

1 **Title:** Systematic Review and Evidence Integration for Literature-Based Environmental Health  
2 Science Assessments

3 **Authors:** Andrew A. Rooney, Abee L. Boyles, Mary S. Wolfe, John R. Bucher, and Kristina A.  
4 Thayer

5 **Affiliation:** Office of Health Assessment and Translation, Division of the National Toxicology  
6 Program, National Institute of Environmental Health Sciences (NIEHS), National Institutes of  
7 Health (NIH), Research Triangle Park, North Carolina, USA

8 **Corresponding Author:** Andrew A. Rooney, *address:* NIEHS, P.O. Box 12233, Mail Drop  
9 K2-04, RTP, NC 27709; *express mail address:* 530 Davis Drive, Morrisville, NC, 27560; *phone:*  
10 (919) 541-2999; *email:* [Andrew.Rooney@nih.gov](mailto:Andrew.Rooney@nih.gov)

11 **Running Title:** Framework for Systematic Review

12 **Key Words:** hazard identification, systematic review, human health risk assessment, risk of bias,  
13 weight of evidence

14 **Acknowledgements:** We appreciate the valuable advice and comments on the development of  
15 this systematic review framework from a number of technical experts, the public, the NTP  
16 Executive Committee, and the NTP Board of Scientific Counselors.

17 **Competing Financial Interests:** The authors have no competing financial interests.

18

19 **Abstract**

20 **BACKGROUND:** Systematic review methodologies provide objectivity and transparency to the  
21 process of collecting and synthesizing scientific evidence for reaching conclusions on specific  
22 research questions. There is increasing interest in applying these procedures to address  
23 environmental health questions.

24 **OBJECTIVES:** To develop a systematic review framework to address environmental health  
25 questions by extending approaches developed for clinical medicine to handle the breadth of data  
26 relevant to environmental health sciences (e.g., human, animal, and mechanistic studies).

27 **METHODS:** The Office of Health Assessment and Translation (OHAT) adapted guidance from  
28 systematic-review authorities and sought advice during development of the OHAT Approach  
29 through consultation with technical experts in systematic review and human health assessments  
30 as well as scientific advisory groups and the public. The method was refined by considering  
31 expert and public comments and through application to case studies.

32 **RESULTS AND DISCUSSION:** Presented here is a 7-step framework for systematic review and  
33 evidence integration for reaching hazard identification conclusions: problem formulation and  
34 protocol development, search for and select studies for inclusion, extract data from studies,  
35 assess the quality or risk of bias of individual studies, rate the confidence in the body of  
36 evidence, translate the confidence ratings into levels of evidence, and integrate the information  
37 from different evidence streams (human, animal, and “other relevant data” including mechanistic  
38 or *in vitro* studies) to develop hazard identification conclusions.

39 **CONCLUSION:** The principles of systematic review can be successfully applied to environmental  
40 health questions to provide greater objectivity and transparency to the process of developing  
41 conclusions.

## 42 **Introduction**

43 Systematic review methodologies increase the objectivity and transparency in the process of  
44 collecting and synthesizing scientific evidence on specific questions. The product of a systematic  
45 review can then be used to inform decisions, reach conclusions, or identify research needs. There  
46 is increasing interest in applying the principles of systematic review to questions in  
47 environmental health (EFSA 2010; NRC 2011, 2013a; Rhomberg et al. 2013; Woodruff and  
48 Sutton 2011).

49 While systematic review methodologies are well established in clinical medicine to  
50 assess data for reaching health care recommendations (AHRQ 2013; Guyatt et al. 2011a; Higgins  
51 and Green 2011; Viswanathan et al. 2012), these approaches are most developed for human  
52 clinical trials, and therefore, typically consider small datasets of similar study design in  
53 developing conclusions. Questions in environmental health require the evaluation of a broader  
54 range of relevant data including experimental animal and mechanistic studies as well as  
55 observational human studies. Also, there is a need to integrate data from multiple evidence  
56 streams (human, animal, and “other relevant data” including mechanistic or *in vitro* studies) in  
57 order to reach conclusions regarding potential health effects from exposure to substances in our  
58 environment.

59 The National Toxicology Program (NTP) Office of Health Assessment and Translation  
60 (OHAT) conducts literature-based evaluations to assess the evidence that environmental

61 chemicals, physical substances, or mixtures (collectively referred to as "substances") cause  
62 adverse health effects and provides opinions on whether these substances may be of concern  
63 given levels of current human exposure (Bucher et al. 2011). Building on a history of rigorous  
64 and objective scientific review, OHAT has been working to incorporate systematic-review  
65 procedures in its evaluations since 2011 through a process that has included adoption of current  
66 practice as well as methods development (Birnbaum et al. 2013; NTP 2012a, b, 2013c). This  
67 article explains the framework developed by OHAT with procedures to integrate multiple  
68 evidence streams including observational human study findings, experimental animal toxicology  
69 results, and other relevant data in developing hazard identification conclusions or state-of-the-  
70 science evaluations regarding health effects from exposure to environmental substances. The 7-  
71 step framework outlines methods to increase transparency and consistency in the process, but it  
72 also presents opportunities to increase efficiencies in data management and data display that  
73 facilitate the process of reaching and communicating hazard identification conclusions.

