

# 1 Same data, different analysts: variation in 2 effect sizes due to analytical decisions in 3 ecology and evolutionary biology. 4

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

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

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

9 Shinichi Nakagawa, School of Biological, Earth & Environmental Sciences, University of New South  
10 Wales, Australia

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

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

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

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

15 Robin N Abbey-Lee, Länsstyrelsen Östergötland, Sweden

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

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

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

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

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

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

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

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

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

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

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

28 Suleiman Kehinde Bello, Department of Arid Land Agriculture, King Abdulaziz University, Kingdom of  
29 Saudi Arabia

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

31 Bernd J. Berauer, Department of Plant Ecology, University of Hohenheim, Institute of Landscape and  
32 Plant Ecology, Germany

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

35 Peter D. Billman, Department of Ecology and Evolutionary Biology, University of Connecticut, USA

36 Charlie K. Blake, STEM Center, Southern Illinois University Edwardsville, USA

37 Shannon Blake, University of Guelph, Canada

38 Louis Bliard, Department of Evolutionary Biology and Environmental Studies, University of Zurich,  
39 Switzerland

40 Andrea Bonisoli-Alquati, Department of Biological Sciences, California State Polytechnic University,  
41 Pomona, USA

42 Timothée Bonnet, Centre d'Études Biologiques de Chizé, UMR 7372 Université de la Rochelle - Centre  
43 National de la Recherche Scientifique, France

44 Camille Nina Marion Bordes, Faculty of Life Sciences, Bar Ilan University, Israel

45 Aneesh P. H. Bose, Department of Wildlife, Fish, and Environmental Studies, Swedish University of  
46 Agricultural Sciences, Sweden

47 Thomas Botterill-James, School of Natural Sciences, University of Tasmania, Australia

48 Melissa Anna Boyd, Whitebark Institute, USA

49 Sarah A. Boyle, Department of Biology, Rhodes College, USA

50 Tom Bradfer-Lawrence, Centre for Conservation Science, RSPB, UK

51 Jennifer Bradham, Environmental Studies, Wofford College, USA

52 Jack A. Brand, Department of Wildlife, Fish and Environmental Studies, Swedish University of  
53 Agricultural Sciences, Sweden

54 Martin I. Brengdahl, IFM Biology, Linköping University, Sweden

55 Martin Bulla, Faculty of Environmental Sciences, Czech University of Life Sciences Prague, Czech  
56 Republic

57 Luc Bussière, Biological and Environmental Sciences & Gothenburg Global Biodiversity Centre,  
58 University of Gothenburg, Sweden

59 Ettore Camerlenghi, School of Biological Sciences, Monash University, Australia

60 Sara E. Campbell, Ecology and Evolutionary Biology, University of Tennessee Knoxville, USA

61 Leonardo L. F. Campos, Departamento de Ecologia e Zoologia, Universidade Federal de Santa  
62 Catarina, Brazil

63 Anthony Caravaggi, School of Biological and Forensic Sciences, University of South Wales, UK

64 Pedro Cardoso, Centre for Ecology, Evolution and Environmental Changes (cE3c) & CHANGE - Global  
65 Change and Sustainability Institute, Faculdade de Ciências, Universidade de Lisboa, Portugal

66 Charles J.W. Carroll, Forest and Rangeland Stewardship, Colorado State University, USA

67 Therese A. Catanach, Department of Ornithology, Academy of Natural Sciences of Drexel University,  
68 USA

69 Xuan Chen, Biology, Salisbury University, USA

70 Heung Ying Janet Chik, Groningen Institute for Evolutionary Life Sciences, University of Groningen,  
71 Netherlands

72 Emily Sarah Choy, Department of Biology, McMaster University, Canada

73 Alec Philip Christie, Department of Zoology, University of Cambridge, UK

74 Angela Chuang, Entomology and Nematology, University of Florida, USA

75 Amanda J. Chunco, Environmental Studies, Elon University, USA

76 Bethany L. Clark, BirdLife International, UK

77 Andrea Contina, School of Integrative Biological and Chemical Sciences, The University of Texas Rio  
78 Grande Valley, USA

79 Garth A. Covernton, Department of Ecology and Evolutionary Biology, University of Toronto, Canada

80 Murray P. Cox, Department of Statistics, University of Auckland, New Zealand

81 Kimberly A. Cressman, Catbird Stats, LLC, USA

82 Marco Crotti, School of Biodiversity, One Health & Veterinary Medicine, University of Glasgow, UK

83 Connor Davidson Crouch, School of Forestry, Northern Arizona University, USA

84 Pietro B. D'Amelio, Department of Behavioural Neurobiology, Max Planck Institute for Biological  
85 Intelligence, Germany

86 Alexandra Allison de Sousa, School of Sciences: Center for Health and Cognition, Bath Spa University,  
87 UK

88 Timm Fabian Döbert, Department of Biological Sciences, University of Alberta, Canada

89 Ralph Dobler, Applied Zoology, TU Dresden, Germany

90 Adam J. Dobson, School of Molecular Biosciences, College of Medical Veterinary & Life Sciences,  
91 University of Glasgow, UK

92 Tim S. Doherty, School of Life and Environmental Sciences, The University of Sydney, Australia

93 Szymon Marian Drobniak, Institute of Environmental Sciences, Jagiellonian University, Poland

94 Alexandra Grace Duffy, Biology Department, Brigham Young University, USA

95 Alison B. Duncan, Institute of Evolutionary Sciences Montpellier, University of Montpellier, CNRS,  
96 IRD., France

97 Robert P. Dunn, Baruch Marine Field Laboratory, University of South Carolina, USA

98 Jamie Dunning, Department of Life Sciences, Imperial College London, UK

99 Trishna Dutta, European Forest Institute, Germany

100 Luke Eberhart-Hertel, Department of Ornithology, Max Planck Institute for Biological Intelligence,  
101 Germany

102 Jared Alan Elmore, Forestry and Environmental Conservation, National Bobwhite and Grassland  
103 Initiative, Clemson University, USA

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

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

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

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

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

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

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

113 John Fieberg, Department of Fisheries, Wildlife, and Conservation Biology, University of Minnesota,  
114 USA

115 Elizabeth A. Finch, CABI, UK

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

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

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

120 Wolfgang Forstmeier, Department of Ornithology, Max Planck Institute for Biological Intelligence,  
121 Germany

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

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

124 Cathryn A. Freund, Wake Forest University, USA

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

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

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

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

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

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

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

199 Kristen Marianne Lalla, Department of Natural Resource Sciences, McGill University, Canada

200 Daniel J. Larkin, Department of Fisheries, Wildlife and Conservation Biology, University of Minnesota-  
201 Twin Cities, USA

202 Courtney L. Larson, The Nature Conservancy, USA

203 Katherine S. Lauck, Department of Wildlife, Fish, and Conservation Biology, University of California,  
204 Davis, USA

205 M. Elise Lauterbur, Ecology and Evolutionary Biology, University of Arizona, USA

206 Alan Law, Biological and Environmental Sciences, University of Stirling, UK

207 Don-Jean Léandri-Breton, Department of Natural Resource Sciences, McGill University, Canada

208 Jonas J. Lembrechts, Department of Biology, University of Antwerp, Belgium

209 Kiara L'Herpinier, Natural sciences, Macquarie University, Australia

210 Eva J. P. Lievens, Aquatic Ecology and Evolution Group, Limnological Institute, University of Konstanz,  
211 Germany

212 Daniela Oliveira de Lima, Campus Cerro Largo, Universidade Federal da Fronteira Sul, Brazil

213 Shane Lindsay, School of Psychology and Social Work, University of Hull, UK

214 Martin Luquet, UMR 1224 ECOBIOP, Université de Pau et des Pays de l'Adour, France

215 Ross MacLeod, School of Biological & Environmental Sciences, Liverpool John Moores University, UK

216 Kirsty H. Macphie, Institute of Ecology and Evolution, University of Edinburgh, UK

217 Kit Magellan, Cambodia

218 Magdalena M. Mair, Statistical Ecotoxicology, Bayreuth Center of Ecology and Environmental  
219 Research (BayCEER), University of Bayreuth, Germany

220 Lisa E. Malm, Ecology and Environmental Science, Umeå University, Sweden

221 Stefano Mammola, Molecular Ecology Group (MEG), Water Research Institute (IRSA), National  
222 Research Council of Italy (CNR), Italy

223 Caitlin P. Mandeville, Department of Natural History, Norwegian University of Science and  
224 Technology, Norway

225 Michael Manhart, Center for Advanced Biotechnology and Medicine, Rutgers University Robert  
226 Wood Johnson Medical School, USA

227 Laura Milena Manrique-Garzon, Departamento de Ciencias Biológicas, Universidad de los Andes,  
228 Colombia

229 Elina Mäntylä, Department of Biology, University of Turku, Finland

230 Philippe Marchand, Institut de recherche sur les forêts, Université du Québec en Abitibi-  
231 Témiscamingue, Canada

232 Benjamin Michael Marshall, Biological and Environmental Sciences, University of Stirling, UK

233 Charles A. Martin, Université du Québec à Trois-Rivières, Canada

234 Dominic Andreas Martin, Institute of Plant Sciences, University of Bern, Switzerland

235 Jake Mitchell Martin, Department of Wildlife, Fish, and Environmental Studies, Swedish University of  
236 Agricultural Sciences, Sweden

237 April Robin Martinig, School of Biological, Earth and Environmental Sciences, University of New South  
238 Wales, Australia

239 Erin S. McCallum, Department of Wildlife, Fish and Environmental Studies, Swedish University of  
240 Agricultural Sciences, Sweden

241 Mark McCauley, Whitney Laboratory for Marine Bioscience, University of Florida, USA

242 Sabrina M. McNew, Ecology and Evolutionary Biology, University of Arizona, USA

243 Scott J. Meiners, Biological Sciences, Eastern Illinois University, USA

244 Thomas Merkling, Centre d'Investigations Clinique Plurithématique - Institut Lorrain du Coeur et des  
245 Vaisseaux, Université de Lorraine, Inserm1433 CIC-P CHRU de Nancy, France

246 Marcus Michelangeli, Department of Wildlife, Fish and Environmental Studies, Swedish University of  
247 Agricultural Sciences, Sweden

248 Maria Moiron, Evolutionary biology department, Bielefeld University, Germany

249 Bruno Moreira, Department of Ecology and global change, Centro de Investigaciones sobre  
250 Desertificación, Consejo Superior de Investigaciones Científicas (CIDE-CSIC/UV/GV), Spain

251 Jennifer Mortensen, Department of Biological Sciences, University of Arkansas, USA

252 Benjamin Mos, School of the Environment, Faculty of Science, The University of Queensland,  
253 Australia

254 Taofeek Olatunbosun Muraina, Department of Animal Health and Production, Oyo State College of  
255 Agriculture and Technology, Nigeria

256 Penelope Wrenn Murphy, Department of Forest & Wildlife Ecology, University of Wisconsin-Madison,  
257 USA

258 Luca Nelli, School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow, UK

259 Petri Niemelä, Organismal and Evolutionary Biology Research Programme, Faculty of Biological and  
260 Environmental Sciences, University of Helsinki, Finland

261 Josh Nightingale, South Iceland Research Centre, University of Iceland, Iceland

262 Gustav Nilsson, Department of Clinical Neuroscience, Karolinska Institutet, Sweden

263 Sergio Nolzco, School of Biological Sciences, Monash University, Australia

264 Sabine S. Nooten, Animal Ecology and Tropical Biology, University of Würzburg, Germany

265 Jessie Lanterman Novotny, Biology, Hiram College, USA

266 Agnes Birgitta Olin, Department of Aquatic Resources, Swedish University of Agricultural Sciences,  
267 Sweden



268 Chris L. Organ, Department of Earth Sciences, Montana State University, USA

269 Kate L. Ostevik, Department of Evolution, Ecology, and Organismal Biology, University of California,  
270 Riverside, USA

271 Facundo Xavier Palacio, Sección Ornitología, Universidad Nacional de La Plata, Argentina

272 Matthieu Paquet, Department of Ecology, Swedish University of Agricultural Sciences, Sweden

273 Darren James Parker, Bangor University, UK

274 David J. Pascall, MRC Biostatistics Unit, University of Cambridge, UK

275 Valerie J. Pasquarella, Harvard Forest, Harvard University, USA

276 John Harold Paterson, Biological and Environmental Sciences, University of Stirling, Scotland

277 Ana Payo-Payo, Departamento de Biodiversidad, Ecología y Evolución., Universidad Complutense de  
278 Madrid, Spain

279 Karen Marie Pedersen, Biology Department, Technische Universität Darmstadt, Germany

280 Grégoire Perez, UMR 1309 ASTRE, CIRAD, France

281 Kayla I. Perry, Department of Entomology, The Ohio State University, USA

282 Patrice Pottier, Evolution & Ecology Research Centre, School of Biological, Earth and Environmental  
283 Sciences, The University of New South Wales, Australia

284 Michael J. Proulx, Department of Psychology, University of Bath, UK

285 Raphaël Proulx, Chaire de recherche en intégrité écologique, Université du Québec à Trois-Rivières,  
286 Canada

287 Jessica L Pruet, Mississippi Based RESTORE Act Center of Excellence, University of Southern  
288 Mississippi, USA

289 Veronarindra Ramananjato, Department of Integrative Biology, University of California, Berkeley, USA

290 Finaritra Tolotra Randimbarison, Mention Zoologie et Biodiversité Animale, Université  
291 d'Antananarivo, Madagascar

292 Onja H. Razafindratsima, Department of Integrative Biology, University of California, Berkeley, USA

293 Diana J. Rennison, Department of Ecology, Behavior and Evolution, University of California, San  
294 Diego, USA

295 Federico Riva, Institute for Environmental Sciences, VU Amsterdam, The Netherlands

296 Sepand Riyahi, Department of Evolutionary Anthropology, University of Vienna, Austria

297 Michael James Roast, Konrad Lorenz Institute for Ethology, University of Veterinary Medicine, Austria

298 Felipe Pereira Rocha, School of Biological Sciences, The University of Hong Kong, China

299 Dominique G. Roche, Institut de biologie, Université de Neuchâtel, Switzerland

300 Cristian Román-Palacios, School of Information, University of Arizona, USA

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

334 Michael I. Sitvarin, Independent researcher, USA

335 Fabrício Luiz Skupien, Programa de Pós-Graduação em Ecologia, Instituto de Biologia, Centro de  
336 Ciências da Saúde, Universidade Federal do Rio de Janeiro, Brazil

337 Heather Lea Slinn, Vive Crop Protection, Canada

338 Grania Polly Smith, University of Cambridge, UK

339 Jeremy A. Smith, British Trust for Ornithology, UK

340 Rahel Sollmann, Department of Wildlife, Fish, and Conservation Biology, University of California,  
341 Davis, USA

342 Kaitlin Stack Whitney, Science, Technology & Society Department, Rochester Institute of Technology,  
343 USA

344 Shannon Michael Still, Nomad Ecology, USA

345 Erica F. Stuber, Wildland Resources Department, Utah State University, USA

346 Guy F. Sutton, Center for Biological Control, Department of Zoology and Entomology, Rhodes  
347 University, South Africa

348 Ben Swallow, School of Mathematics and Statistics and Centre for Research in Ecological and  
349 Environmental Modelling, University of St Andrews, UK

350 Conor Claverie Taff, Department of Ecology and Evolutionary Biology, Cornell University, USA

351 Elina Takola, Department of Computational Landscape Ecology, Helmholtz Centre for Environmental  
352 Research – UFZ, Germany

353 Andrew J. Tanentzap, Ecosystems and Global Change Group, School of the Environment, Trent  
354 University, Canada

355 Rocío Tarjuelo, Instituto Universitario de Investigación en Gestión Forestal Sostenible (iuFOR),  
356 Universidad de Valladolid, Spain

357 Richard J. Telford, Department of Biological Sciences, University of Bergen, Norway

358 Christopher J. Thawley, Department of Biological Science, University of Rhode Island, USA

359 Hugo Thierry, Department of Geography, McGill University, Canada

360 Jacqueline Thomson, Integrative Biology, University of Guelph, Canada

361 Svenja Tidau, School of Biological and Marine Sciences, University of Plymouth, UK

362 Emily M. Tompkins, Biology Department, Wake Forest University, USA

363 Claire Marie Tortorelli, Plant Sciences, University of California, Davis, USA

364 Andrew Trlica, College of Natural Resources, North Carolina State University, USA

365 Biz R. Turnell, Institute of Zoology, Technische Universität Dresden, Germany

366 Lara Urban, Helmholtz AI, Helmholtz Zentrum Muenchen, Germany

367 Stijn Van de Vondel, Department of Biology, University of Antwerp, Belgium

368 Jessica Eva Megan van der Wal, FitzPatrick Institute of African Ornithology, University of Cape Town,  
369 South Africa

370 Jens Van Eeckhoven, Department of Cell & Developmental Biology, Division of Biosciences, University  
371 College London, UK

372 Francis van Oordt, Natural Resource Sciences, McGill University, Canada

373 K. Michelle Vanderwel, Biology, University of Saskatchewan, Canada

374 Mark C. Vanderwel, Department of Biology, University of Regina, Canada

375 Karen J. Vanderwolf, Biology, University of Waterloo, Canada

376 Juliana Vélez, Department of Fisheries, Wildlife and Conservation Biology, University of Minnesota,  
377 USA

378 Diana Carolina Vergara-Florez, Department of Ecology & Evolutionary Biology, University of Michigan,  
379 USA

380 Brian C. Verrelli, Center for Biological Data Science, Virginia Commonwealth University, USA

381 Marcus Vinícius Vieira, Dept. Ecologia, Instituto de Biologia, Universidade Federal do Rio de Janeiro,  
382 Brazil

383 Nora Villamil, Lothian Analytical Services, Public Health Scotland, UK

384 Valerio Vitali, Institute for Evolution and Biodiversity, University of Muenster, Germany

385 Julien Vollering, Department of Environmental Sciences, Western Norway University of Applied  
386 Sciences, Norway

387 Jeffrey Walker, Department of Biological Sciences, University of Southern Maine, USA

388 Xanthe J. Walker, Center for Ecosystem Science and Society, Northern Arizona University, USA

389 Jonathan A. Walter, Center for Watershed Sciences, University of California, Davis, USA

390 Pawel Waryszak, School of Agriculture and Environmental Science, University of Southern  
391 Queensland, Australia

392 Ryan J. Weaver, Department of Ecology, Evolution, and Organismal Biology, Iowa State University,  
393 USA

394 Ronja E. M. Wedegärtner, Fram Project AS, Norway

395 Daniel L. Weller, Department of Food Science & Technology, Virginia Polytechnic Institute and State  
396 University, USA

397 Shannon Whelan, Department of Natural Resource Sciences, McGill University, Canada

398 Rachel Louise White, School of Applied Sciences, University of Brighton, UK

399 David William Wolfson, Department of Fisheries, Wildlife and Conservation Biology, University of  
400 Minnesota, USA

401 Andrew Wood, Department of Biology, University of Oxford, UK

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

## 411 Abstract

412 Although variation in effect sizes and predicted values among studies of similar phenomena is  
413 inevitable, such variation far exceeds what might be produced by sampling error alone. One possible  
414 explanation for variation among results is differences among researchers in the decisions they make  
415 regarding statistical analyses. A growing array of studies has explored this analytical variability in  
416 different fields and has found substantial variability among results despite analysts having the same  
417 data and research question. Many of these studies have been in the social sciences, but one small  
418 ‘many analyst’ study found similar variability in ecology. We expanded the scope of this prior work by  
419 implementing a large-scale empirical exploration of the variation in effect sizes and model  
420 predictions generated by the analytical decisions of different researchers in ecology and evolutionary  
421 biology. We used two unpublished datasets, one from evolutionary ecology (blue tit, *Cyanistes*  
422 *caeruleus*, to compare sibling number and nestling growth) and one from conservation ecology  
423 (*Eucalyptus*, to compare grass cover and tree seedling recruitment). The project leaders recruited  
424 174 analyst teams, comprising 246 analysts, to investigate the answers to prespecified research  
425 questions. Analyses conducted by these teams yielded 141 usable effects (compatible with our meta-  
426 analyses and with all necessary information provided) for the blue tit dataset, and 85 usable effects  
427 for the *Eucalyptus* dataset. We found substantial heterogeneity among results for both datasets,  
428 although the patterns of variation differed between them. For the blue tit analyses, the average  
429 effect was convincingly negative, with less growth for nestlings living with more siblings, but there  
430 was near continuous variation in effect size from large negative effects to effects near zero, and even  
431 effects crossing the traditional threshold of statistical significance in the opposite direction. In  
432 contrast, the average relationship between grass cover and *Eucalyptus* seedling number was only  
433 slightly negative and not convincingly different from zero, and most effects ranged from weakly  
434 negative to weakly positive, with about a third of effects crossing the traditional threshold of  
435 significance in one direction or the other. However, there were also several striking outliers in  
436 the *Eucalyptus* dataset, with effects far from zero. For both datasets, we found substantial variation  
437 in the variable selection and random effects structures among analyses, as well as in the ratings of  
438 the analytical methods by peer reviewers, but we found no strong relationship between any of these  
439 and deviation from the meta-analytic mean. In other words, analyses with results that were far from  
440 the mean were no more or less likely to have dissimilar variable sets, use random effects in their  
441 models, or receive poor peer reviews than those analyses that found results that were close to the  
442 mean. The existence of substantial variability among analysis outcomes raises important questions  
443 about how ecologists and evolutionary biologists should interpret published results, and how they  
444 should conduct analyses in the future.

## 445 Introduction

446 One value of science derives from its production of replicable, and thus reliable, results. When we  
447 repeat a study using the original methods, we should be able to expect a similar result. However,  
448 perfect replicability is not a reasonable goal. Effect sizes will vary, and even reverse in sign, by chance  
449 alone (Gelman and Weakliem 2009). Observed patterns can differ for other reasons as well. It could  
450 be that we do not sufficiently understand the conditions that led to the original result so when we  
451 seek to replicate it, the conditions differ due to some ‘hidden moderator’. This hidden moderator  
452 hypothesis is described by meta-analysts in ecology and evolutionary biology as ‘true biological  
453 heterogeneity’ (Senior et al. 2016). This idea of true heterogeneity is popular in ecology and  
454 evolutionary biology, and there are good reasons to expect it in the complex systems in which we

455 work (Shavit and Ellison 2017). However, despite similar expectations in psychology, recent evidence  
456 in that discipline contradicts the hypothesis that moderators are common obstacles to replicability,  
457 as variability in results in a large ‘many labs’ collaboration was mostly unrelated to commonly  
458 hypothesized moderators such as the conditions under which the studies were administered (Klein et  
459 al. 2018). Another possible explanation for variation in effect sizes is that researchers often present  
460 biased samples of results, thus reducing the likelihood that later studies will produce similar effect  
461 sizes (Open Science Collaboration 2015; Parker et al. 2016; Forstmeier, Wagenmakers, and Parker  
462 2017; Fraser et al. 2018; Parker and Yang 2023). It also may be that although researchers did  
463 successfully replicate the conditions, the experiment, and measured variables, analytical decisions  
464 differed sufficiently among studies to create divergent results (Simonsohn, Simmons, and Nelson  
465 2015; Silberzahn et al. 2018).

