# Guidelines for Systematic Review in Conservation and Environmental Management.

**CENTRE FOR EVIDENCE-BASED CONSERVATION**

**SCHOOL OF BIOSCIENCES**

**UNIVERSITY OF BIRMINGHAM**

**EDGBASTON**

**BIRMINGHAM B15 2TT, UK**

**VERSION 2.0**

**AUGUST 2006**

<u>**Preamble**</u>

In response to problems of accessing scientific information to support decision-making, many applied disciplines are utilising an evidence-based framework for knowledge transfer involving systematic review and dissemination of evidence on effectiveness of interventions at the practical and policy levels (Stevens & Milne 1997; Khan 2003). The framework is most fully developed in the health services sector, where global review and dissemination units have been established and are linked by networks such as the Cochrane Collaboration (www.cochrane.org). Within these networks systematic reviews are undertaken following set guidelines that include peer review to ensure that they meet required standards before dissemination. The need for such a framework in conservation has been argued elsewhere (Pullin & Knight 2001; Fazey et al. 2004; Pullin et al. 2004; Sutherland et al. 2004). Here we present the latest guidelines for systematic review and dissemination in conservation and environmental management.

In the following we have used established guidelines from the health services sector (NHS CRD 2001; Khan 2003; Higgins & Green 2005) as our models, undertaken our own systematic reviews to test these models, and modified the guidelines through analysis of procedures and outcomes for their application to conservation and environmental management. Although the basic ethos of systematic review remains unchanged, ecological data are often fundamentally different in nature from data on human health (Fazey et al. 2004; Pullin et al. 2004), and this is reflected in our guidelines. At first glance many of the guidelines may seem routine and common sense, but the rigour and objectivity applied at key stages, and the underlying philosophy of transparency and independence, sets them apart from the majority of traditional reviews recently published in the field of applied ecology (Roberts et al. 2006). Pullin and Knight (2001), Fazey et al. (2004), Pullin et al. (2004), and Sutherland et al. (2004) argue that, once established, systematic review methodology will significantly improve the identification and provision of evidence to support practice and policy in conservation and environmental management. For this methodology to have an impact on conservation effectiveness, more conservation biologists need to undertake reviews, and we encourage this community to use (and improve) these guidelines and help establish an evidence-based framework for our discipline.


<u>**Systematic Review Guidelines**</u>

For clarity the guidelines are split into three stages and key phases within each. We use examples of our own reviews to highlight key issues.


**Stage 1 - Planning the review**

**Question formulation**
A systematic review starts with a specific question, clearly defined with subject, intervention and outcome elements (Table 1), that is answerable in scientific terms (Jackson 1980; Cooper 1984; Hedges 1994). The question is critical to the process because it generates the literature search terms and determines relevance criteria

(NHS CRD 2001). Finding the right question is a compromise (probably more so in ecology than in medicine) between taking a holistic approach, thus increasing realism by involving a large number of variables but limiting the number of relevant studies, and a reductionist approach that may limit the review's relevance, utility, and value (Stewart et al. 2005). The question should be practice or policy relevant and should therefore be generated by, or at least in collaboration with, relevant decision-makers (or organizations) for whom the question is real. It may also be important for the question to be seen as neutral to stakeholder groups. Ideally meetings should be held with key stakeholders to try and reach consensus on the nature of the question. This may be more critical for ecological review than medical review because, unlike the benefit of improving human health, the benefit of conserving biodiversity is often contested (Fazey et al. 2004).

Table 1. Elements of a reviewable question; normally a permutation of 'does intervention x on subject y produce outcome z'.

| Question element | Definition |
|---|---|
| *Subject* | unit of study (e.g., ecosystem, habitat, species) that should be defined in terms of the subject(s) to which the intervention will be applied |
| *Intervention* | proposed management regime, policy, or action |
| *Outcome* | all relevant objectives of the proposed management intervention that can be reliably measured. Particular consideration must be given to the most important outcome and to any outcome critical to whether the proposed intervention has greater benefits or disadvantages than any alternatives (i.e., the outcome desired) |
| *Comparator* | Either a control with no intervention or an alternative intervention |

**Example of question formulation**

English Nature, a UK statutory conservation agency, was concerned about the ecological impacts of burning management carried out by landowners in upland areas of England. Discussion with English Nature personnel enabled this general concern to be "unpacked," allowing definition of subject, intervention, and outcome elements of two specific review questions (Stewart *et al.* 2005): "Does burning of U.K. submontane, dry dwarf-shrub heath maintain vegetation diversity?" and "Does burning degrade blanket bog?" Identification of these two related questions allowed specific hypotheses to be tested whilst retaining broader policy relevance. These also provided examples of habitat-based reviews.

Although discussions with review proposers have proven effective in the formulation of a review question, other stakeholders may disagree. In the above example, a key stakeholder disagreed with the outcome measure (a measure of favourable ecological condition based on the relative abundance of key species) used in the "blanket bog" review. To avoid post-review problems such as this we advocate involvement of multiple stakeholders early in the review process.

**Review scoping**

Conducting an initial scoping search enables an assessment of the likely form the review will take. Specifically it may enable:

1. An assessment of whether the review will identify a knowledge gap or if it has the potential to develop into a meta-analyzable review. This has implications in terms of the time and resources required to complete the review and, in some cases, deciding whether it is worthwhile proceeding with a particular review topic.
2. Identification of relevant reviews on the same or similar topics. These will assist in determining the amount of relevant information available on the review topic in relation to the review question.
3. The development of a search strategy using a set of search terms that are both thorough and specific.
4. Identification of potential sources of data, whether individuals or organizations, and also potential lead reviewers if commissioning a review.
5. Identification of stakeholders and experts who should be contacted during the question and protocol development stages.

Scoping searches should be performed using a basic set of search terms in a small number of larger electronic databases (e.g., Web of Science or Scopus) or in topic-relevant web search engines. As such, a number of search terms can be developed and tested quickly, and the number of hits returned can be readily assessed to progress a topic to the full review stage. All scoping searches should be saved so that they may be accessed during the search phase of the review (see Stage 2 – Searching for data), removing duplication of effort where possible. However, if the scoping search is conducted well in advance of the actual review search, it would be prudent to conduct the search again in order to ensure all recent literature has been identified.

**Developing a review protocol**

The review protocol acts as a document that all stakeholders agree upon, after which the review itself can be conducted. The protocol makes clear what the review relates to, and is useful for getting the engagement of experts who may have data to contribute. Anyone reading the protocol should clearly see what the question is and what data are required. The draft review protocol should be made openly available for comment (e.g. on the CEBC website for a one month period), enabling others who have not been contacted during the development stage to provide comments on the direction of the review. The CEBC protocols are made available on our website to show which reviews are in progress, enabling others to see if a review is being conducted that they may be interested in, or to prevent starting a review on a topic that is already underway (see www.cebc.bham.ac.uk for examples). A review protocol can also be used to determine the amount of resources required to conduct a review, whether people, time or money required, and to allocate activities to different members of a review team.