## 74 **Methods**

75 In 2011, OHAT began exploring systematic-review methodology as a means to enhance  
76 transparency and increase efficiency in summarizing and synthesizing findings from studies in its  
77 literature-based health assessments. OHAT used a multi-pronged strategy to develop the OHAT  
78 Approach working with advisors to adapt and extend existing methods from clinical medicine  
79 and obtaining input from technical experts and the public on early drafts (Supplemental Material,  
80 Table S1). The methods development process is described in detail in Supplemental Material. In  
81 brief, OHAT reviewed guidance from authoritative systematic-review groups (AHRQ 2013;  
82 Guyatt et al. 2011a; Higgins and Green 2011) in developing an initial draft and sought additional  
83 advice through web-based discussions and consultation with technical experts, the NTP

84 Executive Committee, the NTP Board of Scientific Counselors, and the public (NTP 2012a, b,  
85 2013a, b, c, g). The resulting OHAT Approach has been refined based on the input received and  
86 through application to case studies.

## 87 **Results**

88 The OHAT framework is a 7-step process (**Figure 1**). It includes all of the recommended  
89 elements for conducting and reporting a systematic review (outlined in the PRISMA statement  
90 on Preferred Reporting Items for Systematic Reviews and Meta-Analyses)(Moher et al. 2009).  
91 The specific procedures for performance of each step are described in a detailed protocol that is  
92 developed for each evaluation (NTP 2013e, f).

### 93 ***Step 1: Problem Formulation and Protocol Development***

94 Prior to conducting an evaluation, the scope and focus of the topic is defined through  
95 consultation with subject-matter experts. For OHAT, objective(s) are typically to identify a  
96 potential health hazard or to assess the state of the science in order to identify research needs on  
97 topics of importance to environmental health. The objectives of the evaluation must be clearly  
98 stated including the key questions to be addressed. The evaluation is structured to answer these  
99 key questions that guide the systematic-review process for the literature search, study selection,  
100 data extraction, and synthesis. The questions define the Populations, Exposures, Comparators,  
101 Outcomes, Timings, and Settings of interest (PECOTS) eligibility criteria for the evaluation  
102 (e.g., see discussion in AHRQ 2013). PECOTS is the environmental equivalent of AHRQ's  
103 PICOTS expansion of the original PICO approach developed for clinical evaluations that focuses  
104 on Interventions rather than Exposures, and did not initially include Timing or Setting in the  
105 inclusion criteria (Whitlock et al. 2010).

106 A concept document, or brief proposal, and a specific, detailed protocol for OHAT  
107 evaluations are developed through an iterative process in which information is obtained by  
108 outreach to federal partners, technical experts, the public, and through consultation with the NTP  
109 Board of Scientific Counselors (NTP 2013d). Through this process, the protocol is developed *a*  
110 *priori* and guidance in the protocol forms the basis for scientific judgments throughout the  
111 evaluation. However, it is important to acknowledge that the protocol can be modified to address  
112 unanticipated issues that might arise while conducting the review (e.g., see FDA 2010; Khan et  
113 al. 2001). Revisions to the protocol are documented and justified with notation of when in the  
114 process the revisions were made.

115 ***Step 2: Search for and Select Studies for Inclusion***

116 **Search for Studies:** A comprehensive search of the primary scientific literature is performed.  
117 The search covers multiple databases (including, but not limited to, PubMed, TOXNET, Scopus,  
118 Embase, etc.) with sufficient details of the search strategy documented in the protocol such that it  
119 could be reproduced. The protocol also lists the dates of the search, frequency of updates, and  
120 any limits placed on the search (e.g., language or date of publication). The protocol establishes  
121 requirements for consideration of data from meeting abstracts or other unpublished sources. If a  
122 study that may be critical to the evaluation has not been peer reviewed, and the authors agree to  
123 make all study materials available, the NTP will have it peer reviewed by independent scientists  
124 with relevant expertise. The peer review requirement assures that studies considered in the  
125 evaluation have been reviewed by subject matter experts and the information from this review  
126 would be available in Step 4 when evaluating individual study quality.

127 **Select Studies for Inclusion:** All references identified in the search are screened for relevance to  
128 the key question(s) of the evaluation based on the PECOTS eligibility criteria established when

129 formulating the problem in Step 1. The protocol establishes criteria for including or excluding  
130 references based on, for example, applicable outcomes, relevant exposures, and types of studies.  
131 These criteria contain sufficient detail to develop an inclusion and exclusion checklist such that  
132 use of scientific judgment during the literature-selection process is limited. If major limitations in  
133 a specific study type or design for addressing the question are known in advance (e.g., unreliable  
134 methods to assess exposure or health outcome), the basis for excluding those studies must be  
135 described *a priori* in the protocol.

136         The protocol also outlines the specific plans for reviewing studies for inclusion, resolving  
137 conflicts between reviewers, and documenting the reasons that studies were excluded. Two  
138 reviewers independently screen all references at the title and abstract level and resolve  
139 differences by reaching agreement through discussion. References that meet the inclusion criteria  
140 are retrieved for full text review, as are those with insufficient information to determine  
141 eligibility from just the title and abstract. Procedures for full text review are tailored to the scope  
142 of the review and follow procedures established in the protocol. Reporting the number of  
143 references retrieved, duplicates removed, and studies excluded as references move through the  
144 screening process by creating a flow diagram is one of several required elements for reporting  
145 based on the PRISMA statement (Liberati et al. 2009; Moher et al. 2009) that we have include in  
146 this framework.