466 Analytical decisions vary among studies because researchers have many options. Researchers need  
467 to decide how to exclude possibly anomalous or unreliable data, how to construct variables, which  
468 variables to include in their models, and which statistical methods to use. Depending on the dataset,  
469 this short list of choices could encompass thousands or millions of possible alternative  
470 specifications (Simonsohn, Simmons, and Nelson 2015). However, researchers making these  
471 decisions presumably do so with the goal of doing the best possible analysis, or at least the best  
472 analysis within their current skill set. Thus, it seems likely that some specification options are more  
473 probable than others, possibly because they have previously been shown (or claimed) to be better,  
474 or because they are more well known. Of course, some of these different analyses (maybe many of  
475 them) may be equally valid alternatives. Regardless, on probably any topic in ecology and  
476 evolutionary biology, we can encounter differences in choices of data analysis. The extent of these  
477 differences in analyses and the degree to which these differences influence the outcomes of analyses  
478 and therefore studies’ conclusions are important empirical questions. These questions are especially  
479 important given that many papers draw conclusions after applying a single method, or even a single  
480 statistical model, to analyze a dataset.

481 The possibility that different analytical choices could lead to different outcomes has long been  
482 recognized (Gelman and Loken 2013), and various efforts to address this possibility have been  
483 pursued in the literature. For instance, one common method in ecology and evolutionary biology  
484 involves creating a set of candidate models, each consisting of a different (though often similar) set  
485 of predictor variables, and then, for the predictor variable of interest, averaging the slope across all  
486 models (i.e. model averaging) (Burnham and Anderson 2002; Grueber et al. 2011). This method  
487 reduces the chance that a conclusion is contingent upon a single model specification, though use and  
488 interpretation of this method is not without challenges (Grueber et al. 2011). Further, the models  
489 compared to each other typically differ only in the inclusion or exclusion of certain predictor  
490 variables and not in other important ways, such as methods of parameter estimation. More explicit  
491 examination of outcomes of differences in model structure, model type, data exclusion, or other  
492 analytical choices can be implemented through sensitivity analyses (e.g., Noble et al. 2017).  
493 Sensitivity analyses, however, are typically rather narrow in scope, and are designed to assess the  
494 sensitivity of analytical outcomes to a particular analytical choice rather than to a large universe of  
495 choices. Recently, however, analysts in the social sciences have proposed extremely thorough  
496 sensitivity analysis, including ‘multiverse analysis’ (Steege et al. 2016) and the ‘specification  
497 curve’ (Simonsohn, Simmons, and Nelson 2015), as a means of increasing the reliability of results.  
498 With these methods, researchers identify relevant decision points encountered during analysis and  
499 conduct the analysis many times to incorporate many plausible decisions made at each of these  
500 points. The study’s conclusions are then based on a broad set of the possible analyses and so allow  
501 the analyst to distinguish between robust conclusions and those that are highly contingent on

502 particular model specifications. These are useful outcomes, but specifying a universe of possible  
503 modelling decisions is not a trivial undertaking. Further, the analyst's knowledge and biases will  
504 influence decisions about the boundaries of that universe, and so there will always be room for  
505 disagreement among analysts about what to include. Including more specifications is not necessarily  
506 better. Some analytical decisions are better justified than others, and including biologically  
507 implausible specifications may undermine this process. Regardless, these powerful methods have yet  
508 to be adopted, and even the more limited forms of sensitivity analyses are not particularly  
509 widespread. Most studies publish a small set of analyses and so the existing literature does not  
510 provide much insight into the degree to which published results are contingent on analytical  
511 decisions.

512 Despite the potential major impacts of analytical decisions on variance in results, the outcomes of  
513 different individuals' data analysis choices have only recently begun to receive much empirical  
514 attention. The only formal exploration of this that we were aware of when we submitted our Stage 1  
515 manuscript were (1) an analysis in social science that asked whether male professional football  
516 (soccer) players with darker skin tone were more likely to be issued red cards (ejection from the  
517 game for rule violation) than players with lighter skin tone (Silberzahn et al. 2018) and (2) an analysis  
518 in neuroimaging which evaluated nine separate hypotheses involving the neurological responses  
519 detected with fMRI in 108 participants divided between two treatments in a decision making  
520 task (Botvinik-Nezer et al. 2020). Several others have been published since (e.g., Huntington-Klein et  
521 al. 2021; Schweinsberg et al. 2021; Breznau et al. 2022; Coretta et al. 2023), and we recently learned  
522 of an earlier small study in ecology (Stanton-Geddes, Freitas, and Sales Dambros 2014). In the red  
523 card study, 29 teams designed and implemented analyses of a dataset provided by the study  
524 coordinators (Silberzahn et al. 2018). Analyses were peer reviewed (results blind) by at least two  
525 other participating analysts; a level of scrutiny consistent with standard pre-publication peer review.  
526 Among the final 29 analyses, odds-ratios varied from 0.89 to 2.93, meaning point estimates varied  
527 from having players with lighter skin tones receive more red cards (odds ratio < 1) to a strong effect  
528 of players with darker skin tones receiving more red cards (odds ratio > 1). Twenty of the 29 teams  
529 found a statistically-significant effect in the predicted direction of players with darker skin tones  
530 being issued more red cards. This degree of variation in peer-reviewed analyses from identical data is  
531 striking, but the generality of this finding has only just begun to be formally  
532 investigated (e.g., Huntington-Klein et al. 2021; Schweinsberg et al. 2021; Breznau et al.  
533 2022; Coretta et al. 2023).

534 In the neuroimaging study, 70 teams evaluated each of the nine different hypotheses with the  
535 available fMRI data (Botvinik-Nezer et al. 2020). These 70 teams followed a divergent set of  
536 workflows that produced a wide range of results. The rate of reporting of statistically significant  
537 support for the nine hypotheses ranged from 21% to 84%, and for each hypothesis on average, 20%  
538 of research teams observed effects that differed substantially from the majority of other teams.  
539 Some of the variability in results among studies could be explained by analytical decisions such as  
540 choice of software package, smoothing function, and parametric versus non-parametric corrections  
541 for multiple comparisons. However, substantial variability among analyses remained unexplained,  
542 and presumably emerged from the many different decisions each analyst made in their long  
543 workflows. Such variability in results among analyses from this dataset and from the very different  
544 red-card dataset suggests that sensitivity of analytical outcome to analytical choices may characterize  
545 many distinct fields, as several more recent many-analyst studies also suggest (Huntington-Klein et al.  
546 2021; Schweinsberg et al. 2021; Breznau et al. 2022).



547 To further develop the empirical understanding of the effects of analytical decisions on study  
548 outcomes, we chose to estimate the extent to which researchers' data analysis choices drive  
549 differences in effect sizes, model predictions, and qualitative conclusions in ecology and evolutionary  
550 biology. This is an important extension of the meta-research agenda of evaluating factors influencing  
551 replicability in ecology, evolutionary biology, and beyond (Fidler et al. 2017). To examine the effects  
552 of analytical decisions, we used two different datasets and recruited researchers to analyze one or  
553 the other of these datasets to answer a question we defined. The first question was "To what extent  
554 is the growth of nestling blue tits (*Cyanistes caeruleus*) influenced by competition with siblings?" To  
555 answer this question, we provided a dataset that includes brood size manipulations from 332 broods  
556 conducted over three years at Wytham Wood, UK. The second question was "How does grass cover  
557 influence *Eucalyptus* spp. seedling recruitment?" For this question, analysts used a dataset that  
558 includes, among other variables, number of seedlings in different size classes, percentage cover of  
559 different life forms, tree canopy cover, and distance from canopy edge from 351 quadrats spread  
560 among 18 sites in Victoria, Australia.

561 We explored the impacts of data analysts' choices with descriptive statistics and with a series of tests  
562 to attempt to explain the variation among effect sizes and predicted values of the dependent variable  
563 produced by the different analysis teams for both datasets separately. To describe the variability, we  
564 present forest plots of the standardized effect sizes and predicted values produced by each of the  
565 analysis teams, estimate heterogeneity (both absolute,  $\tau^2$ , and proportional,  $I^2$ ) in effect size and  
566 predicted values among the results produced by these different teams, and calculate a similarity  
567 index that quantifies variability among the predictor variables selected for the different statistical  
568 models constructed by the different analysis teams. These descriptive statistics provide the first  
569 estimates of the extent to which explanatory statistical models and their outcomes in ecology and  
570 evolutionary biology vary based on the decisions of different data analysts. We then quantified the  
571 degree to which the variability in effect size and predicted values could be explained by (1) variation  
572 in the quality of analyses as rated by peer reviewers and (2) the similarity of the choices of predictor  
573 variables between individual analyses.

## 574 Methods

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

### 580 Step 1: Select datasets

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

#### 583 Evolutionary ecology

584 Our evolutionary ecology dataset is relevant to a sub-discipline of life-history research which focuses  
585 on identifying costs and trade-offs associated with different phenotypic conditions. These data were  
586 derived from a brood-size manipulation experiment imposed on wild birds nesting in boxes provided  
587 by researchers in an intensively studied population. Understanding how the growth of nestlings is  
588 influenced by the numbers of siblings in the nest can give researchers insights into factors such as the

589 evolution of clutch size, determination of provisioning rates by parents, and optimal levels of sibling  
590 competition (Vander Werf 1992; DeKogel 1997; Royle et al. 1999; Verhulst, Holveck, and Riebel  
591 2006; Nicolaus et al. 2009). Data analysts were provided this dataset and instructed to answer the  
592 following question: “To what extent is the growth of nestling blue tits (*Cyanistes caeruleus*)  
593 influenced by competition with siblings?”

594 Researchers conducted brood size manipulations and population monitoring of blue tits at Wytham  
595 Wood, a 380 ha woodland in Oxfordshire, U.K (1° 20'W, 51° 47'N). Researchers regularly checked  
596 approximately 1100 artificial nest boxes at the site and monitored the 330 to 450 blue tit pairs  
597 occupying those boxes in 2001-2003 during the experiment. Nearly all birds made only one breeding  
598 attempt during the April to June study period in a given year. At each blue tit nest, researchers  
599 recorded the date the first egg appeared, clutch size, and hatching date. For all chicks alive at age 14  
600 days, researchers measured mass and tarsus length and fitted a uniquely numbered, British Trust for  
601 Ornithology (BTO) aluminium leg ring. Researchers attempted to capture all adults at their nests  
602 between day 6 and day 14 of the chick-rearing period. For these captured adults, researchers  
603 measured mass, tarsus length, and wing length and fitted a uniquely numbered BTO leg ring. During  
604 the 2001-2003 breeding seasons, researchers manipulated brood sizes using cross fostering. They  
605 matched broods for hatching date and brood size and moved chicks between these paired nests one  
606 or two days after hatching. They sought to either enlarge or reduce all manipulated broods by  
607 approximately one fourth. To control for effects of being moved, each reduced brood had a portion  
608 of its brood replaced by chicks from the paired increased brood, and vice versa. Net manipulations  
609 varied from plus or minus four chicks in broods of 12 to 16 to plus or minus one chick in broods of 4  
610 or 5. Researchers left approximately one third of all broods unmanipulated. These unmanipulated  
611 broods were not selected systematically to match manipulated broods in clutch size or laying date.  
612 We have mass and tarsus length data from 3720 individual chicks divided among 167 experimentally  
613 enlarged broods, 165 experimentally reduced broods, and 120 unmanipulated broods. The full list of  
614 variables included in the dataset is publicly available (<https://osf.io/hdv8m>), along with the data  
615 (<https://osf.io/qjzby>).

## 616 Ecology and conservation

### **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.

617 Our ecology and conservation dataset is relevant to a sub-discipline of conservation research which  
618 focuses on investigating how best to revegetate private land in agricultural landscapes. These data  
619 were collected on private land under the Bush Returns program, an incentive system where  
620 participants entered into a contract with the Goulburn Broken Catchment Management Authority  
621 and received annual payments if they executed predetermined restoration activities. This particular  
622 dataset is based on a passive regeneration initiative, where livestock grazing was removed from the  
623 property in the hopes that the *Eucalyptus* spp. overstorey would regenerate without active (and  
624 expensive) planting. Analyses of some related data have been published (Miles 2008; Vesk et al.  
625 2016) but those analyses do not address the question analysts answered in our study. Data analysts  
626 were provided this dataset and instructed to answer the following question: “How does grass cover  
627 influence *Eucalyptus* spp. seedling recruitment?”

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

641

## 642 Step 2: Recruitment and initial survey of analysts

643 The lead team (TP, HF, SN, EG, SG, PV, DH, FF) created a publicly available document providing a

### 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.

644 general description of the project (<https://osf.io/mn5aj/>). The project was advertised at conferences,  
645 via Twitter, using mailing lists for ecological societies (including *Ecolog*, *Evoldir*, and lists for *the*  
646 *Environmental Decisions Group*, and *Transparency in Ecology and Evolution*), and via word of mouth.  
647 The target population was active ecology, conservation, or evolutionary biology researchers with a  
648 graduate degree (or currently studying for a graduate degree) in a relevant discipline. Researchers  
649 could choose to work independently or in a small team. For the sake of simplicity, we refer to these  
650 as ‘analysis teams’ though some comprised one individual. We aimed for a minimum of 12 analysis  
651 teams independently evaluating each dataset (see sample size justification below). We  
652 simultaneously recruited volunteers to peer review the analyses conducted by the other volunteers  
653 through the same channels. Our goal was to recruit a similar number of peer reviewers and analysts,  
654 and to ask each peer reviewer to review a minimum of four analyses. If we were unable to recruit at  
655 least half the number of reviewers as analysis teams, we planned to ask analysts to serve also as  
656 reviewers (after they had completed their analyses), but this was unnecessary. Therefore, no data  
657 analysts peer reviewed analyses of the dataset they had analyzed. All analysts and reviewers were  
658 offered the opportunity to share co-authorship on this manuscript and we planned to invite them to  
659 participate in the collaborative process of producing the final manuscript. All analysts signed  
660 [digitally] a consent (ethics) document (<https://osf.io/xyp68/>) approved by the Whitman College  
661 Institutional Review Board prior to being allowed to participate.

662 We identified our minimum number of analysts per dataset by considering the number of effects  
663 needed in a meta-analysis to generate an estimate of heterogeneity ( $\tau^2$ ) with a 95% confidence  
664 interval that does not encompass zero. This minimum sample size is invariant regardless of  $\tau^2$ . This is  
665 because the same t-statistic value will be obtained by the same sample size regardless of variance

666 ( $\tau^2$ ). We see this by first examining the formula for the standard error, SE for variance, ( $\tau^2$ ) or ( $SE\tau^2$ )  
667 assuming normality in an underlying distribution of effect sizes (Knight 2000):

668 
$$SE(\tau^2) = \sqrt{\frac{2\tau^4}{n-1}}$$

669 and then rearranging the above formula to show how the t-statistic is independent of  $\tau^2$ , as seen  
670 below.

671 
$$t = \frac{\tau^2}{SE(\tau^2)} = \sqrt{\frac{n-1}{2}}$$

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

### 673 Step 3: Primary data analyses

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

684

#### **Additional Information:**

As is common in many studies in ecology and evolutionary biology, the datasets we provided contained many variables, and the research questions we provided could be addressed by our datasets in many different ways. For instance, volunteer analysts had to choose the dependent (response) variable and the independent variable, and make numerous other decisions about which variables and data to use and how to structure their model.

685

#### **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>).

## 686 Step 4: Peer reviews of analyses

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

699 The process for a single reviewer was as follows. First, the reviewer received a description of the  
700 methods of a single analysis. This included the narrative methods section, the analysis team’s  
701 answers to our survey questions regarding their methods, including analysis code, and the dataset.  
702 The reviewer was then asked, in an online survey (<https://osf.io/4t36u/>), to rate that analysis on a  
703 scale of 0-100 based on this prompt: “Rate the overall appropriateness of this analysis to answer the  
704 research question (*one of the two research questions inserted here*) with the available data. To help  
705 you calibrate your rating, please consider the following guidelines:

- 706 • 100. A perfect analysis with no conceivable improvements from the reviewer
- 707 • 75. An imperfect analysis but the needed changes are unlikely to dramatically alter outcomes
- 708 • 50. A flawed analysis likely to produce either an unreliable estimate of the relationship or an  
709 over-precise estimate of uncertainty
- 710 • 25. A flawed analysis likely to produce an unreliable estimate of the relationship and an over-  
711 precise estimate of uncertainty
- 712 • 0. A dangerously misleading analysis, certain to produce both an estimate that is wrong and  
713 a substantially over-precise estimate of uncertainty that places undue confidence in the  
714 incorrect estimate.

715 \*Please note that these values are meant to calibrate your ratings. We welcome ratings of any  
716 number between 0 and 100.

717 After providing this rating, the reviewer was presented with this prompt, in multiple-choice format:  
718 “Would the analytical methods presented produce an analysis that is (a) publishable as is, (b)  
719 publishable with minor revision, (c) publishable with major revision, (d) deeply flawed and  
720 unpublishable?” The reviewer was then provided with a series of text boxes and the following  
721 prompts: “Please explain your ratings of this analysis. Please evaluate the choice of statistical analysis  
722 type. Please evaluate the process of choosing variables for and structuring the statistical model.  
723 Please evaluate the suitability of the variables included in (or excluded from) the statistical model.  
724 Please evaluate the suitability of the structure of the statistical model. Please evaluate choices to  
725 exclude or not exclude subsets of the data. Please evaluate any choices to transform data (or, if there  
726 were no transformations, but you think there should have been, please discuss that choice).” After  
727 submitting this review, a methods section from a second analysis was then made available to the  
728 reviewer. This same sequence was followed until all analyses allocated to a given reviewer were  
729 provided and reviewed. After providing the final review, the reviewer was simultaneously provided  
730 with all four (or more) methods sections the reviewer had just completed reviewing, the option to

731 revise their original ratings, and a text box to provide an explanation. The invitation to revise the  
 732 original ratings was as follows: “If, now that you have seen all the analyses you are reviewing, you  
 733 wish to revise your ratings of any of these analyses, you may do so now.” The text box was prefaced  
 734 with this prompt: “Please explain your choice to revise (or not to revise) your ratings.”

735

**Additional Information: unregistered analysis**

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).

736 **Step 5: Evaluate variation**

737

**Additional Information: analysis schematic**

The lead team conducted a range of preregistered and exploratory analyses to understand variation between analyses and their results. Figure 1 is intended to clarify the analyses described below.

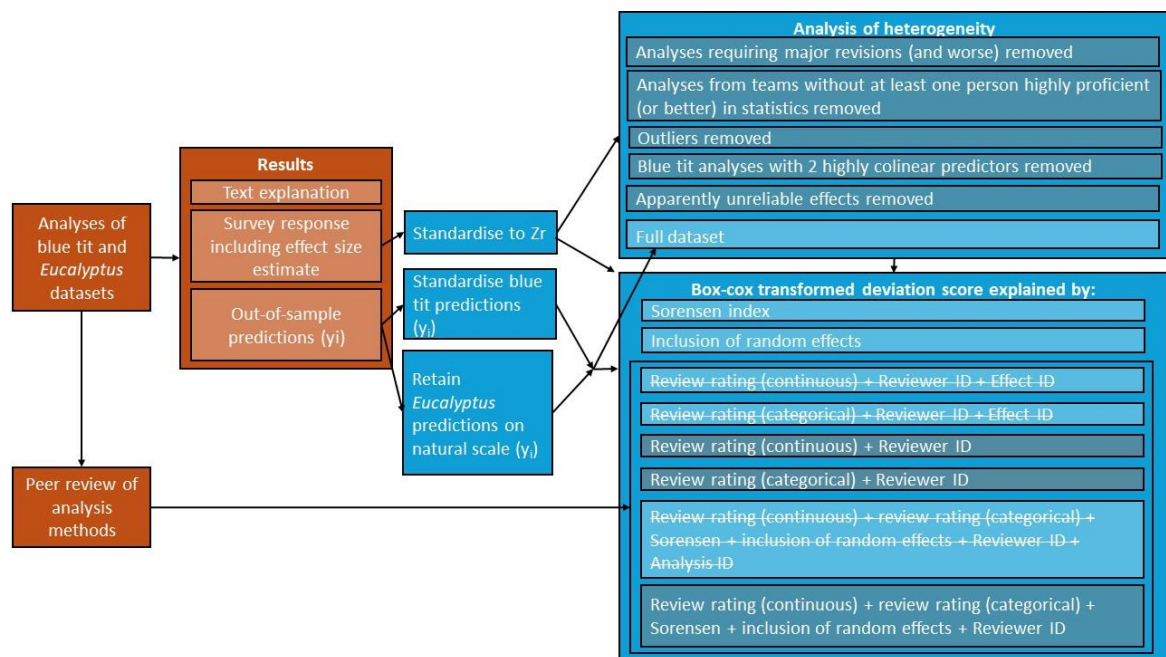


Figure 1: Schematic of research process showing recruited analyst and reviewer contributions in orange and core team contributions in blue. Items that are crossed out were preregistered but could not be conducted. Items with a greyed background were added as exploratory analyses after preregistration.

738

739 The lead team conducted the analyses outlined in this section. We described the variation in model  
740 specification in several ways. We calculated summary statistics describing variation among analyses,  
741 including mean, SD, and range of number of variables per model included as fixed effects, the  
742 number of interaction terms, the number of random effects, and the mean, SD, and range of sample  
743 sizes. We also present the number of analyses in which each variable was included. We summarized  
744 the variability in standardized effect sizes and predicted values of dependent variables among the  
745 individual analyses using standard random effects meta-analytic techniques. First, we derived  
746 standardized effect sizes from each individual analysis. We did this for all linear models or  
747 generalized linear models by converting the  $t$  value and the degree of freedom ( $df$ ) associated with  
748 regression coefficients (e.g. the effect of the number of siblings [predictor] on growth [response] or  
749 the effect of grass cover [predictor] on seedling recruitment [response]) to the correlation  
750 coefficient,  $r$ , using the following:

751 
$$r = \frac{t^2}{(t^2 + df)}$$

752 This formula can only be applied if  $t$  and  $df$  values originate from linear or generalized linear models  
753 [GLMs; Nakagawa and Cuthill (2007)]. If, instead, linear mixed-effects models (LMMs) or generalized  
754 linear mixed-effects models (GLMMs) were used by a given analysis, the exact  $df$  cannot be  
755 estimated. However, adjusted  $df$  can be estimated, for example, using the Satterthwaite  
756 approximation of  $df$ ,  $df_s$ , [note that SAS uses this approximation to obtain  $df$  for LMMs and  
757 GLMMs; Luke (2017)]. For analyses using either LMMs or GLMMs that do not produce  $df_s$  we  
758 planned to obtain  $df_s$  by rerunning the same (G)LMMs using the `lmer()` or `glmer()` function in  
759 the `lmerTest` package in R (Kuznetsova, Brockhoff, and Christensen 2017; R Core Team 2024).

760

**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>).

761 We then used the  $t$  values and  $df_s$  from the models to obtain  $r$  as per the formula above. All  $r$  and  
762 accompanying  $df$  (or  $df_s$ ) were converted to Fisher’s  $Z_r$ .

