Barrett, Biology, Indiana University Purdue University Indianapolis, USA

25 Kristian Bell, School of Life and Environmental Sciences, Deakin University, Australia

26 Suleiman Kehinde Bello, Department of Arid Land Agriculture, King Abdulaziz University, Kingdom of
27 Saudi Arabia

28 Iván Beltrán, Department of Biological Sciences, Macquarie University, Australia

29 Bernd J. Berauer, Department of Plant Ecology, University of Hohenheim, Institute of Landscape and
30 Plant Ecology, Germany

31 Michael Grant Bertram, Department of Wildlife, Fish, and Environmental Studies, Swedish University
32 of Agricultural Sciences, Sweden

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 33 Peter D. Billman, Department of Ecology and Evolutionary Biology, University of Connecticut, USA
34 Charlie K. Blake, STEM Center, Southern Illinois University Edwardsville, USA
35 Shannon Blake, University of Guelph, Canada
36 Louis Bliard, Department of Evolutionary Biology and Environmental Studies, University of Zurich,
37 Switzerland
38 Andrea Bonisoli-Alquati, Department of Biological Sciences, California State Polytechnic University,
39 Pomona, USA
40 Timothée Bonnet, Centre d'Études Biologiques de Chizé, UMR 7372 Université de la Rochelle - Centre
41 National de la Recherche Scientifique, France
42 Camille Nina Marion Bordes, Faculty of Life Sciences, Bar Ilan University, Israel
43 Aneesh P. H. Bose, Department of Wildlife, Fish, and Environmental Studies, Swedish University of
44 Agricultural Sciences, Sweden
45 Thomas Botterill-James, School of Natural Sciences, University of Tasmania, Australia
46 Melissa Anna Boyd, Whitebark Institute, USA
47 Sarah A. Boyle, Department of Biology, Rhodes College, USA
48 Tom Bradfer-Lawrence, Centre for Conservation Science, RSPB, UK
49 Jennifer Bradham, Environmental Studies, Wofford College, USA
50 Jack A. Brand, Department of Wildlife, Fish and Environmental Studies, Swedish University of
51 Agricultural Sciences, Sweden
52 Martin I. Brengdahl, IFM Biology, Linköping University, Sweden
53 Martin Bulla, Faculty of Environmental Sciences, Czech University of Life Sciences Prague, Czech
54 Republic
55 Luc Bussière, Biological and Environmental Sciences & Gothenburg Global Biodiversity Centre,
56 University of Gothenburg, Sweden
57 Ettore Camerlenghi, School of Biological Sciences, Monash University, Australia
58 Sara E. Campbell, Ecology and Evolutionary Biology, University of Tennessee Knoxville, USA
59 Leonardo L. F. Campos, Departamento de Ecologia e Zoologia, Universidade Federal de Santa
60 Catarina, Brazil
61 Anthony Caravaggi, School of Biological and Forensic Sciences, University of South Wales, UK
62 Pedro Cardoso, Centre for Ecology, Evolution and Environmental Changes (CE3c) & CHANGE - Global
63 Change and Sustainability Institute, Faculdade de Ciências, Universidade de Lisboa, Portugal
64 Charles J.W. Carroll, Forest and Rangeland Stewardship, Colorado State University, USA
65 Therese A. Catanach, Department of Ornithology, Academy of Natural Sciences of Drexel University,
66 USA

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 67 Xuan Chen, Biology, Salisbury University, USA
- 68 Heung Ying Janet Chik, Groningen Institute for Evolutionary Life Sciences, University of Groningen,
69 Netherlands
- 70 Emily Sarah Choy, Department of Biology, McMaster University, Canada
- 71 Alec Philip Christie, Department of Zoology, University of Cambridge, UK
- 72 Angela Chuang, Entomology and Nematology, University of Florida, USA
- 73 Amanda J. Chunco, Environmental Studies, Elon University, USA
- 74 Bethany L. Clark, BirdLife International, UK
- 75 Andrea Contina, School of Integrative Biological and Chemical Sciences, The University of Texas Rio
76 Grande Valley, USA
- 77 Garth A. Covernton, Department of Ecology and Evolutionary Biology, University of Toronto, Canada
- 78 Murray P. Cox, Department of Statistics, University of Auckland, New Zealand
- 79 Kimberly A. Cressman, Catbird Stats, LLC, USA
- 80 Marco Crotti, School of Biodiversity, One Health & Veterinary Medicine, University of Glasgow, UK
- 81 Connor Davidson Crouch, School of Forestry, Northern Arizona University, USA
- 82 Pietro B. D'Amelio, Department of Behavioural Neurobiology, Max Planck Institute for Biological
83 Intelligence, Germany
- 84 Alexandra Allison de Sousa, School of Sciences: Center for Health and Cognition, Bath Spa University,
85 UK
- 86 Timm Fabian Döbert, Department of Biological Sciences, University of Alberta, Canada
- 87 Ralph Dobler, Applied Zoology, TU Dresden, Germany
- 88 Adam J. Dobson, School of Molecular Biosciences, College of Medical Veterinary & Life Sciences,
89 University of Glasgow, UK
- 90 Tim S. Doherty, School of Life and Environmental Sciences, The University of Sydney, Australia
- 91 Szymon Marian Drobniak, Institute of Environmental Sciences, Jagiellonian University, Poland
- 92 Alexandra Grace Duffy, Biology Department, Brigham Young University, USA
- 93 Alison B. Duncan, Institute of Evolutionary Sciences Montpellier, University of Montpellier, CNRS,
94 IRD., France
- 95 Robert P. Dunn, Baruch Marine Field Laboratory, University of South Carolina, USA
- 96 Jamie Dunning, Department of Life Sciences, Imperial College London, UK
- 97 Trishna Dutta, European Forest Institute, Germany
- 98 Luke Eberhart-Hertel, Department of Ornithology, Max Planck Institute for Biological Intelligence,
99 Germany

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

100 Jared Alan Elmore, Forestry and Environmental Conservation, National Bobwhite and Grassland
101 Initiative, Clemson University, USA

102 Mahmoud Medhat Elsherif, Department of Psychology and Vision Science, University of Birmingham,
103 Baily Thomas Grant, UK

104 Holly M. English, School of Biology and Environmental Science, University College Dublin, Ireland

105 David C. Ensminger, Department of Biological Sciences, San José State University, USA

106 Ulrich Rainer Ernst, Apicultural State Institute, University of Hohenheim, Germany

107 Stephen M. Ferguson, Department of Biology, St. Norbert College, USA

108 Esteban Fernandez-Juricic, Department of Biological Sciences, Purdue University, USA

109 Thalita Ferreira-Arruda, Biodiversity, Macroecology & Biogeography, Faculty of Forest Sciences and
110 Forest Ecology, University of Göttingen, Germany

111 John Fieberg, Department of Fisheries, Wildlife, and Conservation Biology, University of Minnesota,
112 USA

113 Elizabeth A. Finch, CABI, UK

114 Evan A. Fiorenza, Department of Ecology and Evolutionary Biology, School of Biological Sciences,
115 University of California, Irvine, USA

116 David N. Fisher, School of Biological Sciences, University of Aberdeen, UK

117 Amélie Fontaine, Department of Natural Resource Sciences, McGill University, Canada

118 Wolfgang Forstmeier, Department of Ornithology, Max Planck Institute for Biological Intelligence,
119 Germany

120 Yoan Fourcade, Institute of Ecology and Environmental Sciences (iEES), Univ. Paris-Est Creteil, France

121 Graham S. Frank, Department of Forest Ecosystems and Society, Oregon State University, USA

122 Cathryn A. Freund, Wake Forest University, USA

123 Eduardo Fuentes-Lillo, Laboratorio de Invasiones Biológicas (LIB), Instituto de Ecología y
124 Biodiversidad, Chile

125 Sara L. Gandy, Institute for Biodiversity, Animal Health and Comparative Medicine, University of
126 Glasgow, UK

127 Dustin G. Gannon, Department of Forest Ecosystems and Society, College of Forestry, Oregon State
128 University, USA

129 Ana I. García-Cervigón, Biodiversity and Conservation Area, Rey Juan Carlos University, Spain

130 Alexis C. Garretson, Graduate School of Biomedical Sciences, Tufts University, USA

131 Xuezheng Ge, Department of Integrative Biology, University of Guelph, Canada

132 William L. Geary, School of Life and Environmental Sciences (Burwood Campus), Deakin University,
133 Australia

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 134 Charly Géron, CNRS, University of Rennes, France
- 135 Marc Gilles, Department of Behavioural Ecology, Bielefeld University, Germany
- 136 Antje Girndt, Fakultät für Biologie, Arbeitsgruppe Evolutionsbiologie, Universität Bielefeld, Germany
- 137 Daniel Gliksman, Chair of Meteorology, Institute for Hydrology and Meteorology, Faculty of
138 Environmental Sciences, Technische Universität Dresden, Germany
- 139 Harrison B. Goldspiel, Department of Wildlife, Fisheries, and Conservation Biology, University of
140 Maine, USA
- 141 Dylan G. E. Gomes, Department of Biological Sciences, Boise State University, USA
- 142 Megan Kate Good, School of Agriculture, Food and Ecosystem Sciences, The University of Melbourne,
143 Australia
- 144 Sarah C. Goslee, Pastures Systems and Watershed Management Research Unit, USDA Agricultural
145 Research Service, USA
- 146 J. Stephen Gosnell, Department of Natural Sciences, Baruch College, City University of New York, USA
- 147 Eliza M. Grames, Department of Biological Sciences, Binghamton University, USA
- 148 Paolo Gratton, Dipartimento di Biologia, Università di Roma "Tor Vergata", Italy
- 149 Nicholas M. Grebe, Department of Anthropology, University of Michigan, USA
- 150 Skye M. Greenler, College of Forestry, Oregon State University, USA
- 151 Maaïke Griffioen, University of Antwerp, Belgium
- 152 Daniel M. Griffith, Earth & Environmental Sciences, Wesleyan University, USA
- 153 Frances J. Griffith, Yale School of Medicine, Department of Psychiatry, Yale University, USA
- 154 Jake J. Grossman, Biology Department and Environmental Studies Department, St. Olaf College, USA
- 155 Ali Güncan, Department of Plant Protection, Faculty of Agriculture, Ordu University, Turkey
- 156 Stef Haesen, Department of Earth and Environmental Sciences, KU Leuven, Belgium
- 157 James G. Hagan, Department of Marine Sciences, University of Gothenburg, Sweden
- 158 Heather A. Hager, Department of Biology, Wilfrid Laurier University, Canada
- 159 Jonathan Philo Harris, Natural Resource Ecology and Management, Iowa State University, USA
- 160 Natasha Dean Harrison, School of Biological Sciences, University of Western Australia, Australia
- 161 Sarah Syedia Hasnain, Department of Biological Sciences, Middle East Technical University, Turkey
- 162 Justin Chase Havird, Dept. of Integrative Biology, University of Texas at Austin, USA
- 163 Andrew J. Heaton, Grand Bay National Estuarine Research Reserve, USA
- 164 María Laura Herrera-Chaustre, Universidad de los Andes, Colombia
- 165 Tanner J. Howard

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 166 Bin-Yan Hsu, Department of Biology, University of Turku, Finland
- 167 Fabiola Iannarilli, Dept of Fisheries, Wildlife and Conservation Biology, University of Minnesota, USA
- 168 Esperanza C. Iranzo, Instituto de Ciencia Animal. Facultad de Ciencias Veterinarias, Universidad
169 Austral de Chile, Chile
- 170 Erik N. K. Iverson, Department of Integrative Biology, The University of Texas at Austin, USA
- 171 Saheed Olaide Jimoh, Department of Botany, University of Wyoming, USA
- 172 Douglas H. Johnson, Department of Fisheries, Wildlife, and Conservation Biology, University of
173 Minnesota, USA
- 174 Martin Johnsson, Department of Animal Breeding and Genetics, Swedish University of Agricultural
175 Sciences, Sweden
- 176 Jesse Jorna, Department of Biology, Brigham Young University, Brigham Young University, USA
- 177 Tommaso Jucker, School of Biological Sciences, University of Bristol, UK
- 178 Martin Jung, International Institute for Applied Systems Analysis (IIASA), Austria
- 179 Ineta Kačergytė, Department of Ecology, Swedish University of Agricultural Sciences, Sweden
- 180 Oliver Kaltz, Université de Montpellier, France
- 181 Alison Ke, Department of Wildlife, Fish, and Conservation Biology, University of California, Davis, USA
- 182 Clint D. Kelly, Département des Sciences biologiques, Université du Québec à Montréal, Canada
- 183 Katharine Keogan, Institute of Evolutionary Biology, University of Edinburgh, UK
- 184 Friedrich Wolfgang Keppeler, Center for Limnology, Center for Limnology, University of Wisconsin -
185 Madison, USA
- 186 Alexander K. Killion, Center for Biodiversity and Global Change, Yale University, USA
- 187 Dongmin Kim, Department of Ecology, Evolution, and Behavior, University of Minnesota, St. Paul, USA
- 188 David P. Kochan, Institute of Environment and Department of Biological Sciences, Florida
189 International University, USA
- 190 Peter Korsten, Department of Life Sciences, Aberystwyth University, UK
- 191 Shan Kothari, Institut de recherche en biologie végétale, Université de Montréal, Canada
- 192 Jonas Kuppler, Institute of Evolutionary Ecology and Conservation Genomics, Ulm University,
193 Germany
- 194 Jillian M. Kusch, Department of Biology, Memorial University of Newfoundland, Canada
- 195 Malgorzata Lagisz, Evolution & Ecology Research Centre and School of Biological, Earth &
196 Environmental Sciences, University of New South Wales, Australia
- 197 Kristen Marianne Lalla, Department of Natural Resource Sciences, McGill University, Canada
- 198 Daniel J. Larkin, Department of Fisheries, Wildlife and Conservation Biology, University of Minnesota-
199 Twin Cities, USA

- 200 Courtney L. Larson, The Nature Conservancy, USA
- 201 Katherine S. Lauck, Department of Wildlife, Fish, and Conservation Biology, University of California,
202 Davis, USA
- 203 M. Elise Lauterbur, Ecology and Evolutionary Biology, University of Arizona, USA
- 204 Alan Law, Biological and Environmental Sciences, University of Stirling, UK
- 205 Don-Jean Léandri-Breton, Department of Natural Resource Sciences, McGill University, Canada
- 206 Jonas J. Lembrechts, Department of Biology, University of Antwerp, Belgium
- 207 Kiara L'Herpinier, Natural sciences, Macquarie University, Australia
- 208 Eva J. P. Lievens, Aquatic Ecology and Evolution Group, Limnological Institute, University of Konstanz,
209 Germany
- 210 Daniela Oliveira de Lima, Campus Cerro Largo, Universidade Federal da Fronteira Sul, Brazil
- 211 Shane Lindsay, School of Psychology and Social Work, University of Hull, UK
- 212 Martin Luquet, UMR 1224 ECOBIOP, Université de Pau et des Pays de l'Adour, France
- 213 Ross MacLeod, School of Biological & Environmental Sciences, Liverpool John Moores University, UK
- 214 Kirsty H. Macphie, Institute of Ecology and Evolution, University of Edinburgh, UK
- 215 Kit Magellan, Cambodia
- 216 Magdalena M. Mair, Statistical Ecotoxicology, Bayreuth Center of Ecology and Environmental
217 Research (BayCEER), University of Bayreuth, Germany
- 218 Lisa E. Malm, Ecology and Environmental Science, Umeå University, Sweden
- 219 Stefano Mammola, Molecular Ecology Group (MEG), Water Research Institute (IRSA), National
220 Research Council of Italy (CNR), Italy
- 221 Caitlin P. Mandeville, Department of Natural History, Norwegian University of Science and
222 Technology, Norway
- 223 Michael Manhart, Center for Advanced Biotechnology and Medicine, Rutgers University Robert
224 Wood Johnson Medical School, USA
- 225 Laura Milena Manrique-Garzon, Departamento de Ciencias Biológicas, Universidad de los Andes,
226 Colombia
- 227 Elina Mäntylä, Department of Biology, University of Turku, Finland
- 228 Philippe Marchand, Institut de recherche sur les forêts, Université du Québec en Abitibi-
229 Témiscamingue, Canada
- 230 Benjamin Michael Marshall, Biological and Environmental Sciences, University of Stirling, UK
- 231 Charles A. Martin, Université du Québec à Trois-Rivières, Canada
- 232 Dominic Andreas Martin, Institute of Plant Sciences, University of Bern, Switzerland

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 233 Jake Mitchell Martin, Department of Wildlife, Fish, and Environmental Studies, Swedish University of
234 Agricultural Sciences, Sweden
- 235 April Robin Martinig, School of Biological, Earth and Environmental Sciences, University of New South
236 Wales, Australia
- 237 Erin S. McCallum, Department of Wildlife, Fish and Environmental Studies, Swedish University of
238 Agricultural Sciences, Sweden
- 239 Mark McCauley, Whitney Laboratory for Marine Bioscience, University of Florida, USA
- 240 Sabrina M. McNew, Ecology and Evolutionary Biology, University of Arizona, USA
- 241 Scott J. Meiners, Biological Sciences, Eastern Illinois University, USA
- 242 Thomas Merklung, Centre d'Investigations Clinique Plurithématique - Institut Lorrain du Coeur et des
243 Vaisseaux, Université de Lorraine, Inserm1433 CIC-P CHRU de Nancy, France
- 244 Marcus Michelangeli, Department of Wildlife, Fish and Environmental Studies, Swedish University of
245 Agricultural Sciences, Sweden
- 246 Maria Moiron, Evolutionary biology department, Bielefeld University, Germany
- 247 Bruno Moreira, Department of Ecology and global change, Centro de Investigaciones sobre
248 Desertificación, Consejo Superior de Investigaciones Científicas (CIDE-CSIC/UV/GV), Spain
- 249 Jennifer Mortensen, Department of Biological Sciences, University of Arkansas, USA
- 250 Benjamin Mos, School of the Environment, Faculty of Science, The University of Queensland,
251 Australia
- 252 Taofeek Olatunbosun Muraina, Department of Animal Health and Production, Oyo State College of
253 Agriculture and Technology, Nigeria
- 254 Penelope Wrenn Murphy, Department of Forest & Wildlife Ecology, University of Wisconsin-Madison,
255 USA
- 256 Luca Nelli, School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow, UK
- 257 Petri Niemelä, Organismal and Evolutionary Biology Research Programme, Faculty of Biological and
258 Environmental Sciences, University of Helsinki, Finland
- 259 Josh Nightingale, South Iceland Research Centre, University of Iceland, Iceland
- 260 Gustav Nilsson, Department of Clinical Neuroscience, Karolinska Institutet, Sweden
- 261 Sergio Nolzco, School of Biological Sciences, Monash University, Australia
- 262 Sabine S. Nooten, Animal Ecology and Tropical Biology, University of Würzburg, Germany
- 263 Jessie Lanterman Novotny, Biology, Hiram College, USA
- 264 Agnes Birgitta Olin, Department of Aquatic Resources, Swedish University of Agricultural Sciences,
265 Sweden
- 266 Chris L. Organ, Department of Earth Sciences, Montana State University, USA