### 147 ***Step 3: Extract Data from Studies***

148 Relevant data from individual studies selected for inclusion are extracted or copied from the  
149 publication to a database to facilitate critical evaluation of the results including data summary  
150 and display using separate data collection forms for human, animal, and *in vitro* studies. For each  
151 study, one member of the evaluation team performs the data extraction and quality assurance

152 procedures are undertaken as specified in the protocol (e.g., review and confirmation by another  
153 team member). Following completion of an evaluation, the data extracted and summarized will  
154 be made publicly available in the NTP Chemical Effects in Biological Systems (CEBS) database  
155 (<http://www.niehs.nih.gov/research/resources/databases/cebs/index.cfm>).

156 ***Step 4: Assess the Quality or Risk of Bias of Individual Studies***

157 Despite the critical importance of assessing the credibility of individual studies when developing  
158 literature-based evaluations, the meaning of the term “quality” varies widely across the fields of  
159 systematic review, toxicology, and public health (see discussion in Viswanathan et al. 2012).  
160 Broadly defined, study quality includes: (1) **reporting quality**—how well or completely a study  
161 was reported, (2) **internal validity** or **risk of bias**—how credible are the findings based on the  
162 design and apparent conduct of a study, and (3) **external validity** or **directness** and  
163 **applicability**—how well a study addresses the topic under review (see Cochrane Collaboration  
164 2013 for detailed definitions). Study quality assessment tools that mix different aspects of study  
165 quality or provide a single summary score are discouraged (Balshem et al. 2011; Higgins and  
166 Green 2011; Liberati et al. 2009; Viswanathan et al. 2012).

167 The OHAT risk-of-bias tool adapts guidance from the Agency for Healthcare Research  
168 and Quality (AHRQ) (Viswanathan et al. 2012). Individual risk-of-bias questions are designated  
169 as only applicable to certain types of study designs (e.g., human controlled trials, experimental  
170 animal studies, cohort studies, case-control studies, cross-sectional studies, and case series or  
171 case reports), with a subset of the questions applying to each study design (**Table 1**).

172 Published tools do not address risk-of-bias criteria for animal studies because risk-of-bias  
173 tools, as with systematic review methods in general, have been focused on guidelines for clinical

174 medicine. OHAT evaluates risk of bias in experimental animal studies using criteria similar to  
175 those applied to human randomized controlled trials, because these study designs are similar in  
176 their ability to control timing and dose of exposure and to minimize the impact of confounding  
177 factors. Using the same set of questions for all study types, including experimental animal  
178 studies, allows for comparison of particular risk-of-bias issues across a body of evidence and  
179 facilitates comparison of the strengths and weaknesses of different bodies of evidence.

180 All references are independently assessed for risk of bias for each outcome of interest by  
181 two reviewers who answer all of the applicable questions with one of four options (definitely  
182 low, probably low, probably high, or definitely high risk of bias (CLARITY Group at McMaster  
183 University 2013) following pre-specified criteria detailed in the protocol. Discrepancies between  
184 the reviewers are resolved by reaching agreement through discussion.

### 185 ***Step 5: Rate the Confidence in the Body of Evidence***

186 For each outcome, the confidence in the body of evidence is rated by considering the strengths  
187 and weaknesses of a collection of studies with similar study design features. Ratings reflect  
188 confidence that the study findings accurately reflect the true association between exposure and  
189 effect including aspects of external validity (or directness and applicability) for the studies. The  
190 OHAT method is based on the Grading of Recommendations Assessment, Development and  
191 Evaluation Working Group (GRADE, <http://www.gradeworkinggroup.org/>) guidelines which  
192 have been adopted by the Cochrane Collaboration (Schünemann et al. 2012) and AHRQ  
193 approaches (Balshem et al. 2011; Lohr 2012), which are conceptually very similar. The method  
194 uses four descriptors to indicate the level of confidence in the separate bodies of evidence  
195 (**Table 2**). In the context of identifying research needs, a conclusion of “High Confidence”  
196 indicates that further research is very unlikely to change the confidence in the apparent

197 relationship between exposure to the substance and the outcome. Conversely, a conclusion of  
198 “Very Low Confidence” suggests that further research is very likely to impact confidence in the  
199 apparent relationship. Human and non-human animal data are considered separately throughout  
200 Steps 5 and 6. Conclusions developed in the subsequent steps of the approach are based on the  
201 evidence with the highest confidence.

202 For each outcome, studies are given an initial confidence rating that reflects the presence  
203 or absence of key study-design features (**Figure 1** for Step 5 schematic). Then studies that have  
204 the same number of features are considered together as a group to begin the process of rating  
205 confidence in a body of evidence for that outcome. The initial rating of each group is  
206 downgraded for factors that decrease confidence and upgraded for factors that increase  
207 confidence in the results. Then, confidence across all studies with the same outcome is assessed  
208 by considering the ratings for all groups of studies with that outcome and the highest rating for  
209 that outcome moves forward.

210 While confidence ratings for each outcome are developed for groups of studies, the  
211 number of studies comprising the group will vary and in some cases this group may be  
212 represented by only one study. Therefore, it is worth noting that a single, well conducted study  
213 may provide evidence of toxicity or a health effect associated with exposure to the substance in  
214 question (e.g., see Germolec (2009) and Foster (2009) for explanation of the NTP levels of  
215 evidence for determination of “toxicity” for individual studies). If a sufficient body of very  
216 similar studies is available, a quantitative meta-analysis may be completed to generate an overall  
217 estimate of effect, but this is not required. Finally, confidence conclusions are developed across  
218 multiple outcomes for those outcomes that are biologically related.