763 
$$Z_r = \frac{1}{2} \ln \left( \frac{1+r}{1-r} \right)$$

764 and its sampling variance;  $1/(n-3)$  where  $n=df+1$ . Any analyses from which we could not derive a  
765 signed  $Z_r$ , for instance one with a quadratic function in which the slope changed sign, were  
766 considered unusable for analyses of  $Z_r$ . We expected such analyses would be rare. In fact, most  
767 submitted analyses excluded from our meta-analysis of  $Z_r$  were excluded because of a lack of  
768 sufficient information provided by the analyst team rather than due to the use of effects that could  
769 not be converted to  $Z_r$ . Regardless, as we describe below, we generated a second set of standardized  
770 effects (predicted values) that could (in principle) be derived from any explanatory model produced  
771 by these data.

772 Besides  $Z_r$ , which describes the strength of a relationship based on the amount of variation in a  
773 dependent variable explained by variation in an independent variable, we also examined differences

774 in the shape of the relationship between the independent and dependent variables. To accomplish  
775 this, we derived a point estimate (out-of-sample predicted value) for the dependent variable of  
776 interest for each of three values of our primary independent variable. We originally described these  
777 three values as associated with the 25th percentile, median, and 75th percentile of the independent  
778 variable and any covariates.  
779

**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>).

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

787

788 We plotted individual effect size estimates ( $Z_r$ ) and predicted values of the dependent variable ( $y_i$ )  
789 and their corresponding 95% confidence / credible intervals in forest plots to allow visualization of



## Preregistration Deviation:

### 1. Standardizing blue tit out-of-sample predictions ( $y_i$ )

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 ([Supplementary Material B, Equation B.1](#)).

Note that we were unable to standardise some analyst-constructed variables, so these analyses were excluded from the final out-of-sample estimates meta-analysis, see [Supplementary Material B, section B.1.2.1](#) for details and explanation.

### 2. Log-transforming *Eucalyptus* out-of-sample predictions $y_i$

All analyses of the *Eucalyptus* data chose dependent variables that were on the same scale, that is, *Eucalyptus* seedling counts. Although analysts may have used different size-classes of *Eucalyptus* seedlings for their dependent variable, we considered these choices to be akin to subsetting, rather than as different response variables, since changing the size-class of the dependent variable ultimately results in observations being omitted or included. Consequently, we did not standardise *Eucalyptus* out-of-sample predictions.

We were unable to fit quasi-Poisson or Poisson meta-regressions, as desired (O'Hara and Kotze 2010), because available meta-analysis packages (e.g. `metafor::` and `metainc::`) do not provide implementation for outcomes as estimates-only, methods are only provided for outcomes as ratios or rate-differences between two groups. Consequently, we log-transformed the out-of-sample predictions for the *Eucalyptus* data and use the mean estimate for each prediction scenario as the dependent variable in our meta-analysis with the associated SE as the sampling variance in the meta-analysis (Nakagawa et al. 2023, Table 2).

790 the range and precision of effect size and predicted values. Further, we included these estimates in  
791 random effects meta-analyses (Higgins et al. 2003; Borenstein et al. 2017) using the *metafor* package  
792 in R (Viechtbauer 2010; R Core Team 2024):

$$793 \quad Z_r \sim 1 + (1|Effect\ ID)$$

$$794 \quad y_i \sim 1 + (1|Effect\ ID)$$

795 where  $y_i$  is the predicted value for the dependent variable at the 25th percentile, median, or 75th  
796 percentile of the independent variables. The individual  $Z_r$  effect sizes were weighted with the inverse  
797 of sampling variance for  $Z_r$ . The individual predicted values for dependent variable ( $y_i$ ) were weighted  
798 by the inverse of the associated  $SE^2$  (original registration omitted "inverse of the" in error). These  
799 analyses provided an average  $Z_r$  score ( $\bar{Z}_r$ ) or an average  $y_i$  ( $\bar{y}_i$ ) with corresponding 95% confidence  
800 interval and allowed us to estimate two heterogeneity indices,  $\tau^2$  and  $I^2$ . The former,  $\tau^2$ , is the  
801 absolute measure of heterogeneity or the between-study variance (in our case, between-effect  
802 variance) whereas  $I^2$  is a relative measure of heterogeneity. We obtained the estimate of relative  
803 heterogeneity ( $I^2$ ) by dividing the between-effect variance by the sum of between-effect and within-

804 effect variance (sampling error variance).  $I^2$  is thus, in a standard meta-analysis, the proportion of  
805 variance that is due to heterogeneity as opposed to sampling error. When calculating  $I^2$ , within-study  
806 variance is amalgamated across studies to create a “typical” within-study variance which serves as  
807 the sampling error variance (Higgins et al. 2003; Borenstein et al. 2017). Our goal here was to  
808 visualize and quantify the degree of variation among analyses in effect size estimates (Nakagawa and  
809 Cuthill 2007). We did not test for statistical significance.

810

**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 (Rosenberg 2013). 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 at 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  $\tau^2$  also differs somewhat from traditional meta-analysis, and we discuss this further in the Results.

811

812 Finally, we assessed the extent to which deviations from the meta-analytic mean by individual effect  
813 sizes ( $Z_r$ ) or the predicted values of the dependent variable ( $y_i$ ) were explained by the peer rating of  
814 each analysis team’s method section, by a measurement of the distinctiveness of the set of predictor  
815 variables included in each analysis, and by the choice of whether or not to include random effects in  
816 the model. The deviation score, which served as the dependent variable in these analyses, is the  
817 absolute value of the difference between the meta-analytic mean  $\bar{Z}_r$  (or  $\bar{y}_i$ ) and the  
818 individual  $Z_r$  (or  $y_i$ ) estimate for each analysis. We used the Box-Cox transformation on the absolute  
819 values of deviation scores to achieve an approximately normal distribution (c.f. Fanelli and Ioannidis  
820 2013; Fanelli, Costas, and Ioannidis 2017). We described variation in this dependent variable with  
821 both a series of univariate analyses and a multivariate analysis. All these analyses were general linear

822 (mixed) models. These analyses were secondary to our estimation of variation in effect sizes  
823 described above. We wished to quantify relationships among variables, but we had no *a*  
824 *priori* expectation of effect size and made no dichotomous decisions about statistical significance.

825 When examining the extent to which reviewer ratings (on a scale from 0 to 100) explained deviation  
826 from the average effect (or predicted value), each analysis had been rated by multiple peer  
827 reviewers, so for each reviewer score to be included, we include each deviation score in the analysis  
828 multiple times. To account for the non-independence of multiple ratings of the same analysis, we  
829 planned to include analysis identity as a random effect in our general linear mixed model in  
830 the *lme4* package in R (Bates et al. 2015; R Core Team 2024). To account for potential differences  
831 among reviewers in their scoring of analyses, we also planned to include reviewer identity as a

**Additional explanation:**

In our meta-analyses based on Box-Cox transformed deviation scores, we leave these deviation scores unweighted. This is consistent with our registration, which did not mention weighting these scores. However, the fact that we did not mention weighting the scores was actually an error: we had intended to weight them, as is standard in meta-analysis, using the inverse variance of the Box-Cox transformed deviation scores [Supplementary Material C, equation C.1](#). Unfortunately, when we did conduct the weighted analyses, they produced results in which some weighted estimates differed radically from the unweighted estimate because the weights were invalid. Such invalid weights can sometimes occur when the variance (upon which the weights depend) is partly a function of the effect size, as in our Box-Cox transformed deviation scores (Nakagawa et al. 2022). In the case of the *Eucalyptus* analyses, the most extreme outlier was weighted much more heavily (by close to two orders of magnitude) than any other effect sizes because the effect size was, itself, so high. Therefore, we made the decision to avoid weighting by inverse variance in all analyses of the Box-Cox transformed deviation scores. This was further justified because (a) most analyses have at least some moderately unreliable weights, and (b) the sample sizes were mostly very similar to each other across submitted analyses, and so meta-analytic weights are not particularly important here (Buck et al. 2022). We systematically investigated the impact of different weighting schemes and random effects on model convergence and results, see [Supplementary Material C, section C.8](#) for more details.

832 random effect:

$$833 \text{DeviationScore}_j = \text{BoxCox}(\text{DeviationFromMean}_j)$$

$$834 \text{DeviationScore}_{ij} \sim \text{Rating}_{ij} + \text{ReviewerID}_i + \text{EffectID}_j$$

$$835 \text{ReviewerID}_i \sim N(0, \sigma_i^2)$$

$$836 \text{EffectID} \sim N(0, \sigma_j^2)$$

837 Where  $\text{DeviationFromMean}_j$  is the deviation from the meta-analytic mean for the  $j$ th  
838 analysis,  $\text{ReviewerID}_i$  is the random intercept assigned to each  $i$  reviewer, and  $\text{EffectID}_j$  is the  
839 random intercept assigned to each  $j$  analysis, both of which are assumed to be normally distributed  
840 with a mean of 0 and a variance of  $\sigma^2$ . Absolute deviation scores were Box-Cox transformed using  
841 the `step_box_cox()` function from the *timetk* package in R (Dancho and Vaughan 2023; R Core Team  
842 2024).

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

851

**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 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 ([Supplementary Material C, section C.2](#)). Removing one of the random effects was necessary for the models to converge. For both models of deviation from the meta-analytic mean explained by categorical or continuous reviewer ratings, we removed the random effect of Effect ID, leaving Reviewer ID as the only random effect.
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, the initial ratings were over-written, therefore we did not have access to those initial values.

852 The next set of univariate analyses sought to explain deviations from the mean effects based on a  
853 measure of the distinctiveness of the set of variables included in each analysis. As a 'distinctiveness'  
854 score, we used Sorensen's Similarity Index (an index typically used to compare species composition  
855 across sites), treating variables as species and individual analyses as sites. To generate an individual  
856 Sorensen's value for each analysis required calculating the pairwise Sorensen's value for all pairs of  
857 analyses (of the same dataset), and then taking the average across these Sorensen's values for each  
858 analysis. We calculated the Sorensen's index values using the *betapart* package (Baselga et al.  
859 2023) in R:

860

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

861 where a is the number of variables common to both analyses, b is the number of variables that occur  
862 in the first analysis but not in the second and c is the number of variables that occur in the second  
863 analysis. We then used the per-model average Sorensen's index value as an independent variable to  
864 predict the deviation score in a general linear model, and included no random effect since each  
865 analysis is included only once, in R (R Core Team 2024):

$$866 \text{DeviationScore}_j \sim \beta \text{Sorensen}_j$$

**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.

867 Next, we assessed the relationship between the inclusion of random effects in the analysis and the  
868 deviation from the mean effect size. We anticipated that most analysts would use random effects in a  
869 mixed model framework, but if we were wrong, we wanted to evaluate the differences in outcomes  
870 when using random effects versus not using random effects. Thus, if there were at least 5 analyses  
871 that did and 5 analyses that did not include random effects, we would add a binary predictor variable  
872 "random effects included (yes/no)" to our set of univariate analyses and would add this predictor  
873 variable to our multivariate model described below. This standard was only met for  
874 the *Eucalyptus* analyses, and so we only examined inclusion of random effects as a predictor variable  
875 in meta-analysis of this set to analyses.

876 Finally, we conducted a multivariate analysis with the five predictors described above (peer ratings 0-  
877 100 and peer ratings of publishability 1-4; both original and revised and Sorensen's index, plus a sixth  
878 for *Eucalyptus*, presence / absence of random effects) with random effects of analysis identity and  
879 reviewer identity in the *lme4* package in R (Bates et al. 2015; R Core Team 2024). We had stated here  
880 in the text that we would use only the revised (final) peer ratings in this analysis, so the absence of  
881 the initial ratings is not a deviation from our plan:

$$882 \text{DeviationScore}_j = \text{BoxCox}(\text{DeviationFromMean}_j)$$

$$883 \text{DeviationScore}_{ij} \sim \text{RatingContinuous}_{ij} + \text{RatingCategorical}_{ij} + \beta \text{Sorensen}_j + \text{ReviewerID}_i \\ 884 + \text{Effect ID}_j$$

$$885 \text{ReviewerID}_i \sim N(0, \sigma_i^2)$$

$$886 \text{EffectID}_j \sim N(0, \sigma_j^2)$$

887

888 We conducted all the analyses described above eight times; for each of the four responses  
889 ( $Z_r$ ,  $y_{25}$ ,  $y_{50}$ ,  $y_{75}$ ) one time for each of the two datasets.

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

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

**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 provided to teams for analysis.

### **Additional explanation: unregistered analyses**

#### **1. Evaluating model fit.**

We evaluated all fitted models using the [performance::performance\(\)](#) function from the *performance* package (Lüdtke, Ben-Shachar, et al. 2021) and the `glance()` function from the *broom.mixed* package (Bolker et al. 2024). For all models, we calculated the square root of the residual variance (Sigma) and the root mean squared error (RMSE). For GLMMs [performance::performance\(\)](#) calculates the marginal and conditional  $R^2$  values as well as the contribution of random effects (ICC), based on Nakagawa et al. (2017). 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 created forest plots and 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. Additionally, to assess potential impacts of highly collinear predictor variables on estimates of  $Z_r$  in blue tit analyses, we created forest plots ([Supplementary Material B, Figure B.5](#)) and recalculated our heterogeneity estimates after we removed analyses that contained the brood count after manipulation and the highly correlated (correlation of 0.89, [Supplementary Material D, Figure D.2](#)) brood count at day 14. This removal included the one effect based on a model that contained both these variables and a third highly correlated variable, the estimate of number of chicks fledged (the only model that included the estimate of number of chicks fledged). We did not conduct a similar analysis for the *Eucalyptus* dataset because there were no variables highly collinear with the primary predictors (grass cover variables) in that dataset ([Supplementary Material D, Figure D.1](#)).

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

Our meta-analyses of variation in  $Z_r$  required variance estimates derived from estimates of the degrees of freedom in original analyses from which  $Z_r$  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  $Z_r$  estimates with identical estimates of degrees of freedom and we present these analyses in the supplement.

#### **Additional explanation: Best practices in many-analysts research**

After we initiated our project, a paper was published outlining best practices in many-analysts studies (Aczel et al. 2021). Although we did not have access to this document when we implemented our project, our study complies with these practices nearly completely. The one exception is that although we requested analysis code from analysts, we did not require submission of code.

910

## 911 Step 6: Facilitated Discussion and Collaborative Write-Up of 912 Manuscript

913 We planned for analysts and initiating authors to discuss the limitations, results, and implications of  
914 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.

915 We built an R package, `ManyEcoEvo::` to conduct the analyses described in this study (Gould et al.  
916 2023), which can be downloaded from <https://github.com/egouldo/ManyEcoEvo/> to reproduce our  
917 analyses or replicate the analyses described here using alternate datasets. Data cleaning and  
918 preparation of analysis-data, as well as the analysis, is conducted in R (R Core Team  
919 2024) reproducibly using the `targets` package (Landau 2021). This data and analysis pipeline is stored  
920 in the `ManyEcoEvo::` package repository and its outputs are made available to users of the package  
921 when the library is loaded.

922 The full manuscript, including further analysis and presentation of results is written in Quarto (Allaire  
923 et al. 2024). The source code to reproduce the manuscript is hosted at  
924 <https://github.com/egouldo/ManyAnalysts/> (Gould et al. 2024), and the rendered version of the  
925 source code may be viewed at <https://egouldo.github.io/ManyAnalysts/>. All R packages and their  
926 versions used in the production of the manuscript are listed in Table 7 at the end of this paper.

## 927 Results

### 928 Summary Statistics

929 In total, 173 analyst teams, comprising 246 analysts, contributed 182 usable analyses (compatible  
930 with our meta-analyses and provided with all information needed for inclusion) of the two datasets  
931 examined in this study which yielded 215 effects. Analysts produced 134 distinct effects that met our  
932 criteria for inclusion in at least one of our meta-analyses for the blue tit dataset. Analysts produced  
933 81 distinct effects meeting our criteria for inclusion for the *Eucalyptus* dataset. Excluded analyses and  
934 effects either did not answer our specified biological questions, were submitted with insufficient  
935 information for inclusion in our meta-analyses, or were incompatible with production of our effect  
936 size(s). We expected cases of this final scenario (incompatible analyses), for instance we cannot  
937 extract a  $Z_r$  from random forest models, which is why we analyzed two distinct types of



938 effects,  $Z_r$  and out-of-sample predictions. Some effects only provided sufficient information for a  
939 subset of analyses and were only included in that subset. For both datasets, most submitted analyses  
940 incorporated mixed effects. Submitted analyses of the blue tit dataset typically specified normal error  
941 and analyses of the *Eucalyptus* dataset typically specified a non-normal error distribution  
942 ([Supplementary Material A, Table A.1](#)).

943 For both datasets, the composition of models varied substantially in regards to the number of fixed  
944 and random effects, interaction terms, and the number of data points used, and these patterns  
945 differed somewhat between the blue tit and *Eucalyptus* analyses (See [Supplementary Material A,](#)  
946 [Table A.2](#)). Focusing on the models included in the  $Z_r$  analyses (because this is the larger sample),  
947 blue tit models included a similar number of fixed effects on average (mean  $5.2 \pm 2.92$  SD, range: 1 to  
948 19) as *Eucalyptus* models (mean  $5.01 \pm 3.83$  SD, range: 1 to 13), but the standard deviation in  
949 number of fixed effects was somewhat larger in the *Eucalyptus* models. The average number of  
950 interaction terms was much larger for the blue tit models (mean  $0.44 \pm 1.11$  SD, range: 0 to 10) than  
951 for the *Eucalyptus* models (mean  $0.16 \pm 0.65$  SD, range: 0 to 5), but still under 0.5 for both, indicating  
952 that most models did not contain interaction terms. Blue tit models also contained more random  
953 effects (mean  $3.53 \pm 2.08$  SD, range: 0 to 10) than *Eucalyptus* models (mean  $1.41 \pm 1.09$  SD, range: 0  
954 to 4). The maximum possible sample size in the blue tit dataset (3720 nestlings) was an order of  
955 magnitude larger than the maximum possible in the *Eucalyptus* dataset (351 plots), and the means  
956 and standard deviations of the sample size used to derive the effects eligible for our study were also  
957 an order of magnitude greater for the blue tit dataset (mean  $2611.09 \pm 937.48$  SD, range: 76 to 76)  
958 relative to the *Eucalyptus* models (mean  $298.43 \pm 106.25$  SD, range: 18 to 351). However, the  
959 standard deviation in sample size from the *Eucalyptus* models was heavily influenced by a few cases  
960 of dramatic sub-setting (described below). Approximately three quarters of *Eucalyptus* models used  
961 sample sizes within 3% of the maximum. In contrast, fewer than 20% of blue tit models relied on  
962 sample sizes within 3% of the maximum, and approximately 50% of blue tit models relied on sample  
963 sizes 29% or more below the maximum.

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

- 970 • *Mixed*: some evidence supporting a positive effect, some evidence supporting a negative effect
- 971 • *Conclusive negative*: negative relationship described without caveat
- 972 • *Qualified negative*: negative relationship but only in certain circumstances or where analysts  
973 express uncertainty in their result
- 974 • *Conclusive none*: analysts interpret the results as conclusive of no effect
- 975 • *Qualified none*: analysts describe finding no evidence of a relationship but they describe the  
976 potential for an undetected effect
- 977 • *Qualified positive*: positive relationship described but only in certain circumstances or where  
978 analysts express uncertainty in their result
- 979 • *Conclusive positive*: positive relationship described without caveat

980 For the blue tit dataset, most analysts concluded that there was negative relationship between  
981 measures of sibling competition and nestling growth, though half the teams expressed qualifications  
982 or described effects as mixed or absent. No analysts concluded that there was a positive relationship

983 even though some individual effect sizes were positive, apparently because all analysts who  
 984 produced effects indicating positive relationships also produced effects indicating negative  
 985 relationships and therefore described their results as qualified, mixed, or absent. For  
 986 the *Eucalyptus* dataset, there was a broader spread of conclusions with at least one analyst team  
 987 providing conclusions consistent with each conclusion category. The most common conclusion for  
 988 the *Eucalyptus* dataset was that there was no relationship between grass cover  
 989 and *Eucalyptus* recruitment (either conclusive or qualified description of no relationship), but more  
 990 than half the teams concluded that there were effects; negative, positive, or mixed.

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

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

992

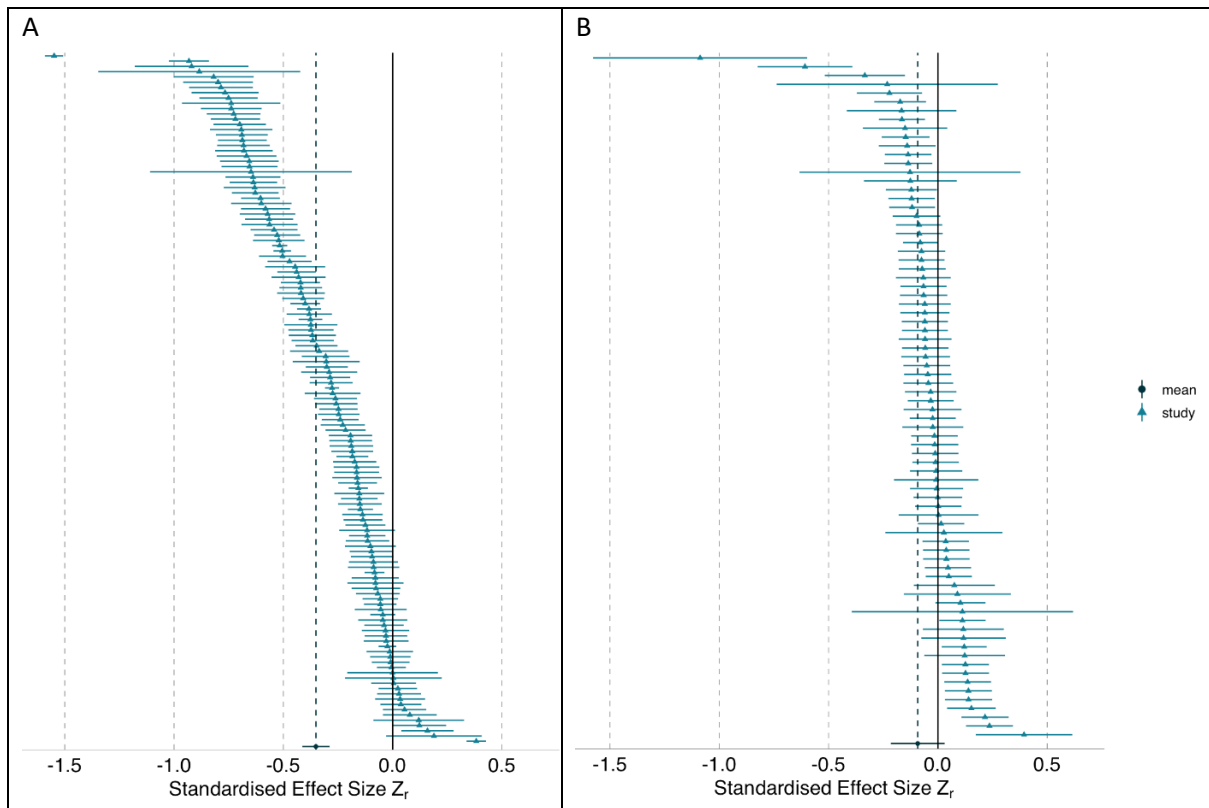
## 993 Distribution of effects

### 994 Effect sizes ( $Z_r$ )

995 Although the majority (118 of 131) of the usable  $Z_r$  effects from the blue tit dataset found nestling  
 996 growth decreased with sibling competition, and the meta-analytic mean  $\bar{Z}_r$  (Fisher's transformation  
 997 of the correlation coefficient) was convincingly negative ( $-0.35 \pm 0.06$  95%CI), there was substantial  
 998 variability in the strength and the direction of this effect.  $Z_r$  ranged from -1.55 to 0.38, and  
 999 approximately continuously from -0.93 to 0.19 (Figure 2a and Table 4), and of the 118 effects with  
 1000 negative slopes, 93 had confidence intervals excluding 0. Of the 13 with positive slopes indicating  
 1001 increased nestling growth in the presence of more siblings, 2 had confidence intervals excluding zero  
 1002 (Figure 2a).

