



ECOLOGICAL SOCIETY OF AMERICA

Ecology/Ecological Monographs/Ecological Applications

PREPRINT

This preprint is a PDF of a manuscript that has been accepted for publication in an ESA journal. It is the final version that was uploaded and approved by the author(s). While the paper has been through the usual rigorous peer review process of ESA journals, it has not been copy-edited, nor have the graphics and tables been modified for final publication. Also note that the paper may refer to online Appendices and/or Supplements that are not yet available. We have posted this preliminary version of the manuscript online in the interest of making the scientific findings available for distribution and citation as quickly as possible following acceptance. However, readers should be aware that the final, published version will look different from this version and may also have some differences in content.

The doi for this manuscript and the correct format for citing the paper are given at the top of the online (html) abstract.

Once the final published version of this paper is posted online, it will replace the preliminary version at the specified doi.

1 **An evidence assessment tool for ecosystem services and**
2 **conservation studies**

3 Anne-Christine Mupepele^{1,2} & Jessica C. Walsh³ & William J. Sutherland³ &
4 Carsten F. Dormann¹

essa
preprint

¹ Department of Biometry and Environmental System Analysis, University of Freiburg, Tennenbacherstr. 4, 79106 Freiburg, Germany

² anne-christine.mupepele@biom.uni-freiburg.de

³ Department of Zoology, University of Cambridge, Downing Street, Cambridge, CB2 3EJ, UK.

Abstract

Reliability of scientific findings is important, especially if they directly impact decision making, such as in environmental management. In the 1990s, assessments of reliability in the medical field resulted in the development of evidence-based practice. Ten years later, evidence-based practice was translated into conservation, but so far no guidelines exist on how to assess the evidence of individual studies. Assessing the evidence of individual studies is essential to appropriately identify and synthesize the confidence in research findings. We develop a tool to assess the strength of evidence of ecosystem services and conservation studies. This tool consists of (1) a hierarchy of evidence, based on the experimental design of studies and (2) a critical-appraisal checklist that identifies the quality of research implementation. The application is illustrated with 13 examples and we suggest further steps required to move towards more evidence-based environmental management.

Keywords: governance - quality checklist - quantification - rigour - valuation

Conservation and ecosystem services studies are important scientific sources for decision-makers seeking advice on environmental management (Daily and Matson, 2008; Kareiva and Marvier, 2012). Their results potentially influence actions and it is therefore crucial to assess transparently the reliability of current research and its recommendations (Pullin and Knight, 2003; Boyd, 2013).

Evidence-based practice was introduced in the medical field aiming to assess the reliability of scientific statements and identify the best available information to answer a question of interest (Sackett *et al.*, 1996; GRADE Working Group,

26 2004; OCEBM Levels of Evidence Working Group, 2011, Cochrane
 27 Collaboration - www.cochrane.org). In conservation, evidence-based practice was
 28 first mentioned 15 years ago (Sutherland, 2000; Pullin and Knight, 2001). Today,
 29 the ‘Collaboration for Environmental Evidence’ (www.environmentalevidence.org)
 30 fosters the creation of systematic reviews to collate the strongest possible
 31 evidence (Petrokofsky *et al.*, 2011; Collaboration for Environmental Evidence,
 32 2013, see also *Journal for Environmental Evidence*), together with ‘Conservation
 33 Evidence’ (www.conservationevidence.org), which focuses on the development of
 34 summaries and guidelines, and the communication of evidence to practitioners
 35 (Sutherland *et al.*, 2012; Dicks *et al.*, 2014). Summaries, contrary to systematic
 36 reviews, do not focus on a specific question but bring together information from a
 37 much broader topic, e.g. from a whole animal group, such as bees (Dicks *et al.*,
 38 2010, 2014; Walsh *et al.*, 2015).

39 Systematic reviews and summaries compile individual studies and therefore
 40 require the evaluation of the evidence at the level of the individual study. In
 41 systematic reviews this is typically mentioned as one step of the critical appraisal.
 42 However, to date such critical appraisal is often implicit, based on criteria varying
 43 for every systematic review (Collaboration for Environmental Evidence, 2013;
 44 Carroll and Booth, 2015; Stewart and Schmid, 2015). We therefore introduce an

45 evidence assessment tool providing a clear appraisal guideline to score the
 46 reliability of individual studies.

47 **Evidence assessment tool**

48 A well-defined terminology is essential for effective communication between
 49 practitioners and scientists. **Evidence** is the ‘ground for belief’ or ‘the available
 50 body of information indicating whether a belief or proposition is true or valid’
 51 (Howick, 2011, OED Online, Oxford University Press, September 2015; Oxford
 52 Dictionaries: www.oxforddictionaries.com/definition/english/evidence). Evidence describes
 53 the knowledge behind a statement and expresses how solid our recommendations
 54 are (see also Higgs and Jones 2000, p.311; Rychetnik *et al.* 2001; Lohr 2004;
 55 Binkley and Menyailo 2005; Pullin and Knight 2005). The **strength of evidence**
 56 reflects the reliability of information and we can identify whether a statement is
 57 based on **strong or weak evidence**, i.e. very reliable or hardly reliable. Hence
 58 **evidence-based practice** means to *identify* the reliability of current knowledge,
 59 based on research integrated with expertise, and to *act* according to this best
 60 available knowledge. The collation and appraisal of the best available evidence
 61 follow strict criteria to ensure transparency and to reduce bias. A goal of
 62 evidence-based practice is to act on best available evidence while being aware of

63 the strength of inference this evidence permits (Howick, 2011, p.15).