A review protocol is developed as a document that guides the review. As in any scientific endeavour, methodology should be established and made available for scrutiny and comment at an early stage. Because reviews are retrospective by nature, the protocol is essential to make the review process as rigorous, transparent, and well defined as possible (Light 1984). Beside a formal presentation of the question and its

background (the "real world" context), a review protocol sets out the strategy for obtaining primary data and defines relevance criteria for data inclusion or exclusion (NHS CRD 2001). The subject, intervention and outcome elements defined in the question-setting stage provide *a priori* inclusion criteria. If the relevant population, intervention, or outcome measures are present, then data are included although data quality thresholds may result in the subsequent exclusion of otherwise relevant material either from quantitative analysis or from the review in entirety (see Stage 2 – Assessing quality of methodology). The protocol should also establish the methods to be used for data extraction and synthesis, and state any conflicts of interest in the review plus sources of funding. For guidance on developing a review protocol go to http://www.cebc.bham.ac.uk/gettinginvolved.htm.

By planning the review in advance, the protocol helps minimise bias within the review. It may become necessary during the course of a review to make changes to the protocol. These changes should be clearly documented within the final review so that transparency and repeatability is maintained.

## Developing a search strategy

The search strategy is constructed from search terms extracted from the subject, intervention, and outcome elements of the question. It is important that the search is sufficiently rigorous and broad so that all studies eligible for inclusion are identified. This may include considering synonyms, alternative spellings, and non-English language terms with the search strategy. Search protocols must balance sensitivity (getting all information of relevance) and specificity (the proportion of "hits" that are relevant) (NHS CRD 2001). A comprehensive search improves the credibility of the review because evidence of a systematic approach is a key factor in judging the validity of review conclusions. In ecology, resource-intensive searches of high sensitivity are required, even though this is at the expense of specificity, because ecology lacks the mesh-heading indexes and integrated databases of medicine and public health. A high-sensitivity and low-specificity approach is necessary to reduce bias and increase repeatability (see below). Typically, large numbers of references are therefore rejected. For example, of 317 articles with relevant titles concerning the impact of burning on blanket bog, only 8 (2.5%) had comparators (Stewart *et al.* 2005). Similarly, reviews regarding burning of dry heath and the impact of windfarms on bird abundance resulted in meta-analysis of 1.7% and 12%, respectively, of material with relevant titles.

### Example of a search strategy

A review of the effectiveness of control methodologies on introduced populations of the American Mink (*Mustela vison*) in Europe (Tyler et al. 2005) searched 14 electronic databases (Agricola, BIOSIS Previews, CAB Abstracts, Copac, Digital Dissertations, Index to Theses Online, ISI Current Contents, ISI Proceedings, ISI Web of Science, JSTOR, ScienceDirect, Scirus, Scopus, English Nature's Wildlink catalogue); the World Wide Web (first 100 "hits" from www.alltheweb.com, www.google.co.uk, U.K. Department for the Environment, Food and Rural Affairs, Scottish Natural Heritage, Oxford University's Wildlife Conservation Research Unit, The Royal Society for the Protection of Birds, The National Trust, British Wildlife, The Mammal Society, Mammals Trust, and The British Trust for Ornithology); and bibliographies of relevant articles. The search terms used were: *Mustela* AND *vison*,

*Mustela* AND *vison* AND trap\*, *Mustela* AND *vison* AND control\*, *Mustela* AND *vison* AND management, *Mustela* AND *vison* AND pest, Mink AND trap\*, Mink AND control\*, Mink AND management, Mink AND pest). The specificity of this search was low with many references identified multiple times. Specificity could have been increased by using the species name as a search term rather than separating it, i.e. *"Mustela vison" and "M. vison"*. The grey literature search was largely U.K.-based due to resource limitations, although the inclusion of non-U.K. theses was possible. The low specificity of the review (only 1% of retrieved material was judged relevant), however, limits the potential for bias notwithstanding the geographical scope of the grey-literature search. The documented search is fully repeatable and transparent; thus, readers can judge its validity.


## Stage 2 - Conducting the review

### Searching for data

It is perhaps self-evident that the widest possible range of sources should be accessed to capture information. The primary method for information retrieval is the systematic literature search, but this is supplemented by the checking of bibliographies, the provision of supplementary data from authors and through contact with subject experts. Different questions may require the use of different resources, and searches may produce different types of results depending on the information available. Many databases do not have full Boolean search capacity and are sensitive to word order and number. It is therefore necessary to ascertain the capabilities of each database prior to use, and to search using individual search terms and different word orders to extract all the relevant information from the database. Obviously resource availability will constrain the numbers of search term permutations, which will also be subject to diminishing return. "Cut-off" points where the effort is judged to be too high for the return should be identified. It is important to record the methods used in all parts of the search so that others can judge the probability that important research has been missed.

The literature search should be comprised of three distinct phases:
1. searching online databases and catalogues
2. searching organizations and professional networks
3. searching the web.


### Searching online databases and catalogues

There are a number of general scientific electronic databases that may be useful for identifying relevant articles and data sets, such as Web of Science and Scopus. Access to most of these depend on Library subscriptions, and so varies between institutions and organisations. Contacting the subject librarian to identify and discuss the resources available within your institution is recommended at an early stage of the protocol development. As well as the general scientific databases, there are also some subject-specific databases that may contain relevant information, and it may be necessary to search local databases for questions with a regional focus. Organisations often have access to different resources, and so the list of resources searched for each review will vary, but checking bibliographies and contact with authors should help to ensure all references are retrieved.

To minimize the problem of publication bias (e.g., Leimu & Koricheva 2005), both published and unpublished data must be included, a standard rarely satisfied in traditional reviews. The next two stages of the literature search help to address this issue.

**Searching organisations and professional networks**

Many organisations and professional networks make documents freely available through their web pages, and many more contain lists of projects, datasets and references. Often, items referred to on a website will be provided if an organisation is contacted. Sometimes, a visit may be necessary when a large number of documents are required. Searching these organisations and networks targets the grey literature which would not come up in a conventional database search. The list of organisations to be searched is dependent on both the subject of the systematic review and any regional focus.

Hand searching of specific sources and visits to institutions (e.g. libraries and museums) may be necessary to extract all relevant material. However, given time and resource constraints, this will not always prove feasible.

**Web searching**

Searching the web can potentially be a time consuming task, with relatively little useful data being returned. For this reason, a pre-defined, structured approach to web searching eliminates time wasted beyond a point of diminishing return. Specialised subject gateways as well as general search engines can help to focus the searching process and ensure relevance, whilst ensuring that relevant grey literature is located. These gateways, such as Intute.ac.uk, ScienceResearch.com or AcademicInfo.net, contain links to hand-selected sites of relevance for a particular topic or subject area.

Specific guidance on how much web searching is acceptable is difficult to give. Search engines rank their results in different ways, and position within the results is often not correlated to the quality or even relevance of the documents retrieved. In medicine, papers sometimes cite a "first 50 hits" rule, whereby the first 50 results for each search are viewed, but this appears to be an arbitrary number. In order to provide a consistent and practical way to limit the web-searching phase of the search, we recommend full viewing of each of the first 50 hits and then checking for any relevant hits in the next fifty. The actual number of hits retrieved is related both to the search terms used and the quantity of information available so, for some searches, there may be a case for modifying the recommended search limits (e.g., if there are particularly large or small numbers of relevant hits, or if time constraints are an issue).