1003 Meta-analysis of the *Eucalyptus* dataset also showed substantial variability in the strength of effects  
 1004 as measured by  $Z_r$ , and unlike with the blue tits, a notable lack of consistency in the direction of  
 1005 effects (Figure 2b, Table 4).  $Z_r$  ranged from -4.47 (Supplementary Material A, Figure A.2), indicating a  
 1006 strong tendency for reduced *Eucalyptus* seedling success as grass cover increased, to 0.39, indicating  
 1007 the opposite. Although the range of reported effects skewed strongly negative, this was due to a  
 1008 small number of substantial outliers. Most values of  $Z_r$  were relatively small with values  $<|0.2|$  and  
 1009 the meta-analytic mean effect size was close to zero ( $-0.09 \pm 0.12$  95%CI). Of the 79 effects, fifty-  
 1010 three had confidence intervals overlapping zero, approximately a quarter (fifteen) crossed the  
 1011 traditional threshold of statistical significance indicating a negative relationship between grass cover  
 1012 and seedling success, and eleven crossed the significance threshold indicating a positive relationship  
 1013 between grass cover and seedling success (Figure 2b).

1014



1015

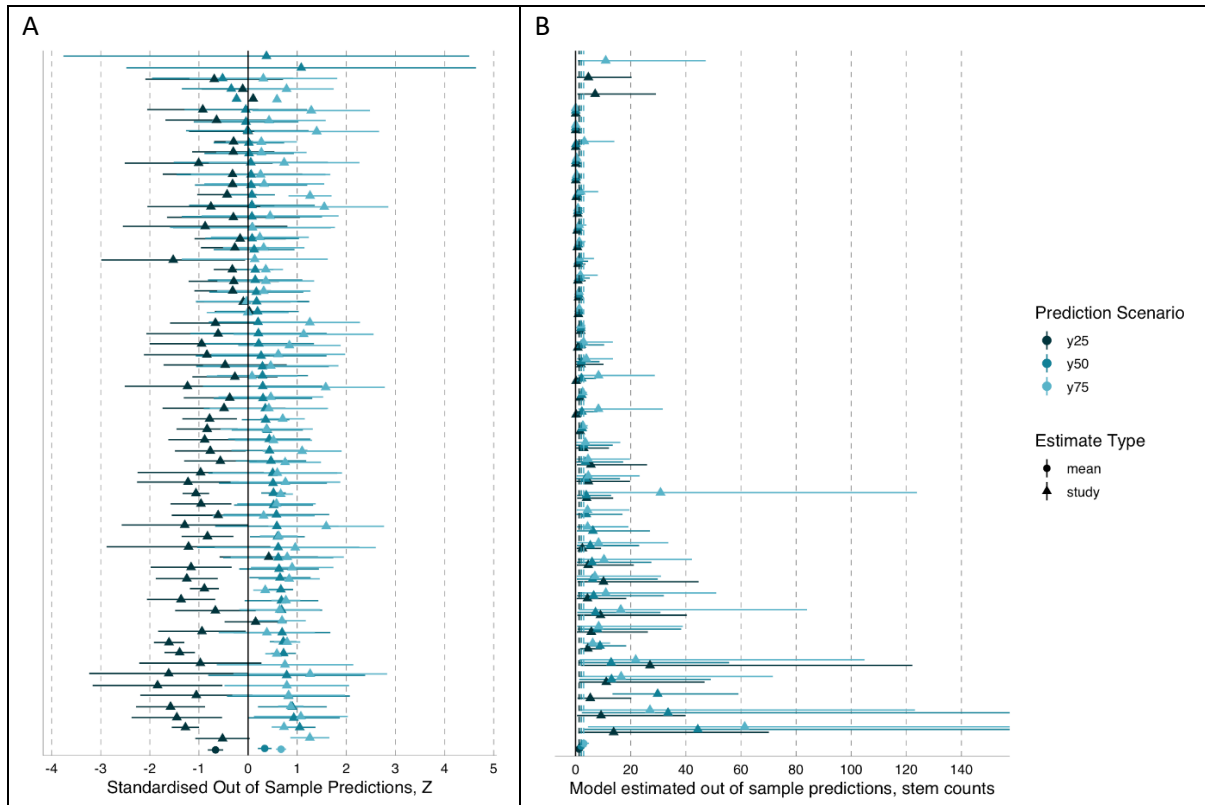
1016 Figure 2: Forest plots of meta-analytic estimated standardized effect sizes ( $Z_r$ , blue triangles) and  
 1017 their 95% confidence intervals for each effect size included in the meta-analysis model. (A) Blue tit  
 1018 analyses: Points where  $Z_r$  are less than 0 indicate analyses that found a negative relationship  
 1019 between sibling number and nestling growth. (B) *Eucalyptus* analyses: Points where  $Z_r$  are less than 0  
 1020 indicate a negative relationship between grass cover and *Eucalyptus* seedling success. The meta-  
 1021 analytic mean effect size is denoted by a black circle and a dashed vertical line, with error bars also  
 1022 representing the 95% confidence interval. The solid black vertical line demarcates effect size of 0,  
 1023 indicating no relationship between the test variable and the response variable. Note that  
 1024 the *Eucalyptus* plot omits one extreme outlier with the value of -4.47 ([Supplementary Material A,](#)  
 1025 [Figure A.2](#)) in order to standardize the x-axes on these two panels.

## 1026 Out-of-sample predictions ( $y_i$ )

1027 As with the effect size  $Z_r$ , we observed substantial variability in the size of out-of-sample predictions  
 1028 derived from the analysts' models. Blue tit predictions (Figure 3a), which were z-score-standardised  
 1029 to accommodate the use of different response variables, always ranged far in excess of one standard  
 1030 deviation. In the  $y_{25}$  scenario, model predictions ranged from -1.84 to 0.42 (a range of 2.68 standard  
 1031 deviations), in the  $y_{50}$  they ranged from -0.52 to 1.08 (a range of 1.63 standard deviations), and in  
 1032 the  $y_{75}$  scenario they ranged from -0.03 to 1.59 (a range of 1.9 standard deviations). As should be  
 1033 expected given the existence of both negative and positive  $Z_r$  values, all three out-of-sample  
 1034 scenarios produced both negative and positive predictions, although as with the  $Z_r$  values, there is a  
 1035 clear trend for scenarios with more siblings to be associated with smaller nestlings. This is supported  
 1036 by the meta-analytic means of these three sets of predictions which were -0.66 (95%CI -0.82–0.5) for  
 1037 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  $y_{75}$ .

1038 *Eucalyptus* out-of-sample predictions also varied substantially (Figure 3b), but because they were not  
 1039 z-score-standardised and are instead on the original count scale, the types of interpretations we can

1040 make differ. The predicted *Eucalyptus* seedling counts per 15 x 15 m plot for the  $y_{25}$  scenario ranged  
 1041 from 0.04 to 26.99, for the  $y_{50}$  scenario ranged from 0.04 to 44.34, and for the  $y_{75}$  scenario they  
 1042 ranged from 0.03 to 61.34. The meta-analytic mean predictions for these three scenarios were  
 1043 similar; 1.27 (95%CI 0.59-2.3) for the  $y_{25}$ , 2.92 (95%CI 0.98-3.89) for the  $y_{50}$ , and 2.92 (95%CI 1.59-  
 1044 4.9) for the  $y_{75}$  scenarios respectively.



1045

1046

1047 Figure 3: Forest plot of meta-analytic estimated out-of-sample predictions. A) Standardized (z-score)  
 1048 blue tit out-of-sample predictions,  $y_i$ . B) response-scale (stem counts) *Eucalyptus* out-of-sample  
 1049 predictions. Triangles represent individual estimates. Circles represent the meta-analytic mean for  
 1050 each prediction scenario. Dark-blue points correspond to  $y_{25}$  scenario, medium-blue points  
 1051 correspond to the  $y_{50}$  scenario, while light blue points correspond to the  $y_{75}$  scenario. Error bars are  
 1052 95% confidence intervals. Note that, for the *Eucalyptus* analysis, outliers (observations more than 3  
 1053 SD above the mean) have been removed prior to model fitting and do not appear on this figure. The  
 1054 x-axis is truncated to approximately 140, and thus some error bars are incomplete.  
 1055 See [Supplementary Material B, Figure B.6](#) for full figure.

1056

## 1057 Quantifying heterogeneity

### 1058 Effect sizes ( $Z_r$ )

1059 We quantified both absolute ( $\tau^2$ ) and relative ( $I^2$ ) heterogeneity resulting from analytical variation.  
 1060 Both measures suggest that substantial variability among effect sizes was attributable to the  
 1061 analytical decisions of analysts.

1062 The total absolute level of variance beyond what would typically be expected due to sampling  
 1063 error,  $\tau^2$  (Table 2), among all usable blue tit effects was 0.08 and for *Eucalyptus* effects was 0.27. This  
 1064 is similar to or exceeding the median value (0.105) of  $\tau^2$  found across 31 recent meta-  
 1065 analyses (calculated from the data in [Yang et al. 2023](#)). The similarity of our observed values to  
 1066 values from meta-analyses of different studies based on different data suggest the potential for a  
 1067 large portion of heterogeneity to arise from analytical decisions. For further discussion of  
 1068 interpretation of  $\tau^2$  in our study, please consult discussion of post hoc analyses below.

1069 Table 2: Heterogeneity in the estimated effects  $Z_r$  for meta-analyses of: the full dataset, as well as  
 1070 from post hoc analyses wherein analyses with outliers are removed, analyses with effects from  
 1071 analysis teams with at least one “unpublishable” rating are excluded, analyses receiving at least one  
 1072 “major revisions” rating or worse excluded, analyses from teams with at least one analyst self-rated  
 1073 as “highly proficient” or “expert” in statistical analysis are included, and (blue tit only) analyses that  
 1074 did not included the pair of highly collinear predictors together.  $\tau^2_{\text{Team}}$  is the absolute heterogeneity  
 1075 for the random effect Team.  $\tau^2_{\text{Effect ID}}$  is the absolute heterogeneity for the random effect Effect  
 1076 ID nested under Team. Effect ID is the unique identifier assigned to each individual statistical effect  
 1077 submitted by an analysis team. We nested Effect ID within analysis team identity (Team) because  
 1078 analysis teams often submitted >1 statistical effect, either because they considered >1 model or  
 1079 because they derived >1 effect per model, especially when a model contained a factor with multiple  
 1080 levels that produced >1 contrast.  $\tau^2_{\text{Total}}$  is the total absolute heterogeneity.  $I^2_{\text{Total}}$  is the proportional  
 1081 heterogeneity; the proportion of the variance among effects not attributable to sampling  
 1082 error,  $I^2_{\text{Team}}$  is the subset of the proportional heterogeneity due to differences  
 1083 among Teams and  $I^2_{\text{Team, Effect ID}}$  is subset of the proportional heterogeneity attributable to among-  
 1084 Effect ID differences.

Dataset	N <sub>Obs</sub>	$\tau^2_{\text{Total}}$	$\tau^2_{\text{Team}}$	$\tau^2_{\text{Effect ID}}$	$I^2_{\text{Total}}$	$I^2_{\text{Team}}$	$I^2_{\text{Team, Effect ID}}$
All Analyses							
<i>Eucalyptus</i>	79	0.27	0.02	0.25	98.59%	6.89%	91.70%
blue tit	131	0.08	0.03	0.05	97.61%	36.71%	60.90%
Blue tit analyses containing highly collinear predictors removed							
blue tit	117	0.07	0.04	0.03	96.92%	58.18%	38.75%
All analyses, outliers removed							
<i>Eucalyptus</i>	75	0.01	0.00	0.01	66.19%	19.25%	46.94%
blue tit	127	0.07	0.04	0.02	96.84%	64.63%	32.21%
Analyses receiving at least one 'Unpublishable' rating removed							
<i>Eucalyptus</i>	55	0.01	0.01	0.01	79.74%	28.31%	51.43%
blue tit	109	0.08	0.03	0.05	97.52%	35.68%	61.84%
Analyses receiving at least one 'Unpublishable' and or 'Major Revisions' rating removed							
<i>Eucalyptus</i>	13	0.03	0.03	0.00	88.91%	88.91%	0.00%
blue tit	32	0.14	0.01	0.13	98.72%	5.17%	93.55%
Analyses from teams with highly proficient or expert data analysts							
<i>Eucalyptus</i>	34	0.58	0.02	0.56	99.41%	3.47%	95.94%
blue tit	89	0.09	0.03	0.06	97.91%	31.43%	66.49%

1085

1086 In our analyses,  $I^2$  is a plausible index of how much more variability among effect sizes we have  
 1087 observed, as a proportion, than we would have observed if sampling error were driving variability.  
 1088 We discuss our interpretation of  $I^2$  further in the methods, but in short, it is a useful metric for  
 1089 comparison to values from published meta-analyses and provides a plausible value for how much

1090 heterogeneity could arise in a normal meta-analysis with similar sample sizes due to analytical  
 1091 variability alone. In our study, total  $I^2$  for the blue tit  $Z_r$  estimates was extremely large, at 97.61%, as  
 1092 was the *Eucalyptus* estimate (98.59% Table 2).

1093 Although the overall  $I^2$  values were similar for both *Eucalyptus* and blue tit analyses, the relative  
 1094 composition of that heterogeneity differed. For both datasets, the majority of heterogeneity  
 1095 in  $Z_r$  was driven by differences among effects as opposed to differences among teams, though this  
 1096 was more prominent for the *Eucalyptus* dataset, where nearly all of the total heterogeneity was  
 1097 driven by differences among effects (91.7%) as opposed to differences among teams (6.89%)  
 1098 (Table 2).

### 1099 Out-of-sample predictions ( $y_i$ )

1100 We observed substantial heterogeneity among out-of-sample estimates, but the pattern differed  
 1101 somewhat from the  $Z_r$  values (Table 3). Among the blue tit predictions,  $I^2$  ranged from medium-high  
 1102 for the  $y_{25}$  scenario (68.54) to low (27.9) for the  $y_{75}$  scenario. Among  
 1103 the *Eucalyptus* predictions,  $I^2$  values were uniformly high (>82%). For both datasets, most of the  
 1104 existing heterogeneity among predicted values was attributable to among-team differences, with the  
 1105 exception of the  $y_{50}$  analysis of the *Eucalyptus* dataset. We are limited in our interpretation of  $\tau^2$  for  
 1106 these estimates because, unlike for the  $Z_r$  estimates, we have no benchmark for comparison with  
 1107 other meta-analyses.

1108 Table 3: Heterogeneity among the out-of-sample predictions  $y_i$  for both blue tit  
 1109 and *Eucalyptus* datasets.  $\tau^2_{\text{Team}}$  is the absolute heterogeneity for the random effect Team.  $T^2_{\text{Effect ID}}$  is  
 1110 the absolute heterogeneity for the random effect Effect ID nested under Team. Effect ID is the unique  
 1111 identifier assigned to each individual statistical effect submitted by an analysis team. We  
 1112 nested Effect ID within analysis team identity (Team) because analysis teams often submitted >1  
 1113 statistical effect, either because they considered >1 model or because they derived >1 effect per  
 1114 model, especially when a model contained a factor with multiple levels that produced >1  
 1115 contrast.  $\tau^2_{\text{Total}}$  is the total absolute heterogeneity.  $I^2_{\text{Total}}$  is the proportional heterogeneity; the  
 1116 proportion of the variance among effects not attributable to sampling error,  $I^2_{\text{Team}}$  is the subset of the  
 1117 proportional heterogeneity due to differences among Teams and  $I^2_{\text{Team, Effect ID}}$  is subset of the  
 1118 proportional heterogeneity attributable to among-Effect ID differences.

Prediction Scenario	N <sub>Obs</sub>	T <sub>Total</sub>	T <sup>2</sup> <sub>Team</sub>	T <sup>2</sup> <sub>Effect ID</sub>	I <sup>2</sup> <sub>Total</sub>	I <sup>2</sup> <sub>Team</sub>	I <sup>2</sup> <sub>Team, Effect ID</sub>
<b>blue tit</b>							
y25	63	0.23	0.11	0.03	68.54%	53.43%	15.11%
y50	60	0.23	0.06	0.00	50%	46.29%	3.71%
y75	63	0.23	0.02	0.00	27.9%	27.89%	0.01%
<b><i>Eucalyptus</i></b>							
y25	38	5.75	1.48	0.68	86.93%	59.54%	27.39%
y50	38	5.75	1.32	0.83	89.63%	55%	34.64%
y75	38	5.75	1.03	0.41	80.19%	57.41%	22.78%

1119

## 1120 *Post-hoc* analysis: Exploring outlier characteristics and the effect of 1121 outlier removal on heterogeneity

### 1122 Effect sizes ( $Z_r$ )

1123 The outlier *Eucalyptus*  $Z_r$  values were striking and merited special examination. The three negative  
1124 outliers had very low sample sizes that were based on either small subsets of the dataset or, in one  
1125 case, extreme aggregation of data. The outliers associated with small subsets had sample sizes  
1126 ( $n= 117, 90, 18$ ) that were less than half of the total possible sample size of 351. The case of extreme  
1127 aggregation involved averaging all values within each of the 351 sites in the dataset.

1128 Surprisingly, both the largest and smallest effect sizes in the blue tit analyses (Figure 2a) come from  
1129 the same analyst (anonymous ID: 'Adelong'), with identical models in terms of the explanatory  
1130 variable structure, but with different response variables. However, the radical change in effect was  
1131 primarily due to collinearity with covariates. The primary predictor variable (brood count after  
1132 manipulation) was accompanied by several collinear variables, including the highly collinear  
1133 (correlation of 0.89 [Supplementary Material D, Figure D.2](#)) covariate (brood count at day 14) in both  
1134 analyses. In the analysis of nestling weight, brood count after manipulation showed a strong positive  
1135 partial correlation with weight after controlling for brood count at day 14 and treatment category  
1136 (increased, decreased, unmanipulated). In that same analysis, the most collinear covariate (the day  
1137 14 count) had a negative partial correlation with weight. In the analysis with tarsus length as the  
1138 response variable, these partial correlations were almost identical in absolute magnitude, but  
1139 reversed in sign and so brood count after manipulation was now the collinear predictor with the  
1140 negative relationship. The two models were therefore very similar, but the two collinear predictors  
1141 simply switched roles, presumably because a subtle difference in the distribution of weight and  
1142 tarsus length data.

1143 When we dropped the *Eucalyptus* outliers,  $I^2$  decreased from high (98.59 %), using Higgins' (Higgins  
1144 et al. 2003) suggested benchmark, to between moderate and high (66.19 %, Table 2). However, more  
1145 notably,  $\tau^2$  dropped from 0.27 to 0.01, indicating that, once outliers were excluded, the observed  
1146 variation in effects was similar to what we would expect if sampling error were driving the  
1147 differences among effects (since  $\tau^2$  is the variance beyond that driven by sampling error). The  
1148 interpretation of this value of  $\tau^2$  in the context of our many-analyst study is somewhat different than  
1149 a typical meta-analysis, however, since in our study (especially for *Eucalyptus*, where most analyses  
1150 used almost exactly the same data points), there is almost no role for sampling error in driving the  
1151 observed differences among the estimates. Thus, rather than concluding that the variability we  
1152 observed among estimates (after removing outliers) was due only to sampling  
1153 error (because  $\tau^2$  became small: 10% of the median from Yang et al. 2023), we instead conclude that  
1154 the observed variability, which must be due to the divergent choices of analysts rather than sampling  
1155 error, is approximately of the same magnitude as what we would have expected if, instead, sampling  
1156 error, and not analytical heterogeneity, were at work. Conversely, dropping outliers from the set of  
1157 blue tit effects did not meaningfully reduce  $I^2$ , and only modestly reduced  $\tau^2$  (Table 2). Thus, effects  
1158 at the extremes of the distribution were much stronger contributors to total heterogeneity for effects  
1159 from analyses of the *Eucalyptus* than for the blue tit dataset.

1160 Table 4: Estimated mean value of the standardised correlation coefficient,  $\bar{Z}_r$ , along with its standard  
1161 error and 95% confidence intervals. We re-computed the meta-analysis for different post hoc subsets  
1162 of the data: All eligible effects, removal of effects from blue tit analyses that contained a pair of  
1163 highly collinear predictor variables, removal of effects from analysis teams that received at least one

1164 peer rating of “deeply flawed and unpublishable”, removal of any effects from analysis teams that  
 1165 received at least one peer rating of either “deeply flawed and unpublishable” or “publishable with  
 1166 major revisions”, inclusion of only effects from analysis teams that included at least one member  
 1167 who rated themselves as “highly proficient” or “expert” at conducting statistical analyses in their  
 1168 research area.

Dataset	$\hat{\mu}$	$SE[\hat{\mu}]$	95% CI	statistic	p
All analyses					
<i>Eucalyptus</i>	-0.09	0.06	[-0.22,0.03]	-1.47	0.14
blue tit	-0.35	0.03	[-0.41,-0.29]	-11.02	<0.001
Blue tit analyses containing highly collinear predictors removed					
blue tit	-0.36	0.03	[-0.42,-0.29]	-10.97	<0.001
All analyses, outliers removed					
<i>Eucalyptus</i>	-0.03	0.01	[-0.06,0.00]	-2.23	0.026
blue tit	-0.36	0.03	[-0.42,-0.30]	-11.48	<0.001
Analyses receiving at least one 'Unpublishable' rating removed					
<i>Eucalyptus</i>	-0.02	0.02	[-0.07,0.02]	-1.15	0.3
blue tit	-0.36	0.03	[-0.43,-0.30]	-10.82	<0.001
Analyses receiving at least one 'Unpublishable' and or 'Major Revisions' rating removed					
<i>Eucalyptus</i>	-0.04	0.05	[-0.15,0.07]	-0.77	0.4
blue tit	-0.37	0.07	[-0.51,-0.23]	-5.34	<0.001
Analyses from teams with highly proficient or expert data analysts					
<i>Eucalyptus</i>	-0.17	0.13	[-0.43,0.10]	-1.24	0.2
blue tit	-0.36	0.04	[-0.44,-0.28]	-8.93	<0.001

1169

## 1170 Out-of-sample predictions ( $y_i$ )

1171 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of  
 1172 eligible effects was smaller and the pattern of outliers differed.