219 It is recognized that the scientific judgments involved in developing these confidence  
220 ratings are inherently subjective. A key advantage of the systematic review process for this step  
221 and throughout an evaluation is that it provides a framework to document and justify the  
222 decisions made, and thereby provides for greater transparency in the scientific basis of judgments  
223 made in reaching conclusions.

224 *Initial confidence set by key features of study design for each outcome*

225 An initial confidence rating is determined by the ability of the study design to address causality  
226 as reflected in the confidence that exposure preceded and was associated with the outcome  
227 (**Figure 1**, Step 5, column 1). This ability is reflected in the presence or absence of four key  
228 study-design features that determine initial confidence ratings, and studies are differentiated  
229 based on whether or not: (1) the exposure to the substance is controlled, (2) the exposure  
230 assessment represents exposures occurring prior to development of the outcome, (3) the outcome  
231 is assessed on the individual level (i.e., not population aggregate data), and (4) a comparison or  
232 control group is used within the study. The first key feature, “controlled exposure” reflects the  
233 ability of experimental studies in humans and animals to largely eliminate confounding by  
234 randomizing allocation of exposure. Therefore, these studies will usually have all four features  
235 and receive an initial rating of “High Confidence.” Observational studies do not have controlled  
236 exposure and are differentiated by the presence or absence of the three remaining study-design  
237 features. For example, prospective cohort studies usually have all three remaining features and  
238 receive an initial rating of “Moderate Confidence,” while a case report may have only one key  
239 feature and receive an initial rating of “Very Low Confidence” (see Supplemental Material,  
240 Table S2 for key features for standard study designs and discussion). The presence or absence of  
241 these study design features capture and discriminate studies on an outcome-specific basis

242 (experimental, prospective, etc.) but do not replace consideration of risk of bias elements or  
243 external validity in other steps.

244 *Downgrade confidence rating*

245 Five properties of the body of evidence (risk of bias, unexplained inconsistency, indirectness,  
246 imprecision, and publication bias) are considered to determine if the initial confidence rating  
247 should be downgraded (**Figure 1**, Step 5, column 2). For each of the five properties, a judgment  
248 is made and documented regarding whether or not there are substantial issues that decrease the  
249 confidence rating in each aspect of the body of evidence for the outcome. Factors that would  
250 downgrade confidence by one versus two levels are specified in the protocol. The reasons for  
251 downgrading confidence may not fit neatly into a single property of the body of evidence. If the  
252 decision to downgrade for two properties is borderline, the body of evidence is downgraded once  
253 to account for both partial concerns. Similarly, the body of evidence is not downgraded twice for  
254 what is essentially the same limitation that could be considered applicable to more than one  
255 property of the body of evidence.

256 **Risk of bias of the body of evidence:** Risk-of-bias criteria were described in Step 4  
257 where study-quality issues for individual studies are evaluated on an outcome-specific  
258 basis. In this step, the previous risk-of-bias assessments for individual studies now serve  
259 as the basis for an overall risk-of-bias conclusion for the entire body of evidence.  
260 Downgrading for risk of bias should reflect the entire body of studies and therefore the  
261 decision to downgrade should be applied conservatively. The decision to downgrade  
262 should be reserved for cases where there is substantial risk of bias across most of the  
263 studies comprising the body of evidence (Guyatt et al. 2011e).

264 **Unexplained inconsistency:** Inconsistency, or large variability in the magnitude or  
265 direction of estimates of effect across studies that cannot be explained, reduces  
266 confidence in the body of evidence. Large inconsistency across studies should be  
267 explored, preferably through *a priori* hypotheses that might explain the heterogeneity.

268 **Indirectness:** Indirectness can refer to external validity or indirect measures of the health  
269 outcome. Indirectness can lower confidence in the body of evidence when the population,  
270 exposure, or outcome(s) measured differs from those that are of most interest. Concerns  
271 about directness could apply to the relationship between: (1) a measured outcome and a  
272 health effect (i.e., upstream biomarker of a health effect), (2) the route of exposure and  
273 the typical human exposure, (3) the study population and the population of interest  
274 (Guyatt et al. 2011c; Lohr 2012), (4) timing of the exposure relative to the appropriate  
275 biological window to affect the outcome, or (5) timing of outcome assessment and the  
276 duration of time required after an exposure for the development of the outcome  
277 (Viswanathan et al. 2012).

278 Note that the administered dose or exposure level is not considered a factor under  
279 indirectness for developing a confidence rating for the purpose of hazard identification.  
280 While exposure level is an important factor when considering the relevance of study  
281 findings to human health effects at known human exposure levels, in the OHAT  
282 evaluation process, this consideration occurs after hazard identification as part of  
283 reaching a “level of concern” conclusion (Jahnke et al. 2005; Medlin 2003; Shelby 2005;  
284 Twombly 1998). The accuracy of an exposure metric (e.g., market basket survey vs.  
285 individual blood levels of a substance) is also not considered a factor under indirectness,

286 and the confidence in the exposure assessment is considered in the risk-of-bias evaluation  
287 of individual studies on an outcome basis in Step 4.

288 **Imprecision:** Imprecision is the lack of certainty in an estimate of effect for a specific  
289 outcome. A precise estimate enables the evaluator to determine whether or not there is an  
290 effect (i.e., it is different from the comparison group). Confidence intervals for the  
291 estimates of effect provide the primary evidence used in considering the imprecision of  
292 the body of evidence (Guyatt et al. 2011b).