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 267 Kate L. Ostevik, Department of Evolution, Ecology, and Organismal Biology, University of California,
268 Riverside, USA
- 269 Facundo Xavier Palacio, Sección Ornitología, Universidad Nacional de La Plata, Argentina
- 270 Matthieu Paquet, Department of Ecology, Swedish University of Agricultural Sciences, Sweden
- 271 Darren James Parker, Bangor University, UK
- 272 David J. Pascall, MRC Biostatistics Unit, University of Cambridge, UK
- 273 Valerie J. Pasquarella, Harvard Forest, Harvard University, USA
- 274 John Harold Paterson, Biological and Environmental Sciences, University of Stirling, Scotland
- 275 Ana Payo-Payo, Departamento de Biodiversidad, Ecología y Evolución., Universidad Complutense de
276 Madrid, Spain
- 277 Karen Marie Pedersen, Biology Department, Technische Universität Darmstadt, Germany
- 278 Grégoire Perez, UMR 1309 ASTRE, CIRAD, France
- 279 Kayla I. Perry, Department of Entomology, The Ohio State University, USA
- 280 Patrice Pottier, Evolution & Ecology Research Centre, School of Biological, Earth and Environmental
281 Sciences, The University of New South Wales, Australia
- 282 Michael J. Proulx, Department of Psychology, University of Bath, UK
- 283 Raphaël Proulx, Chaire de recherche en intégrité écologique, Université du Québec à Trois-Rivières,
284 Canada
- 285 Jessica L Pruet, Mississippi Based RESTORE Act Center of Excellence, University of Southern
286 Mississippi, USA
- 287 Veronarindra Ramananjato, Department of Integrative Biology, University of California, Berkeley, USA
- 288 Finaritra Tolotra Randimbarison, Mention Zoologie et Biodiversité Animale, Université
289 d'Antananarivo, Madagascar
- 290 Onja H. Razafindratsima, Department of Integrative Biology, University of California, Berkeley, USA
- 291 Diana J. Rennison, Department of Ecology, Behavior and Evolution, University of California, San
292 Diego, USA
- 293 Federico Riva, Institute for Environmental Sciences, VU Amsterdam, The Netherlands
- 294 Sepand Riyahi, Department of Evolutionary Anthropology, University of Vienna, Austria
- 295 Michael James Roast, Konrad Lorenz Institute for Ethology, University of Veterinary Medicine, Austria
- 296 Felipe Pereira Rocha, School of Biological Sciences, The University of Hong Kong, China
- 297 Dominique G. Roche, Institut de biologie, Université de Neuchâtel, Switzerland
- 298 Cristian Román-Palacios, School of Information, University of Arizona, USA
- 299 Michael S. Rosenberg, Center for Biological Data Science, Virginia Commonwealth University, USA

- 300 Jessica Ross, University of Wisconsin, USA
- 301 Freya E. Rowland, School of the Environment, Yale University, USA
- 302 Deusedith Rugemalila, Institute of the Environment, Florida International University, USA
- 303 Avery L. Russell, Department of Biology, Missouri State University, USA
- 304 Suvi Ruuskanen, Department of Biological and Environmental Science, University of Jyväskylä,
305 Finland
- 306 Patrick Saccone, Institute for Interdisciplinary Mountain Research, OeAW (Austrian Academy of
307 Sciences), Austria
- 308 Asaf Sadeh, Department of Natural Resources, Newe Ya'ar Research Center, Agricultural Research
309 Organization (Volcani Institute), Israel
- 310 Stephen M. Salazar, Department of Animal Behaviour, Bielefeld University, Germany
- 311 Kris Sales, Office for National Statistics, UK
- 312 Pablo Salmón, Institute of Avian Research "Vogelwarte Helgoland", Germany
- 313 Alfredo Sánchez-Tójar, Department of Evolutionary Biology, Bielefeld University, Germany
- 314 Leticia Pereira Santos, Ecology Department, Universidade Federal de Goiás, Brazil
- 315 Francesca Santostefano, University of Exeter, University of Exeter, UK
- 316 Hayden T. Schilling, New South Wales Department of Primary Industries Fisheries, Australia
- 317 Marcus Schmidt, Research Data Management, Leibniz Centre for Agricultural Landscape Research
318 (ZALF), Germany
- 319 Tim Schmoll, Evolutionary Biology, Bielefeld University, Germany
- 320 Adam C. Schneider, Biology Department, University of Wisconsin-La Crosse, USA
- 321 Allie E. Schrock, Department of Evolutionary Anthropology, Duke University, USA
- 322 Julia Schroeder, Department of Life Sciences, Imperial College London, UK
- 323 Nicolas Schtickzelle, Earth and Life Institute, Ecology and Biodiversity, UCLouvain, Belgium
- 324 Nick L. Schultz, Future Regions Research Centre, Federation University Australia, Australia
- 325 Drew A. Scott, United States Department of Agriculture- Agricultural Research Service-, USA
- 326 Michael Peter Scroggie, Arthur Rylah Insitute for Environmental Research, Australia
- 327 Julie Teresa Shapiro, Epidemiology and Surveillance Support Unit, University of Lyon - French Agency
328 for Food, Environmental and Occupational Health and Safety (ANSES), France
- 329 Nitika Sharma, UCLA Anderson Center for Impact, University of California, Los Angeles, USA
- 330 Caroline L. Shearer, Department of Evolutionary Anthropology, Duke University, USA
- 331 Diego Simón, Facultad de Ciencias, Universidad de la República, Uruguay
- 332 Michael I. Sitvarin, Independent researcher, USA

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 333 Fabrício Luiz Skupien, Programa de Pós-Graduação em Ecologia, Instituto de Biologia, Centro de
334 Ciências da Saúde, Universidade Federal do Rio de Janeiro, Brazil
- 335 Heather Lea Slinn, Vive Crop Protection, Canada
- 336 Grania Polly Smith, University of Cambridge, UK
- 337 Jeremy A. Smith, British Trust for Ornithology, UK
- 338 Rahel Sollmann, Department of Wildlife, Fish, and Conservation Biology, University of California,
339 Davis, USA
- 340 Kaitlin Stack Whitney, Science, Technology & Society Department, Rochester Institute of Technology,
341 USA
- 342 Shannon Michael Still, Nomad Ecology, USA
- 343 Erica F. Stuber, Wildland Resources Department, Utah State University, USA
- 344 Guy F. Sutton, Center for Biological Control, Department of Zoology and Entomology, Rhodes
345 University, South Africa
- 346 Ben Swallow, School of Mathematics and Statistics and Centre for Research in Ecological and
347 Environmental Modelling, University of St Andrews, UK
- 348 Conor Claverie Taff, Department of Ecology and Evolutionary Biology, Cornell University, USA
- 349 Elina Takola, Department of Computational Landscape Ecology, Helmholtz Centre for Environmental
350 Research – UFZ, Germany
- 351 Andrew J. Tanentzap, Ecosystems and Global Change Group, School of the Environment, Trent
352 University, Canada
- 353 Rocío Tarjuelo, Instituto Universitario de Investigación en Gestión Forestal Sostenible (iuFOR),
354 Universidad de Valladolid, Spain
- 355 Richard J. Telford, Department of Biological Sciences, University of Bergen, Norway
- 356 Christopher J. Thawley, Department of Biological Science, University of Rhode Island, USA
- 357 Hugo Thierry, Department of Geography, McGill University, Canada
- 358 Jacqueline Thomson, Integrative Biology, University of Guelph, Canada
- 359 Svenja Tidau, School of Biological and Marine Sciences, University of Plymouth, UK
- 360 Emily M. Tompkins, Biology Department, Wake Forest University, USA
- 361 Claire Marie Tortorelli, Plant Sciences, University of California, Davis, USA
- 362 Andrew Trlica, College of Natural Resources, North Carolina State University, USA
- 363 Biz R. Turnell, Institute of Zoology, Technische Universität Dresden, Germany
- 364 Lara Urban, Helmholtz AI, Helmholtz Zentrum Muenchen, Germany
- 365 Stijn Van de Vondel, Department of Biology, University of Antwerp, Belgium

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 366 Jessica Eva Megan van der Wal, FitzPatrick Institute of African Ornithology, University of Cape Town,
367 South Africa
- 368 Jens Van Eeckhoven, Department of Cell & Developmental Biology, Division of Biosciences, University
369 College London, UK
- 370 Francis van Oordt, Natural Resource Sciences, McGill University, Canada
- 371 K. Michelle Vanderwel, Biology, University of Saskatchewan, Canada
- 372 Mark C. Vanderwel, Department of Biology, University of Regina, Canada
- 373 Karen J. Vanderwolf, Biology, University of Waterloo, Canada
- 374 Juliana Vélez, Department of Fisheries, Wildlife and Conservation Biology, University of Minnesota,
375 USA
- 376 Diana Carolina Vergara-Florez, Department of Ecology & Evolutionary Biology, University of Michigan,
377 USA
- 378 Brian C. Verrelli, Center for Biological Data Science, Virginia Commonwealth University, USA
- 379 Marcus Vinícius Vieira, Dept. Ecologia, Instituto de Biologia, Universidade Federal do Rio de Janeiro,
380 Brazil
- 381 Nora Villamil, Lothian Analytical Services, Public Health Scotland, UK
- 382 Valerio Vitali, Institute for Evolution and Biodiversity, University of Muenster, Germany
- 383 Julien Vollering, Department of Environmental Sciences, Western Norway University of Applied
384 Sciences, Norway
- 385 Jeffrey Walker, Department of Biological Sciences, University of Southern Maine, USA
- 386 Xanthe J. Walker, Center for Ecosystem Science and Society, Northern Arizona University, USA
- 387 Jonathan A. Walter, Center for Watershed Sciences, University of California, Davis, USA
- 388 Pawel Waryszak, School of Agriculture and Environmental Science, University of Southern
389 Queensland, Australia
- 390 Ryan J. Weaver, Department of Ecology, Evolution, and Organismal Biology, Iowa State University,
391 USA
- 392 Ronja E. M. Wedegärtner, Fram Project AS, Norway
- 393 Daniel L. Weller, Department of Food Science & Technology, Virginia Polytechnic Institute and State
394 University, USA
- 395 Shannon Whelan, Department of Natural Resource Sciences, McGill University, Canada
- 396 Rachel Louise White, School of Applied Sciences, University of Brighton, UK
- 397 David William Wolfson, Department of Fisheries, Wildlife and Conservation Biology, University of
398 Minnesota, USA
- 399 Andrew Wood, Department of Biology, University of Oxford, UK

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 400 Scott W. Yanco, Department of Integrative Biology, University of Colorado, Denver, USA
- 401 Jian D. L. Yen, Arthur Rylah Institute for Environmental Research, Australia
- 402 Casey Youngflesh, Ecology, Evolution, and Behavior Program, Michigan State University, USA
- 403 Giacomo Zilio, ISEM, University of Montpellier, CNRS, France
- 404 Cédric Zimmer, Laboratoire d’Ethologie Expérimentale et Comparée, LEEC, UR4443, Université
405 Sorbonne Paris Nord, USA
- 406 Gregory Mark Zimmerman, Department of Science and Environment, Lake Superior State University,
407 USA
- 408 Rachel A. Zitomer, Department of Forest Ecosystems and Society, Oregon State University, USA

409 Abstract

410 Although variation in effect sizes and predicted values among studies of similar phenomena is
411 inevitable, such variation far exceeds what might be produced by sampling error alone. One possible
412 explanation for variation among results is differences among researchers in the decisions they make
413 regarding statistical analyses. A growing array of studies has explored this analytical variability in
414 different (mostly social science) fields, and has found substantial variability among results, despite
415 analysts having the same data and research question. We implemented an analogous study in
416 ecology and evolutionary biology, fields in which there have been no empirical exploration of the
417 variation in effect sizes or model predictions generated by the analytical decisions of different
418 researchers. We used two unpublished datasets, one from evolutionary ecology (blue tit, *Cyanistes*
419 *caeruleus*, to compare sibling number and nestling growth) and one from conservation ecology
420 (*Eucalyptus*, to compare grass cover and tree seedling recruitment), and the project leaders recruited
421 174 analyst teams, comprising 246 analysts, to investigate the answers to prespecified research
422 questions. Analyses conducted by these teams yielded 141 usable effects for the blue tit dataset, and
423 85 usable effects for the *Eucalyptus* dataset. We found substantial heterogeneity among results for
424 both datasets, although the patterns of variation differed between them. For the blue tit analyses,
425 the average effect was convincingly negative, with less growth for nestlings living with more siblings,
426 but there was near continuous variation in effect size from large negative effects to effects near zero,
427 and even effects crossing the traditional threshold of statistical significance in the opposite direction.
428 In contrast, the average relationship between grass cover and *Eucalyptus* seedling number was only
429 slightly negative and not convincingly different from zero, and most effects ranged from weakly
430 negative to weakly positive, with about a third of effects crossing the traditional threshold of
431 significance in one direction or the other. However, there were also several striking outliers in
432 the *Eucalyptus* dataset, with effects far from zero. For both datasets, we found substantial variation
433 in the variable selection and random effects structures among analyses, as well as in the ratings of
434 the analytical methods by peer reviewers, but we found no strong relationship between any of these

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

435 and deviation from the meta-analytic mean. In other words, analyses with results that were far from
436 the mean were no more or less likely to have dissimilar variable sets, use random effects in their
437 models, or receive poor peer reviews than those analyses that found results that were close to the
438 mean. The existence of substantial variability among analysis outcomes raises important questions
439 about how ecologists and evolutionary biologists should interpret published results, and how they
440 should conduct analyses in the future.

441 [Key Words](#)

442 credibility revolution, heterogeneity, meta-analysis, metascience, Replicability, reproducibility

443 [Introduction](#)

444 One value of science derives from its production of replicable, and thus reliable, results. When we
445 repeat a study using the original methods we should be able to expect a similar result. However,
446 perfect replicability is not a reasonable goal. Effect sizes will vary, and even reverse in sign, by
447 chance alone [1]. Observed patterns can differ for other reasons as well. It could be that we do not
448 sufficiently understand the conditions that led to the original result so when we seek to replicate it,
449 the conditions differ due to some ‘hidden moderator’. This hidden moderator hypothesis is
450 described by meta-analysts in ecology and evolutionary biology as ‘true biological heterogeneity’ [2].
451 This idea of true heterogeneity is popular in ecology and evolutionary biology, and there are good
452 reasons to expect it in the complex systems in which we work [3]. However, despite similar
453 expectations in psychology, recent evidence in that discipline contradicts the hypothesis that
454 moderators are common obstacles to replicability, as variability in results in a large ‘many labs’
455 collaboration was mostly unrelated to commonly hypothesized moderators such as the conditions
456 under which the studies were administered [4]. Another possible explanation for variation in effect
457 sizes is that researchers often present biased samples of results, thus reducing the likelihood that
458 later studies will produce similar effect sizes [5–9]. It also may be that although researchers did

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

459 successfully replicate the conditions, the experiment, and measured variables, analytical decisions
460 differed sufficiently among studies to create divergent results [10, 11].

461 Analytical decisions vary among studies because researchers have many options. Researchers need
462 to decide how to exclude possibly anomalous or unreliable data, how to construct variables, which
463 variables to include in their models, and which statistical methods to use. Depending on the dataset,
464 this short list of choices could encompass thousands or millions of possible alternative
465 specifications [10]. However, researchers making these decisions presumably do so with the goal of
466 doing the best possible analysis, or at least the best analysis within their current skill set. Thus it
467 seems likely that some specification options are more probable than others, possibly because they
468 have previously been shown (or claimed) to be better, or because they are more well known. Of
469 course, some of these different analyses (maybe many of them) may be equally valid alternatives.
470 Regardless, on probably any topic in ecology and evolutionary biology, we can encounter differences
471 in choices of data analysis. The extent of these differences in analyses and the degree to which these
472 differences influence the outcomes of analyses and therefore studies' conclusions are important
473 empirical questions. These questions are especially important given that many papers draw
474 conclusions after applying a single method, or even a single statistical model, to analyze a dataset.

475 The possibility that different analytical choices could lead to different outcomes has long been
476 recognized [12], and various efforts to address this possibility have been pursued in the literature.
477 For instance, one common method in ecology and evolutionary biology involves creating a set of
478 candidate models, each consisting of a different (though often similar) set of predictor variables, and
479 then, for the predictor variable of interest, averaging the slope across all models (i.e. model
480 averaging) [13, 14]. This method reduces the chance that a conclusion is contingent upon a single
481 model specification, though use and interpretation of this method is not without challenges [14].
482 Further, the models compared to each other typically differ only in the inclusion or exclusion of
483 certain predictor variables and not in other important ways, such as methods of parameter

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

484 estimation. More explicit examination of outcomes of differences in model structure, model type,
485 data exclusion, or other analytical choices can be implemented through sensitivity
486 analyses [e.g., [15](#)]. Sensitivity analyses, however, are typically rather narrow in scope, and are
487 designed to assess the sensitivity of analytical outcomes to a particular analytical choice rather than
488 to a large universe of choices. Recently, however, analysts in the social sciences have proposed
489 extremely thorough sensitivity analysis, including ‘multiverse analysis’ [[16](#)] and the ‘specification
490 curve’ [[10](#)], as a means of increasing the reliability of results. With these methods, researchers
491 identify relevant decision points encountered during analysis and conduct the analysis many times
492 to incorporate many plausible decisions made at each of these points. The study’s conclusions are
493 then based on a broad set of the possible analyses and so allow the analyst to distinguish between
494 robust conclusions and those that are highly contingent on particular model specifications. These are
495 useful outcomes, but specifying a universe of possible modelling decisions is not a trivial
496 undertaking. Further, the analyst’s knowledge and biases will influence decisions about the
497 boundaries of that universe, and so there will always be room for disagreement among analysts
498 about what to include. Including more specifications is not necessarily better. Some analytical
499 decisions are better justified than others, and including biologically implausible specifications may
500 undermine this process. Regardless, these powerful methods have yet to be adopted, and even
501 more limited forms of sensitivity analyses are not particularly widespread. Most studies publish a
502 small set of analyses and so the existing literature does not provide much insight into the degree to
503 which published results are contingent on analytical decisions.

504 Despite the potential major impacts of analytical decisions on variance in results, the outcomes of
505 different individuals’ data analysis choices have received limited empirical attention. The only formal
506 exploration of this that we were aware of when we submitted our Stage 1 manuscript were (1) an
507 analysis in social science that asked whether male professional football (soccer) players with darker
508 skin tone were more likely to be issued red cards (ejection from the game for rule violation) than

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

509 players with lighter skin tone [11] and (2) an analysis in neuroimaging which evaluated nine separate
510 hypotheses involving the neurological responses detected with fMRI in 108 participants divided
511 between two treatments in a decision making task [17]. Several others have been published
512 since [e.g., 18, 19–21]. In the red card study, twenty-nine teams designed and implemented analyses
513 of a dataset provided by the study coordinators [11]. Analyses were peer reviewed (results blind) by
514 at least two other participating analysts; a level of scrutiny consistent with standard pre-publication
515 peer review. Among the final 29 analyses, odds-ratios varied from 0.89 to 2.93, meaning point
516 estimates varied from having players with lighter skin tones receive more red cards (odds ratio < 1)
517 to a strong effect of players with darker skin tones receiving more red cards (odds ratio > 1). Twenty
518 of the 29 teams found a statistically-significant effect in the predicted direction of players with
519 darker skin tones being issued more red cards. This degree of variation in peer-reviewed analyses
520 from identical data is striking, but the generality of this finding has only just begun to be formally
521 investigated.

522 In the neuroimaging study, 70 teams evaluated each of the nine different hypotheses with the
523 available fMRI data [17]. These 70 teams followed a divergent set of workflows that produced a wide
524 range of results. The rate of reporting of statistically significant support for the nine hypotheses
525 ranged from 21% to 84%, and for each hypothesis on average, 20% of research teams observed
526 effects that differed substantially from the majority of other teams. Some of the variability in results
527 among studies could be explained by analytical decisions such as choice of software package,
528 smoothing function, and parametric versus non-parametric corrections for multiple comparisons.
529 However, substantial variability among analyses remained unexplained, and presumably emerged
530 from the many different decisions each analyst made in their long workflows. Such variability in
531 results among analyses from this dataset and from the very different red-card dataset suggests that
532 sensitivity of analytical outcome to analytical choices may characterize many distinct fields, as
533 several more recent many-analyst studies also suggest [18–20].