## 1173 Post hoc analysis: Exploring the effect of removing analyses with poor 1174 peer ratings on heterogeneity

### 1175 Effect sizes ( $Z_r$ )

1176 Removing poorly rated analyses had limited impact on the meta-analytic means ([Supplementary](#)  
 1177 [Material B, Figure B.3](#)). For the *Eucalyptus* dataset, the meta-analytic mean shifted from -0.09 to -  
 1178 0.02 when effects from analyses rated as unpublishable were removed, and to -0.04 when effects  
 1179 from analyses rated, at least once, as unpublishable or requiring major revisions were removed.  
 1180 Further, the confidence intervals for all of these means overlapped each of the other means  
 1181 (Table 4). We saw similar patterns for the blue tit dataset, with only small shifts in the meta-analytic  
 1182 mean, and confidence intervals of all three means overlapping each other mean (Table 4). Refitting  
 1183 the meta-analysis with a fixed effect for categorical ratings also showed no indication of differences  
 1184 in group meta-analytic means due to peer ratings ([Supplementary Material B, Figure B.1](#)).

1185 For the blue tit dataset, removing poorly-rated analyses led to only negligible changes in  $I^2_{\text{Total}}$  and  
 1186 relatively minor impacts on  $\tau^2$ . However, for the *Eucalyptus* dataset, removing poorly-rated analyses  
 1187 led to notable reductions in  $I^2_{\text{Total}}$  and substantial reductions in  $\tau^2$ . When including all analyses,  
 1188 the *Eucalyptus*  $I^2_{\text{Total}}$  was 98.59% and  $\tau^2$  was 0.27, but eliminating analyses with ratings of



1189 “unpublishable” reduced  $I^2_{\text{Total}}$  to 79.74% and  $\tau^2$  to 0.01, and removing also those analyses “needing  
1190 major revisions” left  $I^2_{\text{Total}}$  at 88.91% and  $\tau^2$  at 0.03 (Table 2). Additionally, the allocations of  $I^2$  to the  
1191 team versus individual effect were altered for both blue tit and *Eucalyptus* meta-analyses by  
1192 removing poorly-rated analyses, but in different ways. For blue tit meta-analysis, between a third and  
1193 two-thirds of the total  $I^2$  was attributable to among-team variance in most analyses until both  
1194 analyses rated “unpublishable” and analyses rated in need of “major revision” were eliminated, in  
1195 which case almost all remaining heterogeneity was attributable to among-effect differences. In  
1196 contrast, for *Eucalyptus* meta-analysis, the among-team component of  $I^2$  was less than third until  
1197 both analyses rated “unpublishable” and analyses rated in need of “major revision” were eliminated,  
1198 in which case almost 90% of heterogeneity was attributable to differences among teams.

### 1199 Out-of-sample predictions ( $y_i$ )

1200 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of  
1201 eligible effects was smaller and our ability to interpret heterogeneity values for these analyses was  
1202 limited

### 1203 Post hoc analysis: Exploring the effect of including only analyses 1204 conducted by analysis teams with at least one member self-rated as 1205 “highly proficient” or “expert” in conducting statistical analyses in 1206 their research area

### 1207 Effect sizes ( $Z_r$ )

1208 Including only analyses conducted by teams that contained at least one member who rated  
1209 themselves as “highly proficient” or “expert” in conducting the relevant statistical methods had  
1210 negligible impacts on the meta-analytic means (Table 4), the distribution of  $Z_r$  effects  
1211 ([Supplementary Material B, Figure B.4](#)), or heterogeneity estimates (Table 2), which remained  
1212 extremely high.

### 1213 Out-of-sample predictions ( $y_i$ )

1214 We did not conduct these post hoc analyses on the out-of-sample predictions as the number of  
1215 eligible effects was smaller.

### 1216 Post hoc analysis: Exploring the effect of excluding estimates of $Z_r$ in 1217 which we had reduced confidence

1218 As described in our addendum to the methods, we identified a subset of estimates of  $Z_r$  in which we  
1219 had less confidence because of features of the submitted degrees of freedom. Excluding these effects  
1220 in which we had lower confidence had minimal impact on the meta-analytic mean and the estimates  
1221 of total  $I^2$  and  $\tau^2$  for both blue tit and *Eucalyptus* meta-analyses, regardless of whether outliers were  
1222 also excluded ([Supplementary Material B, Table B.1](#)).

1223 **Post hoc analysis: Exploring the effect of excluding effects from blue**  
 1224 **tit models that contained two highly collinear predictors**

1225 **Effect sizes ( $Z_r$ )**

1226 Excluding effects from blue tit models that contained the two highly collinear predictors (brood count  
 1227 after manipulation and brood count at day 14) had negligible impacts on the meta-analytic means  
 1228 (Table 4), the distribution of  $Z_r$  effects ([Supplementary Material B, Figure B.5](#)), or heterogeneity  
 1229 estimates (Table 2), which remained high.

1230 **Out-of-sample predictions**

1231 Inclusion of collinear predictors does not harm model prediction, and so we did not conduct these  
 1232 post hoc analyses.

1233 **Explaining Variation in Deviation Scores**

1234 None of the pre-registered predictors explained substantial variation in deviation among submitted  
 1235 statistical effects from the meta-analytic mean (Table 5, Table 6).

1236 Table 5: Summary metrics for registered models seeking to explain deviation (Box-Cox transformed  
 1237 absolute deviation scores) from  $\bar{Z}_r$  as a function of Sorensen’s Index, categorical peer ratings, and  
 1238 continuous peer ratings for blue tit and *Eucalyptus* analyses, and as a function of the presence or  
 1239 absence of random effects (in the analyst’s models) for *Eucalyptus* analyses. We report coefficient of  
 1240 determination,  $R^2$ , for our models including only fixed effects as predictors of deviation, and we  
 1241 report  $R^2_{\text{Conditional}}$ ,  $R^2_{\text{Marginal}}$  and the intra-class correlation (ICC) from our models that included both  
 1242 fixed and random effects. For all our models, we calculated the residual standard deviation  $\sigma$  and  
 1243 root mean squared error (RMSE).

Dataset	NObs	$R^2$	$R^2_{\text{Conditional}}$	$R^2_{\text{Marginal}}$	ICC	$\sigma$	RMSE
Deviation explained by categorical ratings							
<i>Eucalyptus</i>	346		0.13	0.01	0.12	1.06	1.02
blue tit	473		0.09	$7.47 \times 10^{-3}$	0.08	0.5	0.48
Deviation explained by continuous ratings							
<i>Eucalyptus</i>	346		0.12	$7.44 \times 10^{-3}$	0.11	1.06	1.03
blue tit	473		0.09	$3.44 \times 10^{-3}$	0.09	0.5	0.48
Deviation explained by Sorensen's index							
<i>Eucalyptus</i>	79	$1.84 \times 10^{-4}$				1.12	1.1
blue tit	131	$6.32 \times 10^{-3}$				0.51	0.51
Deviation explained by inclusion of random effects							
<i>Eucalyptus</i>	79	$8.75 \times 10^{-8}$				1.12	1.1

1244

1245 Table 6: Parameter estimates from models of Box-Cox transformed deviation scores from  $\bar{Z}_r$  as a  
 1246 function of continuous and categorical peer ratings, Sorensen scores, and the inclusion of random  
 1247 effects. Standard Errors (SE), 95% confidence intervals (95% CI) are reported for all estimates, while  $t$   
 1248 values, degrees of freedom and p-values are presented for fixed-effects. Note that positive  
 1249 parameter estimates mean that as the predictor variable increases, so does the absolute value of the  
 1250 deviation from the meta-analytic mean.

Parameter	Random effect	Coefficient	SE	95% CI	t	df	p
Deviation explained by inclusion of random effects - <i>Eucalyptus</i>							
(Intercept)		-2.53	0.27	[-3.06, -1.99]	-9.31	77	<0.001
Mixed model		0.00	0.31	[-0.60, 0.60]	0.00	77	>0.9
Deviation explained by Sorensen's index - <i>Eucalyptus</i>							
(Intercept)		-2.65	1.05	[-4.70, -0.60]	-2.53	77	0.011
Mean Sorensen's index		0.18	1.51	[-2.78, 3.14]	0.12	77	>0.9
Deviation explained by Sorensen's index - blue tit							
(Intercept)		-1.53	0.28	[-2.08, -0.98]	-5.42	129	<0.001
Mean Sorensen's index		0.42	0.47	[-0.49, 1.34]	0.91	129	0.4
Deviation explained by continuous ratings - <i>Eucalyptus</i>							
(Intercept)		-2.23	0.23	[-2.69, -1.78]	-9.65	342	<0.001
RateAnalysis		-0.004	0	[-0.011, 0]	-1.44	342	0.15
SD (Intercept)	Reviewer ID	0.37	0.09	[ 0.24, 0.60]			
SD (Observations)	Residual	1.06	0.04	[0.98, 1.15]			
Deviation explained by continuous ratings - blue tit							
(Intercept)		-1.16	0.11	[-1.37, -0.94]	-10.60	469	<0.001
RateAnalysis		-0.002	0	[-0.004, 0]	-1.22	469	0.2
SD (Intercept)	Reviewer ID	0.16	0.03	[0.10,0.24]			
SD (Observations)	Residual	0.5	0.02	[0.46,0.53]			
Deviation explained by categorical ratings - <i>Eucalyptus</i>							
(Intercept)		-2.66	0.27	[-3.18, -2.13]	-9.97	340	<0.001
Publishable with major revision		0.29	0.29	[-0.27, 0.85]	1.02	340	0.3
Publishable with minor revision		0.01	0.28	[-0.54, 0.56]	0.04	340	>0.9
Publishable as is		0.05	0.31	[-0.55, 0.66]	0.17	340	0.9
SD (Intercept)	Reviewer ID	0.39	0.09	[ 0.25, 0.61]			
SD (Observations)	Residual	1.06	0.04	[0.98, 1.15]			
Deviation explained by categorical ratings - blue tit							
(Intercept)		-1.11	0.11	[-1.33, -0.89]	-9.91	467	<0.001
Publishable with major revision		-0.19	0.12	[-0.42, 0.04]	-1.62	467	0.10
Publishable with minor revision		-0.19	0.12	[-0.42, 0.04]	-1.65	467	0.10
Publishable as is		-0.13	0.13	[-0.39, 0.12]	-1.02	467	0.3
SD (Intercept)	Reviewer ID	0.15	0.04	[ 0.10, 0.24]			
SD (Observations)	Residual	0.5	0.02	[0.46, 0.53]			

1251

## 1252 Deviation scores as explained by reviewer ratings

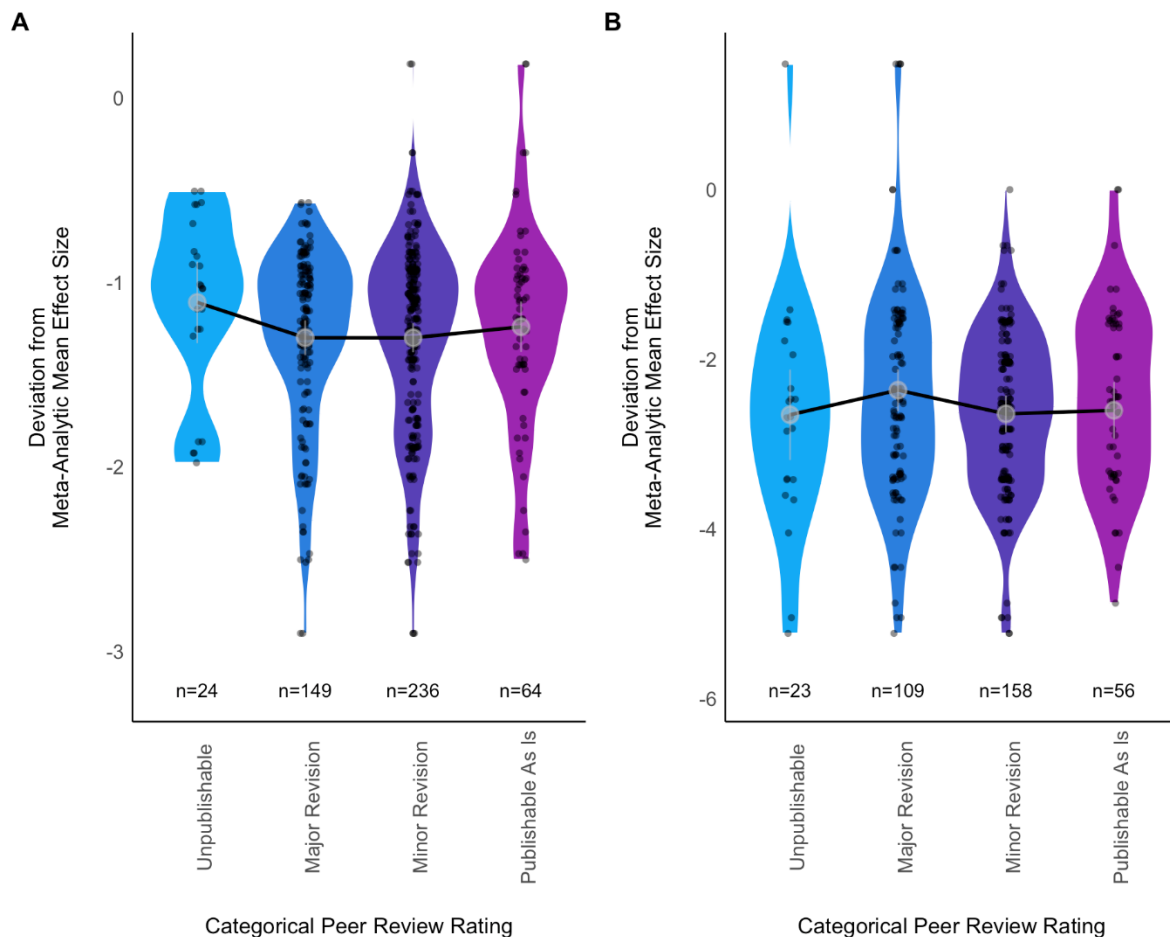
### 1253 Effect sizes ( $Z_r$ )

1254 We obtained reviews from 153 reviewers who reviewed analyses for a mean of 3.27 (range 1 - 11)  
1255 analysis teams. Analyses of the blue tit dataset received a total of 240 reviews, each was reviewed by  
1256 a mean of 3.87 (SD 0.71, range 3-5) reviewers. Analyses of the *Eucalyptus* dataset received a total of  
1257 178 reviews, each was reviewed by a mean of 4.24 (SD 0.79, range 3-6) reviewers. We tested for  
1258 inter-rater-reliability (IRR) to examine how similarly reviewers reviewed each analysis and found

1259 approximately no agreement among reviewers. When considering continuous ratings, IRR was 0.01,  
1260 and for categorical ratings, IRR was -0.14.

1261 Many of the models of deviation as a function of peer ratings faced issues of failure to converge or  
1262 singularity due to sparse design matrices with our pre-registered random effects (Effect ID and  
1263 Reviewer ID) (see [Supplementary Material C](#)). These issues persisted after increasing the tolerance  
1264 and changing the optimizer. For both *Eucalyptus* and blue tit datasets, models with continuous  
1265 ratings as a predictor were singular when both pre-registered random effects were included.

1266 When using both categorical and continuous ratings as predictors, only models converged and  
1267 allowed 95% confidence intervals to be calculated when specifying Reviewer ID as a random effect.  
1268 The categorical ratings model had a  $R^2_C$  of 0.09 and a  $R^2_M$  of 0.01, the continuous ratings model had  
1269 a  $R^2_C$  of 0.09 and a  $R^2_M$  of 0.01 for the blue tit dataset and a  $R^2_C$  of 0.12 and a  $R^2_M$  of 0.01 for the  
1270 *Eucalyptus* dataset. Neither continuous or categorical reviewer ratings of the analyses meaningfully  
1271 predicted deviance from the meta-analytic mean (Table 6, Figure 4). We re-ran the multi-level meta-  
1272 analysis with a fixed effect for the categorical publishability ratings and found no difference in mean  
1273 standardised effect sizes among publishability ratings ([Supplementary Material B, Figure B.1](#)).



1274

1275 Figure 4: Violin plot of Box-Cox transformed deviation from meta-analytic mean  $\bar{Z}_r$  as a function of  
1276 categorical peer rating. Grey points for each rating group denote model-estimated marginal mean  
1277 deviation, and error bars denote 95%CI of the estimate. **A** Blue tit dataset, **B** *Eucalyptus* dataset.

## 1278 Out-of-sample predictions ( $y_i$ )

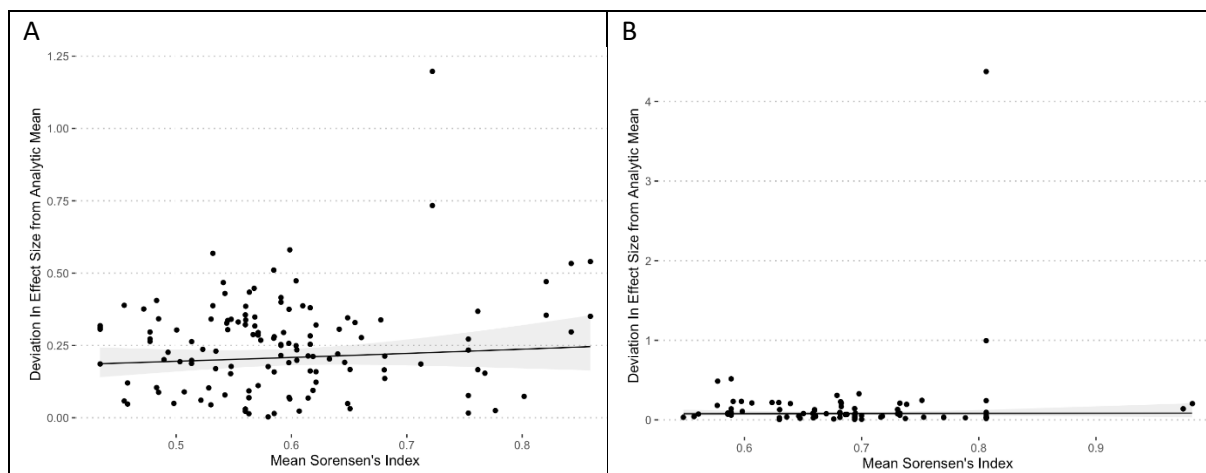
1279 Some models of the influence of reviewer ratings on out-of-sample predictions ( $y_i$ ) had issues with  
1280 convergence and singularity of fit (see [Supplementary Material C, Table C.3](#)) and those models that  
1281 converged and were not singular showed no strong relationship ([Supplementary Material C,](#)  
1282 [Figure C.2, Supplementary Material C, Figure C.3](#)), as with the  $Z_r$  analyses.

## 1283 Deviation scores as explained by the distinctiveness of variables in 1284 each analysis

### 1285 Effect sizes ( $Z_r$ )

1286 We employed Sorensen's index to calculate the distinctiveness of the set of predictor variables used  
1287 in each model (Figure 5). The mean Sorensen's score for blue tit analyses was 0.59 (SD: 0.1, range  
1288 0.43-0.86), and for *Eucalyptus* analyses was 0.69 (SD: 0.08, range 0.55-0.98).

1289 We found no meaningful relationship between distinctiveness of variables selected and deviation  
1290 from the meta-analytic mean (Table 6, Figure 5) for either blue tit (mean 0.42, 95%CI -0.49,1.34)  
1291 or *Eucalyptus* effects (mean 0.18, 95%CI -2.78,3.14).



1292

1293 Figure 5: Fitted model of the Box-Cox-transformed deviation score (deviation in effect size from  
1294 meta-analytic mean) as a function of the mean Sorensen's index showing distinctiveness of the set of  
1295 predictor variables. Grey ribbons on predicted values are 95% CI's. A) blue tit dataset, B) *Eucalyptus*  
1296 dataset.

## 1297 Out-of-sample predictions ( $y_i$ )

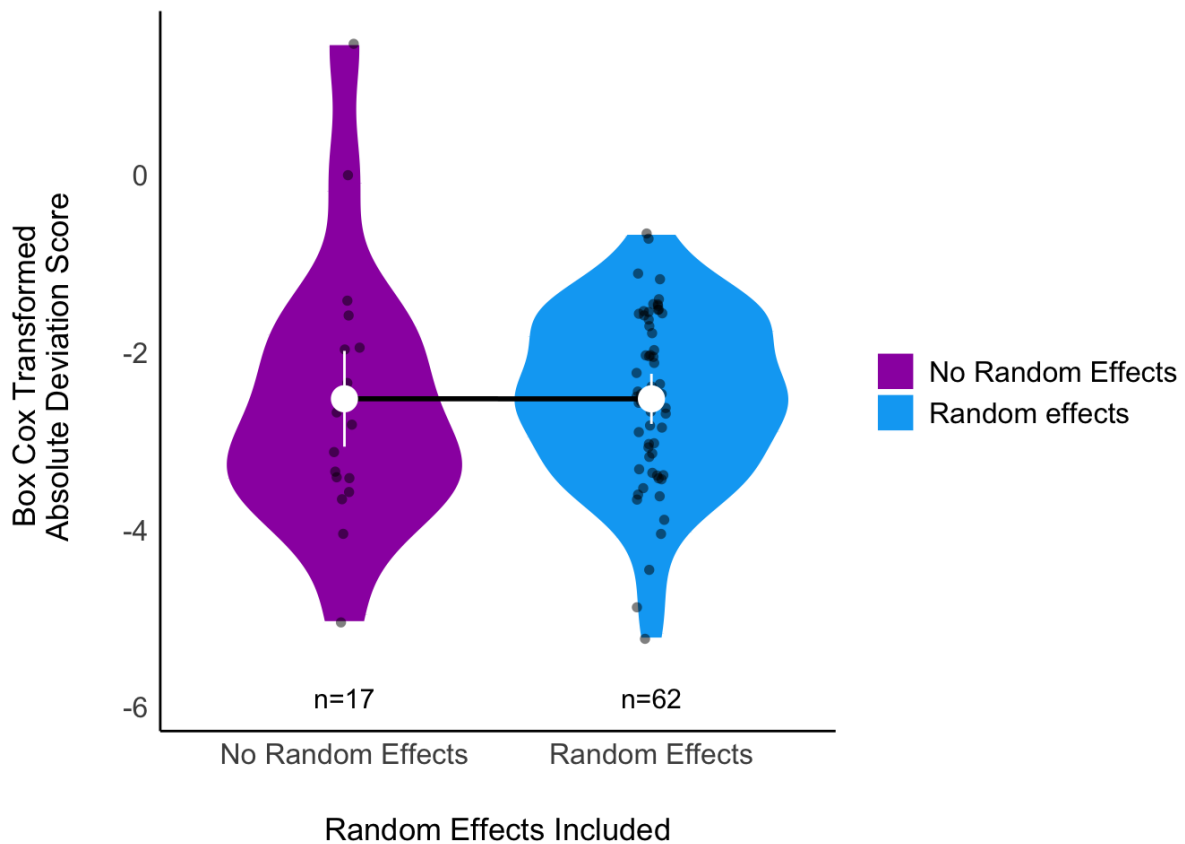
1298 As with the  $Z_r$  estimates, we did not observe any convincing relationships between deviation scores  
1299 of out-of-sample predictions and Sorensen's index values (see [Supplementary Material C4.1](#)).