293 **Publication bias:** Publication bias is addressed specifically in rating the body of  
294 evidence, and selective reporting within a study is covered in the risk-of-bias criteria  
295 addressing these limitations (Guyatt et al. 2011d). Funnel plots provide a useful tool to  
296 visualize asymmetrical or symmetrical patterns of study results for assessing publication  
297 bias when there is a sufficient body of studies for a specific outcome (e.g., Ahmed et al.  
298 2012). There is empirical evidence that studies with negative results (null findings for  
299 clinical trials) are less likely to be in the published literature (Hopewell et al. 2009).  
300 Negative studies may also be affected by “lag bias” or longer time to publication (Stern  
301 and Simes 1997), and therefore it is important to carefully consider data sets limited to  
302 few positive studies with small sample size that might indicate a lag time between early  
303 positive studies and lagging negative studies. While some publication bias is expected,  
304 downgrading is reserved for when serious concern for publication bias significantly  
305 decreases confidence in the body of evidence.

306 *Upgrade confidence rating*

307 Four properties of the body of evidence (large magnitude of effect, dose response, residual  
308 confounding increases confidence, and cross-species/population/study consistency) are

309 considered to determine if the confidence rating should be upgraded (**Figure 1**, Step 5,  
310 column 3). For each of the four properties, a judgment is made and documented regarding  
311 whether or not there are substantial factors that increase the confidence rating in the body of  
312 evidence for the outcome. As discussed in downgrading, two borderline upgrades could be  
313 combined for one upgrade and the body should not be upgraded twice for essentially the same  
314 attribute. Factors that would upgrade confidence by one versus two levels are specified in the  
315 protocol.

316 **Large magnitude of effect:** A large magnitude of effect is defined as an observed effect  
317 that is sufficiently large such that it is unlikely to have occurred as a result of bias from  
318 potential confounding factors.

319 **Dose response:** A plausible dose–response relationship between level of exposure and  
320 the outcome increases confidence in the result because it reduces concern that the result  
321 could be due to chance. In addition to considering dose-response within a study with a  
322 range of exposure levels, multiple studies with varied exposure levels can contribute to an  
323 overall picture of the dose response. It is important to recognize that prior knowledge  
324 may lead to an expectation for a non-monotonic dose response. Therefore, the plausibility  
325 of the observed biological response should be considered in evaluating the dose–response  
326 relationship.

327 **Residual confounding increases confidence:** This element refers to consideration of  
328 residual confounding, healthy worker effect, or effect modification that would bias the  
329 effect estimate towards the null. If a study reports an effect or association despite the  
330 presence of residual confounding that would diminish the association, confidence in the  
331 association is increased. This confounding can push in either direction, and therefore

332 confidence in the results are increased when there is an indication that a body of evidence  
333 is potentially biased by factors counter to the observed effect.

334 **Cross-species/population/study consistency:** Three types of consistency in the body of  
335 evidence can increase confidence in the results: across animal studies—consistent results  
336 reported in multiple experimental animal models or species; across dissimilar  
337 populations—consistent results reported across populations (human or wildlife) that differ  
338 in factors such as time, location, and/or exposure; and across study types—consistent  
339 results reported from studies with different design features.

340 **Other:** Additional factors specific to the topic being evaluated (e.g., particularly rare  
341 outcomes) may result in increasing a confidence rating. These other factors would be  
342 specified and defined in the protocol.

343 *Combine confidence conclusions for all study types and multiple outcomes*

344 Conclusions are based on the evidence with the highest confidence when considering evidence  
345 across study types and multiple outcomes. Confidence ratings are initially set based on key  
346 design features of the available studies for a given outcome (e.g., for experimental studies  
347 separately from observational studies). The studies with the highest confidence rating form the  
348 basis for the confidence conclusion for each evidence stream. As outlined previously, consistent  
349 results across studies with different design features increase confidence in the combined body of  
350 evidence and can result in an upgraded confidence rating moving forward to Step 6. If the only  
351 available body of evidence receives a “Very Low Confidence” rating, then conclusions for those  
352 outcomes will not move on to Step 6.

353 After confidence conclusions are developed for a given outcome, conclusions for  
354 multiple outcomes are developed. The project-specific definition of an outcome and the grouping

355 of biologically related outcomes used in this step follow the definitions developed *a priori* in the  
356 protocol; deviations are taken with care, justified, and documented. When outcomes are  
357 sufficiently biologically related that they may inform confidence on the overall health outcome,  
358 confidence conclusions may be developed in two steps. Each outcome would first be considered  
359 separately. Then, the related outcomes would be considered together and re-evaluated for  
360 properties that relate to downgrading and upgrading the body of evidence.

361 ***Step 6: Translate the Confidence Ratings into Level of Evidence for Health***  
362 ***Effect***

363 The level of evidence is assessed separately within the human, experimental animal, and to the  
364 extent possible and necessary, other relevant data sets. The conclusions for the level of evidence  
365 for health effects reflect the overall confidence in the association between exposure to the  
366 substance and the outcome (effect or no effect); **Figure 1** for Step 6 schematic). The strategy  
367 uses four terms to describe the level of evidence for health effects. These descriptors reflect both  
368 the confidence in the body of evidence for a given outcome and the direction of effect. There are  
369 three descriptors used in Step 6 (“High Level of Evidence,” “Moderate Level of Evidence,” and  
370 “Low Level of Evidence”) that directly translate from the confidence-in-the-evidence ratings that  
371 exposure to the substance is associated with a health effect, and a fourth designation (“Evidence  
372 of No Health Effect”) to indicate confidence that the substance is not associated with a health  
373 effect (see Supplemental Table 3 for definitions of the level of evidence for health effects  
374 descriptors). Because of the inherent difficulty in proving a negative, the conclusion “Evidence  
375 of No Health Effect” is only reached when there is high confidence in the body of evidence. In  
376 the context of evidence potentially supporting a conclusion of no health effect, a low or moderate  
377 level of evidence results in a conclusion of inadequate evidence to reach a conclusion.