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

534 To further develop the empirical understanding of the effects of analytical decisions on study
535 outcomes, we chose to estimate the extent to which researchers' data analysis choices drive
536 differences in effect sizes, model predictions, and qualitative conclusions in ecology and evolutionary
537 biology. This is an important extension of the meta-research agenda of evaluating factors influencing
538 replicability in ecology, evolutionary biology, and beyond [22]. To examine the effects of analytical
539 decisions, we used two different datasets and recruited researchers to analyze one or the other of
540 these datasets to answer a question we defined. The first question was "To what extent is the
541 growth of nestling blue tits (*Cyanistes caeruleus*) influenced by competition with siblings?" To
542 answer this question, we provided a dataset that includes brood size manipulations from 332 broods
543 conducted over three years at Wytham Wood, UK. The second question was "How does grass cover
544 influence *Eucalyptus* spp. seedling recruitment?" For this question, analysts used a dataset that
545 includes, among other variables, number of seedlings in different size classes, percentage cover of
546 different life forms, tree canopy cover, and distance from canopy edge from 351 quadrats spread
547 among 18 sites in Victoria, Australia.

548 We explored the impacts of data analysts' choices with descriptive statistics and with a series of tests
549 to attempt to explain the variation among effect sizes and predicted values of the dependent variable
550 produced by the different analysis teams for both datasets separately. To describe the variability, we
551 present forest plots of the standardized effect sizes and predicted values produced by each of the
552 analysis teams, estimate heterogeneity (both absolute, τ^2 , and proportional, I^2) in effect size and
553 predicted values among the results produced by these different teams, and calculate a similarity
554 index that quantifies variability among the predictor variables selected for the different statistical
555 models constructed by the different analysis teams. These descriptive statistics provide the first
556 estimates of the extent to which explanatory statistical models and their outcomes in ecology and
557 evolutionary biology vary based on the decisions of different data analysts. We then quantified the
558 degree to which the variability in effect size and predicted values could be explained by (1) variation

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

559 in the quality of analyses as rated by peer reviewers and (2) the similarity of the choices of predictor
560 variables between individual analyses.

561 [Methods](#)

562 This project involved a series of steps (1-6) that began with identifying datasets for analyses and
563 continued through recruiting independent groups of scientists to analyze the data, allowing the
564 scientists to analyze the data as they saw fit, generating peer review ratings of the analyses (based
565 on methods, not results), evaluating the variation in effects among the different analyses, and
566 producing the final manuscript.

567 [Step 1: Select Datasets](#)

568 We used two previously unpublished datasets, one from evolutionary ecology and the other from
569 ecology and conservation.

570 [Evolutionary Ecology](#)

571 Our evolutionary ecology dataset is relevant to a sub-discipline of life-history research which focuses
572 on identifying costs and trade-offs associated with different phenotypic conditions.

573 These data were derived from a brood-size manipulation experiment imposed on wild birds nesting
574 in boxes provided by researchers in an intensively studied population.

575 Understanding how the growth of nestlings is influenced by the numbers of siblings in the nest can
576 give researchers insights into factors such as the evolution of clutch size, determination of
577 provisioning rates by parents, and optimal levels of sibling competition (Vander Werf 1992; DeKogel
578 1997; Royle et al. 1999; Verhulst, Holveck, and Riebel 2006; Nicolaus et al. 2009). Data analysts were
579 provided this dataset and instructed to answer the following question: “To what extent is the growth
580 of nestling blue tits (*Cyanistes caeruleus*) influenced by competition with siblings?”

581

582 Researchers conducted brood size manipulations and population monitoring of blue tits at Wytham
583 Wood, a 380ha woodland in Oxfordshire, U.K (1° 20'W, 51° 47'N). Researchers regularly checked

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

584 approximately 1100 artificial nest boxes at the site and monitored the 330 to 450 blue tit pairs
585 occupying those boxes in 2001-2003 during the experiment. Nearly all birds made only one breeding
586 attempt during the April to June study period in a given year. At each blue tit nest, researchers
587 recorded the date the first egg appeared, clutch size, and hatching date. For all chicks alive at age 14
588 days, researchers measured mass and tarsus length and fitted a uniquely numbered, British Trust for
589 Ornithology (BTO) aluminium leg ring. Researchers attempted to capture all adults at their nests
590 between day 6 and day 14 of the chick-rearing period. For these captured adults, researchers
591 measured mass, tarsus length, and wing length and fitted a uniquely numbered BTO leg ring. During
592 the 2001-2003 breeding seasons, researchers manipulated brood sizes using cross fostering. They
593 matched broods for hatching date and brood size and moved chicks between these paired nests one
594 or two days after hatching. They sought to either enlarge or reduce all manipulated broods by
595 approximately one fourth. To control for effects of being moved, each reduced brood had a portion
596 of its brood replaced by chicks from the paired increased brood, and vice versa. Net manipulations
597 varied from plus or minus four chicks in broods of 12 to 16 to plus or minus one chick in broods of 4
598 or 5. Researchers left approximately one third of all broods unmanipulated. These unmanipulated
599 broods were not selected systematically to match manipulated broods in clutch size or laying date.
600 We have mass and tarsus length data from 3720 individual chicks divided among 167 experimentally
601 enlarged broods, 165 experimentally reduced broods, and 120 unmanipulated broods. The full list of
602 variables included in the dataset is publicly available (<https://osf.io/hdv8m>), along with the data
603 (<https://osf.io/qjzby>).

Additional explanation:

Shortly after beginning to recruit analysts, several analysts noted a small set of related errors in the blue tit dataset. We corrected the errors, replaced the dataset on our OSF site, and emailed the analysts on 19 April 2020 to instruct them to use the revised data. The email to analysts is available here (<https://osf.io/4h53z>). The errors are explained in that email.

604

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

605 Ecology and Conservation

606 Our ecology and conservation dataset is relevant to a sub-discipline of conservation research which
607 focuses on investigating how best to revegetate private land in agricultural landscapes. These data
608 were collected on private land under the Bush Returns program, an incentive system where
609 participants entered into a contract with the Goulburn Broken Catchment Management Authority
610 and received annual payments if they executed predetermined restoration activities. This particular
611 dataset is based on a passive regeneration initiative, where livestock grazing was removed from the
612 property in the hopes that the *Eucalyptus* spp. overstorey would regenerate without active (and
613 expensive) planting. Analyses of some related data have been published (Miles 2008; Veski et al.
614 2016) but those analyses do not address the question analysts answered in our study. Data analysts
615 were provided this dataset and instructed to answer the following question: “How does grass cover
616 influence *Eucalyptus* spp. seedling recruitment?”.

617 Researchers conducted three rounds of surveys at 18 sites across the Goulburn Broken catchment in
618 northern Victoria, Australia in winter and spring 2006 and autumn 2007. In each survey period, a
619 different set of 15 x 15 m quadrats were randomly allocated across each site within 60 m of existing
620 tree canopies. The number of quadrats at each site depended on the size of the site, ranging from
621 four at smaller sites to 11 at larger sites. The total number of quadrats surveyed across all sites and
622 seasons was 351. The number of *Eucalyptus* spp. seedlings was recorded in each quadrat along with
623 information on the GPS location, aspect, tree canopy cover, distance to tree canopy, and position in
624 the landscape. Ground layer plant species composition was recorded in three 0.5 x 0.5 m sub-
625 quadrats within each quadrat. Subjective cover estimates of each species as well as bare ground,
626 litter, rock and moss/lichen/soil crusts were recorded. Subsequently, this was augmented with
627 information about the precipitation and solar radiation at each GPS location. The full list of variables
628 included in the dataset is publicly available (<https://osf.io/r5gbn>), along with the data
629 (<https://osf.io/qz5cu>).

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

630 Step 2: Recruitment and initial survey of analysts

631 The lead team (TP, HF, SN, EG, SG, PV, DH, FF) created a publicly available document providing a
632 general description of the project (<https://osf.io/mn5aj/>). The project was advertised at conferences,
633 via Twitter, using mailing lists for ecological societies (including Ecolog, Evoldir, and lists for the
634 Environmental Decisions Group, and Transparency in Ecology and Evolution), and via word of mouth.
635 The target population was active ecology, conservation, or evolutionary biology researchers with a
636 graduate degree (or currently studying for a graduate degree) in a relevant discipline. Researchers
637 could choose to work independently or in a small team. For the sake of simplicity, we refer to these
638 as ‘analysis teams’ though some comprised one individual. We aimed for a minimum of 12 analysis
639 teams independently evaluating each dataset (see sample size justification below). We
640 simultaneously recruited volunteers to peer review the analyses conducted by the other volunteers
641 through the same channels. Our goal was to recruit a similar number of peer reviewers and analysts,
642 and to ask each peer reviewer to review a minimum of four analyses. If we were unable to recruit at
643 least half the number of reviewers as analysis teams, we planned to ask analysts to serve also as
644 reviewers (after they had completed their analyses), but this was unnecessary. All analysts and
645 reviewers were offered the opportunity to share co-authorship on this manuscript and we planned to
646 invite them to participate in the collaborative process of producing the final manuscript. All analysts
647 signed [digitally] a consent (ethics) document (<https://osf.io/xyp68/>) approved by the Whitman
648 College Institutional Review Board prior to being allowed to participate.

Preregistration Deviation:

Due to the large number of recruited analysts and reviewers and the anticipated challenges of receiving and integrating feedback from so many authors, we limited analyst and reviewer participation in the production of the final manuscript to an invitation to call attention to serious problems with the manuscript draft.

649

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

650 We identified our minimum number of analysts per dataset by considering the number of effects
651 needed in a meta-analysis to generate an estimate of heterogeneity (τ^2) with a 95% confidence
652 interval that does not encompass zero. This minimum sample size is invariant regardless of τ^2 . This is
653 because the same t-statistic value will be obtained by the same sample size regardless of variance
654 (τ^2). We see this by first examining the formula for the standard error, SE for variance, (τ^2) or $SE(\tau^2)$
655 assuming normality in an underlying distribution of effect sizes [30]:

$$656 \quad SE(\tau^2) = \sqrt{\frac{t^4}{n-1}}$$

657 and then rearranging the above formula to show how the t-statistic is independent of τ^2 , as seen
658 below.

$$659 \quad t = \frac{\tau^2}{SE \tau^2} = \sqrt{\frac{n-1}{2}}$$

660 We then find a minimum $n = 12$ according to this formula.

661 [Step 3: Primary Data Analysis](#)

662 Analysis teams registered and answered a demographic and expertise survey (<https://osf.io/seqzy/>).
663 We then provided them with the dataset of their choice and requested that they answer a specific
664 research question. For the evolutionary ecology dataset that question was “To what extent is the
665 growth of nestling blue tits (*Cyanistes caeruleus*) influenced by competition with siblings?” and for
666 the conservation ecology dataset it was “How does grass cover influence *Eucalyptus* spp. seedling
667 recruitment?” Once their analysis was complete, they answered a structured survey
668 (<https://osf.io/neyc7/>), providing analysis technique, explanations of their analytical choices,
669 quantitative results, and a statement describing their conclusions. They also were asked to upload
670 their analysis files (including the dataset as they formatted it for analysis and their analysis code [if
671 applicable]) and a detailed journal-ready statistical methods section.

Preregistration Deviation:

We originally planned to have analysts complete a single survey (<https://osf.io/neyc7/>), but after we evaluated the results of that survey, we realized we would need a second survey (<https://osf.io/8w3v5/>) to adequately collect the information we needed to evaluate heterogeneity of results (step 5). We provided a set of detailed instructions with the follow-up survey, and these instructions are publicly available and can be found within the following files (blue tit: <https://osf.io/kr2g9>, *Eucalyptus*: <https://osf.io/dfvym>).

672

673 [Step 4: Peer Review of Analysis](#)

674 At minimum, each analysis was evaluated by four different reviewers, and each volunteer peer
675 reviewer was randomly assigned methods sections from at least four analyst teams (the exact
676 number varied). Each peer reviewer registered and answered a demographic and expertise survey
677 identical to that asked of the analysts, except we did not ask about ‘team name’ since reviewers did
678 not work in teams. Reviewers evaluated the methods of each of their assigned analyses one at a time
679 in a sequence determined by the project leaders. We systematically assigned the sequence so that, if
680 possible, each analysis was allocated to each position in the sequence for at least one reviewer. For
681 instance, if each reviewer were assigned four analyses to review, then each analysis would be the
682 first analysis assigned to at least one reviewer, the second analysis assigned to another reviewer, the
683 third analysis assigned to yet another reviewer, and the fourth analysis assigned to a fourth reviewer.
684 Balancing the order in which reviewers saw the analyses controls for order effects, e.g. a reviewer
685 might be less critical of the first methods section they read than the last.

686 The process for a single reviewer was as follows. First, the reviewer received a description of the
687 methods of a single analysis. This included the narrative methods section, the analysis team’s
688 answers to our survey questions regarding their methods, including analysis code, and the dataset.
689 The reviewer was then asked, in an online survey (<https://osf.io/4t36u/>), to rate that analysis on a

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

690 scale of 0-100 based on this prompt: "Rate the overall appropriateness of this analysis to answer the
691 research question (one of the two research questions inserted here) with the available data. To help
692 you calibrate your rating, please consider the following guidelines:

693

694 100: A perfect analysis with no conceivable improvements from the reviewer

695 75: An imperfect analysis but the needed changes are unlikely to dramatically alter outcomes

696 50: A flawed analysis likely to produce either an unreliable estimate of the relationship or an over-
697 precise estimate of uncertainty

698 25: A flawed analysis likely to produce an unreliable estimate of the relationship and an over-precise
699 estimate of uncertainty

700 0: A dangerously misleading analysis, certain to produce both an estimate that is wrong and a
701 substantially over-precise estimate of uncertainty that places undue confidence in the incorrect
702 estimate.

703 *Please note that these values are meant to calibrate your ratings. We welcome ratings of any
704 number between 0 and 100."

705

706 After providing this rating, the reviewer was presented with this prompt, in multiple-choice format:

707 "Would the analytical methods presented produce an analysis that is (a) publishable as is, (b)

708 publishable with minor revision, (c) publishable with major revision, (d) deeply flawed and

709 unpublishable?" The reviewer was then provided with a series of text boxes and the following

710 prompts: "Please explain your ratings of this analysis. Please evaluate the choice of statistical analysis

711 type. Please evaluate the process of choosing variables for and structuring the statistical model.

712 Please evaluate the suitability of the variables included in (or excluded from) the statistical model.

713 Please evaluate the suitability of the structure of the statistical model. Please evaluate choices to

714 exclude or not exclude subsets of the data. Please evaluate any choices to transform data (or, if there

715 were no transformations, but you think there should have been, please discuss that choice)." After

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

716 submitting this review, a methods section from a second analysis was then made available to the
717 reviewer. This same sequence was followed until all analyses allocated to a given reviewer were
718 provided and reviewed. After providing the final review, the reviewer was simultaneously provided
719 with all four (or more) methods sections the reviewer had just completed reviewing, the option to
720 revise their original ratings, and a text box to provide an explanation. The invitation to revise the
721 original ratings was as follows: “If, now that you have seen all the analyses you are reviewing, you
722 wish to revise your ratings of any of these analyses, you may do so now.” The text box was prefaced
723 with this prompt: “Please explain your choice to revise (or not to revise) your ratings.”

Additional Explanation:

To determine how consistent peer reviewers were in their ratings, we assessed inter-rater reliability among reviewers for both the categorical and quantitative ratings combining blue tit and *Eucalyptus* data using Krippendorff’s alpha for ordinal and continuous data respectively. This provides a value that is between -1 (total disagreement between reviewers) and 1 (total agreement between reviewers).

724

725 [Step 5: Evaluate Variation](#)

726 The lead team conducted the analyses outlined in this section. We described the variation in model
727 specification in several ways. We calculated summary statistics describing variation among analyses,
728 including mean, SD, and range of number of variables per model included as fixed effects, the
729 number of interaction terms, the number of random effects, and the mean, SD, and range of sample
730 sizes. We also present the number of analyses in which each variable was included. We summarized
731 the variability in standardized effect sizes and predicted values of dependent variables among the
732 individual analyses using standard random effects meta-analytic techniques. First, we derived
733 standardized effect sizes from each individual analysis. We did this for all linear models or
734 generalized linear models by converting the t value and the degree of freedom (df) associated with

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

735 regression coefficients (e.g. the effect of the number of siblings [predictor] on growth [response] or
736 the effect of grass cover [predictor] on seedling recruitment [response]) to the correlation coefficient
737 (r), using the following:

738
$$r = \sqrt{\frac{t^2}{t^2 + df}}$$

739 This formula can only be applied if t and df values originate from linear or generalized linear models
740 [GLMs; [31](#)]. If, instead, linear mixed-effects models (LMMs) or generalized linear mixed-effects
741 models (GLMMs) were used by a given analysis, the exact df cannot be estimated. However, adjusted
742 df can be estimated, for example, using the Satterthwaite approximation of df, df_s , [note that SAS
743 uses this approximation to obtain df for LMMs and GLMMs; [32](#)]. For analyses using either LMMs or
744 GLMMs that do not produce df_s , we planned to obtain df_s by rerunning the same (G)LMMs using the
745 `lmer()` or `glmer()` function in the `lmerTest` package in R [33](#), [34](#)].

Preregistration Deviation

Rather than re-run these analyses ourselves, we sent a follow-up survey (referenced above under “Primary data analyses”) to analysts and asked them to follow our instructions for producing this information. The instructions are publicly available and can be found within the following files (blue tit: <https://osf.io/kr2g9>, *Eucalyptus*: <https://osf.io/dfvym>).

746
747 We then used the t values and df_s from the models to obtain r as per the formula above. All r and
748 accompanying df (or df_s) were converted to Zr and its sampling variance $1/(n-3)$ where $n=df+1$. Any
749 analyses from which we could not derive a signed Zr, for instance one with a quadratic function in
750 which the slope changed sign, were excluded from the analyses of Fisher’s Zr. We expected such
751 analyses would be rare. In fact, most submitted analyses excluded from our meta-analysis of Zr were
752 excluded because of a lack of sufficient information provided by the analyst team rather than due to
753 the use of effects that could not be converted to Zr. Regardless, as we describe below, we generated

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

754 a second set of standardized effects (predicted values) that could (in principle) be derived from any
755 explanatory model produced by these data.

756 Besides Z_r , which describes the strength of a relationship based on the amount of variation in a
757 dependent variable explained by variation in an independent variable, we also examined differences
758 in the shape of the relationship between the independent and dependent variables. To accomplish
759 this, we derived a point estimate (out-of-sample predicted value) for the dependent variable of
760 interest for each of three values of our primary independent variable. We originally described these
761 three values as associated with the 25th percentile, median, and 75th percentile of the independent
762 variable and any covariates.

Preregistration Deviation

The original description of the out-of-sample specifications did not account for the facts that (a) some variables are not distributed in a way that allowed division in percentiles and that (b) variables could be either positively or negatively correlated with the dependent variable. We provide a more thorough description here: We derived three point-estimates (out-of-sample predicted values) for the dependent variable of interest; one for each of three values of our primary independent variable that we specified. We also specified values for all other variables that could have been included as independent variables in analysts' models so that we could derive the predicted values from a fully specified version of any model produced by analysts. For all potential independent variables, we selected three values or categories. Of the three we selected, one was associated with small, one with intermediate, and one with large values of one typical dependent variable (day 14 chick weight for the blue tit data and total number of seedlings for the *Eucalyptus* data; analysts could select other variables as their dependent variable, but the others typically correlated with the two identified here). For continuous variables, this means we identified the 25th percentile, median, and 75th percentile and, if the slope of the linear relationship between this variable and the typical dependent variable was positive, we left the quartiles ordered as is. If, instead, the slope was negative, we reversed the order of the independent variable quartiles so that the 'lower' quartile value was the one associated with the lower value for the dependent variable. In the case of categorical variables, we identified categories associated with the 25th percentile, median, and 75th percentile values of the typical dependent variable after averaging the values for each category. However, for some continuous and categorical predictors, we also made selections based on the principle of internal consistency between certain related variables, and we fixed a few categorical variables as identical across all three levels where doing so would simplify the modelling process (specification tables available: blue tit: <https://osf.io/86akx>; *Eucalyptus*: <https://osf.io/jh7g5>).