It is also worth running searches after setting a restriction on the file type, for example by limiting the search to spreadsheets, as these may contain raw data yet rank low in an unrestricted search. This can often be achieved using the "advanced search" function within a search engine. The number of results returned is likely to be low enough that they can all be checked. Searches of this nature should be recorded as part of the search strategy.

**Selection of relevant data**

Once searching is complete, relevant articles must be efficiently selected without wasting resources examining irrelevant articles in detail. Selecting only relevant articles from a potentially large body of initial literature requires the reviewer to use inclusion and exclusion criteria stated *a priori* in the protocol to impose a number of filters of increasing rigor. First, if a long list of articles or data sources is acquired (1000s rather than 100s) and the list of relevant sources is likely to be much shorter, it may be efficient to eliminate some material on title only (especially if obviously spurious hits arise from ambiguity in the use of words in the literature). The second filter should examine title and abstract to determine relevance. The approach should be conservative so as to retain data if there is reasonable doubt over its relevance. It is good practice at this stage to employ a second reviewer to go through the same process on a random sub-sample of articles from the original list (recommended sample is minimum of 25% or a maximum of 2000 references) and to ensure relevance decisions are comparable by performing a kappa analysis, which adjusts the proportion of records for which there was agreement by the amount of agreement expected by chance alone (Cohen 1960; Edwards 2002). A kappa rating of 'substantial' (0.6 or above) is recommended to pass the assessment. If comparability is not achieved, then the criteria should be further developed and the process repeated.

Remaining articles should be viewed in full to determine whether they contain relevant and usable data. Obtaining the full text of all articles can be very time consuming and a realistic deadline may have to be imposed and a record kept of those articles not obtained. The conservative approach and independent checking of a sub-sample by kappa analysis can be repeated at this stage. Short lists of relevant articles and datasets should be made available for scrutiny by stakeholders and subject experts. All should be invited, within a set deadline, to identify relevant data sources they believe are missing from the list. Reviewers should be aware that investigators often cite selectively studies with positive results (Gotzsche 1987; Ravnskov 1992); thus, checking bibliographies and direct contacts must be used only to augment the search.

**Assessing quality of methodology**

The quality of the studies included into a systematic review is of critical importance to the resulting quality of the review; if the data are of poor quality then the conclusions cannot be considered to be robust. Therefore, in an ideal world, each data set included in a systematic review should be of high methodological quality, thus ensuring that the potential for bias is minimized and that any differences in the outcome measure between experimental groups can be attributed to the intervention under scrutiny. To determine the level of confidence that may be placed in selected data sets, each one must be critically appraised to assess the extent to which its research methodology is likely to prevent systematic errors or bias (Moher 1995).

In the health services a hierarchy of research is recognized that scores the value of the data in terms of the scientific rigor of the methodology used (Stevens & Milne 1997). The hierarchy of methodological design can be viewed as generic and has been transferred from medicine to ecology (Pullin & Knight 2003). Where a number of well-designed, high-quality studies are available, others with inferior methodology may be demoted from subsequent quantitative analysis to qualitative tabulation, or

rejected from the systematic review entirely. Alternatively, the effects of individual studies can be weighted according to their position in the "quality hierarchy." However, there are dangers in the rigid application of this hierarchy to ecology as the importance of various methodological dimensions within studies will vary, depending on the study system to which an intervention is being applied. Hypothetically, a rigorous methodology, such as a randomized controlled trial (RCT), could be viewed as superior, even though it was applied over inadequately short time and small spatial scales, to a time series experiment providing data over longer time and larger spatial scales more appropriate to the question. This problem carries with it the threat of misinterpretation of evidence. Potential pitfalls of this kind need to be considered at this stage and addressed by more pragmatic quality weightings (e.g., experimental duration or study area: see Downing et al. 1999 and Côté et al. 2001 respectively) and judicious use of sensitivity analysis (see below).

Four sources of systematic bias that may threaten the internal validity of a study are routinely considered in healthcare (Khan 2003, Moher 1995, Moher 1996, Feinstein 1985). Three of these have, to date, required consideration in ecological systematic reviews. Selection bias results from the way that comparison (e.g., treatment and control) groups are assembled (Kunz 1998) and is a primary reason for randomization in studies. Performance bias refers to systematic differences in the care provided to subjects in the comparison groups and is dealt with by the experimenter being unaware of which are treatments and which controls (blinding) (Shultz 1995). We postulate that the ecological equivalents of performance bias arise from biased baseline comparisons i.e. unequal balancing of heterogeneity in treatment and control arms and failure to consider the impact of covariables that may confound the effectiveness of the intervention. However, it is not possible to account for the influence of potentially confounding variables that are not known or were not measured. Even for those that have been identified, difficulties can arise in extracting standardised information for analysis. Measurement or detection bias refers to systematic differences incurred when knowledge of the intervention influences the assessment of the results in the comparison groups and is also addressed by blinding (Shultz 1995). Blinding is generally not possible in ecology and the extent of detection bias will therefore vary, depending on the rigour and objectivity of sampling methodology (e.g., percent cover assessed by eye is subject to greater potential detection bias than frequency). The fourth, attrition bias (systematic differences between the comparison groups in the loss of samples) has not been an issue in ecological systematic review to date.

Assessing the quality of methodology is a critical part of the systematic review process. It requires a number of subjective decisions about the relative importance of different sources of bias and data quality elements specific to ecology, particularly the appropriateness of variable temporal and spatial scales. It is therefore vital that the assessment process be standardized and as transparent and repeatable as possible. At least 25 scales and 9 checklists have been used to assess the validity of randomized controlled trials in medicine (Moher 1995; Moher 1996). Juni et al. (1999) evaluated 17 health care trials from a meta-analysis, using these 25 different methodological quality scales. For 12 of the scales, the outcomes of the trials were comparable. However, for 6 scales, high quality trials showed little or no benefit of treatment compared to low quality trials, whilst for the remaining 7 scales the opposite trend was observed. Quality scales can therefore give very different results depending on

the data quality items considered and the relevant importance assigned to each one. Similar criteria have also been used to critically appraise the validity of observational studies (Horwitz 1979; Feinstein 1982; Levine 1994; Bero 1999). These checklists do not consider specific ecological criteria. We therefore suggest that review-specific *a priori* assessment forms and two or more assessors should be used to assess study quality in ecological reviewing. The subjective decisions may be a focus of criticism; thus, we advocate consultation with stakeholders to try and reach consensus before moving on to data extraction.

Finally, at this stage, it may be necessary to reject articles that are seemingly relevant but do not present data in extractable format (e.g., if they do not report standard deviations for control and treatment group(s) or the information required to calculate the statistic). If possible, authors of such articles should be contacted and asked whether they can provide data in suitable format.