378           Although the conclusions describe associations, a causal relationship is implied and the  
379 ratings describe the level of evidence for health effects in terms of confidence in the association  
380 or the estimate of effect determined from the body of evidence. **Table 3** outlines how the  
381 Bradford Hill considerations on causality (Hill 1965) are related to the process of evaluating the  
382 confidence in the body of evidence and then integrating the evidence (similar to GRADE  
383 approach as described in Schünemann et al. 2011).

384 ***Step 7: Integrate the Evidence to Develop Hazard Identification Conclusions***

385 The highest level of evidence for a health effect from each of the evidence streams is combined  
386 in the final step of the evidence assessment process to determine the hazard identification  
387 conclusion. Hazard identification conclusions may be reached on individual outcomes (health  
388 effects) or groups of biologically related outcomes, as appropriate, based on the evaluation's  
389 objectives and the available data. The rationale for such conclusions is documented as the  
390 evidence is combined within and across evidence streams, and the conclusions are clearly stated  
391 as to which outcomes are incorporated into each conclusion. The five hazard identification  
392 conclusion categories are:

- 393           •       Known to be a hazard to humans
- 394           •       Presumed to be a hazard to humans
- 395           •       Suspected to be a hazard to humans
- 396           •       Not classifiable as a hazard to humans
- 397           •       Not identified to be a hazard to humans

398 In Step 7, the evidence streams for human studies and non-human animal studies, which  
399 have remained separate through the previous steps, are integrated along with other relevant data.  
400 Hazard identification conclusions are reached by integrating the highest level-of-evidence  
401 conclusion for a health effect(s) from the human and the animal evidence streams. On an  
402 outcome basis, this approach applies to whether the data support a health effect conclusion or  
403 evidence of no health effect.

404 When the data support a health effect, the level-of-evidence conclusion for human data  
405 from Step 6 (“High,” “Moderate,” or “Low”) is considered together with the level of evidence  
406 for non-human animal data to reach one of four hazard identification conclusions (Step 7 in  
407 **Figure 1**). If one evidence stream (either human or animal) has no studies, then conclusions are  
408 based on the remaining evidence stream alone (which is equivalent to treating the missing  
409 evidence stream as “Low” in Step 7 **Figure 1**).

410 Any impact of other relevant data on the hazard identification conclusion derived by  
411 integrating the human and non-human animal streams is considered next (Step 7 in **Figure 1**).  
412 Other relevant data could include, but are not limited to, mechanistic data, *in vitro* data, or data  
413 based on upstream indicators of a health effect. Note that mechanistic data or another type of  
414 other relevant data is not required to reach a final hazard identification conclusion.

- 415 • If other relevant data provide strong support for biological plausibility of the relationship  
416 between exposure and the health effect, the hazard identification conclusion may be  
417 upgraded (indicated by black “up” arrows in Step 7 graphic in **Figure 1**) from that  
418 initially derived by considering the human and non-human animal evidence together. It is  
419 envisioned that strong evidence for a relevant biological process from mechanistic or *in*

420 *vitro* data could result in a conclusion of “suspected” in the absence of human  
421 epidemiology or experimental animal data.

- 422 • If other relevant data provide strong opposition for biological plausibility of the  
423 relationship between exposure and the health effect, the hazard identification conclusion  
424 may be downgraded (indicated by gray “down” arrows in Step 7 graphic in **Figure 1**).

425 When the data provide evidence of no health effect, the level-of-evidence conclusion for  
426 human data from Step 6 is considered together with the level-of-evidence for health effects  
427 conclusion for non-human animal data. And again, any impact of other relevant data on the  
428 hazard identification conclusion is considered.

- 429 • If the human level-of-evidence conclusion of no health effect is supported by animal  
430 evidence of no health effect, the hazard identification conclusion is “not identified.”

431 The outcome of the evaluation includes any hazard identification conclusions reached or  
432 data needs identified along with a detailed rationale outlining how human, animal, and other  
433 relevant data contributed to the conclusions. Draft OHAT evaluations undergo peer review and  
434 public comment as part of the overall process for finalization and publication  
435 (<http://ntp.niehs.nih.gov/go/38138>).

## 436 **Discussion**

437 Aspects of systematic review methodology designed to increase objectivity and transparency  
438 may add to the time and investment required to develop literature-based evaluations, and NTP is  
439 mindful of these concerns. In applying the OHAT Approach to case studies (NTP 2013e, f), NTP  
440 found that Steps 2-4 were the most time intensive: selecting studies, extracting data, and

441 assessing the quality of individual studies. While not formally part of the systematic review  
442 process, data management resources were used to increase transparency and efficiency in  
443 developing the case studies so that time invested in the early steps was recouped in later steps by  
444 entering study information into a database. Summary tables and graphics were readily made from  
445 the database to facilitate decision making in Steps 6 and 7 when evaluating confidence in a body  
446 of studies and integrating evidence streams to develop conclusions. The value of these  
447 efficiencies and further development of these web-based systems for data display, data  
448 management, and data sharing cannot be understated.