## 1300 Deviation scores as explained by the inclusion of random effects

### 1301 Effect sizes ( $Z_r$ )

1302 There were only three blue tit analyses that did not include random effects, which is below the pre-  
1303 registered threshold for fitting a model of the Box-Cox transformed deviation from the meta-analytic  
1304 mean as a function of whether the analysis included random-effects. However,  
1305 17 *Eucalyptus* analyses included only fixed effects, which crossed our pre-registered threshold.

1306 Consequently, we performed this analysis for the *Eucalyptus* dataset only. There was no relationship  
1307 between random-effect inclusion and deviation from meta-analytic mean among  
1308 the *Eucalyptus* analyses (Table 6, Figure 6).



1309

1310 Figure 6: Violin plot of mean Box-Cox transformed deviation from meta-analytic mean as a function  
1311 of random-effects inclusion in *Eucalyptus* analyses. White point for each group of analyses denotes  
1312 model-estimated marginal mean deviation, and error bars denote 95% CI of the estimate.

### 1313 Out-of-sample predictions ( $y_i$ )

1314 As with the  $Z_r$  estimates, we did not examine the possibility of a relationship between the inclusion  
1315 of random effects and the deviation scores of the blue tit out-of-sample predictions. When we  
1316 examined the possibility of this relationship for the *Eucalyptus* effects, we found consistent evidence  
1317 of somewhat higher Box-Cox-transformed deviation values for models including a random effect,  
1318 meaning the models including random effects averaged slightly higher deviation from the meta-  
1319 analytic means ([Supplementary Material C, Figure C.5](#)).

### 1320 Multivariate Analysis Effect size ( $Z_r$ ) and out-of-sample predictions ( $y_i$ )

1321 Like the univariate models, the multivariate models did a poor job of explaining deviations from the  
1322 meta-analytic mean. Because we pre-registered a multivariate model that contained collinear  
1323 predictors that produce results which are not readily interpretable, we present these models in the  
1324 supplement. We also had difficulty with convergence and singularity for multivariate models of out-  
1325 of-sample ( $y_i$ ) result, and had to adjust which random effects we included ([Supplementary Material  
1326 C, Table C.8](#)). However, no multivariate analyses of *Eucalyptus* out-of-sample results avoided  
1327 problems of convergence or singularity, no matter which random effects we included  
1328 ([Supplementary Material C, Table C.8](#)). We therefore present no multivariate *Eucalyptus*  $y_i$  models.

1329 We present parameter estimates from multivariate  $Z_r$  models for both datasets (Supplementary  
1330 Material C, [Table C.6](#), [Table C.7](#)) and from  $y_i$  models from the blue tit dataset (Supplementary  
1331 Material C, [Table C.10](#), [Table C.9](#)). We include interpretation of the results from these models in the  
1332 supplement, but the results do not change the interpretations we present above based on the  
1333 univariate analyses.

## 1334 Discussion

1335 When a large pool of ecologists and evolutionary biologists analyzed the same two datasets to  
1336 answer the corresponding two research questions, they produced substantially heterogeneous sets  
1337 of answers. Although the variability in analytical outcomes was high for both datasets, the patterns  
1338 of this variability differed distinctly between them. For the blue tit dataset, there was nearly  
1339 continuous variability across a wide range of  $Z_r$  values. In contrast, for the *Eucalyptus* dataset, there  
1340 was less variability across most of the range, but more striking outliers at the tails. Among out-of-  
1341 sample predictions, there was again almost continuous variation across a wide range (2 SD) among  
1342 blue tit estimates. For *Eucalyptus*, out-of-sample predictions were also notably variable, with about  
1343 half the predicted stem count values at  $<2$  but the other half being much larger, and ranging to  
1344 nearly 40 stems per 15 m x 15 m plot. We investigated several hypotheses for drivers of this  
1345 variability within datasets, but found little support for any of these. Most notably, even when we  
1346 excluded analyses that had received one or more poor peer reviews, the heterogeneity in results  
1347 largely persisted. Regardless of what drives the variability, the existence of such dramatically  
1348 heterogeneous results when ecologists and evolutionary biologists seek to answer the same  
1349 questions with the same data should trigger conversations about how ecologists and evolutionary  
1350 biologists analyze data and interpret the results of their own analyses and those of others in the  
1351 literature (e.g., Silberzahn et al. 2018; Simonsohn, Simmons, and Nelson 2020; Auspurg and Brüderl  
1352 2021; Breznau et al. 2022).

1353 Our observation of substantial heterogeneity due to analytical decisions is consistent with a small  
1354 earlier study in ecology (Stanton-Geddes, de Freitas and de Sales Dambros 2014) and a growing body  
1355 of work from the quantitative social sciences (e.g., Silberzahn et al. 2018; Botvinik-Nezer et al.  
1356 2020; Huntington-Klein et al. 2021; Schweinsberg et al. 2021; Breznau et al. 2022; Coretta et al.  
1357 2023). In these studies, when volunteers from the discipline analyzed the same data, they produced  
1358 a worryingly diverse set of answers to a pre-set question. This diversity included a wide range of  
1359 effect sizes, and in most cases, even involved effects in opposite directions. Thus, our result should  
1360 not be viewed as an anomalous outcome from two particular datasets, but instead as evidence from  
1361 additional disciplines regarding the heterogeneity that can emerge from analyses of complex  
1362 datasets to answer questions in probabilistic science. Not only is our major observation consistent  
1363 with other studies, it is, itself, robust because it derived primarily from simple forest plots that we  
1364 produced based on a small set of decisions that were mostly registered before data gathering and  
1365 which conform to widely accepted meta-analytic practices.

1366 Unlike the strong pattern we observed in the forest plots, our other analyses, both registered  
1367 and post hoc, produced either inconsistent patterns, weak patterns, or the absence of patterns. Our  
1368 registered analyses found that deviations from the meta-analytic mean by individual effect sizes ( $\overline{Zr}$ )  
1369 or the predicted values of the dependent variable ( $\overline{y}$ ) were poorly explained by our hypothesized  
1370 predictors: peer rating of each analysis team's method section, a measurement of the distinctiveness  
1371 of the set of predictor variables included in each analysis, or whether the model included random  
1372 effects. However, in our post hoc analyses, we found that dropping analyses identified as

1373 unpublishable or in need of major revision by at least one reviewer modestly reduced the observed  
1374 heterogeneity among the  $Z_r$  outcomes, but only for *Eucalyptus* analyses, apparently because this led  
1375 to the dropping of the major outlier. This limited role for peer review in explaining the variability in  
1376 our results should be interpreted cautiously because the inter-rater reliability among peer reviewers  
1377 was extremely low, and at least some analyses that appeared flawed to us were not marked as  
1378 flawed by reviewers. Thus, it seems that the peer reviews we received were of mixed quality, possibly  
1379 due to lack of expertise or lack of care on the part of some reviewers. However, the hypothesis that  
1380 poor quality analyses drove a substantial portion of the heterogeneity we observed was also  
1381 contradicted by our observation that analysts' self-declared statistical expertise appeared unrelated  
1382 to heterogeneity. When we retained only analyses from teams including at least one member with  
1383 high self-declared levels of expertise, heterogeneity among effect sizes remained high. Thus, our  
1384 results suggest lack of statistical expertise is not the primary factor responsible for the heterogeneity  
1385 we observed, although further work is merited before rejecting a role for statistical expertise.  
1386 Besides variability in expertise, it is also possible that the volunteer analysts varied in the effort they  
1387 invested, and low effort presumably drove at least some heterogeneity in results. However, analysts  
1388 often submitted thoughtful and extensive code, tables, figures, and textual explanation and  
1389 interpretations, which is evidence of substantial investment. Further, we are confident that low effort  
1390 alone is an insufficient explanation for the heterogeneity we observed because we have worked with  
1391 these datasets ourselves, and we know from experience that there are countless plausible modeling  
1392 alternatives that can produce a diversity of effects. Additionally, heterogeneity in analytical outcomes  
1393 differed notably between datasets, and there is no reason to expect that one set of analysts took this  
1394 project less seriously than the other. Returning to our exploratory analyses, not surprisingly, simply  
1395 dropping outlier values of  $Z_r$  for *Eucalyptus* analyses, which had more extreme outliers, led to less  
1396 observable heterogeneity in the forest plots, and also reductions in our quantitative measures of  
1397 heterogeneity. We did not observe a similar effect in the blue tit dataset because that dataset had  
1398 outliers that were much less extreme and instead had more variability across the core of the  
1399 distribution.

1400 Our major observations raise two broad questions; why was the variability among results so high,  
1401 and why did the pattern of variability differ between our two datasets. One important and plausible  
1402 answer to the first question is that much of the heterogeneity derives from the lack of a precise  
1403 relationship between the two biological research questions we posed and the data we provided. This  
1404 lack of a precise relationship between data and question creates many opportunities for different  
1405 model specifications, and so may inevitably lead to varied analytical outcomes ([Auspurg and Brüderl  
1406 2021](#)). However, we believe that the research questions we posed are consistent with the kinds of  
1407 research question that ecologists and evolutionary biologists typically work from. When designing  
1408 the two biological research questions, we deliberately sought to represent the level of specificity we  
1409 typically see in these disciplines. This level of specificity is evident when we look at the research  
1410 questions posed by some recent meta-analyses in these fields:

- 1411 • “how [does] urbanisation impact mean phenotypic values and phenotypic variation ... [in]  
1412 paired urban and non-urban comparisons of avian life-history traits” (Capilla-Lasheras et al.  
1413 2022)
- 1414 • “[what are] the effects of ocean acidification on the crustacean exoskeleton, assessing both  
1415 exoskeletal ion content (calcium and magnesium) and functional properties (biomechanical  
1416 resistance and cuticle thickness)” (Siegel et al. 2022)
- 1417 • “[what is] the extent to which restoration affects both the mean and variability of  
1418 biodiversity outcomes ... [in] terrestrial restoration” (Atkinson et al. 2022)



1419 • “[does] drought stress [have] a negative, positive, or null effect on aphid fitness” (Leybourne  
1420 et al. 2021)

1421 • “[what is] the influence of nitrogen-fixing trees on soil nitrous oxide emissions” (Kou-  
1422 Giesbrecht and Menge 2021)

1423 There is not a single precise answer to any of these questions, nor to the questions we posed to  
1424 analysts in our study. And this lack of single clear answers will obviously continue to cause  
1425 uncertainty since ecologists and evolutionary biologists conceive of the different answers from the  
1426 different statistical models as all being answers to the same general question. A possible response  
1427 would be a call to avoid these general questions in favor of much more precise alternatives (Auspurg  
1428 and Brüderl 2021). However, the research community rewards researchers who pose broad  
1429 questions (Simons, Shoda, and Lindsay 2017), and so researchers are unlikely to narrow their scope  
1430 without a change in incentives. Further, we suspect that even if individual studies specified narrow  
1431 research questions, other scientists would group these more narrow questions into broader  
1432 categories, for instance in meta-analyses, because it is these broader and more general questions  
1433 that often interest the research community.

1434 Although variability in statistical outcomes among analysts may be inevitable, our results raise  
1435 questions about why this variability differed between our two datasets. We are particularly  
1436 interested in the differences in the distribution of  $Z_r$ , since the distributions of out-of-sample  
1437 predictions were on different scales for the two datasets, thus limiting the value of comparisons. The  
1438 forest plots of  $Z_r$  from our two datasets showed distinct patterns, and these differences are  
1439 consistent with several alternative hypotheses. The results submitted by analysts of  
1440 the *Eucalyptus* dataset showed a small average (close to zero) with most estimates also close to zero  
1441 ( $\pm 0.2$ ), though about a third far enough above or below zero to cross the traditional threshold of  
1442 statistical significance. There were a small number of striking outliers that were very far from zero. In  
1443 contrast, the results submitted by analysts of the blue tit dataset showed an average much further  
1444 from zero (- 0.35) and a much greater spread in the core distribution of estimates across the range  
1445 of  $Z_r$  values ( $\pm 0.5$  from the mean), with few modest outliers. So, why was there more spread in  
1446 effect sizes (across the estimates that are not outliers) in the blue tit analyses relative to  
1447 the *Eucalyptus* analyses?

1448 One possible explanation for the lower heterogeneity among most *Eucalyptus*  $Z_r$  effects is that weak  
1449 relationships may limit the opportunities for heterogeneity in analytical outcome. Some evidence for  
1450 this idea comes from two sets of “many labs” studies in psychology (Klein et al. 2014, 2018). In these  
1451 studies, many independent lab groups each replicated a large set of studies, including, for each  
1452 study, the experiment, data collection, and statistical analyses. These studies showed that, when the  
1453 meta-analytic mean across the replications from different labs was small, there was much less  
1454 heterogeneity among the outcomes than when the mean effect sizes were large (Klein et al.  
1455 2014, 2018). Of course, a weak average effect size would not prevent divergent effects in all  
1456 circumstances. As we saw with the *Eucalyptus* analyses, taking a radically smaller subset of the data  
1457 can lead to dramatically divergent effect sizes even when the mean with the full dataset is close to  
1458 zero.

1459 Our observation that dramatic sub-setting in the *Eucalyptus* dataset was associated with  
1460 correspondingly dramatic divergence in effect sizes leads us towards another hypothesis to explain  
1461 the differences in heterogeneity between the *Eucalyptus* and blue tit analysis sets. It may be that  
1462 when analysts often divide a dataset into subsets, the result will be greater heterogeneity in  
1463 analytical outcome for that dataset. Although we saw sub-setting associated with dramatic outliers in

1464 the *Eucalyptus* dataset, nearly all other analyses of *Eucalyptus* data used close to the same set of 351  
1465 samples, and as we saw, these effects did not vary substantially. However, analysts often analyzed  
1466 only a subset of the blue tit data, and as we observed, sample sizes were much more variable among  
1467 blue tit effects, and the effects themselves were also much more variable. Important to note here is  
1468 that subsets of data may differ from each other for biological reasons, but they may also differ due to  
1469 sampling error. Sampling error is a function of sample size, and sub-samples are, by definition,  
1470 smaller samples, and so more subject to variability in effects due to sampling error (Jennions et al.  
1471 2013).

1472 Other features of datasets are also plausible candidates for driving heterogeneity in analytical  
1473 outcomes, including features of covariates. In particular, relationships between covariates and the  
1474 response variable as well as relationships between covariates and the primary independent variable  
1475 (collinearity) can strongly influence the modeled relationship between the independent variable of  
1476 interest and the dependent variable (Morrissey and Ruxton 2018; Dormann et al. 2013). Therefore,  
1477 inclusion or exclusion of these covariates can drive heterogeneity in effect sizes ( $Z_r$ ). Also, as we saw  
1478 with the two most extreme  $Z_r$  values from the blue tit analyses, in multivariate models with collinear  
1479 predictors, extreme effects can emerge when estimating partial correlation coefficients due to high  
1480 collinearity, and conclusions can differ dramatically depending on which relationship receives the  
1481 researcher's attention. Therefore, differences between datasets in the presence of strong and/or  
1482 collinear covariates could influence the differences in heterogeneity in results among those datasets.

1483 Although it is too early in the many-analyst research program to conclude which analytical decisions  
1484 or which features of datasets are the most important drivers of heterogeneity in analytical outcomes,  
1485 we must still grapple with the possibility that analytical outcomes may vary substantially based on  
1486 the choices we make as analysts. If we assume that, at least sometimes, different analysts will  
1487 produce dramatically different statistical outcomes, what should we do as ecologists and  
1488 evolutionary biologists? We review some ideas below.

1489 The easiest path forward after learning about this analytical heterogeneity would be simply to  
1490 continue with "business as usual", where researchers report results from a small number of statistical  
1491 models. A case could be made for this path based on our results. For instance, among the blue tit  
1492 analyses, the precise values of the estimated  $Z_r$  effects varied substantially, but the average effect  
1493 was convincingly different from zero, and a majority of individual effects (84%) were in the same  
1494 direction. Arguably, many ecologists and evolutionary biologists appear primarily interested in the  
1495 direction of a given effect and the corresponding p-value (Fidler et al. 2006), and so the variability we  
1496 observed when analyzing the blue tit dataset may not worry these researchers. Similarly, most  
1497 effects from the *Eucalyptus* analyses were relatively close to zero, and about two-thirds of these  
1498 effects did not cross the traditional threshold of statistical significance. Therefore, a large proportion  
1499 of people analyzing these data would conclude that there was no effect, and this is consistent with  
1500 what we might conclude from the meta-analysis.

1501 However, we find the counter arguments to "business as usual" to be compelling. For blue tits, there  
1502 were a substantial minority of calculated effects that would be interpreted by many biologists as  
1503 indicating the absence of an effect (28%), and there were three traditionally 'significant' effects in  
1504 the opposite direction to the average. The qualitative conclusions of analysts also reflected  
1505 substantial variability, with fully half of teams drawing a conclusion distinct from the one we draw  
1506 from the distribution as a whole. These teams with different conclusions were either uncertain about  
1507 the negative relationship between competition and nestling growth, or they concluded that effects  
1508 were mixed or absent. For the *Eucalyptus* analyses, this issue is more concerning. Around two-thirds  
1509 of effects had confidence intervals overlapping zero, and of the third of analyses with confidence

1510 intervals excluding zero, almost half were positive, and the rest were negative. Accordingly, the  
1511 qualitative conclusions of the *Eucalyptus* teams were spread across the full range of possibilities. But,  
1512 as we describe in the next paragraph, even this striking lack of consensus may be much less of a  
1513 problem than what could emerge as scientists select which results to publish.

1514 A potentially larger argument against “business as usual” is that it provides the raw material for  
1515 biasing the literature. When different model specifications readily lead to different results, analysts  
1516 may be tempted to report the result that appears most interesting, or that is most consistent with  
1517 expectation (Gelman and Loken 2013; Forstmeier, Wagenmakers and Parker 2017). There is growing  
1518 evidence that researchers in ecology and evolutionary biology often report a biased subset of the  
1519 results they produce (Deressa et al. 2023; Kimmel, Avolio and Ferraro 2023), and that this bias  
1520 exaggerates the average size of effects in the published literature between 30 and 150% (Yang et al.  
1521 2023; Parker and Yang 2023). The bias then accumulates in meta-analyses, apparently more than  
1522 doubling the rate of conclusions of “statistical significance” in published meta-analyses above what  
1523 would have been found in the absence of bias (Yang et al. 2023). Thus, “business as usual” does not  
1524 just create noisy results, it helps create systematically misleading results.

1525 If we move away from “business as usual”, where do we go? Many obvious options involve multiple  
1526 analyses per dataset. For instance, there is the traditional robustness or sensitivity check (e.g., Pei et  
1527 al. 2020; Briga and Verhulst 2021), in which the researcher presents several alternative versions of an  
1528 analysis to demonstrate that the result is ‘robust’ (Lu and White 2014). Unfortunately, robustness  
1529 checks are at risk of the same potential biases of reporting found in other studies (Silberzahn et al.  
1530 2018), especially given the relatively few models typically presented. However, these risks could be  
1531 minimized by running more models and doing so with a pre-registration or registered report.  
1532 Another option is model averaging. Averages across models often perform well (e.g. Taylor and Taylor  
1533 2023), and in some forms this may be a relatively simple solution. Model averaging, as most often  
1534 practiced in ecology and evolutionary biology, involves first identifying a small suite of candidate  
1535 models (see Burnham and Anderson 2002), then using Akaike weights, based on Akaike’s Information  
1536 Criterion (AIC), to calculate weighted averages for parameter estimates from those models. As with  
1537 typical robustness checks, the small number of models limits the exploration of specification space,  
1538 but examining a larger number of models could become the norm. However, there are more  
1539 concerning limitations. The largest of these limitations is that averaging regression coefficients is  
1540 problematic when models differ in interaction terms or collinear variables (Cade 2015). Additionally,  
1541 weighting by AIC may often be inconsistent with our modelling goals. AIC balances the trade-off  
1542 between model complexity and predictive ability, but penalizing models for complexity may not be  
1543 suited for testing hypotheses about causation (Arif and MacNeil 2022). So, AIC may often not offer  
1544 the weight we want to use, and we may also not wish to just generate an average at all. Instead, if we  
1545 hope to understand an extensive universe of possible modelling outcomes, we could conduct a  
1546 multiverse analysis, possibly with a specification curve (Simonsohn, Simmons, and Nelson  
1547 2015, 2020). This could mean running hundreds or thousands of models (or more!) to examine the  
1548 distribution of possible effects, and to see how different model specification choices map onto these  
1549 effects. However, exploring large areas of specification space may come at the cost of including  
1550 biologically implausible specifications. Thus, we expect a trade-off, and attempts to limit models to  
1551 the most biologically plausible may become increasingly difficult in proportion to the number of  
1552 variables and modeling choices. To make selecting plausible models easier, one could recruit multiple  
1553 analysts to design one or a few plausible specifications each as with our ‘many analyst’  
1554 study (Silberzahn et al. 2018). An alternative that may be more labor intensive for the primary  
1555 analyst, but which may lead to a more plausible set of models, could involve hypothesizing about  
1556 causal pathways with DAGs [directed acyclic graphs; Arif and MacNeil (2023)] to constrain the model

1557 set. As with other options outlined above, generating model specifications with DAGs could be  
1558 partnered with pre-registration to hinder bias from undisclosed data dredging.

1559 Responses to heterogeneity in analysis outcomes need not be limited to simply conducting more  
1560 analyses, especially if it turns out that analysis quality drives some of the observed heterogeneity. As  
1561 we noted above, we cannot yet rule out the possibility that insufficient statistical expertise or poor-  
1562 quality analyses might drive some portion of the heterogeneity we observed. Improving the quality  
1563 of analyses might be accomplished with a deliberate increase in investment in statistical education.  
1564 Many ecology and evolutionary biology students learn their statistical practice informally, with many  
1565 ecology doctoral programs in the USA not requiring a statistics course (Touchon and McCoy 2016),  
1566 and no formal courses of any kind included in doctoral degrees in most other countries. In cases  
1567 where formal investment in statistical education is lacking, informal resources, such as guidelines and  
1568 checklists, may help researchers avoid common mistakes. However, unless following guidelines or  
1569 checklists is enforced for publication, the adherence to guidelines is patchy. For example, despite the  
1570 publication of guidelines for conducting meta-analyses in ecology, the quality of meta-analyses did  
1571 not improve substantially over time (Koricheva and Gurevitch 2014). Even in medical research where  
1572 adherence to guidelines such as the PRISMA standards for systematic reviews and meta-analyses is  
1573 more highly valued, adherence is often poor (Page and Moher 2017).

1574 Although we have reviewed a variety of potential responses to the existence of variability in  
1575 analytical outcomes, we certainly do not wish to imply that this is a comprehensive set of possible  
1576 responses. Nor do we wish to imply that the opinions we have expressed about these options are  
1577 correct. Determining how the disciplines of ecology and evolutionary biology should respond to  
1578 knowledge of the variability in analytical outcome will benefit from the contribution and discussion  
1579 of ideas from across these disciplines. We look forward to learning from these discussions and to  
1580 seeing how these disciplines ultimately respond.