64 **1. Setting question and context**

65 The formulation of a clear research question and the purpose of investigation is
 66 highly emphasized throughout the evidence literature (Higgins and Green, 2011;
 67 Collaboration for Environmental Evidence, 2013, p.20-23). Questions should
 68 specify *which* ecosystem service, species or aspect of biodiversity will be
 69 investigated in *which* system, as this will help to determine the external validity
 70 of the answer provided in a study.

71 We further recommend to determine the focus of the question, as either
 72 ‘quantification’, ‘valuation’, ‘management’ or ‘governance’. **Quantification**
 73 studies measure the amount of an ecosystem service, species abundance,
 74 biodiversity or other conservation targets. Measures can be taken in absolute
 75 units or relative to another system. **Valuation** studies assess the societal value of
 76 ecosystem services. The most common way is monetary valuation. **Management**
 77 is the treatment designed to improve or benefit specific ecosystem services, target
 78 species or other conservation aspects. For example: leaving dead wood in forests
 79 to increase biodiversity, or reducing agricultural fertiliser to decrease nearby lake
 80 eutrophication. **Governance** is seen as the strategy or policy to steer a

81 management intervention, such as REDD (Reducing Emissions from
 82 Deforestation and Forest Degradation), which aims to encourage forest protection
 83 and reforestation (Kenward *et al.*, 2011). The strategies used by policy makers
 84 include incentives (subsidies) or penalties (law/tax; see also Bevir, 2012).
 85 When the effectiveness of management and governance strategies is determined,
 86 evidence-based quantification or valuation is required to measure the outcome of
 87 the management or governance intervention. Acuña *et al.* (2013), for example,
 88 used valuation methods to determine success or failure of a management strategy
 89 while Walsh *et al.* (2012) quantified malleefowl abundance through monitoring
 90 survey data to assess the management impact of fox baiting. The distinction of
 91 four different foci is essential to assess the whole range of environmental
 92 management.

93 We have described how to set the context of questions that can be useful in
 94 environmental management. Once the question has been determined, and the
 95 investigation carried out, the strength of the resulting evidence should be assessed
 96 (Fig. 1).

97 **2. Evidence assessment**

98 The reliability of a study is characterized by its study design and the quality of its
 99 implementation. Both are evaluated in the evidence assessment.

100 **2a. Evidence hierarchy**

101 The study design refers to the set-up of the investigation, e.g. controlled or
 102 observational design (GRADE Working Group, 2004). These study designs are
 103 not equally compelling with respect to inferring causality. Differences in study
 104 designs typically translate into weak or strong evidence. To identify the
 105 reliability of a study, study designs can be ranked hierarchically according to a
 106 level-of-evidence scale, hence forth the evidence hierarchy (Fig. 2).

107 **Systematic reviews (LoE1a)** are at the top of the evidence hierarchy and
 108 provide the most reliable information. They summarize all information collated
 109 in several individual studies, have an *a priori* protocol on design and procedure,
 110 and are conducted according to strict guidelines (e.g. Collaboration for
 111 Environmental Evidence, 2013). If possible, they ideally include quantitative
 112 measures, i.e. a meta-analysis (see Koricheva *et al.*, 2013; Vetter *et al.*, 2013).
 113 All other, non-systematic and more **conventional reviews (LoE1b)** may also
 114 include quantitative analysis or are purely qualitative. Both types of review

115 summarize the findings of several studies, but systematic reviews assess the
 116 completeness and reproducibility more carefully and strive to reduce bias by
 117 having transparent, thorough, pre-defined methods (Freeman *et al.*, 2006;
 118 Higgins and Green, 2011; Collaboration for Environmental Evidence, 2013;
 119 Haddaway and Bayliss, 2015; Haddaway and Bilotta, 2015).

120 The necessary condition for any review is that appropriate individual studies
 121 are available. The most reliable individual study design is **a study with a**
 122 **reference/control (LoE2)**. Typically, these are case-control or before-after
 123 control-impact studies (**LoE2a**) (Smith *et al.*, 2014). Investigations that cannot
 124 follow such a controlled design may alternatively seek to gain strong evidence
 125 through multiple lines of moderate evidence (**LoE2b**). Multiple lines of evidence
 126 require at least two unrelated and consistent arguments to confirm the study
 127 conclusions, thereby forming a non-contradicting picture (see also Smith *et al.*,
 128 2002). Illustrative examples are the valuation of ecosystem services (e.g. Mogas
 129 *et al.*, 2006), or long-term environmental processes that are difficult to control
 130 (e.g. Dorman *et al.*, 2015). Multiple lines of evidence can be collected in
 131 individual studies using different approaches within one study context (LoE2b,
 132 LoE3c) or in reviews (LoE1) including evidence from different studies.