## 449 **Conclusions**

450 Applying systematic-review methodologies to environmental health questions is gaining a  
451 critical mass (EFSA 2010; NRC 2013a, b; Woodruff and Sutton 2011). The OHAT Approach  
452 provides a practical method for applying the principles of systematic review to address  
453 environmental health questions. Moving forward, OHAT will apply this framework in future  
454 evaluations (<http://ntp.niehs.nih.gov/go/evals>). As evaluations are completed and practices in the  
455 field of systematic review evolve, OHAT may refine and amend its “evergreen” approach and  
456 post updates to the framework (NTP 2013b). The protocols and the data compiled as part of an  
457 evaluation (e.g., study-level health effects data and risk-of-bias assessment) will be publicly  
458 available following its completion to increase transparency and facilitate data sharing with  
459 government agencies, scientific community, and the public. The scientifically rigorous and  
460 objective procedures, which have been a hallmark of OHAT literature-based health assessments,  
461 will be strengthened by implementation of the OHAT approach for systematic review and  
462 evidence integration (NTP 2013g).

463           The application of the procedures of systematic review to environmental health questions  
464 has the potential to bring an increase in objectivity and transparency similar to what it has  
465 already done for clinical medicine. Developing evaluations with this approach can improve  
466 communication and clarity about how hazard identification conclusions are reached by  
467 documenting the source of the data considered, the methods of quality assessment used, and the  
468 scientific judgments made during evidence integration.

## 469 **References**

470

471 Ahmed I, Sutton AJ, Riley RD. 2012. Assessment of publication bias, selection bias, and unavailable data in meta-  
472 analyses using individual participant data: a database survey. *Br Med J* 344:d7762.

473 AHRQ (Agency for Healthcare Research and Quality). 2013. AHRQ training modules for the systematic reviews  
474 methods guide. Available: <http://www.effectivehealthcare.ahrq.gov/index.cfm/tools-and-resources/slide-library/>  
475 [accessed 11 October 2013].

476 Balshem H, Helfand M, Schunemann HJ, Oxman AD, Kunz R, Brozek J, et al. 2011. GRADE guidelines: 3. Rating  
477 the quality of evidence. *J Clin Epidemiol* 64:401-406.

478 Birnbaum LS, Thayer KA, Bucher JR, Wolfe MS. 2013. Implementing systematic review at the National  
479 Toxicology Program. *Environ Health Perspect* 121:A108-109.

480 Bucher JR, Thayer K, Birnbaum LS. 2011. The Office of Health Assessment and Translation: a problem-solving  
481 resource for the National Toxicology Program. *Environ Health Perspect* 119:A196-197.

482 CLARITY Group at McMaster University. 2013. Tools to assess risk of bias in cohort studies, case control studies,  
483 randomized controlled trials, and longitudinal symptom research studies aimed at the general population. Available:  
484 <http://www.evidencepartners.com/resources/> [accessed 15 January 2013].

485 Cochrane Collaboration. 2013. Glossary of Cochrane terms. Available: <http://www.cochrane.org/glossary>  
486 [accessed 15 January 2013].

487 EFSA (European Food Safety Authority). 2010. Application of systematic review methodology to food and feed  
488 safety assessments to support decision making. *EFSA Journal*: 1-90. Available:  
489 <http://www.efsa.europa.eu/en/efsajournal/pub/1637.htm>.

490 FDA (Food and Drug Administration). 2010. Guidance for industry: Adaptive design clinical trials for drugs and  
491 biologics: Draft guidance. Silver Spring, MD: 1-50. Available:  
492 <http://www.fda.gov/downloads/Drugs/.../Guidances/ucm201790.pdf> [accessed 29 December 2012].

493 Foster PM. 2009. Explanation of levels of evidence for reproductive system toxicity. Research Triangle Park, NC.  
494 Available: <http://ntp.niehs.nih.gov/go/18711> [accessed 15 January 2013].

495 Germolec D. 2009. Explanation of levels of evidence for immune system toxicity. Research Triangle Park, NC.  
496 Available: <http://ntp.niehs.nih.gov/go/9399> [accessed 15 January 2013].

497 Guyatt GH, Oxman AD, Akl EA, Kunz R, Vist G, Brozek J, et al. 2011a. GRADE guidelines: 1. Introduction-  
498 GRADE evidence profiles and summary of findings tables. *J Clin Epidemiol* 64:383-394.

499 Guyatt GH, Oxman AD, Kunz R, Brozek J, Alonso-Coello P, Rind D, et al. 2011b. GRADE guidelines 6. Rating the  
500 quality of evidence--imprecision. *J Clin Epidemiol* 64:1283-1293.

501 Guyatt GH, Oxman AD, Kunz R, Woodcock J, Brozek J, Helfand M, et al. 2011c. GRADE guidelines: 8. Rating the  
502 quality of evidence--indirectness. *J Clin Epidemiol* 64:1303-1310.

503 Guyatt GH, Oxman AD, Montori V, Vist G, Kunz R, Brozek J, et al. 2011d. GRADE guidelines: 5. Rating the  
504 quality of evidence--publication bias. *J Clin Epidemiol* 64:1277-1282.

505 Guyatt GH, Oxman AD, Vist G, Kunz R, Brozek J, Alonso-Coello P, et al. 2011e. GRADE guidelines: 4. Rating the  
506 quality of evidence--study limitations (risk of bias). *J Clin Epidemiol* 64:407-415.