**Examples of study quality / methodology assessment**
Stewart et al. (2005) used the hierarchy of methodology to separate randomized controlled trials and site comparisons addressing the question: "Does burning degrade blanket bog?" This reflected a major data-quality schism; therefore, further data-quality assessment was inappropriate given the very small number of studies. This approach enabled a simple, but discriminatory, vote count of studies with results showing positive, neutral, or negative effects.

When reviewing the impact of windfarms on bird populations, the standard hierarchy of evidence was considered inadequate by itself due to variation in other critical data-quality elements. This particularly related to the widespread occurrence of confounding factors resulting from variation between treatment and control at baseline or from changes concurrent with windfarm operation (ecological performance bias). The rigour of observations was also variable as measured in terms of replication and objectivity (ecological detection bias). To test for the impact of these factors, data-quality scores, summing the different aspects of data quality outlined above, were added as a meta-regression covariable. Data-quality score was not significant, suggesting that bifurcation of the data into high- and low-quality evidence was not necessary, possibly because the low-quality studies (low replication, imprecise estimates of abundance, high intratreatment variation coupled with confounded baselines) had a high variance and therefore a low weighting in meta-analysis by inverse variance. Sensitivity analyses were used to explore the impact of including low-quality unreplicated data, but the impact of individual data quality elements other than time was not examined because a large number of environmental and windfarm correlates were of interest and the potential for Type II errors would have been increased. Although this pragmatic approach is easy to apply, there is no measure of a studies' "true" validity (Emerson 1990; Schulz 1995; Jüni 1999). Caution should be exercised in interpreting study validity, especially if different quality elements are combined in a single data-quality sum.

A review of the effectiveness of *Rhododendron* control methods considered study hierarchy and potential for bias providing a subjective summary of data quality (Table 2). In this instance the number of environmental variables with sufficient data for analysis was low and sample sizes were sufficient to examine the impact of some individual study quality variables, such as length of experiment and whether results

were generated in the field or a glass house. There were statistically significant differences in effectiveness of control, with greenhouse trials showing greater control than field-based experimentation or monitoring, raising questions about the ecological relevance of greenhouse work and the likely modifying variables. This approach has the merit of objectivity, although there is choice regarding which variables are included in the analysis and caution must be exercised to avoid Type II errors, data mining and overinterpreting results, especially when sample sizes are small.

Table 2. Data quality assessment of an article included in a systematic review of the effectiveness of methods for the control of *Rhododendron ponticum* (Tyler et al. 2004).

| | |
|---|---|
| Methods | site comparison based on sites treated with different interventions, no control, comparison methods only |
| Population | no stand-age detail, site located on lowland heath |
| Intervention and cointerventions | drilled holes filled with herbicide, compared with stumps painted with herbicide |
| Outcomes | painted stumps 30-40% kill<br>drilled holes 95% kill |
| Study design | site comparison |
| Baseline comparison | no information regarding the sites prior to treatment, thus not possible to validate baseline |
| Intratreatment variation | no information describing intratreatment variation |
| Measurement of intervention and cointerventions | no information regarding the sites provided, thus not possible to comment on other management within the area |
| Replication & parameter of abundance | no replication or measure of abundance other than percent kill |
| Notes | study appears to comment on use of techniques rather than providing the reader with scientific evidence, resulting in a high potential for bias and subsequently low data quality |

**Data extraction**
Data extracted from articles should be recorded on carefully designed spreadsheets and undertaken with synthesis in mind. Narrative synthesis requires the construction of tables that provide details of the study or population characteristics, data quality, and relevant outcomes, all of which are defined *a priori*. A summary of methodology *in lieu* of study quality assessment may be sufficient where reviews simply summarise available evidence. However, objective qualitative synthesis requires more formal study quality assessment. In such instances data regarding methodology should be extracted to inform critical appraisal in a standardized, transparent and repeatable manner.

Quantitative analysis follows the same model but care must be taken to extract information pertinent to subsequent analysis (e.g., should binary or continuous outcomes be extracted?). In contrast to medicine, consideration of the appropriate spatial scale(s) and level of replication are necessary prior to extracting the variance measures required to weight meta-analyses. Great care must be taken to standardize and document the process of data extraction, the details of which should be recorded in tables of included studies to increase the transparency of the process. To some extent data extraction can be guided by *a priori* rules, but the complexity of the operation means a degree of flexibility must be maintained. Sensitivity analyses can be used to investigate the impact of extracting data in different ways when there is doubt about the optimum extraction method.

In many cases, the information required is not presented and cannot be obtained from authors. Missing data can be substituted by various methods, most commonly substitution of average or standardized values (Deeks *et al*. 2005), but also calculating bootstrapped confidence limits (Gurevitch & Hedges 2001). Such substitution should be accompanied by sensitivity analyses to assess its impact.

Imputation of data may provide a more robust alternative to substitution, and is especially useful where variance measures are unreported. It is relatively straightforward to impute standard deviation from standard errors, confidence intervals, *t*-values, or a one-way *F*-ratio based on two groups (Lipsey and Wilson 2001, Deeks *et al*. 2005). Variance may also be derived from summary information provided that the study is down-weighted (this is almost invariably the case when weighting is performed by inverse variance). As with substitution, sensitivity analyses should accompany imputation.

It is difficult to perform formal kappa analysis on the repeatability of data extraction, but some attempt to verify repeatability should be made. A second reviewer should check a random subset (recommended sample of minimum 25%) of the included studies to ensure that the *a priori* rules have been applied or the rationale of deviations explained. This also acts as a check on data hygiene and human error (e.g. misinterpretation of a standard error as a standard deviation).

**Example of data extraction**
Reviewing the impact of burning on the ecological condition of blanket bog required extraction of data showing changes in floristic composition and structure. Two reviewers extracted data after reaching a consensus regarding which subsets were relevant within the full data set of each article. *A priori* rules increased the repeatability of data-set formation. For example, sites within an experiment were pooled to prevent pseudoreplication, avoiding post hoc justifications for deriving more than one data-set from an experiment and combining unreplicated, pseudoreplicated and replicated data. Pooled treatment and control sites were included once to maintain independence and avoid bias, with the exception of data on rotational burning, which was scarce and therefore admitted to the review provided there was a comparator irrespective of further potential for bias. Where there was a choice of times since burning, priority was given to the longest time range to maintain independence and maximize predictive power. Similarly, grazed sites received priority over ungrazed sites when the maintenance of independence demanded a choice because grazing and burning are carried out concurrently over most of the

British uplands (Stewart *et al.* 2005). If sample sizes had been larger and a quantitative generic outcome measure identified, the impact of these decisions could have been explored with sensitivity analyses. Given the nature of the data, qualitative discussion of the issues was more appropriate.

**Data synthesis**
This stage includes both qualitative synthesis and quantitative analysis with statistical methods as appropriate. Qualitative synthesis allows informal evaluation of the effect of the intervention and the manner in which it may be influenced by measured study characteristics and data quality. Data from the data-extraction spreadsheet is tabulated to form a summary of the number of data sets providing a yes, no, or neutral answer to each question (vote counting). Where the internal validity of studies varies greatly, reviewers may wish to give greater weight to some studies than others. In these instances it is vital that the studies have been subject to standardized *a priori* critical appraisal with the value judgments regarding internal validity clearly stated. Ideally these will have been subject to stakeholder scrutiny prior to application.