133 **Observational studies (LoE3)** are individual studies without a control. These

134 include studies employing inferential and correlative statistics (**LoE3a**), e.g.
 135 testing for the influence of environmental variables on the quantity of an
 136 ecosystem service. Descriptive studies (**LoE3b**) imply data collection and
 137 representation without statistical testing (e.g. data summaries, ordinations,
 138 histograms, surveys). Multiple lines of weak evidence (**LoE3c**) can increase the
 139 evidence of LoE4 investigations; elicitation of independent expert opinions is a
 140 currently well-known example (Sutherland *et al.*, 2013; Morgan, 2014; Smith
 141 *et al.*, 2015; Sutherland and Burgman, 2015, see also Appendix).

142 The lowest level of evidence are statements **without underlying data (LoE4)**.
 143 These are usually individual expert opinions, often not distinguishable from
 144 randomness (Tetlock, 2005; Drolet *et al.*, 2015). Other statements without
 145 underlying data are reasoning based on mechanism. Mechanism-based reasoning
 146 involves an inferential chain linking an intervention to the outcome (Howick
 147 *et al.*, 2010; Howick, 2011). If this chain of mechanisms is not supported by data,
 148 there is no possibility to assess whether all relevant mechanisms linking the
 149 intervention to the outcome have been included. Mechanism-based reasoning
 150 without corroborative data provides only weak evidence. On the other hand,
 151 mechanism-based reasoning can result in a model that is validated and tested on
 152 real world data. With such a data validation, the model could reach moderate

153 evidence or strong evidence, depending on the underlying study design.

154 It is important to note that ‘method’ and ‘design’ should not be confused.

155 Methods are the means used to collect or analyse data, e.g. remote sensing,
 156 questionnaires, ordination techniques. Design reflects how the study was planned
 157 and conducted, e.g. a case-control or observational design (GRADE Working
 158 Group, 2004). The same methods can be employed for different underlying
 159 designs. Remote sensing for example can be done purely descriptively (LoE3b)
 160 or with a reference such as ground-truthing or in a ‘before-and-after’ design
 161 (LoE2a). Analogously, models can represent theories without supporting data
 162 (LoE4), involve data input to determine parameters (LoE3b) or be tested and
 163 validated (LoE3a). To achieve strong evidence, model predictions have to be
 164 confirmed by several unrelated data sets forming a non-contradicting picture
 165 (LoE2b), or should be built on information derived from controlled studies
 166 unequivocally identifying the underlying causal mechanism (LoE2a; Kirchner,
 167 2006).

168 **2b. Critical appraisal**

169 Study design alone is an inadequate marker of the strength of evidence
 170 (Rychetnik *et al.*, 2001). A study with a strong-evidence design may be poorly

171 conducted. The critical appraisal assesses the implementation of the study
 172 design, specifically the methodological quality, the actual realization of the study
 173 design and its reporting (Higgins and Green, 2011). It identifies the study quality
 174 and may lead to a downgrading in the evidence hierarchy. **Quality**, in this
 175 context, is the extent to which all aspects of conducting a study can be shown to
 176 protect against bias, and inferential error (Lohr, 2004). Quality checklists can be
 177 used to detect bias and inferential error. Combining 30 published quality
 178 checklists, we provide the first quality checklist for conservation and ecosystem
 179 services (Appendix Table 1), that can be used to comprehensively assess the
 180 internal validity of a study, covering questions on data collection, analysis and the
 181 presentation of results. The checklist consists of 43 questions, of which some
 182 apply only to a specific context, e.g. for reviews or only studies focusing on
 183 valuation. All questions answered with ‘yes’ receive one point. In the case of
 184 non-reported issues, we advise the answer ‘no’ to indicate a deficient reporting
 185 quality. The percentage of points received can help to decide whether to
 186 downgrade the level of evidence (Appendix Table 2).

187 Reviews provide information at the highest level of evidence and their critical
 188 appraisal is different from other designs, because they are based on studies with
 189 weaker evidence (see Appendix Table 1: Review). Every single study included in

190 the review can be assessed for its level of evidence, using the evidence hierarchy
 191 and the checklist for quality criteria. If only studies based on weak evidence were
 192 included, then the review should be downgraded, regardless of other quality
 193 criteria. In addition, a review can be assessed for other quality shortcomings
 194 using again the quality checklist.