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

764 We used the 25th and 75th percentiles rather than minimum and maximum values to reduce the
765 chance of occupying unrealistic parameter space. We planned to derive these predicted values from
766 the model information provided by the individual analysts. All values (predictions) were first
767 transformed to the original scale along with their standard errors (SE); we used the delta method
768 (Ver Hoef 2012) for the transformation of SE. We used the square of the SE associated with predicted
769 values as the sampling variance in the meta-analyses described below, and we planned to analyze
770 these predicted values in exactly the same ways as we analyzed Z_r in the following analyses.

Preregistration Deviation

Because analysts of blue tit data chose different dependent variables on different scales, after transforming out-of-sample values to the original scales, we standardized all values as z scores ('standard scores') to put all dependent variables on the same scale and make them comparable. This involved taking each relevant value on the original scale (whether a predicted point estimate or a SE associated with that estimate) and subtracting the value in question from the mean value of that dependent variable derived from the full dataset and then dividing this difference by the standard deviation, SD, corresponding to the mean from the full dataset. Thus, all our out-of-sample prediction values from the blue tit data are from a distribution with the mean of 0 and SD of 1. We did not add this step for the *Eucalyptus* data because (a) all responses were on the same scale (counts of *Eucalyptus* stems) and were thus comparable and (b) these data, with many zeros and high skew, are poorly suited for z scores.

771
772 We plotted individual effect size estimates (Z_r) and predicted values of the dependent variable (y_i)
773 and their corresponding 95% confidence / credible intervals in forest plots to allow visualization of
774 the range and precision of effect size and predicted values. Further, we included these estimates in
775 random effects meta-analyses [36, 37] using the metafor package in R [34, 38]:

776 $Z_r \sim 1 + 1 | \text{analysisId}$

777 $y_i \sim 1 + 1 | \text{analysisId}$

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

778 where y_i is the predicted value for the dependent variable at the 25th percentile, median, or 75th
779 percentile of the independent variables. The individual Z_r effect sizes were weighted with the inverse
780 of sampling variance for Z_r . The individual predicted values for dependent variable (y_i) were weighted
781 by the inverse of the associated SE^2 original registration omitted “inverse of the” in error). These
782 analyses provided an average Z_r score or an average y_i with corresponding 95% confidence interval
783 and allowed us to estimate two heterogeneity indices, τ^2 and I^2 . The former, τ^2 , is the absolute
784 measure of heterogeneity or the between-study variance (in our case, between-effect variance)
785 whereas I^2 is a relative measure of heterogeneity. We obtained the estimate of relative heterogeneity
786 (I^2) by dividing the between-effect variance by the sum of between-effect and within-effect variance
787 (sampling error variance). I^2 is thus, in a standard meta-analysis, the proportion of variance that is
788 due to heterogeneity as opposed to sampling error. When calculating I^2 , within-study variance is
789 amalgamated across studies to create a “typical” within-study variance which serves as the sampling
790 error variance [36, 37]. Our goal here was to visualize and quantify the degree of variation among
791 analyses in effect size estimates [31]. We did not test for statistical significance.

Additional Explanation

Our use of I^2 to quantify heterogeneity violates an important assumption, but this violation does not invalidate our use of I^2 as a metric of how much heterogeneity can derive from analytical decisions. In standard meta-analysis, the statistic I^2 quantifies the proportion of variance that is greater than we would expect if differences among estimates were due to sampling error alone [39]. However, it is clear that this interpretation does not apply to our value of I^2 because I^2 assumes that each estimate is based on an independent sample (although these analyses can account for non-independence via hierarchical modelling), whereas all our effects were derived from largely or entirely overlapping subsets of the same dataset. Despite this, we believe that I^2 remains a useful statistic for our purposes. This is because, in calculating I^2 , we are still setting a benchmark of expected variation due to sampling error based on the variance associated with each separate effect size estimate, and we are assessing how much (if it all) the variability among our effect sizes exceeds what would be expected had our effect sizes been based on independent data. In other words, our estimates can tell us how much proportional heterogeneity is possible from analytical decisions alone when sample sizes (and therefore meta-analytic within-estimate variance) are similar to the ones in our analyses. Among other implications, our violation of the independent sample assumption means that we (dramatically) over-estimate the variance expected due to sampling error, and because I^2 is a proportional estimate, we thus underestimate the actual proportion of variance due to differences among analyses other than sampling error. However, correcting this underestimation would create a trivial value since we designed the study so that much of the variance would derive from analytic decisions as opposed to differences in sampled data. Instead, retaining the I^2 value as typically calculated provides a useful comparison to I^2 values from typical meta-analyses.

Interpretation of τ^2 also differs somewhat from traditional meta-analysis, and we discuss this further in the Results.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

793 Finally, we assessed the extent to which deviations from the meta-analytic mean by individual effect
794 sizes (Z_r) or the predicted values of the dependent variable (y_i) were explained by the peer rating of
795 each analysis team's method section, by a measurement of the distinctiveness of the set of predictor
796 variables included in each analysis, and by the choice of whether or not to include random effects in
797 the model. The deviation score, which served as the dependent variable in these analyses, is the
798 absolute value of the difference between the meta-analytic mean Z_r (or y_i) and the individual Z_r (or
799 y_i) estimate for each analysis. We used the Box-Cox transformation on the absolute values of
800 deviation scores to achieve an approximately normal distribution [c.f. [40](#), [41](#)]. We described variation
801 in this dependent variable with both a series of univariate analyses and a multivariate analysis. All
802 these analyses were general linear (mixed) models. These analyses were secondary to our estimation
803 of variation in effect sizes described above. We wished to quantify relationships among variables, but
804 we had no a priori expectation of effect size and made no dichotomous decisions about statistical
805 significance.

806 When examining the extent to which reviewer ratings (on a scale from 0 to 100) explained deviation
807 from the average effect (or predicted value), each analysis had been rated by multiple peer
808 reviewers, so for each reviewer score to be included, we include each deviation score in the analysis
809 multiple times. To account for the non-independence of multiple ratings of the same analysis, we
810 planned to include analysis identity as a random effect in our general linear mixed model in the lme4
811 package in R [[34](#), [42](#)]. To account for potential differences among reviewers in their scoring of
812 analyses, we also planned to include reviewer identity as a random effect:

813 $\text{DeviationScore}_j = \text{BoxCox}(\text{abs}(\text{DeviationFromMean}_j))$

814 $\text{DeviationScore}_{ij} \sim \text{Rating}_{ij} + \text{ReviewerID}_i + \text{AnalysisID}_j$

815 $\text{ReviewerID}_i \sim \mathcal{N}(0, \sigma^2)$

816 $\text{AnalysisID}_i \sim \mathcal{N}(0, \sigma^2)$

817 Where $\text{DeviationFromMean}_j$ is the deviation from the meta-analytic mean for the j th analysis,

818 ReviewerID_i is the random intercept assigned to each i reviewer, and AnalysisID_j is the random

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

819 intercept assigned to each j analysis, both of which are assumed to be normally distributed with a
820 mean of 0 and a variance of σ^2 Absolute deviation scores were Box-Cox transformed using the
821 `step_box_cox()` function from the `timetk` package in R [34, 43].

822 We conducted a similar analysis with the four categories of reviewer ratings ((1) deeply flawed and
823 unpublishable, (2) publishable with major revision, (3) publishable with minor revision, (4)
824 publishable as is) set as ordinal predictors numbered as shown here. As with the analyses above, we
825 planned for these analyses to also include random effects of analysis identity and reviewer identity.
826 Both of these analyses (1: 1-100 ratings as the fixed effect, 2: categorical ratings as the fixed effects)
827 were planned to be conducted eight times for each dataset. Each of the four responses (Z_r , y_{25} , y_{50} ,
828 y_{75}) were to be compared once to the initial ratings provided by the peer reviewers, and again based
829 on the revised ratings provided by the peer reviewers.

Preregistration Deviation

1. We planned to include random effects of both analysis identity and reviewer identity in these models comparing reviewer ratings with deviation scores. However, after we received the analyses, we discovered that a subset of analyst teams had either conducted multiple analyses and/or identified multiple effects per analysis as answering the target question. We therefore faced an even more complex potential set of random effects. We decided that including team ID, analysis ID, and effect ID along with reviewer ID as random effects in the same model would almost certainly lead to model fit problems, and so we started with simpler models including just effect ID and reviewer ID. However, even with this simpler structure, our dataset was sparse, with reviewers rating a small number of analyses, resulting in models with singular fit (Section C.2). Removing one of the random effects was necessary for the models to converge. The models that included the categorical quality rating converged when including reviewer ID, and the models that included the continuous quality rating converged when including effect ID.
2. We conducted analyses only with the final peer ratings after the opportunity for revision, not with the initial ratings. This was because when we recorded the final ratings, they over-wrote the initial ratings, and so we did not have access to those initial values.

830

831 The next set of univariate analyses sought to explain deviations from the mean effects based on a
832 measure of the distinctiveness of the set of variables included in each analysis. As a 'distinctiveness'
833 score, we used Sorensen's Similarity Index (an index typically used to compare species composition
834 across sites), treating variables as species and individual analyses as sites. To generate an individual
835 Sorensen's value for each analysis required calculating the pairwise Sorensen's value for all pairs of
836 analyses (of the same dataset), and then taking the average across these Sorensen's values for each
837 analysis. We calculated the Sorensen's index values using the betapart package [44] in R:

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

838
$$\beta_{Sorensen} = \frac{b + c}{2a + b + c}$$

839 Where a is the number of variables common to both analyses, b is the number of variables that
840 occur in the first analysis but not in the second and c is the number of variables that occur in the
841 second analysis. We then used the per-model average Sorensen's index value as an independent
842 variable to predict the deviation score in a general linear model, and included no random effect since
843 each analysis is included only once, in R [34]:

844
$$DeviationScore_j \sim \beta_{Sorensen}$$

Additional Explanation

When we planned this analysis, we anticipated that analysts would identify a single primary effect from each model, so that each model would appear in the analysis only once. Our expectation was incorrect because some analysts identified >1 effect per analysis, but we still chose to specify our model as registered and not use a random effect. This is because most models produced only one effect and so we expected that specifying a random effect to account for the few cases where >1 effect was included for a given model would prevent model convergence.

Note that this analysis contrasts with the analyses in which we used reviewer ratings as predictors because in the analyses with reviewer ratings, each effect appeared in the analysis approximately four times due to multiple reviews of each analysis, and so it was much more important to account for that variance through a random effect.

845
846 Finally, we conducted a multivariate analysis with the five predictors described above (peer ratings 0-
847 100 and peer ratings of publishability 1-4; both original and revised and Sorensen's index, plus a
848 sixth, presence /absence of random effects) with random effects of analysis identity and reviewer
849 identity in the *lme4* package in R [34, 42]. We had stated here in the text that we would use only the
850 revised (final) peer ratings in this analysis, so the absence of the initial ratings is not a deviation from
851 our plan:

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

852 $DeviationScore_j \sim RatingsContinuous_{ij} + RatingsCategorical_{ij} + \beta Sorensen_j + AnalysisID_j + ReviewerID_i$

853 $ReviewerID_i \sim \chi(0, \sigma^2)$

854 $AnalysisID_j \sim \chi(0, \sigma^2)$

855 We conducted all the analyses described above eight times; for each of the four responses (Zr , γ_{25} ,
856 γ_{50} , γ_{75}) one time for each of the two datasets.

857 We have publicly archived all relevant data, code, and materials on the Open Science Framework
858 (<https://osf.io/mn5aj/>). Archived data includes the original datasets distributed to all analysts, any
859 edited versions of the data analyzed by individual groups, and the data we analyzed with our meta-
860 analyses, which include the effect sizes derived from separate analyses, the statistics describing
861 variation in model structure among analyst groups, and the anonymized answers to our surveys of
862 analysts and peer reviewers. Similarly, we have archived both the analysis code used for each
863 individual analysis (where available) and the code from our meta-analyses. We have also archived
864 copies of our survey instruments from analysts and peer reviewers.

865 Our rules for excluding data from our study were as follows. We excluded from our synthesis any
866 individual analysis submitted after we had completed peer review or those unaccompanied by
867 analysis files that allow us to understand what the analysts did. We also excluded any individual
868 analysis that did not produce an outcome that could be interpreted as an answer to our primary
869 question (as posed above) for the respective dataset. For instance, this means that in the case of the
870 data on blue tit chick growth, we excluded any analysis that did not include something that can be
871 interpreted as growth or size as a dependent (response) variable, and in the case of the *Eucalyptus*
872 establishment data, we excluded any analysis that did not include a measure of grass cover among
873 the independent (predictor) variables. Also, as described above, any analysis that could not produce
874 an effect that could be converted to a signed Zr was excluded from analyses of Zr.

Preregistration Deviation

Some analysts had difficulty implementing our instructions to derive the out-of-sample predictions, and in some cases (especially for the *Eucalyptus* data), they submitted predictions with implausibly extreme values. We believed these values were incorrect and thus made the conservative decision to exclude out-of-sample predictions where the estimates were > 3 standard deviations from the mean value from the full dataset.

875

Additional Explanation

1. Evaluating model fit.

We evaluated all fitted models using the `performance()` function from the `performance` package [45] and the `glance()` function from the `broom.mixed` package [46]. For all models, we calculated the square root of the residual variance (Sigma) and the root mean squared error (RMSE). For GLMMs `performance()` calculates the marginal and conditional R^2 values as well as the contribution of random effects (ICC), based on Nakagawa et al. [47]. The conditional R^2 accounts for both the fixed and random effects, while the marginal R^2 considers only the variance of the fixed effects. The contribution of random effects is obtained by subtracting the marginal R^2 from the conditional R^2 .

2. Exploring outliers and analysis quality.

After seeing the forest plots of Z_r values and noticing the existence of a small number of extreme outliers, especially from the Eucalyptus analyses, we wanted to understand the degree to which our heterogeneity estimates were influenced by these outliers. To explore this question, we removed the highest two and lowest two values of Z_r in each dataset and re-calculated our heterogeneity estimates.

To help understand the possible role of the quality of analyses in driving the heterogeneity we observed among estimates of Z_r , we recalculated our heterogeneity estimates after removing all effects from analysis teams that had received at least one rating of “deeply flawed and unpublishable” and then again after removing all effects from analysis teams with at least one rating of either “deeply flawed and unpublishable” or “publishable with major revisions”. We also used self-identified levels of statistical expertise to examine heterogeneity when we retained analyses only from analysis teams that contained at least one member who rated themselves as “highly proficient” or “expert” (rather than “novice” or “moderately proficient”) in conducting statistical analyses in their research area in our intake survey.

Additional Explanation

3. Exploring possible impacts of lower quality estimates of degrees of freedom.

Our meta-analyses of variation in Zr required variance estimates derived from estimates of the degrees of freedom in original analyses from which Zr estimates were derived. While processing the estimates of degrees of freedom submitted by analysts, we identified a subset of these estimates in which we had lower confidence because two or more effects from the same analysis were submitted with identical degrees of freedom. We therefore conducted a second set of (more conservative) meta-analyses that excluded these Zr estimates with identical estimates of degrees of freedom and we present these analyses in the supplement.

877

878 [Step 6: Facilitated Discussion and Collaborative Write-Up of Manuscript](#)

879 We planned for analysts and initiating authors to discuss the limitations, results, and implications of
880 the study and collaborate on writing the final manuscript for review as a stage-2 Registered Report.

Preregistration Deviation

As described above, due to the large number of recruited analysts and reviewers and the anticipated challenges of receiving and integrating feedback from so many authors, we limited analyst and reviewer participation in the production of the final manuscript to an invitation to call attention to serious problems with the manuscript draft.

881

882 [Results](#)

883 [Summary Statistics](#)

884 In total, 173 analyst teams, comprising 246 analysts, contributed 182 usable analyses of the two
885 datasets examined in this study which yielded 215 effects. Analysts produced 135 distinct effects that
886 met our criteria for inclusion in at least one of our meta-analyses for the blue tit dataset. Analysts

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

887 produced 81 distinct effects meeting our criteria for inclusion for the Eucalyptus dataset. Excluded
888 analyses and effects either did not answer our specified biological questions, were submitted with
889 insufficient information for inclusion in our meta-analyses, or were incompatible with production of
890 our effect size(s). We expected this final scenario (incompatible analyses), for instance we cannot
891 extract a Zr from random forest models, which is why we analyzed two distinct types of effects, Zr
892 and out-of-sample (y_i). Effects included in only a subset of our meta-analyses provided sufficient
893 information for inclusion in only that subset (see Table A.1). For both datasets, most submitted
894 analyses incorporated mixed effects. Submitted analyses of the blue tit dataset typically specified
895 normal error and analyses of the Eucalyptus dataset typically specified a non-normal error
896 distribution ([Supplementary Table A.1](#)).

897 For both datasets, the composition of models varied substantially in regards to the number of fixed
898 and random effects, interaction terms, and the number of data points used, and these patterns
899 differed somewhat between the blue tit and Eucalyptus analyses (See [Supplementary Table A.2](#)).

900 Focusing on the models included in the Zr analyses (because this is the larger sample), blue tit
901 models included a similar number of fixed effects on average (mean 5.2 ± 2.92 SD) as *Eucalyptus*
902 models (mean 5.01 ± 3.83 SD), but the standard deviation in number of fixed effects was somewhat
903 larger in the Eucalyptus models. The average number of interaction terms was much larger for the
904 blue tit models (mean 0.44 ± 1.11 SD) than for the Eucalyptus models (mean 0.16 ± 0.65 SD), but still
905 under 0.5 for both, indicating that most models did not contain interaction terms. Blue tit models
906 also contained more random effects (mean 3.53 ± 2.08 SD) than Eucalyptus models (mean $1.41 \pm$
907 1.09 SD). The maximum possible sample size in the blue tit dataset (3720 nestlings) was an order of
908 magnitude larger than the maximum possible in the Eucalyptus dataset (351 plots), and the means
909 and standard deviations of the sample size used to derive the effects eligible for our study were also
910 an order of magnitude greater for the blue tit dataset (mean 2622.07 ± 939.28 SD) relative to the
911 *Eucalyptus* models (mean 298.43 ± 106.25 SD). However, the standard deviation in sample size from
912 the *Eucalyptus* models was heavily influenced by a few cases of dramatic sub-setting (described

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

913 below). Approximately three quarters of *Eucalyptus* models used sample sizes within 3% of the
914 maximum. In contrast, fewer than 20% of blue tit models relied on sample sizes within 3% of the
915 maximum, and approximately 50% of blue tit models relied on sample sizes 29% or more below the
916 maximum.

917 Analysts provided qualitative descriptions of the conclusions of their analyses. Each analysis team
918 provided one conclusion per dataset. These conclusions could take into account the results of any
919 formal analyses completed by the team as well as exploratory and visual analyses of the data. Here
920 we summarize all qualitative responses, regardless of whether we had sufficient information to use
921 the corresponding model results in our quantitative analyses below. We classified these conclusions
922 into the categories summarized below (Table 1):

923

924 Mixed: some evidence supporting a positive effect, some evidence supporting a negative effect

925 Conclusive negative: negative relationship described without caveat

926 Qualified negative: negative relationship but only in certain circumstances or where analysts express
927 uncertainty in their result

928 Conclusive none: analysts interpret the results as conclusive of no effect

929 None qualified: analysts describe finding no evidence of a relationship but they describe the
930 potential for an undetected effect

931 Qualified positive: positive relationship described but only in certain circumstances or where analysts
932 express uncertainty in their result

933 Conclusive positive: positive relationship described without caveat

934

935 For the blue tit dataset, most analysts concluded that there was negative relationship between
936 measures of sibling competition and nestling growth, though half the teams expressed qualifications
937 or described effects as mixed or absent. For the *Eucalyptus* dataset, there was a broader spread of
938 conclusions with at least one analyst team providing conclusions consistent with each conclusion

939 category. The most common conclusion for the *Eucalyptus* dataset was that there was no
940 relationship between grass cover and *Eucalyptus* recruitment (either conclusive or qualified
941 description of no relationship), but more than half the teams concluded that there were effects;
942 negative, positive, or mixed.