More formal quantitative analysis can be undertaken to generate overall point estimates of the effect size and to analyze reasons for heterogeneity in the effect of the intervention where appropriate data exist. Meta-analysis is now commonly used in ecology (e.g., Arnqvist & Wooster 1995; Osenberg et al. 1999; Gurevitch & Hedges 2001; Gates 2002), consequently we have not treated it in detail here. Meta-analysis provides summary effect sizes with each data set weighted according to some measure of its importance, with more weight given to large studies with precise effect estimates and less to small studies with imprecise effect estimates. Generally each study is weighted in inverse proportion to the variance of its effect. Pooling of individual effects can be undertaken with fixed-effects or random-effects statistical models. Fixed-effects models estimate the average effect and assume there is a single true underlying effect, whereas random-effects models assume there is a distribution of effects that depend on study characteristics. Random effects models include inter-study variability (assuming a normal distribution); thus, when there is heterogeneity, a random-effects model has wider confidence intervals on its summary effect than a fixed-effect model. In medicine both statistical models are used to assess the robustness of statistical synthesis with an *a priori* decision about which is most germane (NHS CRD 2001; Khan 2003). Results of our initial reviews suggest that random-effects models are most appropriate for the analysis of ecological data because the numerous complex interactions common in ecology are likely to result in heterogeneity between studies.

Relationships between differences in characteristics of individual studies and heterogeneity in results can be investigated as part of the meta-analysis, thus aiding the interpretation of ecological relevance of the findings. Exploration of these differences is facilitated by construction of tables that group studies with similar characteristics and outcomes together. Data sets can be stratified into subgroups based on populations, interventions, outcomes, and methodology. Important factors that could produce variation in effect size should be defined *a priori* (see Stage 1 above) and their relative importance considered prior to data extraction to make the most efficient use of data. Differences in subgroups of studies can then be explored.

If sufficient data exist, meta-analysis can be undertaken on subgroups and the significance of differences assessed (see Box 1.). Such analyses must be interpreted with caution because statistical power may be limited (Type I errors possible) and multiple analyses of numerous subgroups could result in spurious significance (Type II errors possible). Alternatively, a meta-regression approach can be adopted whereby linear regression models are fitted for each covariate, with studies weighted according to the precision of the estimate of treatment effect in a random-effects model (Sharp 1998).

Despite the attempt to achieve objectivity in reviewing scientific data, considerable subjective judgment is required when undertaking meta-analyses. These judgements include decisions about choice of effect measure, how data are combined to form datasets, which data sets are relevant and which are methodologically sound enough to be included, methods of meta-analysis, and the issue of whether and how to investigate sources of heterogeneity (Thompson 1994). Reviewers should state explicitly and distinguish between the *a priori* and *post hoc* rationales behind these decisions to minimize bias and increase transparency.

Quantitative research synthesis is still in its infancy. The biases associated with specific techniques are not generally based on empirical evidence, particularly as applied to ecological research. There is considerable potential to improve the statistical models and to provide robust guidance about which models are most germane under which circumstances. Pending these developments, we advise reviewers to search for the broad patterns contained in accumulated ecological knowledge using *a priori* decisions and a pragmatic design wisdom to build repeatable knowledge structures with as much structural integrity as possible.

**Example of data synthesis**
A review of the impact of wind turbines on bird abundance utilized standardized mean difference meta-analysis with weighting by inverse variance to combine data from 19 globally distributed windfarms. Sensitivity analyses were used to explore the effect of including data from unreplicated studies and to assess bias arising from data extraction of pseudoreplicated or aggregated data. Pooled effect sizes remained negative and statistically significant regardless of how the effect sizes were generated, indicating that the patterns in the data were robust. *A priori* and *post hoc* reasons for heterogeneity were explored with meta-regression. Of the *a priori* variables only bird taxon appeared to modify the result, with relationships between turbine number and power being too weak to have biological significance. *Post hoc* analysis revealed that the impact of windfarms became more pronounced over time, a finding not reported by any of the original research or previously assessed in the literature. This has important implications because declines in local bird abundance are more likely to have deleterious population-level impacts if they worsen over time. It also suggests that current windfarm monitoring programs are of inadequate duration to detect deleterious effects.

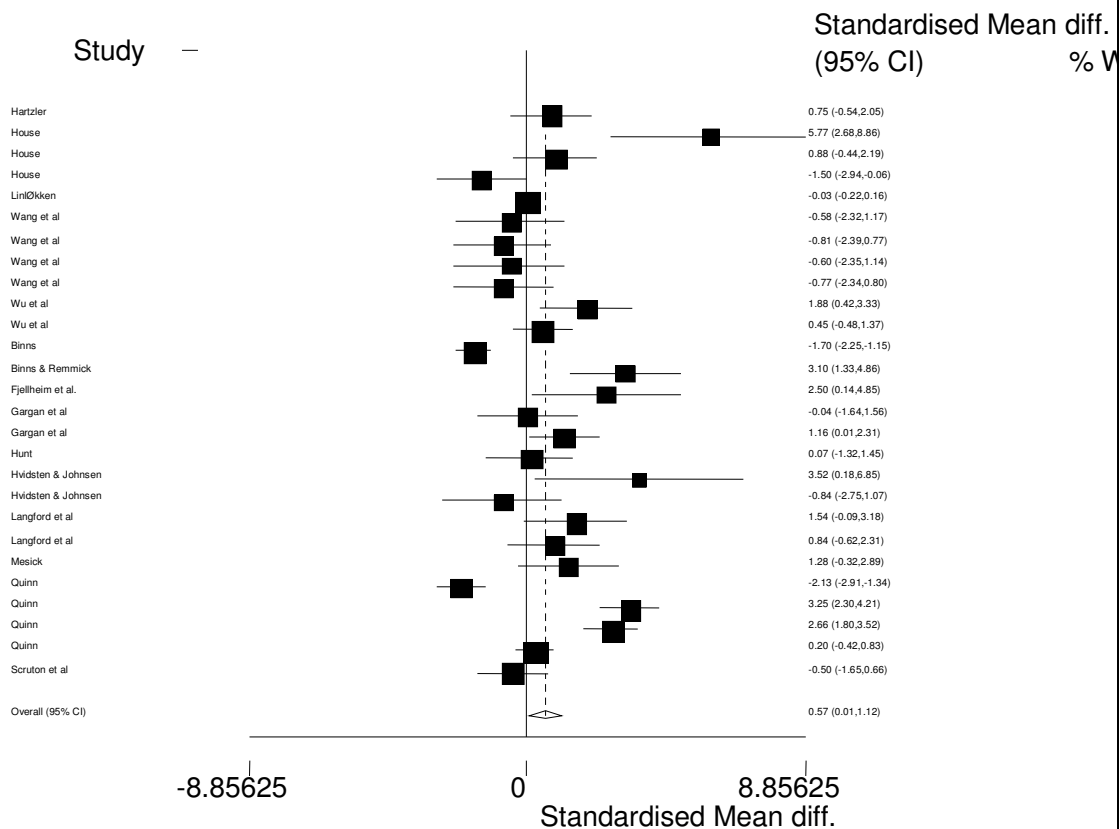## Box 1. Interpretation of Forrest plots-Example using STATA



Figure 1. An example of a Forrest plot generated using STATA and typically included as an outcome in reviews that incorporate a meta-analysis.