195 The checklist should make the assessment more transparent, but we are aware
 196 that the process may not always be straightforward. Questions in the checklist
 197 can be subjective and depend on the judgment of the assessor. Cohen's kappa test
 198 was used to test the agreement in 13 exemplary studies between two different
 199 assessors (Appendix Table 3). It ranges from 0 to 1, representing random to
 200 perfect agreement. Our result revealed a moderate agreement (unweighted
 201 Cohen's kappa = 0.49; p-value < 0.001. Landis and Koch, 1977; Cohen, 1960;
 202 Gamer *et al.*, 2015). Depending on the context, the assessor may decide to give
 203 more weight to particular questions or add questions to the checklist. Although
 204 the procedure cannot be fully standardized, we are not aware of a better
 205 alternative, and we encourage the use of the checklist as a baseline that can be
 206 adapted for specific studies.

207 The combination of study design (Fig. 2) and quality criteria (Appendix Table
 208 1) is the last step and identifies the strength of evidence supporting the study

209 result (schematic representation in Fig. 1). The level of evidence derived by the
 210 study design should be downgraded depending on the quality score calculated
 211 from the quality checklist (Appendix Table 2).

212 **Application of the evidence assessment tool**

213 The suggested method was applied to assess the evidence of 13 studies
 214 (Appendix Table 3). They were selected to serve as examples and illustrate the
 215 applicability of the tool to the whole range of study designs and foci. The first
 216 example was a management-related systematic review of Mant *et al.* (2013),
 217 conducted according to the guidelines of the Collaboration for Environmental
 218 Evidence (2013). They investigated the effect of ‘liming’ rivers or lakes on fish
 219 and invertebrate populations. They found that liming increased fish abundances
 220 and acid-sensitive invertebrates, but may have a negative impact on the
 221 abundance of all invertebrate taxa combined. According to the critical appraisal
 222 the study achieved 21 out of 24 points (88%) and it therefore remained at the
 223 originally assigned LoE1a, the highest level of evidence (Appendix Table 3).

224 A second example tackles the question: ‘How does adding dead wood to rivers
 225 influence the provision of ecosystem services?’ (Acuña *et al.*, 2013). The authors
 226 investigated two ecosystem services (fishing and retention of organic and

227 inorganic matter) in a river-forest ecosystem in Spain and Portugal and studied
 228 the effect of this management intervention. Their study design followed a
 229 before-after control-impact approach, equivalent to LoE2a. The critical appraisal
 230 revealed shortcomings, e.g. no blinding, no randomization and no probability
 231 sampling: only 17 out of 25 points (68%) were achieved. The level of evidence
 232 was downgraded by one level to LoE3a. We therefore conclude that the statement
 233 made by Acuña *et al.* (2013): ‘restoration of natural wood loading in streams
 234 increases the ecosystem service provision’ is based on moderate evidence
 235 (LoE3a).

236 We provide further examples in the Appendix (Appendix Table 3 and 4,
 237 GitHub: <https://github.com/biometry/EvidenceAssessmentTool/blob/master/Examples.xlsx>). All but
 238 one study revealed quality shortcomings and had to be downgraded. Most were
 239 scored as LoE3 or LoE4.

240 **Relevance for different user groups**

241 In the previous section it was elaborated *how* to assess the strength of evidence
 242 for individual studies and reviews. Now we provide a few notes on *who* should
 243 use it:

244 **1. Scientists conducting their own studies** have to be aware of how to achieve

245 strong evidence, particularly during the planning phase. Choosing a study design
 246 that provides strong evidence and respects the quality criteria will substantially
 247 increase the potential contribution to our knowledge.

248 **2. Scientists advising decision-makers** should be explicit about the strength of
 249 evidence of information they include in their recommendations. Weighting all
 250 scientific information equally, or subjectively, runs the risk of overconfidence and
 251 bias.

252 **3. Decision-makers** receiving information from scientists should demand a
 253 level-of-evidence statement for the information provided. Alternatively, they can
 254 assess the strength of evidence themselves. However, this may be difficult as it
 255 takes time and requires some scientific training to identify the study design and
 256 evaluate the quality questions.

257 **4. We further encourage consortia, international panels and learned societies,**
 258 such as the Intergovernmental Platform on Biodiversity & Ecosystem Services
 259 (IPBES), the Ecological Societies (ESA, BES, GFÖ and others), the Society for
 260 Conservation Biology (SCB) and the Ecosystem Services Partnership (ESP) to
 261 support the development of guidelines, that include an evidence assessment
 262 (Graham *et al.*, 2011; Sutherland *et al.*, 2015). These ‘best-practice guides’ are
 263 based on the collection of scientific evidence synthesized and judged by a group

264 of experts. They provide recommendations on how to best quantify, value,
 265 manage or govern a desired ecosystem service or conservation target, giving
 266 decision-makers transparent advice with an emphasis on the strength of the
 267 evidence available (for examples of equivalent Clinical Guidelines see
 268 www.guideline.gov (USA), www.ncgc.ac.uk (UK), www.awmf.org/leitlinien
 269 (Germany)).

270 Discussion

271 We have outlined an evidence assessment tool for ecosystem services and
 272 conservation studies, encompassing a hierarchy to judge the available evidence
 273 based on study design and a quality checklist to facilitate critical appraisal. We
 274 have further illustrated with examples how to apply the tool (see also Appendix
 275 Table 3 and 4).