943 Table 1: Tallies of analysts' qualitative answers to the research questions addressed by their
944 analyses.

Dataset	Mixed	Negative	Negative	None	None	Positive	Positive
		Conclusive	Qualified	Conclusive	Qualified	Qualified	Conclusive
blue tit	5	37	27	4	1	0	0
<i>Eucalyptus</i>	8	6	12	19	12	4	2

945

946 [Distribution of Effects](#)

947 [Effect Size Zr](#)

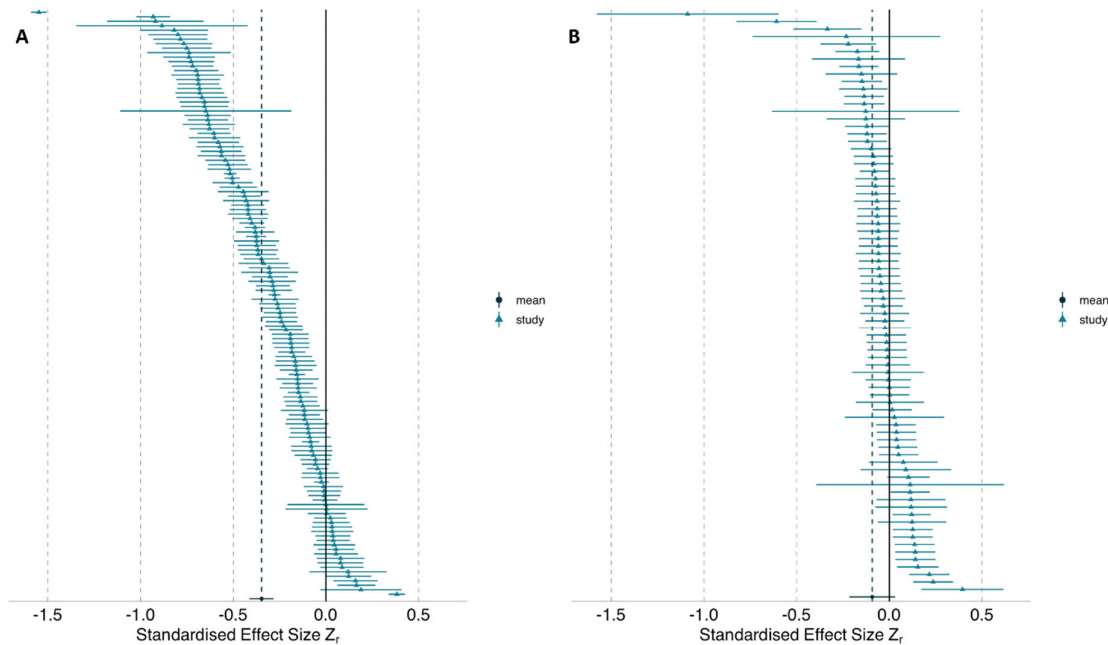
948 Although the majority (111 of 132) of the usable Zr effects from the blue tit dataset found nestling
949 growth decreased with sibling competition, and the meta-analytic mean Zr (Fisher's transformation
950 of the correlation coefficient) was convincingly negative (-0.35 ± 0.06 95% CI), there was substantial
951 variability in the strength and the direction of this effect. Zr ranged approximately continuously from
952 -0.93 to 0.19 , (Figure 1a and Table 4) and of the 111 effects with negative slopes, 92 had confidence
953 intervals excluding 0. Of the 20 with positive slopes indicating increased nestling growth in the
954 presence of more siblings, 3 had confidence intervals excluding zero (Figure 1a).

955 Meta-analysis of the *Eucalyptus* dataset also showed substantial variability in the strength of effects
956 as measured by Zr, and unlike with the blue tits, a notable lack of consistency in the direction of
957 effects (Figure 1b, Table 4). Zr ranged from -4.47 ([Supplementary Figure A.2](#)), indicating a strong
958 tendency for reduced *Eucalyptus* seedling success as grass cover increased, to 0.39 , indicating the
959 opposite. Although the range of reported effects skewed strongly negative, this was due to a small
960 number of substantial outliers. Most values of Zr were relatively small with values < 0.2 and the

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

961 meta-analytic mean effect size was close to zero (-0.09 ± 0.12 95% CI). Of the 79 effects, fifty-three
962 had confidence intervals overlapping zero, approximately a quarter (fifteen) crossed the traditional
963 threshold of statistical significance indicating a negative relationship between grass cover and
964 seedling success, and eleven crossed the significance threshold indicating a positive relationship
965 between grass cover and seedling success (Figure 1b).

966



967

968 Figure 1: Forest plots of meta-analytic estimated standardized effect sizes (Z_r) and their 95%
969 confidence intervals for each effect size included in the meta-analysis model for a) blue tit and b)
970 *Eucalyptus*. The meta-analytic mean effect size is noted in black and as a dashed vertical line, with
971 error bars also representing the 95% confidence interval. The solid black vertical line demarcates
972 effect size of 0, indicating no relationship between the test variable and the response variable. Note
973 that the Eucalyptus plot omits one extreme outlier with the value of -4.47 (Figure A.2) in order to
974 standardize the x-axes on these two panels.

975

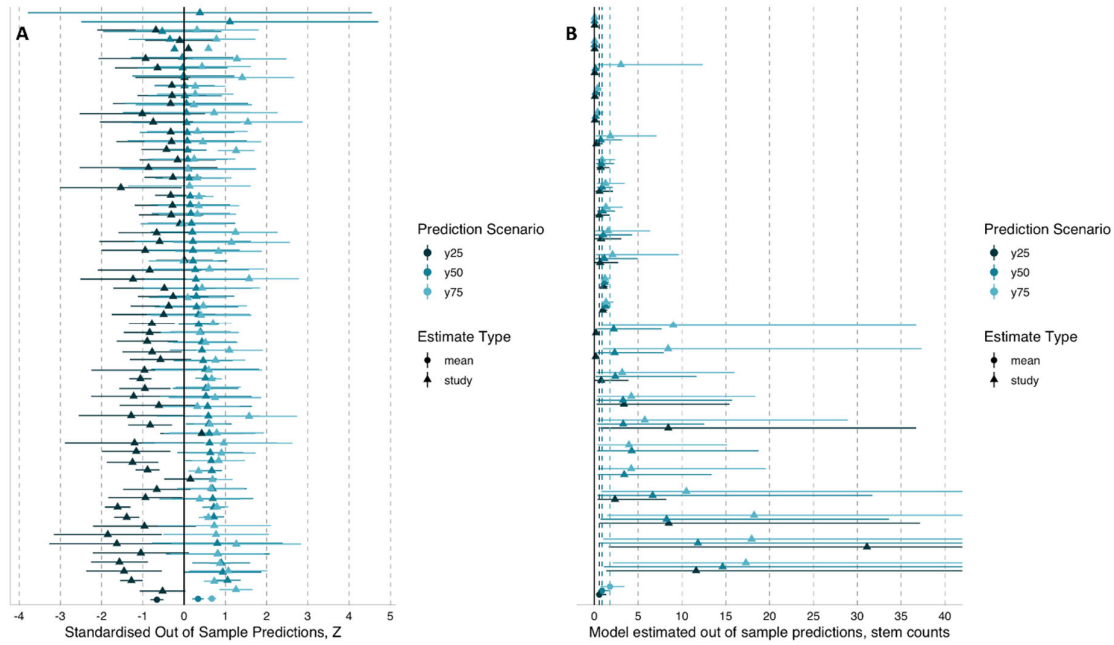
we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

976 Out-of-sample predictions (y_i)

977 As with the effect size Z_r , we observed substantial variability in the size of out-of-sample predictions
978 derived from the analysts' models. Blue tit predictions (Figure 2a), which were z-score-standardised
979 to accommodate the use of different response variables, always ranged far in excess of one standard
980 deviation. In the y_{25} scenario, model predictions ranged from -1.85 to 0.42 (a range of 2.68 standard
981 deviations), in the y_{50} scenario, they ranged from -0.53 to 1.11 (a range of 1.63 standard deviations),
982 and in the y_{75} scenario they ranged from -0.03 to 1.58 (a range of 1.9 standard deviations). As should
983 be expected given the existence of both negative and positive Z_r values, all three out-of-sample
984 scenarios produced both negative and positive predictions, although as with the Z_r values, there is a
985 clear trend for scenarios with more siblings to be associated with smaller nestlings. This is supported
986 by the meta-analytic means of these three sets of predictions which were -0.66 (95% CI -0.82,-0.5)
987 for the y_{25} , 0.34 (95% CI 0.2-0.48) for the y_{50} , and 0.67 (95% CI 0.57-0.77) for the y_{75} .

988 *Eucalyptus* out-of-sample predictions also varied substantially (Figure 2b), but because they were not
989 z-score-standardised and are instead on the original count scale, the types of interpretations we can
990 make differ. The predicted *Eucalyptus* seedling counts per 15 x 15 m plot for the y_{25} scenario ranged
991 from 0.04 to 33.66, for the y_{50} scenario ranged from 0.03 to 13.02, and for the y_{75} scenario they
992 ranged from 0.05 to 21.93. The meta-analytic mean predictions for these three scenarios were
993 similar; 0.58 (95% CI, 0.21,-1.37) for the y_{25} , 0.92 (95% CI 0.36-1.65) for the y_{50} , and 1.67 (95% CI 0.8-
994 2.83) for the y_{75} scenarios respectively.

995



996

997 Figure 2: Forest plot of meta-analytic estimated standardized (z-score) blue tit out-of-sample
998 predictions, y_i , for a) blue tit and b) Eucalyptus. Triangles represent individual estimates, circles
999 represent the meta-analytic mean for each prediction scenario. Error bars are 95% confidence
1000 intervals.

1001 Quantifying Heterogeneity

1002 Effect Size (Z_r)

1003 We quantified both absolute (τ^2) and relative (I^2) heterogeneity resulting from analytical variation.

1004 Both measures suggest that substantial variability among effect sizes was attributable to the
1005 analytical decisions of analysts.

1006 The total absolute level of variance beyond what would typically be expected due to sampling error,

1007 τ^2 (Table 2), among all usable blue tit effects was 0.088 and for Eucalyptus effects was 0.267. This is

1008 similar to or exceeding the median value (0.105) of τ^2 found across 31 recent meta-analyses

1009 (calculated from the data in [48](#)). The similarity of our observed values to values from meta-analyses

1010 of different studies based on different data suggest the potential for a large portion of heterogeneity

1011 to arise from analytical decisions. For further discussion of interpretation of τ^2 in our study, please

1012 consult discussion of post hoc analyses below.

1013 Table 2: Heterogeneity in the estimated effects Z_r for meta-analyses of the full dataset, as well as
 1014 from post hoc analyses including the dataset with outliers removed, the dataset excluding effects
 1015 from analysis teams with at least one “unpublishable” rating, the dataset excluding effects from
 1016 analysis teams with at least one “major revisions” rating or worse, or the dataset including only
 1017 analyses from teams in which at least one analyst rated themselves as "highly proficient" or "expert"
 1018 in statistical analysis. τ_{Team}^2 is the absolute heterogeneity for the random effect Team, $\tau_{EffectID}^2$ is
 1019 the absolute heterogeneity for the random effect *EffectID*, nested under Team, and τ_{total}^2 is the total
 1020 absolute heterogeneity. I^2_{Total} is the proportional heterogeneity; the proportion of the variance
 1021 among effects not attributable to sampling error, I^2_{Team} is the subset of the proportional
 1022 heterogeneity due to differences among Teams and $I^2_{Team, EffectID}$ is subset of the proportional
 1023 heterogeneity attributable to among-*EffectID* differences.

Dataset	τ^2_{Total}	τ^2_{Team}	$\tau^2_{EffectID}$	I^2_{Total}	I^2_{Team}	$I^2_{Team, EffectID}$	N. Obs
All Analyses							
blue tit	0.09	0.04	0.05	97.732%	40.11%	57.63%	131
<i>Eucalyptus</i>	0.27	0.02	0.25	98.589%	6.88%	91.71%	79
All analyses, outliers Removed							
blue tit	0.07	0.05	0.02	97.030%	66.90%	30.13%	127
<i>Eucalyptus</i>	0.01	0.00	0.01	66.193%	19.27%	46.93%	75
Analyses receiving at least one 'Unpublishable' rating removed							
blue tit	0.08	0.03	0.05	97.601%	38.10%	59.50%	109
<i>Eucalyptus</i>	0.01	0.01	0.01	79.741%	28.32%	51.42%	55
Analyses receiving at least one 'Unpublishable' and or 'Major Revisions' rating removed							
blue tit	0.14	0.01	0.13	98.718%	5.17%	93.55%	32
<i>Eucalyptus</i>	0.03	0.03	0.00	88.915%	88.91%	0.00%	13
Analyses from teams that include highly proficient or expert data analysts							
blue tit	0.10	0.04	0.06	98.058%	36.27%	61.78%	89
<i>Eucalyptus</i>	0.58	0.02	0.56	99.412%	3.49%	95.93%	34

1024
 1025 In our analyses, I^2 is a plausible index of how much more variability among effect sizes we have
 1026 observed, as a proportion, than we would have observed if sampling error were driving variability.
 1027 We discuss our interpretation of I^2 further in the methods, but in short, it is a useful metric for
 1028 comparison to values from published meta-analyses and provides a plausible value for how much
 1029 heterogeneity could arise in a normal meta-analysis with similar sample sizes due to analytical

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1030 variability alone. In our study, total I^2 for the blue tit Zr estimates was extremely large, at 97.73%, as
1031 was the Eucalyptus estimate (98.59% Table 2).

1032 Although the overall I^2 values were similar for both Eucalyptus and blue tit analyses, the relative
1033 composition of that heterogeneity differed. For both datasets, the majority of heterogeneity in Zr
1034 was driven by differences among effects as opposed to differences among teams, though this was
1035 more prominent for the Eucalyptus dataset, where nearly all of the total heterogeneity was driven by
1036 differences among effects (91.71%) as opposed to differences among teams (6.88%) (Table 2).

1037 Out-of-sample predictions (y_i)

1038 We observed substantial heterogeneity among out-of-sample estimates, but the pattern differed
1039 somewhat from the Zr values (Table 3). Among the blue tit predictions, I^2 ranged from medium-high
1040 for the y_{25} scenario (68.36) to low (27.02) for the y_{75} scenario. Among the Eucalyptus predictions, I^2
1041 values were uniformly high (>82%). For both datasets, most of the existing heterogeneity among
1042 predicted values was attributable to among-team differences, with the exception of the y_{50} analysis
1043 of the Eucalyptus dataset. We are limited in our interpretation of τ^2 for these estimates because,
1044 unlike for the Zr estimates, we have no benchmark for comparison with other meta-analyses.

1045

1046 Table 3: Heterogeneity among the out-of-sample predictions y_i for both blue tit and *Eucalyptus*
1047 datasets. τ_{Team}^2 is the absolute heterogeneity for the random effect *Team*, $\tau_{EffectID}^2$ is the absolute
1048 heterogeneity for the random effect *EffectID*, nested under *Team*, and τ_{total}^2 is the total absolute
1049 heterogeneity. I^2_{Total} is the proportional heterogeneity; the proportion of the variance among
1050 effects not attributable to sampling error, I^2_{Team} is the subset of the proportional heterogeneity due
1051 to differences among Teams and $I^2_{Team, EffectID}$ is subset of the proportional heterogeneity
1052 attributable to among-*EffectID* differences.

Dataset	Scenario	N. Obs	τ^2_{Total}	τ^2_{Team}	$\tau^2_{EffectID}$	I^2_{Total}	I^2_{Team}	$I^2_{Team, EffectID}$
blue tit	y_{25}	62	0.14	0.11	0.03	68.36%	51.82%	16.54%
	y_{50}	59	0.07	0.06	0.01	50.37%	45.66%	4.71%
	y_{75}	62	0.02	0.02	0.00	27.02%	25.57%	1.45%

<i>Eucalyptus</i>	γ_{25}	22	3.05	1.95	1.10	88.76%	56.76%	32.00%
	γ_{50}	24	1.61	0.53	1.08	83.26%	27.52%	55.73%
	γ_{75}	24	1.69	1.41	0.28	79.76%	66.52%	13.25%

1053

1054 [Post-hoc Analysis: Exploring outlier characteristics and the effect of outlier removal on](#)
1055 [heterogeneity](#)

1056 Effect Sizes (Z_r)

1057 The outlier *Eucalyptus* Z_r values were striking and merited special examination. The three negative
1058 outliers had very low sample sizes were based on either small subsets of the dataset or, in one case,
1059 extreme aggregation of data. The outliers associated with small subsets had sample sizes ($n= 117, 90$)
1060 that were less than half of the total possible sample size of 351. The case of extreme aggregation
1061 involved averaging all values within each of the 18 sites in the dataset.

1062 Surprisingly, both the largest and smallest effect sizes in the blue tit analyses (Figure 1a) come from
1063 the same analyst (anonymous ID: Adelong), with identical models in terms of the explanatory
1064 variable structure, but with different response variables. However, the radical change in effect was
1065 primarily due to collinearity with covariates. The primary predictor variable (brood count after
1066 manipulation) was accompanied by several collinear variables, including the highly collinear
1067 (correlation of approximately 0.9 ([Supplementary Figure D.2](#))) covariate (brood size at day 14) in both
1068 analyses. In the analysis of nestling weight, brood count after manipulation showed a strong positive
1069 partial correlation with weight after controlling for brood count at day 14 and treatment category
1070 (increased, decreased, unmanipulated). In that same analysis, the most collinear covariate (the day
1071 14 count) had a negative partial correlation with weight. In the analysis with tarsus length as the
1072 response variable, these partial correlations were almost identical in absolute magnitude, but
1073 reversed in sign and so brood count after manipulation was now the collinear predictor with the
1074 negative relationship. The two models were therefore very similar, but the two collinear predictors
1075 simply switched roles, presumably because a subtle difference in the distribution of weight and
1076 tarsus length data.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1077 When we dropped the *Eucalyptus* outliers, I^2 decreased from high (98.59%), using Higgins' [36]
1078 suggested benchmark, to between moderate and high (66.19%, Table 2). However, more notably, τ^2
1079 dropped from 0.27 to 0.01, indicating that, once outliers were excluded, the observed variation in
1080 effects was similar to what we would expect if sampling error were driving the differences among
1081 effects (since τ^2 is the variance in addition to that driven by sampling error). The interpretation of this
1082 value of τ^2 in the context of our many-analyst study is somewhat different than a typical meta-
1083 analysis, however, since in our study (especially for *Eucalyptus*, where most analyses used almost
1084 exactly the same data points), there is almost no role for sampling error in driving the observed
1085 differences among the estimates. Thus, rather than concluding that the variability we observed
1086 among estimates (after removing outliers) was due only to sampling error (because τ^2 became small:
1087 10% of the median from 48), we instead conclude that the observed variability, which must be due to
1088 the divergent choices of analysts rather than sampling error, is approximately of the same magnitude
1089 as what we would have expected if, instead, sampling error, and not analytical heterogeneity, were at
1090 work. Presumably, if sampling error had actually also been at work, it would have acted as an
1091 additional source of variability and would have led total variability among estimates to be higher.
1092 With total variability higher and thus greater than expected due to sampling error alone, τ^2 would
1093 have been noticeably larger. Conversely, dropping outliers from the set of blue tit effects did not
1094 meaningfully reduce I^2 , and only modestly reduced τ^2 (Table 2). Thus, effects at the extremes of the
1095 distribution were much stronger contributors to total heterogeneity for effects from analyses of the
1096 *Eucalyptus* than for the blue tit dataset.

1097 Table 4: Estimated mean value of the standardised correlation coefficient, Z_r , along with its standard
1098 error and 95% confidence intervals. We re-computed the meta-analysis for different post-hoc subsets
1099 of the data: All eligible effects, removal of effects from analysis teams that received at least one peer
1100 rating of 'deeply flawed and unpublishable', removal of any effects from analysis teams that received
1101 at least one peer rating of either 'deeply flawed and unpublishable' or 'publishable with major

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1102 revisions', inclusion of only effects from analysis teams that included at least one member who rated
 1103 themselves as "highly proficient" or "expert" at conducting statistical analyses in their research area..