## 1581 Conclusions

1582 Overall, our results suggest to us that, where there is a diverse set of plausible analysis options, no  
1583 single analysis should be considered a complete or reliable answer to a research question. Further,  
1584 because of the evidence that ecologists and evolutionary biologists often present a biased subset of  
1585 the analyses they conduct (Deressa et al. 2023; Yang et al. 2023; Kimmel, Avolio and Ferraro 2023),  
1586 we do not expect that even a collection of different effect sizes from different studies will accurately  
1587 represent the true distribution of effects (Yang et al. 2023). Therefore, we believe that an increased  
1588 level of skepticism of the outcomes of single analyses, or even single meta-analyses, is warranted  
1589 going forward. We recognize that some researchers have long maintained a healthy level of  
1590 skepticism of individual studies as part of sound and practical scientific practice, and it is possible  
1591 that those researchers will be neither surprised nor concerned by our results. However, we doubt  
1592 that many researchers are sufficiently aware of the potential problems of analytical flexibility to be  
1593 appropriately skeptical. We hope that our work leads to conversations in ecology, evolutionary  
1594 biology, and other disciplines about how best to contend with heterogeneity in results that is  
1595 attributable to analytical decisions.

## 1596 Declarations

### 1597 Ethics, consent and permissions

1598 We obtained permission to conduct this research from the Whitman College Institutional Review  
1599 Board (IRB). As part of this permission, the IRB approved the consent form (<https://osf.io/xyp68/>)  
1600 that all participants completed prior to joining the study. The authors declare that they have no  
1601 competing interests.

### 1602 Availability of data and materials

1603 All materials and data are archived and hosted on the OSF at <https://osf.io/mn5aj/>, including survey  
1604 instruments and analyst / reviewer consent forms. The Evolutionary Ecology Data and Ecology and  
1605 Conservation Data provided to analysts are available  
1606 at <https://osf.io/34fzc/> and <https://osf.io/t76uy/> respectively. Data has been anonymised, and the  
1607 non-anonymised data is stored on the project OSF within private components accessible to the lead  
1608 authors.

1609 We built an R package, ManyEcoEvo to conduct the analyses described in this study ([Gould et al.](#)  
1610 [2023](#)), which can be downloaded from <https://github.com/egouldo/ManyEcoEvo/> to reproduce our  
1611 analyses or replicate the analyses described here using alternate datasets. Data cleaning and  
1612 preparation of analysis-data, as well as the analysis, is conducted in R ([R Core Team](#)  
1613 [2024](#)) reproducibly using the targets package ([Landau 2021](#)). This data and analysis pipeline is stored  
1614 in the ManyEcoEvo package repository and its outputs are made available to users of the package  
1615 when the library is loaded.

1616 The full manuscript, including further analysis and presentation of results is written in Quarto ([J. J.](#)  
1617 [Allaire et al. 2024](#)). The source code to reproduce the manuscript is hosted  
1618 at <https://github.com/egouldo/ManyAnalysts/>, and the rendered version of the source code may be  
1619 viewed at <https://egouldo.github.io/ManyAnalysts/>. All R packages and their versions used in the  
1620 production of this manuscript are listed in the session info at [Section 6.6](#).

### 1621 Competing interests

1622 The authors declare that they have no competing interests

### 1623 Funding

1624 EG's contributions were supported by an Australian Government Research Training Program  
1625 Scholarship, AIMOS top-up scholarship (2022) and Melbourne Centre of Data Science Doctoral  
1626 Academy Fellowship (2021). FF's contributions were supported by ARC Future Fellowship  
1627 FT150100297.

### 1628 Author's contributions

1629 HF, THP and FF conceptualized the project. PV provided raw data for *Eucalyptus* analyses and SG and  
1630 THP provided raw data for blue tit analyses. DGH, HF and THP prepared surveys for collecting  
1631 participating analysts and reviewer's data. EG, HF, THP, PV, SN and FF planned the analyses of the  
1632 data provided by our analysts and reviewers, EG, HF, and THP curated the data, EG and HF wrote the  
1633 software code to implement the analyses and prepare data visualisations. EG ensured that analyses  
1634 were documented and reproducible. THP and HF administered the project, including coordinating

1635 with analysts and reviewers. FF provided funding for the project. THP, HF, and EG wrote the  
1636 manuscript. Authors listed alphabetically contributed analyses of the primary datasets or reviews of  
1637 analyses. All authors read and approved the final manuscript.

## 1638 References

- 1639 Aczel, Balazs, Barnabas Szaszi, Gustav Nilsson, Olmo R van den Akker, Casper J Albers, Marcel ALM  
1640 van Assen, Jojanneke A Bastiaansen, et al. 2021. "Consensus-Based Guidance for Conducting and  
1641 Reporting Multi-Analyst Studies." *eLife* 10 (November). <https://doi.org/10.7554/elife.72185>.
- 1642 Allaire, J. J., Charles Teague, Carlos Scheidegger, Yihui Xie, and Christophe Dervieux.  
1643 2024. "Quarto." <https://doi.org/10.5281/zenodo.5960048>.
- 1644 Allaire, JJ, Yihui Xie, Christophe Dervieux, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron  
1645 Atkins, et al. 2024. *rmarkdown: Dynamic Documents for r*. <https://github.com/rstudio/rmarkdown>.
- 1646 Arif, S., & M. Aaron MacNeil. 2022. "Predictive models aren't for causal inference." *Ecology Letters*  
1647 25(8), 1741–1745. <https://doi.org/10.1111/ele.14033>
- 1648 Arif, Suchinta, and M. Aaron MacNeil. 2023. "Applying the Structural Causal Model Framework for  
1649 Observational Causal Inference in Ecology." *Ecological Monographs* 93 (1): e1554.  
1650 <https://doi.org/https://doi.org/10.1002/ecm.1554>.
- 1651 Arnold, Jeffrey B. 2024. *ggthemes: Extra Themes, Scales and Geoms*  
1652 *for "ggplot2"*. <https://jrnold.github.io/ggthemes/>.
- 1653 Atkinson, Joe, Lars A. Brudvig, Max Mallen-Cooper, Shinichi Nakagawa, Angela T. Moles, and Stephen  
1654 P. Bonser. 2022. "Terrestrial Ecosystem Restoration Increases Biodiversity and Reduces Its Variability,  
1655 but Not to Reference Levels: A Global Meta-Analysis." *Ecology Letters* 25 (7): 1725–37.  
1656 <https://doi.org/https://doi.org/10.1111/ele.14025>.
- 1657 Auspurg, Katrin, and Josef Brüderl. 2021. "Has the Credibility of the Social Sciences Been Credibly  
1658 Destroyed? Reanalyzing the 'Many Analysts, One Data Set' Project." *Socius* 7:  
1659 23780231211024421. <https://doi.org/10.1177/23780231211024421>.
- 1660 Bartoń, Kamil. 2023. *MuMIn: Multi-Model Inference*.
- 1661 Baselga, Andres, David Orme, Sebastien Villegier, Julien De Bortoli, Fabien Leprieur, Maxime Logez,  
1662 Sara Martinez-Santalla, et al. 2023. *betapart: Partitioning Beta Diversity into Turnover and*  
1663 *Nestedness Components*. <https://CRAN.R-project.org/package=betapart>.
- 1664 Bates, Douglas, Martin Mächler, Ben Bolker, and Steve Walker. 2015. "Fitting Linear Mixed-Effects  
1665 Models Using lme4." *2015* 67 (1): 48. <https://doi.org/10.18637/jss.v067.i01>.
- 1666 Blake, Kevin. 2022. *NatParksPalettes: Color Palettes Inspired by National*  
1667 *Parks*. <https://github.com/kevinsblake/NatParksPalettes>.
- 1668 Bolker, Ben, David Robinson, Dieter Menne, Jonah Gabry, Paul Buerkner, Christopher Hau, William  
1669 Petry, et al. 2024. *broom.mixed: Tidying Methods for Mixed*  
1670 *Models*. <https://github.com/bbolker/broom.mixed>.
- 1671 Borenstein, Michael, Julian P. T. Higgins, Larry Hedges, and Hannah Rothstein. 2017. "Basics of Meta-  
1672 Analysis:  $I^2$  Is Not an Absolute Measure of Heterogeneity." *Research Synthesis Methods* 8: 5–  
1673 18. <https://doi.org/10.1002/jrsm.1230>.

- 1674 Botvinik-Nezer, Rotem, Felix Holzmeister, Colin F. Camerer, Anna Dreber, Juergen Huber, Magnus  
1675 Johannesson, Michael Kirchler, et al. 2020. "Variability in the Analysis of a Single Neuroimaging  
1676 Dataset by Many Teams." *Nature* 582 (7810): 84–88.
- 1677 Breznau, Nate, Eike Mark Rinke, Alexander Wuttke, Hung H. V. Nguyen, Muna Adem, Jule Adriaans,  
1678 Amalia Alvarez-Benjumea, et al. 2022. "Observing Many Researchers Using the Same Data and  
1679 Hypothesis Reveals a Hidden Universe of Uncertainty." *Proceedings of the National Academy of  
1680 Sciences* 119 (44): e2203150119. <https://doi.org/10.1073/pnas.2203150119>.
- 1681 Briga, Michael, and Simon Verhulst. 2021. "Mosaic Metabolic Ageing: Basal and Standard Metabolic  
1682 Rates Age in Opposite Directions and Independent of Environmental Quality, Sex and Life Span in a  
1683 Passerine." *Functional Ecology* 35 (5): 1055–68. [https://doi.org/https://doi.org/10.1111/1365-  
1684 2435.13785](https://doi.org/https://doi.org/10.1111/1365-2435.13785).
- 1685 Brooks, Mollie E., Kasper Kristensen, Koen J. van Benthem, Arni Magnusson, Casper W. Berg, Anders  
1686 Nielsen, Hans J. Skaug, Martin Maechler, and Benjamin M. Bolker. 2017. "glmmTMB Balances Speed  
1687 and Flexibility Among Packages for Zero-Inflated Generalized Linear Mixed Modeling." *The R Journal* 9  
1688 (2): 378–400. <https://doi.org/10.32614/RJ-2017-066>.
- 1689 Buck, Robert J., John Fieberg, and Daniel J. Larkin. 2022. "The use of weighted averages of Hedges' d  
1690 in meta-analysis: Is it worth it?" *Methods in Ecology and Evolution* 13 (5): 1093–1105.  
1691 <https://doi.org/10.1111/2041-210X.13818>.
- 1692 Burnham, K. P., and D. R. Anderson. 2002. *Model Selection and Multimodel Inference: A Practical  
1693 Information-Theoretical Approach*. Book. 2nd ed. New York: Springer-  
1694 Verlag. <https://doi.org/10.1007/b97636>.
- 1695 Cade, Brian S. 2015. "Model Averaging and Muddled Multimodel Inferences." *Ecology* 96 (9): 2370–  
1696 82. <http://www.jstor.org.ezproxy.whitman.edu/stable/24702343>.
- 1697 Capilla-Lasheras, Pablo, Megan J. Thompson, Alfredo Sánchez-Tójar, Yacob Haddou, Claire J. Branston,  
1698 Denis Réale, Anne Charmantier, and Davide M. Dominoni. 2022. "A Global Meta-Analysis Reveals  
1699 Higher Variation in Breeding Phenology in Urban Birds Than in Their Non-Urban Neighbours." *Ecology  
1700 Letters* 25 (11): 2552–70. <https://doi.org/10.1111/ele.14099>.
- 1701 Coretta, Stefano, Joseph V. Casillas, Simon Roessig, Michael Franke, Byron Ahn, Ali H. Al-Hoorie, Jalal  
1702 Al-Tamimi, et al. 2023. "Multidimensional Signals and Analytic Flexibility: Estimating Degrees of  
1703 Freedom in Human-Speech Analyses." *Advances in Methods and Practices in Psychological Science* 6  
1704 (3): 25152459231162567. <https://doi.org/10.1177/25152459231162567>.
- 1705 Dancho, Matt, and Davis Vaughan. 2023. *Timetk: A Tool Kit for Working with Time  
1706 Series*. <https://CRAN.R-project.org/package=timetk>.
- 1707 DeKogel, C. H. 1997. "Long-Term Effects of Brood Size Manipulation on Morphological Development  
1708 and Sex-Specific Mortality of Offspring." *Journal of Animal Ecology* 66 (2): 167–78. [Go to  
1709 ISI>://WOS:A1997WQ19600003](https://doi.org/10.1111/j.1365-2656.1997.00003.x).
- 1710 Deressa, Teshome, David Stern, Jaco Vangronsveld, Jan Minx, Sebastien Lizin, Robert Malina, and  
1711 Stephan Bruns. 2023. "More Than Half of Statistically Significant Research Findings in the  
1712 Environmental Sciences Are Actually Not." *EcoEvoRxiv*.  
1713 <https://doi.org/https://doi.org/10.32942/X24G6Z>.

1714 Dormann, Carsten F., Jane Elith, Sven Bacher, Carsten Buchmann, Gudrun Carl, Gabriel Carré, Jaime  
1715 R. García Marquéz, et al. 2013. "Collinearity: A Review of Methods to Deal with It and a Simulation  
1716 Study Evaluating Their Performance." *Ecography* 36 (1): 27–46.  
1717 <https://doi.org/https://doi.org/10.1111/j.1600-0587.2012.07348.x>.

1718 Fanelli, Daniele, Rodrigo Costas, and John P. A. Ioannidis. 2017. "Meta-Assessment of Bias in  
1719 Science." *Proceedings of the National Academy of Sciences* 114: 3714–  
1720 19. <https://doi.org/10.1073/pnas.1618569114>.

1721 Fanelli, Daniele, and John P. A. Ioannidis. 2013. "US Studies May Overestimate Effect Sizes in Softer  
1722 Research." *Proceedings of the National Academy of Sciences* 110 (37): 15031–  
1723 36. <https://doi.org/10.1073/pnas.1302997110>.

1724 Fidler, Fiona, Mark A. Burgman, Geoff Cumming, Robert Buttrose, and Neil Thomason. 2006. "Impact  
1725 of Criticism of Null-Hypothesis Significance Testing on Statistical Reporting Practices in Conservation  
1726 Biology." *Conservation Biology* 20 (5): 1539–44. <https://doi.org/10.1111/j.1523-1739.2006.00525.x>.

1727 Fidler, Fiona, Yung En Chee, Bonnie C. Wintle, Mark A. Burgman, Michael A. McCarthy, and Ascelin  
1728 Gordon. 2017. "Metaresearch for Evaluating Reproducibility in Ecology and Evolution." *BioScience* 67  
1729 (3): 282–89. <https://doi.org/10.1093/biosci/biw159>.

1730 Firke, Sam. 2023. *janitor: Simple Tools for Examining and Cleaning Dirty*  
1731 *Data*. <https://github.com/sfirke/janitor>.

1732 Forstmeier, Wolfgang, Eric-Jan Wagenmakers, and T. H. Parker. 2017. "Detecting and Avoiding Likely  
1733 False-Positive Findings – a Practical Guide." *Biological Reviews* 92: 1941–  
1734 68. <https://doi.org/10.1111/brv.12315>.

1735 Fraser, Hannah, Tim Parker, Shinichi Nakagawa, Ashley Barnett, and Fiona Fidler. 2018. "Questionable  
1736 Research Practices in Ecology and Evolution." *PLOS ONE* 13 (7):  
1737 e0200303. <https://doi.org/10.1371/journal.pone.0200303>.

1738 Gamer, Matthias, Jim Lemon, and Ian Fellows Puspendra Singh. 2019. *irr: Various Coefficients of*  
1739 *Interrater Reliability and Agreement*. <https://www.r-project.org>.

1740 Gelman, Andrew, and Eric Loken. 2013. "The Garden of Forking Paths: Why Multiple Comparisons  
1741 Can Be a Problem, Even When There Is No 'Fishing Expedition' or 'p-Hacking' and the Research  
1742 Hypothesis Was Posited Ahead of Time." *Department of Statistics, Columbia University*.

1743 Gelman, Andrew, and David Weakliem. 2009. "Of Beauty, Sex, and Power." *American Scientist* 97:  
1744 310–16.

1745 Gould, Elliot, Hannah S. Fraser, Shinichi Nakagawa, and Timothy H. Parker. 2023. "ManyEcoEvo:  
1746 Meta-Analyse Data from ManyAnalyst Style  
1747 Studies." Zenodo. <https://doi.org/10.5281/zenodo.10046153>.

1748 Gould, Elliot, Hannah S. Fraser, Shinichi Nakagawa, Timothy H. Parker. 2024. *egouldo/ManyAnalysts:*  
1749 *Manuscript Source Code for 'Same data, different analysts: variation in effect sizes due to analytical*  
1750 *decisions in ecology and evolutionary biology.'* Zenodo. <https://doi.org/10.5281/zenodo.13850927>.  
1751 Version 2.0.2.

1752 Grueber, C. E., S. Nakagawa, R. J. Laws, and I. G. Jamieson. 2011. "Multimodel Inference in Ecology  
1753 and Evolution: Challenges and Solutions." *Journal of Evolutionary Biology* 24 (4): 699–  
1754 711. <https://doi.org/doi:10.1111/j.1420-9101.2010.02210.x>.



- 1755 Harrell Jr, Frank E. 2024. *Hmisc: Harrell Miscellaneous*. <https://hbiostat.org/R/Hmisc/>.
- 1756 Hester, Jim, Lionel Henry, Kirill Müller, Kevin Ushey, Hadley Wickham, and Winston Chang.  
1757 2024. *withr: Run Code "With" Temporarily Modified Global State*. <https://withr.r-lib.org>.
- 1758 Higgins, Julian P T, Simon G Thompson, Jonathan J Deeks, and Douglas G Altman. 2003. "Measuring  
1759 Inconsistency in Meta-Analyses." *BMJ* 327 (7414): 557–  
1760 60. <https://doi.org/10.1136/bmj.327.7414.557>.
- 1761 Huntington-Klein, Nick, Andreu Arenas, Emily Beam, Marco Bertoni, Jeffrey R. Bloem, Pralhad Burli,  
1762 Naibin Chen, et al. 2021. "The Influence of Hidden Researcher Decisions in Applied  
1763 Microeconomics." *Economic Inquiry* 59 (3): 944–60.  
1764 <https://doi.org/https://doi.org/10.1111/ecin.12992>.
- 1765 Iannone, Richard, Joe Cheng, Barret Schloerke, Ellis Hughes, Alexandra Lauer, and JooYoung Seo.  
1766 2024. *gt: Easily Create Presentation-Ready Display Tables*. <https://gt.rstudio.com>.
- 1767 Jennions, M. D., C. J. Lortie, M. S. Rosenberg, and H. R. Rothstein. 2013. "Publication and Related  
1768 Biases." Book Section. In *Handbook of Meta-Analysis in Ecology and Evolution*, edited by J. Koricheva,  
1769 J. Gurevitch, and K. Mengersen, 207–36. Princeton, USA: Princeton University Press.
- 1770 Kassambara, Alboukadel. 2023. *ggpubr: "ggplot2" Based Publication Ready  
1771 Plots*. <https://rpkgs.datanovia.com/ggpubr/>.
- 1772 Kimmel, Kaitlin, Meghan L. Avolio, and Paul J. Ferraro. 2023. "Empirical Evidence of Widespread  
1773 Exaggeration Bias and Selective Reporting in Ecology." *Nature Ecology &  
1774 Evolution*. <https://doi.org/10.1038/s41559-023-02144-3>.
- 1775 Klein, Richard A., Kate A. Ratliff, Michelangelo Vianello, Reginald B. Adams Jr., Štěpán Bahník, Michael  
1776 J. Bernstein, Konrad Bocian, et al. 2014. "Investigating Variation in Replicability: A "Many Labs"  
1777 Replication Project." *Social Psychology* 45 (3): 142–52. <https://doi.org/10.1027/1864-9335/a000178>.
- 1778 Klein, Richard A., Michelangelo Vianello, Fred Hasselman, Byron G. Adams, Reginald B. Adams, Sinan  
1779 Alper, Mark Aveyard, et al. 2018. "Many Labs 2: Investigating Variation in Replicability Across Samples  
1780 and Settings." *Advances in Methods and Practices in Psychological Science* 1 (4): 443–  
1781 90. <https://doi.org/10.1177/2515245918810225>.
- 1782 Knight, K. 2000. *Mathematical Statistics*. Book. New York: Chapman; Hall.
- 1783 Koricheva, Julia, and Jessica Gurevitch. 2014. "Uses and Misuses of Meta-Analysis in Plant  
1784 Ecology." *Journal of Ecology* 102 (4): 828–44. [https://doi.org/https://doi.org/10.1111/1365-  
1785 2745.12224](https://doi.org/https://doi.org/10.1111/1365-2745.12224).
- 1786 Kou-Giesbrecht, Sian, and Duncan N. L. Menge. 2021. "Nitrogen-Fixing Trees Increase Soil Nitrous  
1787 Oxide Emissions: A Meta-Analysis." *Ecology* 102 (8): e03415.  
1788 <https://doi.org/https://doi.org/10.1002/ecy.3415>.
- 1789 Kuhn, Max, and Hannah Frick. 2022. *multilevelmod: Model Wrappers for Multi-Level  
1790 Models*. <https://github.com/tidymodels/multilevelmod>.
- 1791 Kuhn, Max, and Hadley Wickham. 2020. *Tidymodels: A Collection of Packages for Modeling and  
1792 Machine Learning Using Tidyverse Principles*. <https://www.tidymodels.org>.

1793 Kuznetsova, Alexandra, Per B. Brockhoff, and Rune H. B. Christensen. 2017. "lmerTest Package: Tests  
1794 in Linear Mixed Effects Models." *Journal of Statistical Software* 82 (13): 1–  
1795 26. <https://doi.org/10.18637/jss.v082.i13>.

1796 Landau, William Michael. 2021. "The Targets r Package: A Dynamic Make-Like Function-Oriented  
1797 Pipeline Toolkit for Reproducibility and High-Performance Computing." *Journal of Open Source*  
1798 *Software* 6 (57): 2959. <https://doi.org/10.21105/joss.02959>.

1799 Leybourne, Daniel J., Katharine F. Preedy, Tracy A. Valentine, Jorunn I. B. Bos, and Alison J. Karley.  
1800 2021. "Drought Has Negative Consequences on Aphid Fitness and Plant Vigor: Insights from a Meta-  
1801 Analysis." *Ecology and Evolution* 11 (17): 11915–29.  
1802 <https://doi.org/https://doi.org/10.1002/ece3.7957>.

1803 Lu, Xun, and Halbert White. 2014. "Robustness Checks and Robustness Tests in Applied  
1804 Economics." *Journal of Econometrics* 178: 194–206.  
1805 <https://doi.org/https://doi.org/10.1016/j.jeconom.2013.08.016>.

1806 Lüdecke, Daniel, Mattan S. Ben-Shachar, Indrajeet Patil, and Dominique Makowski. 2020. "Extracting,  
1807 Computing and Exploring the Parameters of Statistical Models Using R." *Journal of Open Source*  
1808 *Software* 5 (53): 2445. <https://doi.org/10.21105/joss.02445>.