276 Evidence-based practice seeks to complement existing management
 277 frameworks, by emphasizing the importance of systematically collating the
 278 existing scientific evidence and assessing it for its reliability and relevance. The
 279 IPCC report, for example, uses a combined measure of evidence and level of
 280 agreement (Mastrandrea *et al.*, 2010; Spiegelhalter and Riesch, 2011). Our
 281 suggested approach is more detailed, describing *how* one can actually assess the

282 evidence.

283 Evidence-based practice has faced criticism of its evidence hierarchies,
 284 claiming that controlled trials are not always more reliable than observational
 285 studies. A main argument against hierarchies is that they are rigid and only
 286 consider the study design to assign a level of evidence (Petticrew and Roberts,
 287 2003; Adams and Sandbrook, 2013; Stegenga, 2014). With our quality checklist
 288 we emphasize the critical appraisal to check for an appropriate implementation
 289 and methodological quality of study designs. The proposed assessment therefore
 290 does not overestimate the results of deficiently implemented meta-analyses and
 291 controlled studies. Some science sectors have to rely on observational studies,
 292 because their study units cannot be controlled. This usually applies to
 293 environmental governance, conservation biology of rare species, or global
 294 theories that lack a second earth as a control. Multiple lines of evidence can lead
 295 to strong evidence using only observational study designs (Fig. 2, LoE2b).
 296 However, a central task of natural science is to determine causal relationships,
 297 and observational studies do not have the same strength to determine causal
 298 relationships than replicated and randomized case-control studies (Holland, 1986;
 299 Grimes and Schulz, 2002; Illari *et al.*, 2011). We should acknowledge that in
 300 some areas of science causality cannot be established, and hence the reliability

301 achieved remains lower than in areas where it can.

302 Other criticism has been directed towards the fact that every system is unique
 303 and the external validity of studies is low. We are aware that generalizability of
 304 results is problematic in ecosystems, where many different drivers take influence
 305 at the same time and hence the general evidence may not apply due to particular
 306 circumstances. At this point the judgment of experts on the external validity of
 307 the currently best available evidence is irreplaceable (Karanicolas *et al.*, 2008;
 308 Howick, 2011). Evidence-based practice means integrating individual expertise
 309 with the best available evidence from systematic research (Sackett *et al.*, 1996;
 310 Straus *et al.*, 2010). More reflection and responses to criticism of evidence-based
 311 practice can be found in Mullen and Streiner (2004), Sutherland *et al.* (2004,
 312 2005) and Haddaway and Pullin (2013).

313 Despite the criticism raised against evidence-based practice the benefits are
 314 clear (Gilbert *et al.*, 2005; Howick, 2011; Walsh *et al.*, 2014, 2015). Rating the
 315 strength of evidence matters as it clarifies the reliability of research results and,
 316 thus, the strength of conclusions, decisions, or recommendations drawn from that
 317 research (Lohr, 2004).

318 Reliable scientific evidence in environmental management is pivotal, and its
 319 use (or misuse) can have immense impacts on environmental outcomes and the

320 society. It is essential that scientists and decision makers consider the strength of
321 evidence when conducting studies, providing advice and taking decisions. In the
322 interest of responsible use of environmental resources and processes, we strongly
323 encourage embracing evidence-based practice as a paradigm for all research
324 contributing to environmental management.

325 **Acknowledgements**

326 We thank Andrew Pullin, Sven Lautenbach and Ian Bateman for valuable
327 comments on earlier versions of the manuscript. This work was supported by the
328 7th framework programme of the European Commission in the project OPERAs
329 (grant number 308393, www.operas-project.eu).

330 **References**

- 331 Acuña V, Díez JR, Flores L, *et al.* 2013. Does it make economic sense to restore
 332 rivers for their ecosystem services? *Journal of Applied Ecology* **50**: 988–997.
- 333 Adams WM and Sandbrook C. 2013. Conservation, evidence and policy. *Oryx*
 334 **47**: 329–335.
- 335 Bevir M. 2012. Governance: A Very Short Introduction. Oxford University Press.
- 336 Binkley D and Menyailo O. 2005. Gaining insights on the effects of tree species
 337 on soils. In: Tree Species Effects on Soils: Implications for Global Change,
 338 chapter 1, 1–16. Dordrecht: Kluwer Academic Publishers.
- 339 Boyd I. 2013. A standard for policy-relevant science. *Nature* **501**: 159–160.
- 340 Carroll C and Booth A. 2015. Quality assessment of qualitative evidence for
 341 systematic review and synthesis: Is it meaningful, and if so, how should it be
 342 performed? *Research Synthesis Methods* **6**: 149–154.
- 343 Cohen J. 1960. A coefficient of agreement for nominal scales. *Educational and*
 344 *psychological measurement* **20**: 37–46.
- 345 Collaboration for Environmental Evidence. 2013. Guidelines for Systematic
 346 Review and Evidence Synthesis in Environmental Management. Version 4.2.
 347 Environmental Evidence. URL
 348 www.environmentalevidence.org/Documents/Guidelines/Guidelines4.2.pdf.
- 349 Daily GC and Matson PA. 2008. Ecosystem services: from theory to
 350 implementation. *Proceedings of the National Academy of Sciences* **105**:
 351 9455–9456.