Dataset	$\hat{\mu}$	SE[$\hat{\mu}$]	95% CI	statistic	p-value
All Analyses					
blue tit	-0.35	0.03	[-0.41,-0.28]	-10.49	<0.001
<i>Eucalyptus</i>	-0.09	0.06	[-0.22,0.03]	-1.47	0.14
Analyses receiving at least one 'Unpublishable' rating removed					
blue tit	-0.36	0.03	[-0.43,-0.29]	-10.49	<0.001
<i>Eucalyptus</i>	-0.02	0.02	[-0.07,0.02]	-1.15	0.3
Analyses receiving at least one 'Unpublishable' and or 'Major Revisions' rating removed					
blue tit	-0.37	0.07	[-0.51,-0.23]	-5.34	<0.001
<i>Eucalyptus</i>	-0.04	0.05	[-0.15,0.07]	-0.77	0.4
All analyses - outliers removed					
blue tit	-0.35	0.03	[-0.42,-0.29]	-10.95	<0.001
<i>Eucalyptus</i>	-0.03	0.01	[-0.06,0.00]	-2.23	0.026
Analyses from teams with highly proficient or expert data analysts					
blue tit	-0.35	0.04	[-0.44,-0.27]	-8.31	<0.001
<i>Eucalyptus</i>	-0.17	0.13	[-0.43,0.10]	-1.24	0.2

1104

1105 Out-of-sample predictions (y_i)

1106 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of
 1107 eligible effects was smaller and the pattern of outliers differed.

1108 Post-hoc analysis: Exploring the effect of removing analyses with poor peer ratings on
 1109 heterogeneity

1110 Effect Size (Z_r)

1111 Removing poorly rated analyses had limited impact on the meta-analytic means ([Supplementary](#)
 1112 [Figure B.3](#)). For the *Eucalyptus* dataset, the meta-analytic mean shifted from -0.09 to -0.02 when
 1113 effects from analyses rated as unpublishable were removed, and to -0.04 when effects from analyses
 1114 rated, at least once, as unpublishable or requiring major revisions were removed. Further, the
 1115 confidence intervals for all of these means overlapped each of the other means (Table 4). We saw
 1116 similar patterns for the blue tit dataset, with only small shifts in the meta-analytic mean, and
 1117 confidence intervals of all three means overlapping each other mean (Table 4). Refitting the meta-

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1118 analysis with a fixed effect for categorical ratings also showed no indication of differences in group
1119 meta-analytic means due to peer ratings ([Supplementary Figure B.1](#)).

1120 For the blue tit dataset, removing poorly-rated analyses led to only negligible changes in I^2 Total and
1121 relatively minor impacts on τ^2 . However, for the *Eucalyptus* dataset, removing poorly-rated analyses
1122 led to notable reductions in I^2 Total and substantial reductions in τ^2 . When including all analyses, the
1123 *Eucalyptus* I^2 Total was 98.59% and τ^2 was 0.27, but eliminating analyses with ratings of
1124 “unpublishable” reduced I^2 Total to 79.74% and τ^2 to 0.01, and removing also those analyses “needing
1125 major revisions” left I^2 Total at 88.91% and τ^2 at 0.03 (Table 2). Additionally, the allocations of I^2 to the
1126 team versus individual effect were altered for both blue tit and *Eucalyptus* meta-analyses by
1127 removing poorly rated analyses, but in different ways. For blue tit meta-analysis, between a third and
1128 two-thirds of the total I^2 was attributable to among-team variance in most analyses until both
1129 analyses rated “unpublishable” and analyses rated in need of “major revision” were eliminated, in
1130 which case almost all remaining heterogeneity was attributable to among-effect differences. In
1131 contrast, for *Eucalyptus* meta-analysis, the among-team component of I^2 was less than third until
1132 both analyses rated “unpublishable” and analyses rated in need of “major revision” were eliminated,
1133 in which case almost 90% of heterogeneity was attributable to differences among teams.

1134 [Out-of-sample predictions \(\$y_i\$ \)](#)

1135 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of
1136 eligible effects was smaller and our ability to interpret heterogeneity values for these analyses was
1137 limited.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1138 [Post-hoc analysis: Exploring the effect of including only analyses conducted by analysis](#)

1139 [teams with at least one member self-rated as “highly proficient” or “expert” in](#)

1140 [conducting statistical analyses in their research area](#)

1141 [Effect sizes \(\$Z_r\$ \)](#)

1142 Including only analyses conducted by teams that contained at least one member who rated

1143 themselves as “highly proficient” or “expert” in conducting the relevant statistical methods had

1144 negligible impacts on the meta-analytic means (Table 4), the distribution of Z_r effects

1145 ([Supplementary Figure B.4](#)), or heterogeneity estimates (Table 2), which remained extremely high.

1146 [Out-of-sample predictions \(\$y_i\$ \)](#)

1147 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of

1148 eligible effects was smaller.

1149 [Post-hoc analysis: Exploring the effect of excluding estimates of \$Z_r\$ in which we had](#)

1150 [reduced confidence](#)

1151 As described in our addendum to the methods, we identified a subset of estimates of Z_r in which we

1152 had less confidence because of features of the submitted degrees of freedom. Excluding these effects

1153 in which we had lower confidence had minimal impact on the meta-analytic mean and the estimates

1154 of total I^2 and τ^2 for both blue tit and Eucalyptus meta-analyses, regardless of whether outliers were

1155 also excluded ([Supplementary Table B.1](#)).

1156 [Explaining Variation in Deviation Scores](#)

1157 None of the pre-registered predictors explained substantial variation in deviation among submitted

1158 statistical effects from the meta-analytic mean (Table 5, Table 6). Note that the extremely high

1159 $R^2_{\text{Conditional}}$ values from the analyses of continuous peer ratings as predictors of deviation scores are a

1160 function of the random effects, not the fixed effect of interest. These high values of $R^2_{\text{Conditional}}$ result

1161 from the fact that each effect size was included in the analysis multiple times, to allow comparison

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1162 with ratings from the multiple peer reviewers who reviewed each analysis, and therefore when we
 1163 included Effect ID as a random effect, the observations within each random effect category were
 1164 identical.

1165 Table 5: Summary metrics for registered models seeking to explain deviation (Box-Cox transformed
 1166 absolute deviation scores) from the mean Zr as a function of Sorensen’s Index, categorical peer
 1167 ratings, and continuous peer ratings for blue tit and Eucalyptus analyses, and as a function of the
 1168 presence or absence of random effects (in the analyst’s models) for Eucalyptus analyses. We report
 1169 coefficient of determination, R^2 , for our models including only fixed effects as predictors of deviation,
 1170 and we report $R^2_{\text{Conditional}}$, R^2_{Marginal} and the intra-class correlation (ICC) from our models that included
 1171 both fixed and random effects. For all our models, we calculated the residual standard deviation σ
 1172 and root mean squared error (RMSE).

Dataset	R^2	$R^2_{\text{Conditional}}$	R^2_{Marginal}	ICC	σ	RMSE	N. Obs.
Deviation explained by categorical ratings							
blue tit		0.0903	0.0067	0.0842	6.52e-01	6.32e-01	473
<i>Eucalyptus</i>		0.1319	0.0124	0.1209	1.06e+00	1.02e+00	346
Deviation explained by continuous ratings							
blue tit		1.0000	2.00e-26	1.0000	1.63e-05	1.56e-12	473
<i>Eucalyptus</i>		0.9998	6.57e-30	0.9998	7.93e-03	7.09e-14	346
Deviation explained by Sorensen's index							
blue tit	0.0011				0.681	0.676	124
<i>Eucalyptus</i>	0.0005				1.14	1.120	72
Deviation explained by inclusion of random effects							
blue tit	0.0268				0.658	0.653	131
<i>Eucalyptus</i>	8.67e-08				1.12	1.100	79

1173

1174 Table 6: Parameter estimates from models of Box-Cox transformed deviation scores as a function of
 1175 continuous and categorical peer ratings, Sorensen scores, and the inclusion of random effects.
 1176 Standard Errors (SE), 95% confidence intervals (95%CI) are reported for all estimates, while t values,
 1177 degrees of freedom and p-values are presented for fixed-effects. Note that positive parameter
 1178 estimates mean that as the predictor variable increases, so does the absolute value of the deviation
 1179 from the meta-analytic mean.

Dataset	Parameter	Effect	Coeff.	SE	95% CI	t	df	p-value
Deviation explained by inclusion of random effects								

<i>Eucalyptus</i>	Intercept		2.53	0.27	-3.06,-1.99	-9.31	77	<0.001
	Random effects		0.00	0.31	-0.60, 0.60	0.00	77	>0.9
Deviation explained by mean Sorensen's index								
<i>Eucalyptus</i>	Intercept		-2.75	1.07	-4.85,-0.65	-2.57	70	0.010
	Sorensen Index		0.29	1.54	-2.74, 3.32	0.19	70	0.9
blue tit	Intercept		-1.56	0.38	-2.30,-0.82	-4.12	122	<0.001
	Mean Sorensen Index		0.23	0.63	-1.00, 1.46	0.37	122	0.7
Deviation explained by continuous ratings								
<i>Eucalyptus</i>	Intercept	Fixed	-2.52	0.06	-2.63,-2.40	-42.58	342	<0.001
	Continuous Rating	Fixed	6e-17	2e-10	-4e-10, 4e-10	-3e-07	342	>0.9
	SD (Intercept)	Random (EffectID)	0.53	0.04	0.45, 0.62			
	SD (Observations)	Random (Residual)	0.01	3e-04	0.01,0.01			
blue tit	Intercept	Fixed	-1.41	0.03	-1.47,-1.35	-46.54	469	<0.001
	Continuous Rating	Fixed	-3e-15	1e-09	-2e-09,2e-09	-2e-06	469	>0.9
	SD (Intercept)	Random (EffectID)	0.34	0.02	0.30, 0.39			
	SD (Observations)	Random (Residual)	2e-05	6e-07	2e-05,2e-05			
Deviation explained by categorical ratings								
<i>Eucalyptus</i>	Intercept	Fixed	-2.66	0.27	-3.18,-2.13	-9.97	340	<0.001
	Publishable with major revisions	Fixed	0.29	0.29	-0.27, 0.85	1.02	340	0.3
	Publishable with minor revisions	Fixed	0.01	0.28	-0.54, 0.56	0.04	340	>0.9
	Publishable as is	Fixed	0.05	0.31	-0.55, 0.66	0.17	340	0.9
	SD (Intercept)	Random (ReviewerID)	0.39	0.09	0.25, 0.61			
	SD (Observations)	Random (Residual)	1.06	0.04	0.98,1.15			
blue tit	Intercept	Fixed	-1.21	0.15	-1.50,-0.93	-8.29	467	<0.001
	Publishable with major revisions	Fixed	-0.23	0.15	-0.53, 0.07	-1.50	467	0.13
	Publishable with minor revisions	Fixed	-0.23	0.15	-0.53, 0.07	-1.52	467	0.13
	Publishable as is	Fixed	-0.15	0.17	-0.48, 0.18	-0.89	467	0.4
	SD (Intercept)	Random (ReviewerID)	0.20	0.05	0.13, 0.31			
	SD (Observations)	Random (Residual)	0.65	0.02	0.61,0.7			

1180

1181 [Deviation Scores as explained by reviewer ratings](#)

1182 [Effect Sizes \(Zr\)](#)

1183 We obtained reviews from 128 reviewers who reviewed analyses for a mean of 3.27 (range 1 - 11)

1184 analysis teams. Analyses of the blue tit dataset received a total of 240 reviews, each was reviewed by

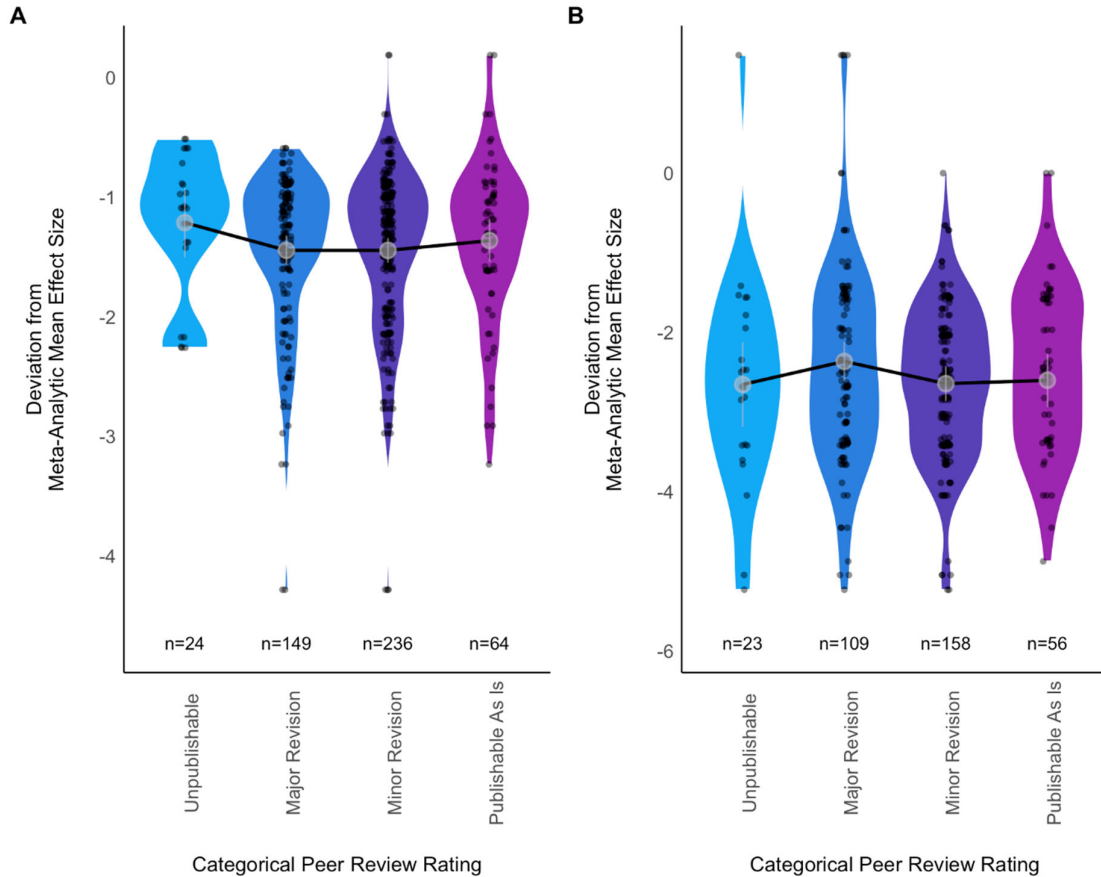
we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1185 a mean of 3.87 (SD 0.71, range 3-5) reviewers. Analyses of the *Eucalyptus* dataset received a total of
1186 178 reviews, each was reviewed by a mean of 4.24 (SD 0.79, range 3-6) reviewers. We tested for
1187 inter-rater reliability to examine how similarly reviewers reviewed each analysis and found
1188 approximately no agreement among reviewers. When considering continuous ratings, IRR was 0.01,
1189 and for categorical ratings, IRR was -0.14.

1190 Many of the models of deviance as a function of peer ratings faced issues of failure to converge or
1191 singularity due to sparse design matrices with our pre-registered random effects (*EffectID* and
1192 *ReviewerID*) (see [Supplementary Table C.1](#)). These issues persisted after increasing the tolerance and
1193 changing the optimizer. For both *Eucalyptus* and blue tit datasets, models with continuous ratings as
1194 a predictor were singular when both pre-registered random effects were included.

1195 When using only categorical ratings as predictors, models converged only when specifying reviewer
1196 ID as a random effect. That model had a R^2_C of 0.09 and a R^2_M of 0.01. The model using the
1197 continuous ratings converged for both random effects (in isolation), but not both. We present results
1198 for the model using study ID as a random effect because we expected it would be a more important
1199 driver of variation in deviation scores. That model had a R^2_C of 1 and a R^2_M of 0.01 for the blue tit
1200 dataset and a R^2_C of 1 and a R^2_M of 0.01 for the *Eucalyptus* dataset. Neither continuous or categorical
1201 reviewer ratings of the analyses meaningfully predicted deviance from the meta-analytic mean
1202 (Table 6, Figure 3). We re-ran the multi-level meta-analysis with a fixed-effect for the categorical
1203 publishability ratings and found no difference in mean standardised effect sizes among publishability
1204 ratings ([Supplementary Figure B.1](#)).

1205



1206

1207 Figure 3: Violin plot of Box-Cox transformed deviation from meta-analytic mean as a function of
1208 categorical peer rating for a) blue tit and b) *Eucalyptus*. Grey points for each rating group denote
1209 model-estimated marginal mean deviation, and error bars denote 95% CI of the estimate.

1210

1211 Out-of-sample predictions (y_i)

1212 Some models of the influence of reviewer ratings on out-of-sample predictions (y_i) had issues with

1213 convergence and singularity of fit (see [Supplementary Table C.2](#)) and those models that converged

1214 and were not singular showed no strong relationship ([Supplementary Figures C.2, Figure C.3](#)), as with

1215 the Zr analyses.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

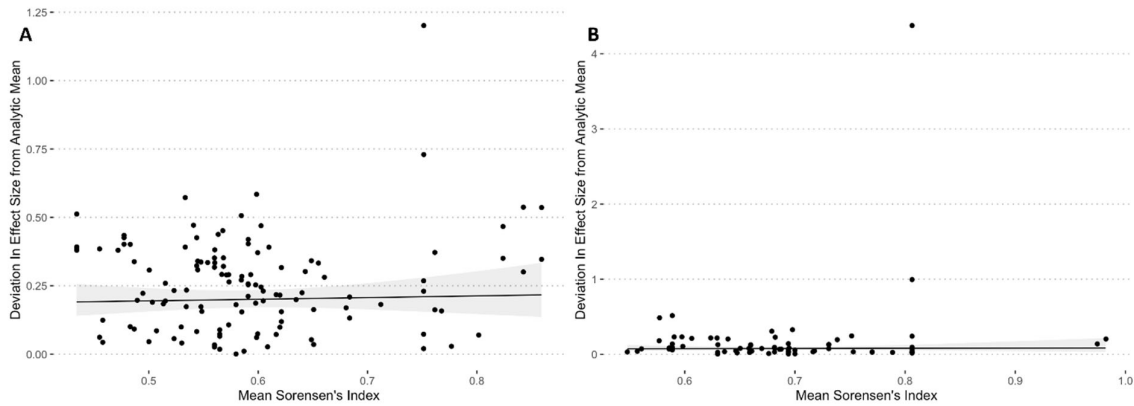
1216 Deviation scores as explained by the distinctiveness of variables in each analysis

1217 Effect Size (Z_r)

1218 We employed Sorensen's index to calculate the distinctiveness of the set of predictor variables used
1219 in each model (Figure 5). The mean Sorensen's score for blue tit analyses was 0.69 (range 0.55-0.98),
1220 and for *Eucalyptus* analyses was 0.59 (range 0.43-0.86).

1221 We found no meaningful relationship between distinctiveness of variables selected and deviation
1222 from the meta-analytic mean (Table 6, Figure 5) for either blue tit (mean 0.23, 95% CI -1,1.46) or
1223 *Eucalyptus* effects (mean 0.29, 95% CI -2.74,3.32).

1224



1225

1226 Figure 4: Fitted model of the Box-Cox-transformed deviation score (deviation in effect size from
1227 meta-analytic mean) as a function of the mean Sorensen's index showing distinctiveness of the set of
1228 predictor variables for a) blue tit, and b) *Eucalyptus*. Grey ribbons on predicted values are 95% CI's.