1809 Lüdecke, Daniel, Mattan S. Ben-Shachar, Indrajeet Patil, Philip Waggoner, and Dominique Makowski.  
1810 2021. "performance: An R Package for Assessment, Comparison and Testing of Statistical  
1811 Models." *Journal of Open Source Software* 6 (60): 3139. <https://doi.org/10.21105/joss.03139>.

1812 Lüdecke, Daniel, Indrajeet Patil, Mattan S. Ben-Shachar, Brenton M. Wiernik, Philip Waggoner, and  
1813 Dominique Makowski. 2021. "see: An R Package for Visualizing Statistical Models." *Journal of Open*  
1814 *Source Software* 6 (64): 3393. <https://doi.org/10.21105/joss.03393>.

1815 Luke, S. G. 2017. "Evaluating Significance in Linear Mixed-Effects Models in r." *Behavior Research*  
1816 *Methods* 49 (4): 1494–1502.

1817 Makowski, Dominique, Mattan S. Ben-Shachar, Indrajeet Patil, and Daniel Lüdecke. 2020. "Estimation  
1818 of Model-Based Predictions, Contrasts and  
1819 Means." CRAN. <https://github.com/easystats/modelbased>.

1820 Masur, Philipp K., and Michael Scharrow. 2020. "specr: Conducting and Visualizing Specification  
1821 Curve Analyses (Version 1.0.0)." <https://CRAN.R-project.org/package=specr>.

1822 Meschiari, Stefano. 2022. *Latex2exp: Use LaTeX Expressions in*  
1823 *Plots*. <https://www.stefanom.io/latex2exp/>.

1824 Miles, C. 2008. "Testing Market-Based Instruments for Conservation in Northern Victoria." Book  
1825 Section. In *Biodiversity: Integrating Conservation and Production: Case Studies from Australian*  
1826 *Farms, Forests and Fisheries*, edited by T. Norton, T. Lefroy, K. Bailey, and G. Unwin, 133–46.  
1827 Melbourne, Australia: CSIRO Publishing.

1828 Millard, Steven P. 2013. *EnvStats: An r Package for Environmental Statistics*. New York:  
1829 Springer. <https://www.springer.com>.

1830 Molina, Isabel, and Yolanda Marhuenda. 2015. "sae: An R Package for Small Area Estimation." *The R*  
1831 *Journal* 7 (1): 81–98. <https://journal.r-project.org/archive/2015/RJ-2015-007/RJ-2015-007.pdf>.

1832 Morrissey, Michael B., and Graeme D. Ruxton. 2018. "Multiple Regression Is Not Multiple  
1833 Regressions: The Meaning of Multiple Regression and the Non-Problem of Collinearity." *Philosophy,*  
1834 *Theory, and Practice in Biology* 10 (3). <https://doi.org/10.3998/ptpbio.16039257.0010.003>.

1835 Müller, Kirill. 2020. *here: A Simpler Way to Find Your Files*. <https://here.r-lib.org/>.

1836 Nakagawa, Shinichi, and Innes C. Cuthill. 2007. "Effect Size, Confidence Interval and Statistical  
1837 Significance: A Practical Guide for Biologists." *Biological Reviews* 82 (4): 591–  
1838 605. <https://doi.org/10.1111/j.1469-185X.2007.00027.x>.

1839 Nakagawa, Shinichi, Malgorzata Lagisz, Michael D. Jennions, Julia Koricheva, Daniel W. A. Noble,  
1840 Timothy H. Parker, Alfredo Sánchez-Tójar, Yefeng Yang, and Rose E. O’Dea. 2022. "Methods for  
1841 Testing Publication Bias in Ecological and Evolutionary Meta-Analyses." *Methods in Ecology and*  
1842 *Evolution* 13 (1): 4–21. <https://doi.org/https://doi.org/10.1111/2041-210X.13724>.

1843 Nakagawa, Shinichi, Malgorzata Lagisz, Rose E. O’Dea, Patrice Pottier, Joanna Rutkowska, Alistair M.  
1844 Senior, Yefeng Yang, and Daniel W. A. Noble. 2023. "orchaRd 2.0: An r Package for Visualizing Meta-  
1845 Analyses with Orchard Plots." *EcoEvoRxiv* 12: 4–12.  
1846 <https://doi.org/https://doi.org/10.32942/X2QC7K>.

1847 Nakagawa, Shinichi, Yefeng Yang, Erin L. Macartney, Rebecca Spake, and Malgorzata Lagisz.  
1848 2023. "Quantitative Evidence Synthesis: A Practical Guide on Meta-Analysis, Meta-Regression, and  
1849 Publication Bias Tests for Environmental Sciences." *Environmental Evidence* 12 (1):  
1850 8. <https://doi.org/10.1186/s13750-023-00301-6>.

1851 Nakagawa, S., D. W. Noble, A. M. Senior, and M. Lagisz. 2017. "Meta-Evaluation of Meta-Analysis: Ten  
1852 Appraisal Questions for Biologists." *BMC Biology* 15 (1): 18. [https://doi.org/10.1186/s12915-017-](https://doi.org/10.1186/s12915-017-0357-7)  
1853 [0357-7](https://doi.org/10.1186/s12915-017-0357-7).

1854 Nicolaus, M., S. P. M. Michler, R. Ubels, M. van der Velde, J. Komdeur, C. Both, and J. M. Tinbergen.  
1855 2009. "Sex-Specific Effects of Altered Competition on Nestling Growth and Survival: An Experimental  
1856 Manipulation of Brood Size and Sex Ratio." *Journal of Animal Ecology* 78 (2): 414–  
1857 26. <https://doi.org/10.1111/j.1365-2656.2008.01505.x>.

1858 Noble, Daniel W. A., Malgorzata Lagisz, Rose E. O’Dea, and Shinichi Nakagawa.  
1859 2017. "Nonindependence and Sensitivity Analyses in Ecological and Evolutionary Meta-  
1860 Analyses." *Molecular Ecology* 26 (9): 2410–25. <https://doi.org/10.1111/mec.14031>.

1861 O’Hara, Robert B., and D. Johan Kotze. 2010. "Do Not Log-Transform Count Data." *Methods in*  
1862 *Ecology and Evolution* 1 (2): 118–22. <https://doi.org/10.1111/j.2041-210x.2010.00021.x>.

1863 Open Science Collaboration. 2015. "Estimating the Reproducibility of Psychological  
1864 Science." *Science* 349 (6251): aac4716. <https://doi.org/10.1126/science.aac4716>.

1865 Page, Matthew J., and David Moher. 2017. "Evaluations of the Uptake and Impact of the Preferred  
1866 Reporting Items for Systematic Reviews and Meta-Analyses (PRISMA) Statement and Extensions: A  
1867 Scoping Review." *Systematic Reviews* 6 (1): 263. <https://doi.org/10.1186/s13643-017-0663-8>.

1868 Parker, Timothy H., Wolfgang Forstmeier, Julia Koricheva, Fiona Fidler, Jarrod D. Hadfield, Yung En  
1869 Chee, Clint D. Kelly, Jessica Gurevitch, and Shinichi Nakagawa. 2016. "Transparency in Ecology and  
1870 Evolution: Real Problems, Real Solutions." *Trends in Ecology & Evolution* 31 (9): 711–  
1871 19. <https://doi.org/10.1016/j.tree.2016.07.002>.

- 1872 Parker, Timothy H., and Yefeng Yang. 2023. "Exaggerated Effects in Ecology." *Nature Ecology &*  
1873 *Evolution*. <https://doi.org/10.1038/s41559-023-02156-z>.
- 1874 Pedersen, Thomas Lin. 2024. *patchwork: The Composer of Plots*. [https://patchwork.data-](https://patchwork.data-imaginist.com)  
1875 [imaginist.com](https://patchwork.data-imaginist.com).
- 1876 Pei, Yifan, Wolfgang Forstmeier, Daiping Wang, Katrin Martin, Joanna Rutkowska, and Bart  
1877 Kempenaers. 2020. "Proximate Causes of Infertility and Embryo Mortality in Captive Zebra  
1878 Finches." *The American Naturalist* 196 (5): 577–96. <https://doi.org/10.1086/710956>.
- 1879 Qiu, Yixuan. 2024. *showtext: Using Fonts More Easily in r*  
1880 *Graphs*. <https://github.com/yixuan/showtext>.
- 1881 R Core Team. 2024. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R  
1882 Foundation for Statistical Computing. <https://www.R-project.org/>.
- 1883 Rosenberg, M. S. 2013. "Moment and Least-Squares Based Approaches to Metaanalytic  
1884 Inference." Book Section. In *Handbook of Meta-Analysis in Ecology and Evolution*, edited by J.  
1885 Koricheva, J. Gurevitch, and K. Mengersen, 108–24. Princeton, USA: Princeton University Press.
- 1886 Royle, N. J., I. R. Hartley, I. P. F. Owens, and G. A. Parker. 1999. "Sibling Competition and the Evolution  
1887 of Growth Rates in Birds." *Proceedings of the Royal Society B-Biological Sciences* 266 (1422): 923–  
1888 32. <https://doi.org/10.1098/rspb.1999.0725>.
- 1889 Scheinin, Ilari, Maria Kalimeri, Vilma Jagerroos, Juuso Parkkinen, Emmi Tikkanen, Peter Würtz, and  
1890 Antti Kangas. 2020. *ggforestplot: Forestplots of Measures of Effects and Their Confidence*  
1891 *Intervals*. <https://github.com/NightingaleHealth/ggforestplot>.
- 1892 Schloerke, Barret, Di Cook, Joseph Larmarange, Francois Briatte, Moritz Marbach, Edwin Thoen,  
1893 Amos Elberg, and Jason Crowley. 2024. *GGally: Extension*  
1894 *to "ggplot2"*. <https://ggobi.github.io/ggally/>.
- 1895 Schweinsberg, M., M. Feldman, N. Staub, O. R. van den Akker, R. C. M. van Aert, Malm van Assen, Y.  
1896 Liu, et al. 2021. "Same Data, Different Conclusions: Radical Dispersion in Empirical Results When  
1897 Independent Analysts Operationalize and Test the Same Hypothesis." *Organizational Behavior and*  
1898 *Human Decision Processes* 165: 228–49. <https://doi.org/10.1016/j.obhdp.2021.02.003>.
- 1899 Senior, Alistair M., Catherine E. Grueber, Tsukushi Kamiya, Malgorzata Lagisz, Katie O'Dwyer, Eduardo  
1900 S. A. Santos, and Shinichi Nakagawa. 2016. "Heterogeneity in Ecological and Evolutionary Meta-  
1901 Analyses: Its Magnitude and Implications." *Ecology* 97 (12): 3293–  
1902 99. <https://doi.org/10.1002/ecy.1591>.
- 1903 Shavit, A., and Aaron M. Ellison. 2017. *Stepping in the Same River Twice: Replication in Biological*  
1904 *Research*. Edited Book. New Haven, Connecticut, USA: Yale University Press.
- 1905 Siegel, Kyle R., Muskanjot Kaur, A. Calvin Grigal, Rebecca A. Metzler, and Gary H. Dickinson.  
1906 2022. "Meta-Analysis Suggests Negative, but pCO<sub>2</sub>-Specific, Effects of Ocean Acidification on the  
1907 Structural and Functional Properties of Crustacean Biomaterials." *Ecology and Evolution* 12 (6):  
1908 e8922. <https://doi.org/https://doi.org/10.1002/ece3.8922>.
- 1909 Silberzahn, R., E. L. Uhlmann, D. P. Martin, P. Anselmi, F. Aust, E. Awtrey, Š. Bahník, et al. 2018. "Many  
1910 Analysts, One Data Set: Making Transparent How Variations in Analytic Choices Affect  
1911 Results." *Advances in Methods and Practices in Psychological Science* 1 (3): 337–  
1912 56. <https://doi.org/10.1177/2515245917747646>.

- 1913 Silge, Julia, and David Robinson. 2016. "tidytext: Text Mining and Analysis Using Tidy Data Principles  
1914 in r." *JOSS* 1 (3). <https://doi.org/10.21105/joss.00037>.
- 1915 Simons, Daniel J., Yuichi Shoda, and D. Stephen Lindsay. 2017. "Constraints on Generality (COG): A  
1916 Proposed Addition to All Empirical Papers." *Perspectives on Psychological  
1917 Science*. <https://doi.org/10.1177/174569161770863>.
- 1918 Simonsohn, Uri, Joseph P. Simmons, and Leif D. Nelson. 2015. "Specification Curve: Descriptive and  
1919 Inferential Statistics on All Reasonable Specifications." Manuscript. *SSRN Electronic  
1920 Journal*. <https://doi.org/10.2139/ssrn.2694998>.
- 1921 ———. 2020. "Specification Curve Analysis." *Nature Human Behaviour* 4 (11): 1208–  
1922 14. <https://doi.org/10.1038/s41562-020-0912-z>.
- 1923 Sjoberg, Daniel D., Karissa Whiting, Michael Curry, Jessica A. Lavery, and Joseph Larmarange.  
1924 2021. "Reproducible Summary Tables with the Gtsummary Package." *The R Journal* 13: 570–  
1925 80. <https://doi.org/10.32614/RJ-2021-053>.
- 1926 Slowikowski, Kamil. 2024. *ggrepel: Automatically Position Non-Overlapping Text Labels  
1927 with "ggplot2"*. <https://ggrepel.slowkow.com/>.
- 1928 Stanton-Geddes, John, Cintia Gomes de Freitas, and Cristian de Sales Dambros. 2014. "In Defense of  
1929 p Values: Comment on the Statistical Methods Actually Used by Ecologists." *Ecology* 95 (3): 637–42.  
1930 <https://doi.org/https://doi.org/10.1890/13-1156.1>.
- 1931 Steegen, Sara, Francis Tuerlinckx, Andrew Gelman, and Wolf Vanpaemel. 2016. "Increasing  
1932 Transparency Through a Multiverse Analysis." *Perspectives on Psychological Science* 11 (5): 702–  
1933 12. <https://doi.org/10.1177/1745691616658637>.
- 1934 Taylor, James W., and Kathryn S. Taylor. 2023. "Combining Probabilistic Forecasts of COVID-19  
1935 Mortality in the United States." *European Journal of Operational Research* 304 (1): 25–41.  
1936 <https://doi.org/https://doi.org/10.1016/j.ejor.2021.06.044>.
- 1937 Tierney, Nicholas, and Dianne Cook. 2023. "Expanding Tidy Data Principles to Facilitate Missing Data  
1938 Exploration, Visualization and Assessment of Imputations." *Journal of Statistical Software* 105 (7): 1–  
1939 31. <https://doi.org/10.18637/jss.v105.i07>.
- 1940 Touchon, Justin C., and Michael W. McCoy. 2016. "The Mismatch Between Current Statistical Practice  
1941 and Doctoral Training in Ecology." *Ecosphere* 7 (8): e01394.  
1942 <https://doi.org/https://doi.org/10.1002/ecs2.1394>.
- 1943 Ushey, Kevin, and Hadley Wickham. 2023. *renv: Project  
1944 Environments*. <https://rstudio.github.io/renv/>.
- 1945 van den Brand, Teun. 2024. *Ggh4x: Hacks for "ggplot2"*. <https://github.com/teunbrand/ggh4x>.
- 1946 Vander Werf, Eric. 1992. "Lack's Clutch Size Hypothesis: An Examination of the Evidence Using Meta-  
1947 Analysis." *Ecology* 73 (5): 1699–1705. <https://doi.org/10.2307/1940021>.
- 1948 Ver Hoef, Jay M. 2012. "Who Invented the Delta Method?" *The American Statistician* 66 (2): 124–  
1949 27. <https://doi.org/10.1080/00031305.2012.687494>.

1950 Verhulst, S., M. J. Holveck, and K. Riebel. 2006. “Long-Term Effects of Manipulated Natal Brood Size  
1951 on Metabolic Rate in Zebra Finches.” *Biology Letters* 2 (3): 478–  
1952 80. <https://doi.org/10.1098/rsbl.2006.0496>.

1953 Vesk, P. A., W. K. Morris, W. McCallum, R. Apted, and C. Miles. 2016. “Processes of Woodland  
1954 Eucalypt Regeneration: Lessons from the Bush Returns Trial.” *Proceedings of the Royal Society of  
1955 Victoria* 128: 54–63.

1956 Viechtbauer, Wolfgang. 2010. “Conducting Meta-Analyses in R with the metafor Package.” *Journal of  
1957 Statistical Software* 36 (3): 1–48. <https://doi.org/10.18637/jss.v036.i03>.

1958 Wickham, Hadley, Mara Averick, Jennifer Bryan, Winston Chang, Lucy D’Agostino McGowan, Romain  
1959 François, Garrett Grolemond, et al. 2019. “Welcome to the tidyverse.” *Journal of Open Source  
1960 Software* 4 (43): 1686. <https://doi.org/10.21105/joss.01686>.

1961 Wickham, Hadley, Jim Hester, Winston Chang, and Jennifer Bryan. 2022. *devtools: Tools to Make  
1962 Developing r Packages Easier*. <https://devtools.r-lib.org/>.

1963 Wickham, Hadley, Thomas Lin Pedersen, and Dana Seidel. 2023. *scales: Scale Functions for  
1964 Visualization*. <https://scales.r-lib.org>.

1965 Wilke, Claus O. 2024. *cowplot: Streamlined Plot Theme and Plot Annotations  
1966 for “ggplot2”*. <https://wilkelab.org/cowplot/>.

1967 Xie, Yihui. 2024a. *knitr: A General-Purpose Package for Dynamic Report Generation in  
1968 r*. <https://yihui.org/knitr/>.

1969 ———. 2024b. *xfun: Supporting Functions for Packages Maintained by “Yihui  
1970 Xie”*. <https://github.com/yihui/xfun>.

1971 Yang, Yefeng, Alfredo Sánchez-Tójar, Rose E. O’Dea, Daniel W. A. Noble, Julia Koricheva, Michael D.  
1972 Jennions, Timothy H. Parker, Malgorzata Lagisz, and Shinichi Nakagawa. 2023. “Publication Bias  
1973 Impacts on Effect Size, Statistical Power, and Magnitude (Type m) and Sign (Type s) Errors in Ecology  
1974 and Evolutionary Biology.” *BMC Biology* 21 (1): 71. <https://doi.org/10.1186/s12915-022-01485-y>.

1975 Zeileis, Achim, Jason C. Fisher, Kurt Hornik, Ross Ihaka, Claire D. McWhite, Paul Murrell, Reto Stauffer,  
1976 and Claus O. Wilke. 2020. “colorspace: A Toolbox for Manipulating and Assessing Colors and  
1977 Palettes.” *Journal of Statistical Software* 96 (1): 1–49. <https://doi.org/10.18637/jss.v096.i01>.

1978

## 1979 R Package References and Session Information

1980 Table 7: R packages used to generate this manuscript. Please see the ManyEcoEvo:: package for a full  
1981 list of packages used in the analysis pipeline.

Package	Version	Citation
base	4.4.0	<a href="#">R Core Team (2024)</a>
betapart	1.6	<a href="#">Baselga et al. (2023)</a>
broom.mixed	0.2.9.5	<a href="#">Bolker et al. (2024)</a>
colorspace	2.1.0	<a href="#">Zeileis et al. (2020)</a>
cowplot	1.1.3	<a href="#">Wilke (2024)</a>

devtools	2.4.5	<a href="#">Wickham et al. (2022)</a>
EnvStats	2.8.1	<a href="#">Millard (2013)</a>
GGally	2.2.1	<a href="#">Schloerke et al. (2024)</a>
ggforestplot	0.1.0	<a href="#">Scheinin et al. (2020)</a>
gg4x	0.2.8	<a href="#">van den Brand (2024)</a>
ggpubr	0.6.0	<a href="#">Kassambara (2023)</a>
ggrepel	0.9.5	<a href="#">Slowikowski (2024)</a>
ggthemes	5.1.0	<a href="#">Arnold (2024)</a>
glmmTMB	1.1.8	<a href="#">Brooks et al. (2017)</a>
gt	0.10.1	<a href="#">Iannone et al. (2024)</a>
gtsummary	1.7.2	<a href="#">Sjoberg et al. (2021)</a>
here	1.0.1	<a href="#">Müller (2020)</a>
Hmisc	5.1.2	<a href="#">Harrell Jr (2024)</a>
irr	0.84.1	<a href="#">Gamer, Lemon, and Singh (2019)</a>
janitor	2.2.0	<a href="#">Firke (2023)</a>
knitr	1.46	<a href="#">Xie (2024a)</a>
latex2exp	0.9.6	<a href="#">Meschiari (2022)</a>
lme4	1.1.35.3	<a href="#">Bates et al. (2015)</a>
ManyEcoEvo	2.7.6	<a href="#">Gould et al. (2023)</a>
metafor	4.6.0	<a href="#">Viechtbauer (2010)</a>
modelbased	0.8.7	<a href="#">Makowski et al. (2020)</a>
multilevelmod	1.0.0	<a href="#">Kuhn and Frick (2022)</a>
MuMIn	1.47.5	<a href="#">Bartoń (2023)</a>
naniar	1.1.0	<a href="#">Tierney and Cook (2023)</a>
NatParksPalettes	0.2.0	<a href="#">Blake (2022)</a>
orchaRd	2	<a href="#">Nakagawa, Lagisz, et al. (2023)</a>
parameters	0.21.7	<a href="#">Lüdecke et al. (2020)</a>
patchwork	1.2.0	<a href="#">Pedersen (2024)</a>
performance	0.11.0	<a href="#">Lüdecke, Ben-Shachar, et al. (2021)</a>
renv	1.0.2	<a href="#">Ushey and Wickham (2023)</a>
rmarkdown	2.27	<a href="#">Allaire et al. (2024)</a>
sae	1.3	<a href="#">Molina and Marhuenda (2015)</a>
scales	1.3.0	<a href="#">Wickham, Pedersen, and Seidel (2023)</a>
see	0.8.4	<a href="#">Lüdecke, Patil, et al. (2021)</a>
showtext	0.9.7	<a href="#">Qiu (2024)</a>
specr	1.0.0	<a href="#">Masur and Scharrow (2020)</a>
targets	1.7.0	<a href="#">Landau (2021)</a>
tidymodels	1.1.1	<a href="#">Kuhn and Wickham (2020)</a>
tidytext	0.4.2	<a href="#">Silge and Robinson (2016)</a>
tidyverse	2.0.0	<a href="#">Wickham et al. (2019)</a>
withr	3.0.0	<a href="#">Hester et al. (2024)</a>
xfun	0.44	<a href="#">Xie (2024b)</a>

1983

1984 — **Session info** —————

1985   setting   value

1986   version   R version 4.4.0 (2024-04-24)

1987   os           macOS Ventura 13.6.9

1988   system   aarch64, darwin20

1989   ui          X11

1990   language   (EN)

1991   collate   en\_US.UTF-8

1992   ctype      en\_US.UTF-8

1993   tz          Australia/Melbourne

1994   date       2024-09-17

1995   pandoc     3.1.12.2 @ /opt/homebrew/bin/ (via rmarkdown)

1996