- 352 Dicks LV, Showler DA, and Sutherland WJ. 2010. Bee conservation: evidence
 353 for the effects of interventions. Pelagic Publishing.
- 354 Dicks LV, Walsh JC, and Sutherland WJ. 2014. Organising evidence for
 355 environmental management decisions: a '4S' hierarchy. *Trends in Ecology &*
 356 *Evolution* **29**: 1–7.
- 357 Dorman M, Svoray T, Perevolotsky A, *et al.* 2015. What determines tree
 358 mortality in dry environments? a multi-perspective approach. *Ecological*
 359 *Applications* **25**: 1054–1071.
- 360 Drolet D, Locke A, Lewis MA, and Davidson J. 2015. Evidence-based tool
 361 surpasses expert opinion in predicting probability of eradication of aquatic
 362 nonindigenous species. *Ecological Applications* **25**: 441–450.
- 363 Freeman SR, Williams HC, and Dellavalle RP. 2006. The increasing importance
 364 of systematic reviews in clinical dermatology research and publication. *The*
 365 *Journal of investigative dermatology* **126**: 2357–2360.
- 366 Gamer M, Lemon J, Fellows I, and Singh P. 2015. The irr package for R: various
 367 coefficients of interrater reliability and agreement.
- 368 Gilbert R, Salanti G, Harden M, and See S. 2005. Infant sleeping position and the
 369 sudden infant death syndrome: systematic review of observational studies and
 370 historical review of recommendations from 1940 to 2002. *International*
 371 *Journal of Epidemiology* **34**: 874–87.
- 372 GRADE Working Group. 2004. Grading quality of evidence and strength of
 373 recommendations. *BMJ* **328**: 1–8.

- 374 Graham R, Mancher M, Wolmann DM, *et al.* 2011. Clinical Practice Guidelines
 375 We Can Trust. Washington, DC: The National Academies Press.
- 376 Grimes DA and Schulz KF. 2002. Descriptive studies: what they can and cannot
 377 do. *Lancet* **359**: 145–149.
- 378 Haddaway NR and Bayliss HR. 2015. Clarification on the applicability of
 379 systematic reviews. *Frontiers in Ecology and the Environment* **13**: 129–129.
- 380 Haddaway NR and Bilotta GS. 2015. Systematic reviews: Separating fact from
 381 fiction. *Environment International* (in press).
- 382 Haddaway NR and Pullin AS. 2013. Evidence-based conservation and
 383 evidence-informed policy: a response to Adams & Sandbrook. *Oryx* **47**:
 384 336–338.
- 385 Higgins JPT and Green S. 2011. Cochrane Handbook for Systematic Reviews of
 386 Interventions. Version 5.1.0. [updated March 2011]. The Cochrane
 387 Collaboration.
- 388 Higgs J and Jones M. 2000. Will evidence-based practice take the reasoning out
 389 of practice? In: Higgs J and Jones M (Eds.) *Clinical Reasoning in the Health*
 390 *Professionals*, 307–315. Oxford: Butterworth Heineman, 2 edition.
- 391 Holland PW. 1986. Statistics and causal inference. *Journal of the American*
 392 *Statistical Association* **81**: 945–960.
- 393 Howick J. 2011. *The Philosophy of Evidence-Based Medicine*. Oxford, UK:
 394 Wiley-Blackwell.

- 395 Howick J, Glasziou P, and Aronson JK. 2010. Evidence-based mechanistic
 396 reasoning. *Journal of the Royal Society of Medicine* **103**: 433–441.
- 397 Illari PM, Russo F, and Williamson J. 2011. *Causality in the Sciences*. Oxford
 398 University Press.
- 399 Karanickolas PJ, Kunz R, and Guyatt GH. 2008. Evidence-based medicine has a
 400 sound scientific base. *CHEST* **133**: 1067.
- 401 Kareiva P and Marvier M. 2012. What is conservation science? *BioScience* **62**:
 402 962–969.
- 403 Kenward RE, Whittingham MJ, Arampatzis S, *et al.* 2011. Identifying
 404 governance strategies that effectively support ecosystem services, resource
 405 sustainability, and biodiversity. *Proceedings of the National Academy of
 406 Sciences* **108**: 5308–5312.
- 407 Kirchner JW. 2006. Getting the right answers for the right reasons: linking
 408 measurements, analyses, and models to advance the science of hydrology.
 409 *Water Resources Research* **42**: 1–5.
- 410 Koricheva J, Gurevitch J, and Mengersen K. 2013. *Handbook of Meta-analysis in
 411 Ecology and Evolution*. Princeton University Press.
- 412 Landis JR and Koch GG. 1977. The measurement of observer agreement for
 413 categorical data. *Biometrics* **33**: 159–174.
- 414 Lohr KN. 2004. Rating the strength of scientific evidence: relevance for quality
 415 improvement programs. *International Journal for Quality in Health Care* **16**:
 416 9–18.