The individual data points included in the meta-analysis are listed down the left side of the diagram. In this example multiple independent points have been extracted from the same references. Individual studies are typically identified by author name and year, with multiple points numbered. Full details of each study can be found in the references at the end of the systematic review and the tables of included studies and data extraction appendices should make it clear how multiple points were derived from individual studies.

Each data point extracted from a study is represented by a square. The size of the square represents the sample size of that data point whilst the error bar typically represents the 95% confidence interval. The position of the square on the x axis denotes the effect size (in this example Cohens D). This example also lists the effect size and confidence interval for each study to the right of the diagram, along with the weight which that study contributes to the overall synthesis (in this example weighting is by inverse variance).

Underneath the studies, there is a pooled estimate of effect represented by an open diamond. This is a graphical representation of the combined outcome for all of the included data points. The width of this diamond represents the confidence interval.

The "line of no effect" where the effect size is zero is represented by a solid vertical line, and anything that crosses this line is not statistically significant (including those studies where only the confidence interval crosses the line). Anything that falls to the left of the line of no effect has less of the outcome; whereas anything that falls to the right has more of the outcome- whether this is a positive or negative result depends on what the outcome of the meta-analysis is. Therefore a beneficial result for a negative outcome (such as habitat loss) has a significant effect size to the left of the vertical line and a beneficial result for a positive outcome (such as increase in suitable habitat) has a significant effect size to the right of the vertical line. Overall interpretation of the Forrest plot relies on consideration of the position and significance of individual points as well as the pooled estimate, because the pooled estimate can be misleading when heterogeneity is high (see above).

**The interpretation of meta-analysis and systematic review evidence**
Systematic reviews synthesise and present evidence but the strength of this evidence and the applicability of the results require careful consideration and interpretation. The discussion and conclusions may consider the implications of the evidence in relationship to practical decisions, but the decision-making context may vary, leading to different decisions based on the same evidence. Authors should, where appropriate, explicitly acknowledge the variation in possible interpretation and simply present the evidence rather than offer advice. Recommendations that depend on assumptions about resources and values should be avoided (Khan 2003, Deeks *et al.* 2005).

Deeks et al (2005) offer the following advice of relevance here. Authors and end-users should be wary of the pitfalls surrounding inconclusive evidence and should beware of unwittingly introducing bias in their desire to draw conclusions rather than pointing out the limits of current knowledge. Where reviews are inconclusive because there is insufficient evidence, it is important not to confuse 'no evidence of an effect' with 'evidence of no effect'. The former results in no change to existing guidelines, but has an important bearing on future research, whereas the latter could have considerable ramifications for current practice or policy.

Review authors, and to a lesser extent end-users, may be tempted to reach conclusions that go beyond the evidence that is reviewed or to present only some of the results. Authors must be careful to be balanced when reporting on and interpreting results. For example if a 'positive' but statistically non-significant trend is described as 'promising', then a 'negative' effect of the same magnitude should be described as a 'warning sign'. Other examples of unbalanced reporting include one-sided reporting of sensitivity analyses or explaining non-significant positive results (e.g. the included studies were too small to detect a reduction in mortality for a statistically non-significant increase in mortality) but not negative ones. If the confidence interval for the estimate of difference in the effects of interventions overlaps the null value, the analysis is compatible with both a true beneficial effect and a true harmful effect. If one of the possibilities is mentioned in the conclusion, the other possibility should be mentioned as well and both should be given equal consideration in discussion of results. One-sided attempts to explain results with reference to indirect evidence external to the review should be avoided. Medical guidance suggests that considering results in a blinded manner can avoid these pitfalls (Deeks *et al.* 2005). Authors

should consider how the results would be presented and framed in the conclusions and discussion if the direction of the results was reversed.

**Evidence of effectiveness**

Medical systematic reviews assess the strength of inferences about the effectiveness of an intervention using guidelines that consider the strength of a causal inference (Hill 1971). Areas for consideration include:

1. The quality of the included studies
2. The size and significance of the observed effects
3. The consistency of the effects across studies or sites?
4. The clarity of the relationship between the intensity of the intervention and the outcome?
5. The existence of any indirect evidence that supports or refutes the inference
6. The lack of other plausible competing explanations of the observed effects (bias or confounding)

There are a range of approaches to grading the strength of evidence presented in medical reviews but there is no universal approach (Deeks *et al.* 2005). We suggest that authors of ecological reviews explicitly state weaknesses associated with each of the areas above, but the overall impact they make on conclusions can only be considered subjectively.

**Applicability of results**

End-users must decide, either implicitly or explicitly, how applicable the evidence presented in a systematic review is to their particular circumstances (Deeks *et al.* 2005). Authors should highlight where the evidence is likely to be applicable and equally importantly where it may not be applicable with reference to variation between studies and study characteristics.

Clearly, variation in the ecological context and geographical location of studies can limit the applicability of results. Authors should be aware of the timescale of included studies which may be insufficiently short to make long-term predictions. Variation in application of the intervention may also be important (and difficult to predict); but authors should be aware of differences between *ex situ* and *in situ* treatments (measuring efficacy versus effectiveness respectively) where they are combined, and should also consider the implications of applying the same intervention at different scales. Variation in baseline risk may also be an important consideration in determining the applicability of results as the net benefit of any intervention depends on the risk of adverse outcomes without intervention, as well as on the effectiveness of the intervention (Deeks *et al.* 2005). Given the myriad factors involved in nature-conservation decision making, consideration of baseline risk is probably best left to end-users. However, reviewers should point out any clear discrepancies between high and low baseline risk groups where there is *a priori* rationale for the split.

Where reviewers identify predictable variation in the relative effect of the intervention in relation to the specified reasons for heterogeneity these should be highlighted. However, these relationships require cautious interpretation (because they are only

correlations) particularly where sample sizes are small, data points are not fully independent, and where multiple confounding occurs.

## Stage 3 - Reporting and dissemination of results

Wide dissemination and open access are key requirements of the evidence-based framework. For systematic reviews to have a real impact in terms of knowledge transfer from the science to the practitioner and policy communities they need to be readily accessible from a recognized central source. To this end the CEBC has established a library of systematic reviews with the intent of managing and servicing the library on a non-profit basis in a similar format to the Cochrane Collaboration Library in medicine (see www.cochrane.org) with its emphasis on transparency of the review process and independence from bias (Fazey et al. 2004).We urge reviewers to submit their reviews to the library and contact us at the earliest possible stage of the review process. The following sets out the principles and conditions of inclusion of reviews in the library.