1229

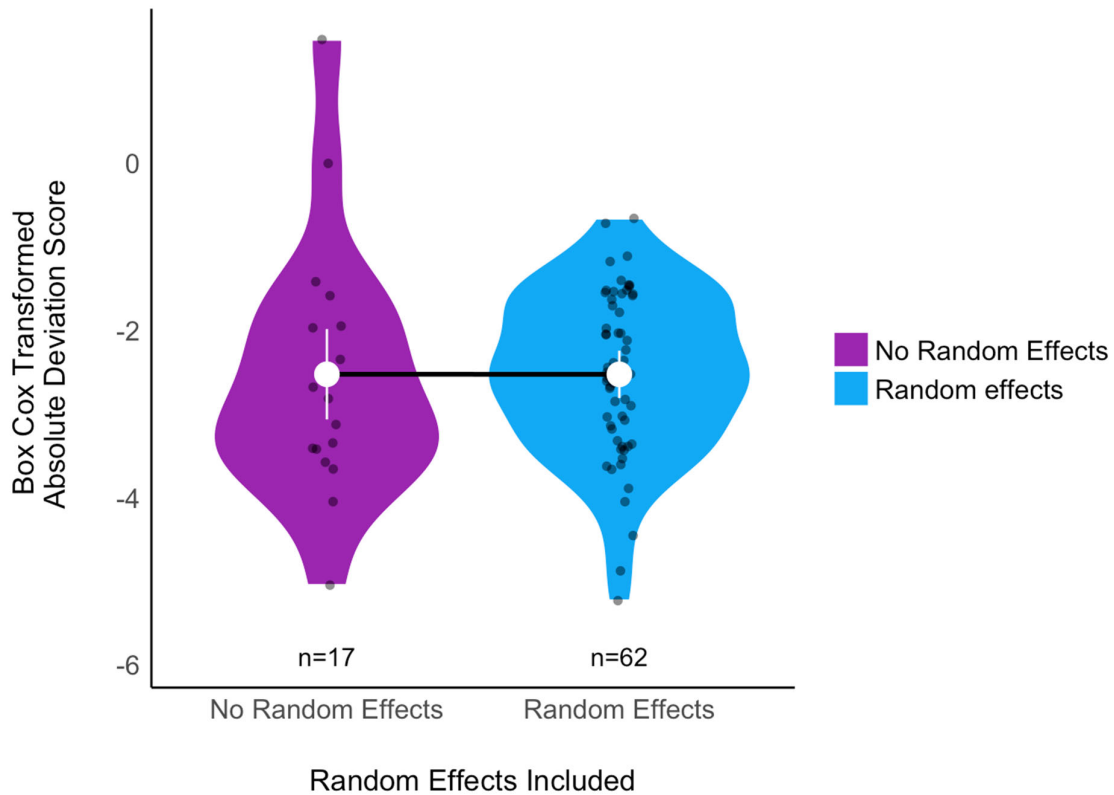
1230 Out-of-sample predictions

1231 As with the Z_r estimates, we did not observe any convincing relationships between deviation scores
1232 of out-of-sample predictions and Sorensen's index values. Please see [Supplementary Material C.4.2](#).

1233 Deviation scores as explained by the inclusion of random effects

1234 Effect Size (Z_r)

1235 There were only three blue tit analyses that did not include random effects, which is below the pre-
1236 registered threshold for fitting a model of the Box-Cox transformed deviation from the meta-analytic
1237 mean as a function of whether the analysis included random-effects. However, 17 Eucalyptus
1238 analyses included only fixed effects, which crossed our pre-registered threshold. Consequently, we
1239 performed this analysis for the Eucalyptus dataset only. There was no relationship between random-
1240 effect inclusion and deviation from meta-analytic mean among the Eucalyptus analyses (Table 6,
1241 Figure 5).



1242

1243 Figure 5: Violin plot of mean Box-Cox transformed deviation from meta-analytic mean as a function
1244 of random-effects inclusion in *Eucalyptus* analyses. '1' indicates random-effects were included in
1245 analyst's model, while 0 indicates no random-effects were included. White points for each group of

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1246 analyses denote model-estimated marginal mean deviation, and error bars denote 95% CI of the
1247 estimate.

1248 [Out-of-sample predictions](#)

1249 As with the Z_r estimates, we did not examine the possibility of a relationship between the inclusion
1250 of random effects and the deviation scores of the blue tit out-of-sample predictions. When we
1251 examined the possibility of this relationship for the Eucalyptus effects, we found consistent evidence
1252 of somewhat higher Box-Cox-transformed deviation values for models including a random effect,
1253 meaning the models including random effects averaged slightly higher deviation from the meta-
1254 analytic means ([Supplementary Figure C.5](#)).

1255 [Multivariate Analysis Effect size \(\$Z_r\$ \) and out-of-sample predictions \(\$y_i\$ \)](#)

1256 Like the univariate models, the multivariate models did a poor job of explaining deviations from the
1257 meta-analytic mean. Because we pre-registered a multivariate model that contained collinear
1258 predictors that produce results which are not readily interpretable, we present these models in the
1259 supplement. We also had difficulty with convergence and singularity for multivariate models of out-
1260 of-sample (y_i) result, and had to adjust which random effects we included ([Supplementary Table C.7](#)).
1261 However, no multivariate analyses of Eucalyptus out-of-sample results avoided problems of
1262 convergence or singularity, no matter which random effects we included ([Supplementary Table C.7](#)).
1263 We therefore present no multivariate Eucalyptus y_i models. We present parameter estimates from
1264 multivariate Z_r models for both datasets ([Supplementary Tables C.5, C.6](#)) and from y_i models from
1265 the blue tit dataset ([Supplementary Tables C.8, C.9](#)). We include interpretation of the results from
1266 these models in the supplement, but the results do not change the interpretations we present above
1267 based on the univariate analyses.

1268 [Discussion](#)

1269 When a large pool of ecologists and evolutionary biologists analyzed the same two datasets to
1270 answer the corresponding two research questions, they produced substantially heterogeneous sets

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1271 of answers. Although the variability in analytical outcomes was high for both datasets, the patterns
1272 of this variability differed distinctly between them. For the blue tit dataset, there was nearly
1273 continuous variability across a wide range of Zr values. In contrast, for the Eucalyptus dataset, there
1274 was less variability across most of the range, but more striking outliers at the tails. Among out-of-
1275 sample predictions, there was again almost continuous variation across a wide range (2 SD) among
1276 blue tit estimates. For Eucalyptus, out-of-sample predictions were also notably variable, with about
1277 half the predicted stem count values at <2 but the other half being much larger, and ranging to
1278 nearly 40 stems per 15 m x 15 m plot. We investigated several hypotheses for drivers of this
1279 variability within datasets, but found little support for any of these. Most notably, even when we
1280 excluded analyses that had received one or more poor peer reviews, the heterogeneity in results
1281 largely persisted. Regardless of what drives the variability, the existence of such dramatically
1282 heterogeneous results when ecologists and evolutionary biologists seek to answer the same
1283 questions with the same data should trigger conversations about how ecologists and evolutionary
1284 biologists analyze data and interpret the results of their own analyses and those of others in the
1285 literature [e.g., [11](#), [20](#), [49](#), [50](#)].

1286 Our observation of substantial heterogeneity due to analytical decisions is consistent with a growing
1287 body of work, much of it from the quantitative social sciences [e.g., [11](#), [17–21](#)]. In all of these
1288 studies, when volunteers from the discipline analyzed the same data, they produced a worryingly
1289 diverse set of answers to a pre-set question. This diversity always included a wide range of effect
1290 sizes, and in most cases, even involved effects in opposite directions. Thus, our result should not be
1291 viewed as an anomalous outcome from two particular datasets, but instead as evidence from
1292 additional disciplines regarding the heterogeneity that can emerge from analyses of complex
1293 datasets to answer questions in probabilistic science. Not only is our major observation consistent
1294 with other studies, it is, itself, robust because it derived primarily from simple forest plots that we
1295 produced based on a small set of decisions that were mostly registered before data gathering and
1296 which conform to widely accepted meta-analytic practices.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1297 Unlike the strong pattern we observed in the forest plots, our other analyses, both registered and
1298 post hoc, produced either inconsistent patterns, weak patterns, or the absence of patterns. Our
1299 registered analyses found that deviations from the meta-analytic mean by individual effect sizes (Z_r)
1300 or the predicted values of the dependent variable (y_i) were poorly explained by our hypothesized
1301 predictors: peer rating of each analysis team's method section, a measurement of the distinctiveness
1302 of the set of predictor variables included in each analysis, or whether the model included random
1303 effects. However, in our post hoc analyses, we found that dropping analyses identified as
1304 unpublishable or in need of major revision by at least one reviewer modestly reduced the observed
1305 heterogeneity among the Z_r outcomes, but only for *Eucalyptus* analyses, apparently because this led
1306 to the dropping of the major outlier. This limited role for peer review in explaining the variability in
1307 our results should be interpreted cautiously because the inter-rater reliability among peer reviewers
1308 was extremely low, and at least some analyses that appeared flawed to us were not marked as
1309 flawed by reviewers. However, the hypothesis that poor quality analyses drove the heterogeneity we
1310 observed was also contradicted by our observation that analysts' self-declared statistical expertise
1311 appeared unrelated to heterogeneity. When we retained only analyses from teams including at least
1312 one member with high self-declared levels of expertise, heterogeneity among effect sizes remained
1313 high. Thus, our results suggest lack of statistical expertise is not the primary factor responsible for
1314 the heterogeneity we observed, although further work is merited before rejecting a role for
1315 statistical expertise. Not surprisingly, simply dropping outlier values of Z_r for *Eucalyptus* analyses,
1316 which had more extreme outliers, led to less observable heterogeneity in the forest plots, and also
1317 reductions in our quantitative measures of heterogeneity. We did not observe a similar effect in the
1318 blue tit dataset because that dataset had outliers that were much less extreme and instead had more
1319 variability across the core of the distribution.

1320 Our major observations raise two broad questions; why was the variability among results so high,
1321 and why did the pattern of variability differ between our two datasets. One important and plausible
1322 answer to the first question is that much of the heterogeneity derives from the lack of a precise

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1323 relationship between the two biological research questions we posed and the data we provided. This
1324 lack of a precise relationship between data and question creates many opportunities for different
1325 model specifications, and so may inevitably lead to varied analytical outcomes [50]. However, we
1326 believe that the research questions we posed are consistent with the kinds of research question that
1327 ecologists and evolutionary biologists typically work from. When designing the two biological
1328 research questions, we deliberately sought to represent the level of specificity we typically see in
1329 these disciplines. This level of specificity is evident when we look at the research questions posed by
1330 some recent meta-analyses in these fields:

1331 “how [does] urbanisation impact mean phenotypic values and phenotypic variation ... [in] paired
1332 urban and non-urban comparisons of avian life-history traits” [51]

1333 “[what are] the effects of ocean acidification on the crustacean exoskeleton, assessing both
1334 exoskeletal ion content (calcium and magnesium) and functional properties (biomechanical
1335 resistance and cuticle thickness)” [52]

1336 “[what is] the extent to which restoration affects both the mean and variability of biodiversity
1337 outcomes ... [in] terrestrial restoration” [53]

1338 “[does] drought stress [have] a negative, positive, or null effect on aphid fitness” [54]

1339 “[what is] the influence of nitrogen-fixing trees on soil nitrous oxide emissions” [55]

1340 There is not a single precise answer to any of these questions, nor to the questions we posed to
1341 analysts in our study. And this lack of single clear answers will obviously continue to cause
1342 uncertainty since ecologists and evolutionary biologists conceive of the different answers from the
1343 different statistical models as all being answers to the same general question. A possible response
1344 would be a call to avoid these general questions in favor of much more precise alternatives [50].

1345 However, the research community rewards researchers who pose broad questions [56], and so
1346 researchers are unlikely to narrow their scope without a change in incentives. Further, we suspect
1347 that even if individual studies specified narrow research questions, other scientists would group

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1348 these more narrow questions into broader categories, for instance in meta-analyses, because it is
1349 these broader and more general questions that often interest the research community.

1350 Although variability in statistical outcomes among analysts may be inevitable, our results raise
1351 questions about why this variability differed between our two datasets. We are particularly
1352 interested in the differences in the distribution of Zr since the distributions of out-of-sample
1353 predictions were on different scales for the two datasets, thus limiting the value of comparisons. The
1354 forest plots of Zr from our two datasets showed distinct patterns, and these differences are
1355 consistent with several alternative hypotheses. The results submitted by analysts of the Eucalyptus
1356 dataset showed a small average (close to zero) with most estimates also close to zero (± 0.2), though
1357 about a third far enough above or below zero to cross the traditional threshold of statistical
1358 significance. There were a small number of striking outliers that were very far from zero. In contrast,
1359 the results submitted by analysts of the blue tit dataset showed an average much further from zero (-
1360 0.35) and a much greater spread in the core distribution of estimates across the range of Zr values (\pm
1361 0.5 from the mean), with few modest outliers. So, why was there more spread in effect sizes (across
1362 the estimates that are not outliers) in the blue tit analyses relative to the Eucalyptus analyses?

1363 One possible explanation for the lower heterogeneity among most Eucalyptus Zr effects is that weak
1364 relationships may limit the opportunities for heterogeneity in analytical outcome. Some evidence for
1365 this idea comes from two sets of “many labs” studies in psychology [4, 57]. In these studies, many
1366 independent lab groups each replicated a large set of studies, including, for each study, the
1367 experiment, data collection, and statistical analyses. These studies showed that, when the meta-
1368 analytic mean across the replications from different labs was small, there was much less
1369 heterogeneity among the outcomes than when the mean effect sizes were large [4, 57]. Of course, a
1370 weak average effect size would not prevent divergent effects in all circumstances. As we saw with the
1371 Eucalyptus analyses, taking a radically smaller subset of the data can lead to dramatically divergent
1372 effect sizes even when the mean with the full dataset is close to zero.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1373 Our observation that dramatic sub-setting in the Eucalyptus dataset was associated with
1374 correspondingly dramatic divergence in effect sizes leads us towards another hypothesis to explain
1375 the differences in heterogeneity between the Eucalyptus and blue tit analysis sets. It may be that
1376 when analysts often divide a dataset into subsets, the result will be greater heterogeneity in
1377 analytical outcome for that dataset. Although we saw sub-setting associated with dramatic outliers in
1378 the Eucalyptus dataset, nearly all other analyses of Eucalyptus data used very close to the same set
1379 of 351 samples, and as we saw, these effects did not vary substantially. However, analysts often
1380 analyzed only a subset of the blue tit data, and as we observed, sample sizes were much more
1381 variable among blue tit effects, and the effects themselves were also much more variable. Important
1382 to note here is that subsets of data may differ from each other for biological reasons, but they may
1383 also differ due to sampling error. Sampling error is a function of sample size, and sub-samples are, by
1384 definition, smaller samples, and so more subject to variability in effects due to sampling error [58].
1385 Other features of datasets are also plausible candidates for driving heterogeneity in analytical
1386 outcomes, including features of covariates. In particular, relationships between covariates and the
1387 response variable as well as relationships between covariates and the primary independent variable
1388 (collinearity) can strongly influence the modeled relationship between the independent variable of
1389 interest and the dependent variable [59, 60]. Therefore, inclusion or exclusion of these covariates
1390 can drive heterogeneity in effect sizes (Z_r). Also, as we saw with the two most extreme Z_r values from
1391 the blue tit analyses, in multivariate models with collinear predictors, extreme effects can emerge
1392 when estimating partial correlation coefficients due to high collinearity, and conclusions can differ
1393 dramatically depending on which relationship receives the researcher's attention. Therefore,
1394 differences between datasets in the presence of strong and/or collinear covariates could influence
1395 the differences in heterogeneity in results among those datasets.
1396 Although it is too early in the many-analyst research program to conclude which analytical decisions
1397 or which features of datasets are the most important drivers of heterogeneity in analytical outcomes,
1398 we must still grapple with the possibility that analytical outcomes may vary substantially based on

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1399 the choices we make as analysts. If we assume that, at least sometimes, different analysts will
1400 produce dramatically different statistical outcomes, what should we do as ecologists and
1401 evolutionary biologists? We review some ideas below.

1402 The easiest path forward after learning about this analytical heterogeneity would be simply to
1403 continue with “business as usual”, where researchers report results from a small number of statistical
1404 models. A case could be made for this path based on our results. For instance, among the blue tit
1405 analyses, the precise values of the estimated Zr effects varied substantially, but the average effect
1406 was convincingly different from zero, and a majority of individual effects (84%) were in the same
1407 direction. Arguably, many ecologists and evolutionary biologists appear primarily interested in the
1408 direction of a given effect and the corresponding p-value[61], and so the variability we observed
1409 when analyzing the blue tit dataset may not worry these researchers. Similarly, most effects from the
1410 Eucalyptus analyses were relatively close to zero, and about two-thirds of these effects did not cross
1411 the traditional threshold of statistical significance. Therefore, a large proportion of people analyzing
1412 these data would conclude that there was no effect, and this is consistent with what we might
1413 conclude from the meta-analysis.

1414 However, we find the counter arguments to “business as usual” to be compelling. For blue tits, there
1415 were a substantial minority of calculated effects that would be interpreted by many biologists as
1416 indicating the absence of an effect (28%), and there were three traditionally ‘significant’ effects in
1417 the opposite direction to the average. The qualitative conclusions of analysts also reflected
1418 substantial variability, with fully half of teams drawing a conclusion distinct from the one we draw
1419 from the distribution as a whole. These teams with different conclusion were either uncertain about
1420 the negative relationship between competition and nestling growth, or they concluded that effects
1421 were mixed or absent. For the Eucalyptus analyses, this issue is more concerning. Around two-thirds
1422 of effects had confidence intervals overlapping zero, and of the third of analyses with confidence
1423 intervals excluding zero, almost half were positive, and the rest were negative. Accordingly, the

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1424 qualitative conclusions of the Eucalyptus teams were spread across the full range of possibilities. But
1425 even these problems are optimistic.

1426 A potentially larger argument against “business as usual” is that it provides the raw material for
1427 biasing the literature. When different model specifications readily lead to different results, analysts
1428 may be tempted to report the result that appears most interesting, or that is most consistent with
1429 expectation [7, 12]. There is growing evidence that researchers in ecology and evolutionary biology
1430 often report a biased subset of the results they produce [62, 63], and that this bias exaggerates the
1431 average size of effects in the published literature between 30 and 150% [9, 48]. The bias then
1432 accumulates in meta-analyses, apparently more than doubling the rate of conclusions of “statistical
1433 significance” in published meta-analyses above what would have been found in the absence of bias
1434 [48]. Thus, “business as usual” does not just create noisy results, it helps create systematically
1435 misleading results.

1436 Conclusions

1437 Overall, our results suggest to us that, where there is a diverse set of plausible analysis options, no
1438 single analysis should be considered a complete or reliable answer to a research question. We
1439 contend that ecologists and evolutionary biologists typically do multiple analyses (as many of our
1440 analyst teams did) however, some of these analyses don't make it into the published manuscript.
1441 Further, because of the evidence that ecologists and evolutionary biologists often present a biased
1442 subset of the analyses they conduct [48, 62, 63], we do not expect that even a collection of different
1443 effect sizes from different studies will accurately represent the true distribution of effects
1444 [48]. Therefore, we believe that an increased level of skepticism of the outcomes of single analyses,
1445 or even single meta-analyses, is warranted going forward. We recognize that some researchers have
1446 long maintained a healthy level of skepticism of individual studies as part of sound and practical
1447 scientific practice, and it is possible that those researchers will be neither surprised nor concerned by

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1448 our results. However, we doubt that many researchers are sufficiently aware of the potential
1449 problems of analytical flexibility to be appropriately skeptical.

1450 If we are skeptical of single analyses, the path forward may be multiple analyses per dataset. One
1451 possibility is the traditional robustness or sensitivity check [e.g., [64](#), [65](#)], in which the researcher
1452 presents several alternative versions of an analysis to demonstrate that the result is ‘robust’ [[66](#)].
1453 Unfortunately, robustness checks are at risk of the same potential biases of reporting found in other
1454 studies [[11](#)], especially given the relatively few models typically presented. However, these risks
1455 could be minimized by running more models and doing so with pre-registration or registered report.
1456 Another option is model averaging. Averages across models often perform well [e.g., [67](#)], and in
1457 some forms this may be a relatively simple solution. As most often practiced in ecology and
1458 evolutionary biology, model averaging involves first identifying a small suite of candidate models
1459 [see [13](#)], then using Akaike weights, based on Akaike’s Information Criterion (AIC), to calculate
1460 weighted averages for parameter estimates from those models. Again, the small number of models
1461 limits the exploration of specification space, but we can examine a larger number of models.