- 417 Mant RC, Jones DL, Reynolds B, *et al.* 2013. A systematic review of the
 418 effectiveness of liming to mitigate impacts of river acidification on fish and
 419 macro-invertebrates. *Environmental Pollution* **179**: 285–293.
- 420 Mastrandrea M, Field C, Stocker TF, *et al.* 2010. Guidance note for lead authors
 421 of the IPCC fifth assessment report on consistent treatment of uncertainties.
- 422 Mogas J, Riera P, and Bennett J. 2006. A comparison of contingent valuation and
 423 choice modelling with second-order interactions. *Journal of Forest Economics*
 424 **12**: 5–30.
- 425 Morgan MG. 2014. Use (and abuse) of expert elicitation in support of decision
 426 making for public policy. *Proceedings of the National Academy of Sciences*
 427 **111**: 7176–7184.
- 428 Mullen EJ and Streiner DL. 2004. The evidence for and against evidence-based
 429 practice. *Brief Treatment and Crisis Intervention* **4**: 111–121.
- 430 OCEBM Levels of Evidence Working Group. 2011. The Oxford Levels of
 431 Evidence 1. URL <http://www.cebm.net/index.aspx?o=5653>.
- 432 Petrokofsky G, Holmgren P, and Brown ND. 2011. Reliable forest carbon
 433 monitoring-systematic reviews as a tool for validating the knowledge base.
 434 *International Forestry Review* **13**: 56–66.
- 435 Petticrew M and Roberts H. 2003. Evidence, hierarchies, and typologies: horses
 436 for courses. *Theory and Methods* **57**: 527–529.
- 437 Pullin AS and Knight TM. 2001. Effectiveness in conservation practice: pointers
 438 from medicine and public health. *Conservation Biology* **15**: 50–54.

- 439 Pullin AS and Knight TM. 2003. Support for decision making in conservation
 440 practice: an evidence-based approach. *Journal for Nature Conservation* **11**:
 441 83–90.
- 442 Pullin AS and Knight TM. 2005. Assessing conservation management’s evidence
 443 base: a survey of management-plan compilers in the United Kingdom and
 444 Australia. *Conservation Biology* **19**: 1989–1996.
- 445 Rychetnik L, Frommer M, Hawe P, and Shiell A. 2001. Criteria for evaluation
 446 evidence on public health interventions. *Journal of Epidemiology and*
 447 *Community Health* **56**: 119–127.
- 448 Sackett DL, Rosenberg WMC, Gray JAM, *et al.* 1996. Evidence based medicine:
 449 what it is and what it isn’t. *Clinical Orthopaedics and Related Research* **455**:
 450 3–5.
- 451 Smith EP, Lipkovich I, and Ye K. 2002. Weight-of-Evidence (WOE): quantitative
 452 estimation of probability of impairment for individual and multiple lines of
 453 evidence. *Human and Ecological Risk Assessment: An International Journal* **8**:
 454 1585–1596.
- 455 Smith RK, Dicks LV, Mitchell R, and Sutherland WJ. 2014. Comparative
 456 effectiveness research: the missing link in conservation. *Conservation*
 457 *Evidence* **11**: 2–6.
- 458 Smith SDP, McIntyre PB, Halpern BS, *et al.* 2015. Rating impacts in a
 459 multi-stressor world: a quantitative assessment of 50 stressors affecting the
 460 Great Lakes. *Ecological Applications* **25**: 717–728.

- 461 Spiegelhalter DJ and Riesch H. 2011. Don't know, can't know: embracing deeper
 462 uncertainties when analysing risks. *Philosophical Transactions of the Royal*
 463 *Society A* **369**: 4730–50.
- 464 Stegenga J. 2014. Down with the hierarchies. *Topoi* **33**: 313–322.
- 465 Stewart GB and Schmid CH. 2015. Lessons from meta-analysis in ecology and
 466 evolution: the need for trans-disciplinary evidence synthesis methodologies.
 467 *Research Synthesis Methods* **6**: 109–110.
- 468 Straus SE, Glasziou P, Richardson WS, and Haynes RB. 2010. Evidence-Based
 469 Medicine: How to Practice and Teach It, 4e (Straus, Evidence-Based
 470 Medicine). Churchill Livingstone.
- 471 Sutherland WJ. 2000. The Conservation Handbook: Research, Management and
 472 Policy. Blackwell Science Ltd.
- 473 Sutherland WJ and Burgman MA. 2015. Use experts wisely. *Nature* **526**: 317.
- 474 Sutherland WJ, Dicks LV, and Smith RK. 2015. What Works in Conservation?
 475 Lessons from Conservation Evidence. OpenBooks, Cambridge.
- 476 Sutherland WJ, Gardner TA, Haider LJ, and Dicks LV. 2013. How can local and
 477 traditional knowledge be effectively incorporated into international
 478 assessments? *Oryx* **48**: 1–2.
- 479 Sutherland WJ, Mitchell R, and Prior SV. 2012. The role of 'Conservation
 480 Evidence' in improving conservation management. *Conservation Evidence* **9**:
 481 1–2.