Before reports are disseminated they should be subjected to expert scrutiny or peer review, including assessment of scientific quality and completeness. This process is organized by the CEBC and is equivalent to that of a journal or grant board, but with a more supportive role in helping reviewers achieve the necessary quality rather than accepting/rejecting outright. If the CEBC is contacted at an early stage and open consultation is undertaken as set out above then the chances of meeting the required standard should be significantly improved.

The full review should be submitted to CEBC in the standard format (see Review Presentation and Formatting Guidelines at
 www.cebc.bham.ac.uk/gettinginvolved.htm). The format for reporting on the CEBC website is a short summary that highlights the main review outcomes. This should be written so as to enable effective communication with managers and policy formers. A full review will normally include too much detail for wider dissemination but will nevertheless be made available, along with the summary, to all who want more information on the conduct of the review process. By mutual agreement, other formats such as policy briefs and guidance notes may also be posted.

Submitting and posting a review on the CEBC website DOES NOT prevent further publication and the review may also be submitted, at the author's discretion, for publication in a peer-reviewed journal.

## Requirement for further work

Systematic review in conservation and environmental management is in its infancy and these guidelines will need updating on a regular basis as well develop methodology and learn from undertaking more reviews on a wider range of subjects. For example, all reviews to date have incorporated comparators, although work in progress involves synthesising experience and evidence employing Bayesian methodologies (Morris 1992; Louis 1993). It could be argued that this is an excessively reductionist approach, applying a narrow definition of evidence (Fox

2005) and that further methodological development might be necessary to integrate different types of evidence (Dixon-Woods et al. 2004) or to assess ecological information of types beyond the experience of the authors.

Other issues require consideration to strengthen the ecological guidelines presented above. Medical systematic review methodology is developing rapidly, with new techniques being developed to handle diverse types of variable quality data in fields such as diagnostic testing. The utility of these techniques for ecological purposes requires further investigation. Likewise, techniques for economic cost-benefit evaluation and disseminating evidence to different audiences (policy, scientific, practitioner and stakeholder groups) (NHMRC 2000) warrant consideration. Addressing all these issues is beyond the scope of these guidelines, but require further development if an ecological evidence base is to be fully established. The ecological guidelines presented evolved from the existing medical model. Table 3 highlights key differences between ecological and medical guidelines at present. As was the experience in the medical field, it will take time for systematic reviews to be recognized and valued as equivalent to other scientific papers in conservation. We hope these guidelines will set standards and facilitate key steps forward in encouraging more systematic reviews (e.g. journals encouraging their submission and publication and funders seeing systematic reviews as a valid form of research). We call on the conservation and environmental management community to engage with the CEBC to further develop the library of systematic reviews and create the accessible evidence base that conservation and environmental management urgently requires.

Table 3. Differences between the medical systematic review guidelines and the ecological review guidelines advocated by the authors

| Review stage | Medical guidelines | Ecological guidelines |
|---|---|---|
| Question formulation | Question formulation generally not limited by complexity and study numbers | Question formulation usually limited by information availability and complexity requiring a balance between holism (more realistic) and reductionism (more studies) |
| | Stakeholder engagement useful but not generally critical | Stakeholder engagement may be critical because conservation actions often result in conflicts in objectives |
| Developing review protocol: Search strategy | Complex searches balancing sensitivity and specificity are possible and recommended | High sensitivity, low specificity searches are recommended to reduce bias and increase repeatability because ecology lacks the sophisticated search infra-structure of medicine |
| Assessing quality of methodology | Clear hierarchy of evidence generally applicable and often used to define a minimum quality threshold | Pragmatic quality weightings and sensitivity analyses must augment data quality hierarchies to avoid misinterpretation, particularly when combining data across the hierarchy to increase sample sizes |
| | Performance bias and detection bias addressed by blinding and easy to assess using published quality weightings. Attrition bias common | Performance bias and detection bias addressed by experimental design and hard to assess especially in a standardized manner, necessitating the use of review specific quality weightings. Attrition bias rare |
| | Numerous "off the shelf" checklists to assess the validity of medical research | No "off the shelf" checklists hence the need for *a priori* review specific criteria preferably validated by consensus with stakeholders |
| Data extraction | Data extraction often relatively straightforward except for missing data and data hygiene problems | Data extraction complex especially with respect to variance measures for weighting. *A priori* rules must be developed in order to extract data in a repeatable standardized manner with independence and (pseudo)replication commonly problematic. |
| Data synthesis: Meta-analysis | Fixed and random effects models applicable | Random effects models generally more useful than fixed effect models because the complex interactions in ecology generally result in ecologically important heterogeneity between studies. |

**Literature Cited**

Arnqvist, G., and D. Wooster 1995. Metaanalysis – synthesizing research findings in ecology and evolution. Trends in Ecology & Evolution **10**: 236–240.

Bero L., R. Grilli, J. Grimshaw, G. Mowatt, A. Oxman, and M. Zwarenstein eds. 1999. Effective Practice and Organisation of Care Module of The Cochrane Database of Systematic Reviews. Issue 2. Update software. The Cochrane Library, Oxford, United Kingdom.

Cohen, J. 1960. A coefficient of agreement for nominal scales. Educational and Psychological Measurement **20:** 37-46.

Cooper, H.M. 1984. Integrating Research. A Guide for Literature Reviews. Sage Publications, Newbury Park.

Côté, I.M., Mosqueira, I. & Reynolds, J.D. (2001). Effects of marive reserve characteristics on the protection of fish populations: a meta-analysis. *Journal of Fish Biology* **59**: SA178-189

Deeks, J.J., J.P.T. Higgins and D.G. Altman 2005. Analysing and presenting results. In: Cochrane Handbook for Systematic Reviews of Interventions 4.2.5 [updated May 2005]; Section 8. (ed by J.P.T. Higgins and S. Green.) http://www.cochrane.org/resources/handbook/hbook.htm (accessed 12th July 2006).

Dixon-Woods, M., S. Agarwal, B. Young. D. Jones, and A. Sutton 2004. Integrative approaches to qualitative and quantitative evidence. National Health Services Health Development Agency, London.

Downing, J.A., Osenberg, C.W. & Sarnelle, O. (1999). Meta-analysis of marine nutrient-enrichment experiements: variation in the magnitude of nutrient limitation. *Ecology* **80**: 1157-1167.

Edwards, P., M. Clarke, C. DiGuiseppi, S. Pratap, I. Roberts, and R. Wentz. 2002. Identification of randomized controlled trials in systematic reviews: accuracy and reliability of screening records. Statistics in Medicine **21**: 1635-1640.

Emerson, J.D., E. Burdick, D.C. Hoaglin, F. Mosteller, and T.C. Chalmers. 1990. An empirical study of the possible relation of treatment differences to quality scores in controlled randomized clinical trials. Controlled Clinical Trials **11**: 339-352.

Fazey, I., J.G. Salisbury, D.B. Lindenmayer, J. Maindonald, and R. Douglas. 2004. Can methods applied in medicine be used to summarize and disseminate conservation research? Environmental Conservation **31**: 190-198.

Feinstein, A.R. 1985. Clinical Epidemiology: The Architecture of Clinical Research. Saunders, Philadelphia.