1462 However, there are more concerning limitations. The largest of these limitations is that averaging
1463 regression coefficients is problematic when models differ in interaction terms or collinear variables
1464 [[68](#)]. Additionally, weighting by AIC may often be inconsistent with our modelling goals. AIC balances
1465 the trade-off between model complexity and predictive ability, but penalizing models for complexity
1466 may not be suited for testing hypotheses about causation. So, AIC may often not offer the weight we
1467 want to use for an average, and we may also not wish to just generate an average. Instead, if we
1468 hope to understand an extensive universe of possible modelling outcomes, we could conduct a
1469 multiverse analysis, possibly with a specification [[10](#), [49](#)]. This could mean running hundreds or
1470 thousands of models (or more!) to examine the distribution of possible effects, and to see how
1471 different specification choices map onto these effects. However, there is a trade-off between
1472 efficiently exploring large areas of specification space and limiting the analyses to biologically
1473 plausible specifications. Instead of simply identifying modelling decisions and creating all possible

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1474 combinations for the multiverse, a researcher could attempt to prevent implausible combinations,
1475 though the more variables in the dataset, the more difficult this becomes. To make this easier, one
1476 could recruit many analysts to each designate one or a few plausible specifications, as with our
1477 ‘many analyst’ study [11]. An alternative that may be more labor intensive for the primary analyst,
1478 but which may lead to a more plausible set of models, could involve hypothesizing about causal
1479 pathways with DAGs [directed acyclic graphs; [69]] to constrain the model set. Devoting this effort to
1480 thoughtful multiverse specifications, possibly combined with pre-registration to hinder undisclosed
1481 data dredging, seems worthy of consideration.

1482 Although we have reviewed a variety of potential responses to the existence of variability in
1483 analytical outcomes, we certainly do not wish to imply that this is a comprehensive set of possible
1484 responses. Nor do we wish to imply that the opinions we have expressed about these options are
1485 correct. Determining how the disciplines of ecology and evolutionary biology should respond to
1486 knowledge of the variability in analytical outcome will benefit from the contribution and discussion
1487 of ideas from across these disciplines. We look forward to learning from these discussions and to
1488 seeing how these disciplines ultimately respond.

1489 [Declarations](#)

1490 [Ethics approval and consent to participate](#)

1491 We obtained permission to conduct this research from the Whitman College Institutional Review
1492 Board (IRB). As part of this permission, the IRB approved the consent form (<https://osf.io/xyp68/>)
1493 that all participants completed prior to joining the study.

1494 [Consent for publication](#)

1495 Not applicable

1496 [Availability of data and materials](#)

1497 All data cleaning and preparation for our analyses was conducted in R (R Core Team 2022) and is
1498 publicly archived at (<https://zenodo.org/doi/10.5281/zenodo.10046152>). Please see session info for

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

1499 the full list of packages and their citations used in our analysis pipeline. We built an R package,
1500 *ManyEcoEvo* to conduct the analyses described in this chapter. This same package can be used to
1501 reproduce our analyses or replicate the analyses described here using alternate datasets.

1502 [Competing interests](#)

1503 The authors declare that they have no competing interests

1504 [Funding](#)

1505 EG's contributions were supported by an Australian Government Research Training Program
1506 Scholarship, AIMOS top-up scholarship (2022) and Melbourne Centre of Data Science Doctoral
1507 Academy Fellowship (2021). FF's contributions were supported by ARC Future Fellowship
1508 FT150100297.

1509 [Author's contributions](#)

1510 HF, THP and FF conceptualized the project. PV provided raw data for Eucalyptus analyses and SG and
1511 THP provided raw data for blue tit analyses. DGH, HF and THP prepared surveys for collecting
1512 participating analysts and reviewer's data. EG, HF, THP, PV, SN and FF planned the analyses of the
1513 data provided by our analysts and reviewers, EG, HF, and THP curated the data, EG and HF wrote the
1514 software code to implement the analyses and prepare data visualisations. EG ensured that analyses
1515 were documented and reproducible. THP and HF administered the project, including coordinating
1516 with analysts and reviewers. FF provided funding for the project. THP, HF, and EG wrote the
1517 manuscript. Authors listed alphabetically contributed analyses of the primary datasets or reviews of
1518 analyses. All authors read and approved the final manuscript.

1519 [Acknowledgements](#)

1520 Not applicable

1521 [References](#)

1522 1. Arif S, MacNeil MA. Applying the structural causal model framework for observational causal
1523 inference in ecology. *Ecological Monographs*. 2023;93:e1554.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 1524 2. Atkinson J, Brudvig LA, Mallen-Cooper M, Nakagawa S, Moles AT, Bonser SP. Terrestrial ecosystem
1525 restoration increases biodiversity and reduces its variability, but not to reference levels: A global
1526 meta-analysis. *Ecology Letters*. 2022;25:1725–37.
- 1527 3. Auspurg K, Brüderl J. Has the credibility of the social sciences been credibly destroyed?
1528 Reanalyzing the “many analysts, one data set” project. *Socius*. 2021;7:23780231211024421.
- 1529 4. Schloerke B, Cook D, Larmarange J, Briatte F, Marbach M, Thoen E, et al. GGally: Extension to
1530 'ggplot2'. 2022.
- 1531 5. Baselga A, Orme D, Villeger S, De Bortoli J, Leprieur F, Logez M, et al. Package “betapart”. 2023.
- 1532 6. Bates D, Mächler M, Bolker B, Walker S. Fitting linear mixed-effects models using lme4. 2015.
1533 2015;67:48.
- 1534 7. Bolker B, Robinson D, Menne D, Gabry J, Buerkner P, Hau C, et al. Package “broom.mixed”. 2022.
- 1535 8. Borenstein M, Higgins JPT, Hedges L, Rothstein H. Basics of meta-analysis: I2 is not an absolute
1536 measure of heterogeneity. *Research Synthesis Methods*. 2017;8:5–18.
- 1537 9. Botvinik-Nezer R, Holzmeister F, Camerer CF, Dreber A, Huber J, Johannesson M, et al. Variability in
1538 the analysis of a single neuroimaging dataset by many teams. *Nature*. 2020;582:84–8.
- 1539 10. Breznau N, Rinke EM, Wuttke A, Nguyen HHV, Adem M, Adriaans J, et al. Observing many
1540 researchers using the same data and hypothesis reveals a hidden universe of uncertainty.
1541 *Proceedings of the National Academy of Sciences*. 2022;119:e2203150119.
- 1542 11. Briga M, Verhulst S. Mosaic metabolic ageing: Basal and standard metabolic rates age in opposite
1543 directions and independent of environmental quality, sex and life span in a passerine. *Functional
1544 Ecology*. 2021;35:1055–68.
- 1545 12. Burnham KP, Anderson DR. Model selection and multimodel inference: A practical information-
1546 theoretical approach. 2nd edition. Book. New York: Springer-Verlag; 2002.
- 1547 13. Cade BS. Model averaging and muddled multimodel inferences. *Ecology*. 2015;96:2370–82.
- 1548 14. Capilla-Lasheras P, Thompson MJ, Sánchez-Tójar A, Haddou Y, Branston CJ, Réale D, et al. A global
1549 meta-analysis reveals higher variation in breeding phenology in urban birds than in their non-urban
1550 neighbours. *Ecology Letters*. 2022;25:2552–70.
- 1551 15. Coretta S, Casillas JV, Roessig S, Franke M, Ahn B, Al-Hoorie AH, et al. Multidimensional signals
1552 and analytic flexibility: Estimating degrees of freedom in human-speech analyses. *Advances in
1553 Methods and Practices in Psychological Science*. 2023;6:25152459231162567.
- 1554 16. DeKogel CH. Long-term effects of brood size manipulation on morphological development and
1555 sex-specific mortality of offspring. *Journal of Animal Ecology*. 1997;66:167–78.
- 1556 17. Deressa T, Stern D, Vangronsveld J, Minx J, Lizin S, Malina R, et al. More than half of statistically
1557 significant research findings in the environmental sciences are actually not. *EcoEvoRxiv*. 2023.
1558 <https://doi.org/https://doi.org/10.32942/X24G6Z>.
- 1559 18. Dormann CF, Elith J, Bacher S, Buchmann C, Carl G, Carré G, et al. Collinearity: A review of
1560 methods to deal with it and a simulation study evaluating their performance. *Ecography*.
1561 2013;36:27–46.

- 1562 19. Fanelli D, Costas R, Ioannidis JPA. Meta-assessment of bias in science. *Proceedings of the National*
1563 *Academy of Sciences*. 2017;114:3714–9.
- 1564 20. Fanelli D, Ioannidis JPA. US studies may overestimate effect sizes in softer research. *Proceedings*
1565 *of the National Academy of Sciences*. 2013;110:15031–6.
- 1566 21. Fidler F, Burgman MA, Cumming G, Buttrose R, Thomason N. Impact of criticism of null-
1567 hypothesis significance testing on statistical reporting practices in conservation biology. *Conservation*
1568 *Biology*. 2006;20:1539–44.
- 1569 22. Fidler F, Chee YE, Wintle BC, Burgman MA, McCarthy MA, Gordon A. Metaresearch for evaluating
1570 reproducibility in ecology and evolution. *BioScience*. 2017;67:282–9.
- 1571 23. Forstmeier W, Wagenmakers E-J, Parker TH. Detecting and avoiding likely false-positive findings –
1572 a practical guide. *Biological Reviews*. 2017;92:1941–68.
- 1573 24. Fraser H, Parker T, Nakagawa S, Barnett A, Fidler F. Questionable research practices in ecology
1574 and evolution. *PLOS ONE*. 2018;13:e0200303.
- 1575 25. Gelman A, Weakliem D. Of beauty, sex, and power. *American Scientist*. 2009;97:310–6.
- 1576 26. Gelman A, Loken E. The garden of forking paths: Why multiple comparisons can be a problem,
1577 even when there is no “fishing expedition” or “p-hacking” and the research hypothesis was posited
1578 ahead of time. Department of Statistics, Columbia University. 2013.
- 1579 27. Grueber CE, Nakagawa S, Laws RJ, Jamieson IG. Multimodel inference in ecology and evolution:
1580 Challenges and solutions. *Journal of Evolutionary Biology*. 2011;24:699–711.
- 1581 28. Higgins JPT, Thompson SG, Deeks JJ, Altman DG. Measuring inconsistency in meta-analyses. *BMJ*.
1582 2003;327:557–60.
- 1583 29. Huntington-Klein N, Arenas A, Beam E, Bertoni M, Bloem JR, Burli P, et al. The influence of hidden
1584 researcher decisions in applied microeconomics. *Economic Inquiry*. 2021;59:944–60.
- 1585 30. Jennions MD, Lortie CJ, Rosenberg MS, Rothstein HR. Publication and related biases. In: Koricheva
1586 J, Gurevitch J, Mengersen K, editors. *Handbook of meta-analysis in ecology and evolution*. Princeton,
1587 USA: Princeton University Press; 2013. p. 207–36.
- 1588 31. Kimmel K, Avolio ML, Ferraro PJ. Empirical evidence of widespread exaggeration bias and
1589 selective reporting in ecology. *Nature Ecology & Evolution*. 2023. <https://doi.org/10.1038/s41559-023-02144-3>.
- 1591 32. Klein RA, Ratliff KA, Vianello M, Jr. RBA, Bahník Š, Bernstein MJ, et al. Investigating variation in
1592 replicability: A “many labs” replication project. *Social Psychology*. 2014;45:142–52.
- 1593 33. Klein RA, Vianello M, Hasselman F, Adams BG, Adams RB, Alper S, et al. Many labs 2: Investigating
1594 variation in replicability across samples and settings. *Advances in Methods and Practices in*
1595 *Psychological Science*. 2018;1:443–90.
- 1596 34. Knight K. *Mathematical statistics*. Book. New York: Chapman; Hall; 2000.
- 1597 35. Kou-Giesbrecht S, Menge DNL. Nitrogen-fixing trees increase soil nitrous oxide emissions: A
1598 meta-analysis. *Ecology*. 2021;102:e03415.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 1599 36. Kuznetsova A, Brockhoff PB, Christensen RHB. lmerTest package: Tests in linear mixed effects
1600 models. *Journal of Statistical Software*. 2017;82:1–26.
- 1601 37. Leybourne DJ, Preedy KF, Valentine TA, Bos JIB, Karley AJ. Drought has negative consequences on
1602 aphid fitness and plant vigor: Insights from a meta-analysis. *Ecology and Evolution*. 2021;11:11915–
1603 29.
- 1604 38. Lu X, White H. Robustness checks and robustness tests in applied economics. *Journal of*
1605 *Econometrics*. 2014;178:194–206.
- 1606 39. Lüdtke D, Ben-Shachar MS, Patil I, Waggoner P, Makowski D. Performance: An r package for
1607 assessment, comparison and testing of statistical models. *Journal of Open Source Software*.
1608 2021;6:3139.
- 1609 40. Luke SG. Evaluating significance in linear mixed-effects models in r. *Behavior Research Methods*.
1610 2017;49:1494–502.
- 1611 41. Miles C. Testing market-based instruments for conservation in northern victoria. In: Norton T,
1612 Lefroy T, Bailey K, Unwin G, editors. *Biodiversity: Integrating conservation and production: Case*
1613 *studies from australian farms, forests and fisheries*. Melbourne, Australia: CSIRO Publishing; 2008. p.
1614 133–46.
- 1615 42. Morrissey MB, Ruxton GD. Multiple regression is not multiple regressions: The meaning of
1616 multiple regression and the non-problem of collinearity. *Philosophy, Theory, and Practice in Biology*.
1617 2018;10.
- 1618 43. Nakagawa S, Cuthill IC. Effect size, confidence interval and statistical significance: A practical
1619 guide for biologists. *Biological Reviews*. 2007;82:591–605.
- 1620 44. Nakagawa S, Noble DW, Senior AM, Lagisz M. Meta-evaluation of meta-analysis: Ten appraisal
1621 questions for biologists. *BMC Biology*. 2017;15:18.
- 1622 45. Nicolaus M, Michler SPM, Ubels R, Velde M van der, Komdeur J, Both C, et al. Sex-specific effects
1623 of altered competition on nestling growth and survival: An experimental manipulation of brood size
1624 and sex ratio. *Journal of Animal Ecology*. 2009;78:414–26.
- 1625 46. Noble DWA, Lagisz M, O’Dea RE, Nakagawa S. Nonindependence and sensitivity analyses in
1626 ecological and evolutionary meta-analyses. *Molecular Ecology*. 2017;26:2410–25.
- 1627 47. Open Science Collaboration. Estimating the reproducibility of psychological science. *Science*.
1628 2015;349:aac4716.
- 1629 48. Parker TH, Forstmeier W, Koricheva J, Fidler F, Hadfield JD, Chee YE, et al. Transparency in ecology
1630 and evolution: Real problems, real solutions. *Trends in Ecology & Evolution*. 2016;31:711–9.
- 1631 49. Parker TH, Yang Y. Exaggerated effects in ecology. *Nature Ecology & Evolution*. 2023.
1632 <https://doi.org/10.1038/s41559-023-02156-z>.
- 1633 50. Pei Y, Forstmeier W, Wang D, Martin K, Rutkowska J, Kempnaers B. Proximate causes of infertility
1634 and embryo mortality in captive zebra finches. *The American Naturalist*. 2020;196:577–96.
- 1635 51. R Core Team. R: A language and environment for statistical computing. Vienna, Austria: R
1636 Foundation for Statistical Computing; 2022.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 1637 52. Rosenberg MS. Moment and least-squares based approaches to metaanalytic inference. In:
1638 Koricheva J, Gurevitch J, Mengersen K, editors. Handbook of meta-analysis in ecology and evolution.
1639 Princeton, USA: Princeton University Press; 2013. p. 108–24.
- 1640 53. Royle NJ, Hartley IR, Owens IPF, Parker GA. Sibling competition and the evolution of growth rates
1641 in birds. *Proceedings of the Royal Society B-Biological Sciences*. 1999;266:923–32.
- 1642 54. Schweinsberg M, Feldman M, Staub N, Akker OR van den, Aert RCM van, Assen M van, et al.
1643 Same data, different conclusions: Radical dispersion in empirical results when independent analysts
1644 operationalize and test the same hypothesis. *Organizational Behavior and Human Decision*
1645 *Processes*. 2021;165:228–49.
- 1646 55. Senior AM, Grueber CE, Kamiya T, Lagisz M, O’Dwyer K, Santos ESA, et al. Heterogeneity in
1647 ecological and evolutionary meta-analyses: Its magnitude and implications. *Ecology*. 2016;97:3293–9.
- 1648 56. Shavit A, Ellison AM. Stepping in the same river twice: Replication in biological research. Edited
1649 Book. New Haven, Connecticut, USA: Yale University Press; 2017.
- 1650 57. Siegel KR, Kaur M, Grigal AC, Metzler RA, Dickinson GH. Meta-analysis suggests negative, but
1651 pCO₂-specific, effects of ocean acidification on the structural and functional properties of crustacean
1652 biomaterials. *Ecology and Evolution*. 2022;12:e8922.
- 1653 58. Silberzahn R, Uhlmann EL, Martin DP, Anselmi P, Aust F, Awtrey E, et al. Many analysts, one data
1654 set: Making transparent how variations in analytic choices affect results. *Advances in Methods and*
1655 *Practices in Psychological Science*. 2018;1:337–56.
- 1656 59. Simons DJ, Shoda Y, Lindsay DS. Constraints on generality (COG): A proposed addition to all
1657 empirical papers. *Perspectives on Psychological Science*. 2017.
1658 <https://doi.org/10.1177/174569161770863>.
- 1659 60. Simonsohn U, Simmons JP, Nelson LD. Specification curve: descriptive and inferential statistics on
1660 all reasonable specifications. *SSRN Electronic Journal*. 2015. <https://doi.org/10.2139/ssrn.2694998>.
- 1661 61. Simonsohn U, Simmons JP, Nelson LD. Specification curve analysis. *Nature Human Behaviour*.
1662 2020;4:1208–14.
- 1663 62. Steegen S, Tuerlinckx F, Gelman A, Vanpaemel W. Increasing transparency through a multiverse
1664 analysis. *Perspectives on Psychological Science*. 2016;11:702–12.
- 1665 63. Taylor JW, Taylor KS. Combining probabilistic forecasts of COVID-19 mortality in the united states.
1666 *European Journal of Operational Research*. 2023;304:25–41.
- 1667 64. Dancho M, Vaughan D. Timetk: A tool kit for working with time series. 2023.
- 1668 65. Vander Werf E. Lack’s clutch size hypothesis: An examination of the evidence using meta-analysis.
1669 *Ecology*. 1992;73:1699–705.
- 1670 66. Ver Hoef JM. Who invented the delta method? *The American Statistician*. 2012;66:124–7.
- 1671 67. Verhulst S, Holveck MJ, Riebel K. Long-term effects of manipulated natal brood size on metabolic
1672 rate in zebra finches. *Biology Letters*. 2006;2:478–80.
- 1673 68. Vesk PA, Morris WK, McCallum W, Apted R, Miles C. Processes of woodland eucalypt
1674 regeneration: Lessons from the bush returns trial. *Proceedings of the Royal Society of Victoria*.
1675 2016;128:54–63.

we recommend viewing this manuscript in html format at <https://egouldo.github.io/ManyAnalysts/>

- 1676 69. Viechtbauer W. Conducting meta-analyses in r with the metafor package. 2010. 2010;36:48.
- 1677 70. Yang Y, Sánchez-Tójar A, O’Dea RE, Noble DWA, Koricheva J, Jennions MD, et al. Publication bias
1678 impacts on effect size, statistical power, and magnitude (type m) and sign (type s) errors in ecology
1679 and evolutionary biology. BMC Biology. 2023;21:71.