- 482 Sutherland WJ, Pullin AS, Dolman PM, and Knight TM. 2004. Response to
 483 Griffiths. Mismatches between conservation science and practice. *Trends in*
 484 *Ecology & Evolution* **19**: 565–566.
- 485 Sutherland WJ, Pullin AS, Dolman PM, and Knight TM. 2005. Response to
 486 Mathevet and Mauchamp: Evidence-based conservation: dealing with social
 487 issues. *Trends in Ecology & Evolution* **20**: 424–425.
- 488 Tetlock PE. 2005. *Expert Political Judgment: How Good Is It? How Can We*
 489 *Know?* Princeton University Press.
- 490 Vetter D, Rucker G, and Storch I. 2013. Meta-analysis: A need for well-defined
 491 usage in ecology and conservation biology. *Ecosphere* **4**.
- 492 Walsh JC, Dicks LV, and Sutherland WJ. 2014. The effect of scientific evidence
 493 on conservation practitioners' management decisions. *Conservation Biology*
 494 **00**: 1–11.
- 495 Walsh JC, Dicks LV, and Sutherland WJ. 2015. The effect of scientific evidence
 496 on conservation practitioners management decisions. *Conservation Biology* **29**:
 497 88–98.
- 498 Walsh JC, Wilson KA, Benshemesh J, and Possingham HP. 2012. Integrating
 499 research, monitoring and management into an adaptive management framework
 500 to achieve effective conservation outcomes. *Animal Conservation* **15**: 334–336.

501 **Appendix A**

502 The appendix provides details and examples for the application of the evidence
503 assessment tool. The quality checklist is given in Table 1. Table 2 guides the
504 downgrading of the level of evidence according to the quality score. We further
505 present the evidence assessment of all 13 examples, together with the detailed
506 quality checklist filled in for each study (also available on GitHub:
507 <https://github.com/biometry/EvidenceAssessmentTool/blob/master/Examples.xlsx>).

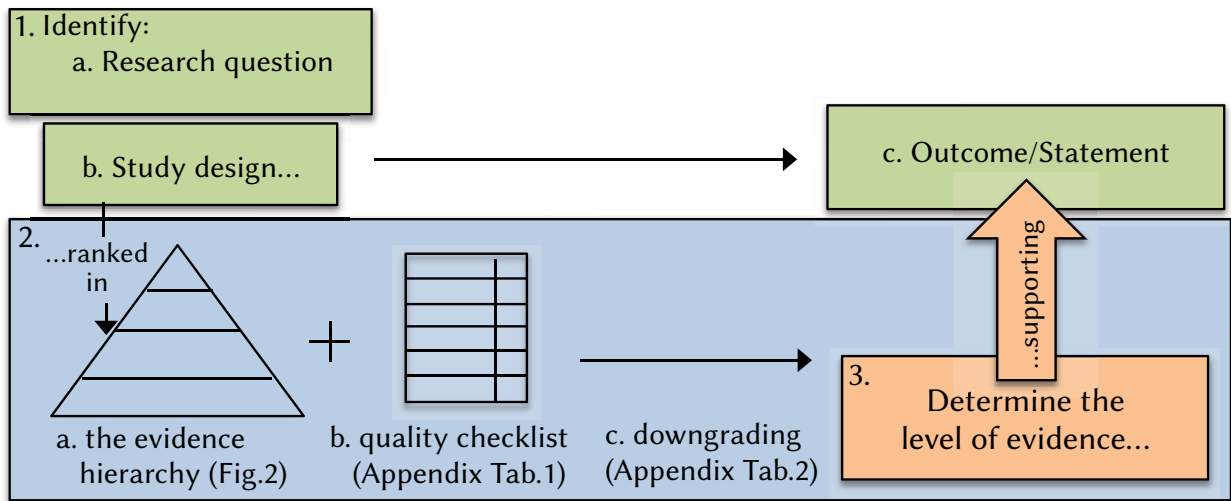
preprint

508 **Figure Legends**

509 **Figure 1:** Schematic representation of the evidence assessment tool: 1.
 510 Identification of study question, design and outcome. 2. Assessing a level of
 511 evidence based on the underlying study design and calculating the quality score
 512 based on the quality checklist. 3. Determine the final level of evidence supporting
 513 the outcome by downgrading the originally assigned level of evidence according
 514 to the quality score.

515 **Figure 2:** Level-of-evidence (LoE) hierarchy ranking study designs according to
 516 their evidence. Very strong evidence (LoE1) to weak evidence (LoE4) with
 517 internally ranked sublevels a, b and c.

518 **Figure 1**



preprint

519 **Figure 2**