Feinstein, A.R., and R.I. Horwitz. 1982. Double standards, scientific methods, and epidemiological research. New England Journal of Medicine **307**: 1611-1617.

Fox, D.M. 2005. Evidence of evidence-based health policy: The politics of systematic reviews in coverage decisions. Health Affairs **24**: 114-122.

Gates, S. 2002. Review of methodology of quantitative reviews using meta-analysis in

ecology. *Journal of Animal Ecology* **71**: 547–557.

Gotzsche, P. C. 1987. Reference bias in reports of drug trials. British Medical Journal **295**: 654-656.

Gurevitch, J. and L.V. Hedges 2001. Meta-analysis Combining the results of independent experiments. In: Design and Analysis of Ecological Experiments (ed by S.M. Scheiner and J. Gurevitch) pp. 347-369. Oxford University Press, New York.

Hedges, L.V. 1994. Statistical considerations. In: The Handbook of Research Synthesis. (ed by H. Cooper and L.V. Hedges) pp. 30-33. Russell Sage Foundation, New York.

Higgins, J.P.T. and S. Green. (eds) 2005. Cochrane Handbook for Systematic Reviews of Interventions 4.2.5. John Wiley & Sons, Ltd, Chichester, UK.

Hill, A.B. 1971 Principles of Medical Statistics. Lancet **9**: 312-20.

Horwitz, R.I., and A.R. Feinstein. 1979. Methodological standards and contradictory results in case-control research. American Journal of Medicine **66**: 556-564.

Jackson, G.B. 1980. Methods for integrative reviews. Review Education Research **50**: 438-460.

Jüni, P., A. Witschi, R. Bloch, and M. Egger. 1999. The hazards of scoring the quality of clinical trials for meta-analysis. Journal of American Medical Association **282**: 1054-60.

Khan, K.S., R. Kunz, J. Kleijnen and G. Antes. 2003. Systematic reviews to support evidence-based medicine: how to apply findings of healthcare research. Royal Society of Medicine Press Ltd, London.

Kunz, R., and A.D. Oxman. 1998. The unpredictability paradox: review of empirical comparisons of randomised and non-randomised trials. British Medical Journal **317**: 1185-1190.

Leimu, R., and J. Koricheva 2005. What determines the citation frequency of ecological papers? Trends in Ecology and Evolution **20:** 28-32.

Levine, M., S. Walter, H. Lee, T. Haines, A. Holbrook, and V. Moyer. 1994. The Evidence-Based Medicine Working Group. Users' guides to the medical literature IV: how to use an article about harm. Journal of American Medical Association **271**:1615-1619.

Light R.J., and D.B. Pillemer. 1984. Summing Up: The Science of Reviewing Research. Harvard University Press, Massachusetts.

Lipsey, M.W. and D.B. Wilson 2001. Practical Meta-analysis. Applied Social Research Methods Series. Volume 49. Sage Publications, Thousand Oaks, California.

Louis. T., and D. Zelterman. 1993. Bayesian approaches to research synthesis. In: The Handbook of Research Synthesis. (ed. by H. Cooper and L.V. Hedges), Pp 411-422. Russell Sage Foundation, New York.

Moher, D., A.R. Jadad, G. Nichol, M. Penman, P. Tugwell, and S. Walsh. 1995. Assessing the quality of randomized controlled trials: an annotated bibliography of scales and checklists. Controlled Clinical Trials **16**: 62-73.

Moher, D., A.R. Jadad, and P. Tugwell. 1996. Assessing the quality of randomized controlled trials: current issues and future directions. International Journal of Technology Assessment in Health Care **12**: 195-208.

Morris, C.N., and S.L. Normand. 1992. Hierarchical models for combining information and for meta-analyses. In: Bayesian Statistics (ed. By J.M. Bernardo, J.O. Berger, A.P. Dawid, and A.F.M. Smith), pp. 321-344. 4th Ed. Oxford University Press, New York.

NHS Centre for Reviews and Dissemination. 2001. Undertaking systematic review of research on effectiveness. NHS CRD, University of York.

Osenberg, C.W., O. Sarnelle, S.D. Cooper and R.D. Holt. 1999. Resolving ecological questions through meta-analysis: goals, metrics and models. Ecology **80**: 1105–1117.

Pullin, A., and T. Knight. 2001. Effectiveness in conservation practice: pointers from medicine and public health. Conservation Biology **15**: 50-54.

Pullin, A., and T. Knight. 2003. Support for decision-making in conservation practice: an evidence-based approach. Journal for Nature Conservation **11**: 83-90.

Pullin, A., T. Knight, D. Stone, and K. Charman. 2004. Do conservation managers use scientific evidence to support their decision-making? Biological Conservation **119**: 245-252

Ravnskov, U. 1992. Cholesterol lowering trials in coronary heart disease: frequency of citation and outcome. British Medical Journal **305**: 9-15**.**

Roberts, P.D., Stewart, G.B. & Pullin, A.S. (in press) Are review articles a reliable source of evidence to support conservation and environmental management? A comparison with medicine. *Biological Conservation.*

Sharp, S. 1998. Meta-analysis regression: statistics, biostatistics, and epidemiology. Stata Technical Bulletin **42**: 16-22.

Schulz, K.F., I. Chalmers, R.J. Hayes, and D.G. Altman. 1995. Empirical evidence of bias: dimensions of methodological quality associated with estimates of treatment effects in controlled trials. Journal of the American Medical Association **273**: 408-412.

Stevens, A., and R. Milne. 1997. The effectiveness revolution and public health. In: Progress in Public Health (ed. By G. Scally), Pp. 197-225. Royal Society of Medicine Press, London.

Stewart, G. B., C. F. Coles, and A. S. Pullin. 2005. Applying evidence-based practice in conservation management: lessons from the first systematic review and dissemination projects. Biological Conservation **126**: 270-278.

Stewart, G.B., A.S. Pullin, and C.F. Coles. 2005. Effects of wind turbines on bird abundance. Systematic Review 4. Centre for Evidence-Based Conservation, Birmingham, UK. Available from www.cebc.bham.ac.uk/completedreviews.htm

Sutherland, W., A. Pullin, P. Dolman, and T. Knight. 2004. The need for evidence-based conservation. Trends in Ecology and Evolution **19**: 305-308.

Thompson, S. 1994. Systematic review: why sources of heterogeneity in meta-analysis should be investigated. British Medical Journal **309:** 1351-1355.

Tyler, C., E. Clark, and A. S. Pullin. 2005. Do management interventions effectively reduce or eradicate populations of the American Mink, Mustela vison? Systematic Review 7. Centre for Evidence-Based Conservation, Birmingham, UK. Available from www.cebc.bham.ac.uk/completedreviews.htm.

Tyler, C., and A.S. Pullin. 2004. Do commonly used interventions effectively control *Rhododendron ponticum*? Systematic Review 6. Centre for Evidence-Based Conservation, Birmingham, UK. Available from www.cebc.bham.ac.uk/completedreviews.htm

Tyler, C., A.S. Pullin, and G. B. Stewart. 2006. Effectiveness of management interventions to control invasion by *Rhododendron ponticum.* Environmental Management 37, 513- 522.