507 Higgins J, Green S. 2011. *Cochrane Handbook for Systematic Reviews of Interventions*. Available: [www.cochrane-](http://www.cochrane-handbook.org)  
508 [handbook.org](http://www.cochrane-handbook.org) [accessed 3 February 2013].

509 Hill AB. 1965. The environment and disease: Association or causation? *Proc Roy Soc Med* 58:295-300.

510 Hopewell S, Loudon K, Clarke MJ, Oxman AD, Dickersin K. 2009. Publication bias in clinical trials due to  
511 statistical significance or direction of trial results. *Cochrane Database Syst Rev*:1-26.

- 512 Jahnke GD, Iannucci AR, Scialli AR, Shelby MD. 2005. Center for the evaluation of risks to human reproduction--  
513 the first five years. *Birth Defects Res B Dev Reprod Toxicol* 74:1-8.
- 514 Khan K, ter Riet G, Glanville J, Sowden A, Kleijnen J, eds. 2001. *Undertaking Systematic Reviews of Research on*  
515 *Effectiveness: CRD's Guidance for those Carrying Out or Commissioning Reviews (CRD Report Number 4) (2nd*  
516 *edition)*. York (UK): NHS Centre for Reviews and Dissemination:University of York.
- 517 Liberati A, Altman DG, Tetzlaff J, Mulrow C, Gotzsche PC, Ioannidis JP, et al. 2009. The PRISMA statement for  
518 reporting systematic reviews and meta-analyses of studies that evaluate health care interventions: explanation and  
519 elaboration. *J Clin Epidemiol* 62:e1-34.
- 520 Lohr KN. 2012. Grading the strength of evidence. In: *The Agency for Healthcare Research and Quality (AHRQ)*  
521 *Training Modules for Systematic Reviews Methods Guide*. Available:  
522 <http://www.effectivehealthcare.ahrq.gov/index.cfm/slides/?pageAction=displaySlides&tk=18> [accessed 13 July  
523 2012].
- 524 Medlin J. 2003. New arrival: CERHR monograph series on reproductive toxicants. *Environ Health Perspect*  
525 111:A696-698.
- 526 Moher D, Liberati A, Tetzlaff J, Altman DG. 2009. Preferred reporting items for systematic reviews and meta-  
527 analyses: The PRISMA statement. *J Clin Epidemiol* 62:1006-1012.
- 528 NRC (National Research Council). 2011. Committee to Review of the Environmental Protection Agency's Draft  
529 IRIS Assessment of Formaldehyde. The National Academies Press: 1-190. Available:  
530 [http://www.nap.edu/openbook.php?record\\_id=13142](http://www.nap.edu/openbook.php?record_id=13142) [accessed 30 July 2012].
- 531 NRC (National Research Council). 2013a. Critical Aspects of EPA's IRIS Assessment of Inorganic Arsenic -  
532 Interim Report. The National Academies Press: 1-117. Available: [http://www.nap.edu/catalog.php?record\\_id=18594](http://www.nap.edu/catalog.php?record_id=18594)  
533 [accessed 7 November 2013].
- 534 NRC (National Research Council). 2013b. Review of the IRIS Process. Available:  
535 <http://www8.nationalacademies.org/cp/projectview.aspx?key=49458> [accessed 1 November 2013].
- 536 NTP (National Toxicology Program). 2012a. Board of Scientific Counselors December 11, 2012 meeting. Meeting  
537 materials. Available: <http://ntp.niehs.nih.gov/go/9741> [accessed 21 February 2013].
- 538 NTP (National Toxicology Program). 2012b. Board of Scientific Counselors June 21-22, 2012 meeting. Meeting  
539 materials. Available: <http://ntp.niehs.nih.gov/go/9741> [accessed 16 June 2013].
- 540 NTP (National Toxicology Program). 2013a. Webinar on the assessment of data quality in animal studies. March 20,  
541 2013. Available: <http://ntp.niehs.nih.gov/go/38752> [accessed 7 April 2013].
- 542 NTP (National Toxicology Program). 2013b. Draft OHAT Approach for Systematic Review and Evidence  
543 Integration for Literature-based Health Assessments – February 2013. RTP, NC: Office of Health Assessment and  
544 Translation. Available: <http://ntp.niehs.nih.gov/go/38138> [accessed 10 March 2013].
- 545 NTP (National Toxicology Program). 2013c. Board of Scientific Counselors June 25, 2013 meeting. Meeting  
546 materials. Available: <http://ntp.niehs.nih.gov/go/9741> [accessed 1 November 2013].
- 547 NTP (National Toxicology Program). 2013d. OHAT Evaluation Process. Available:  
548 <http://ntp.niehs.nih.gov/go/38138> [accessed 16 June 2013].
- 549 NTP (National Toxicology Program). 2013e. Draft Protocol for Systematic Review to Evaluate the Evidence for an  
550 Association Between Perfluorooctanoic Acid (PFOA) or Perfluorooctane Sulfonate (PFOS) Exposure and  
551 Immunotoxicity. RTP, NC: Office of Health Assessment and Translation. Available:  
552 <http://ntp.niehs.nih.gov/go/38673> [accessed 9 April 2013].
- 553 NTP (National Toxicology Program). 2013f. Draft Protocol for Systematic Review to Evaluate the Evidence for an  
554 Association Between Bisphenol A (BPA) and Obesity. RTP, NC: Office of Health Assessment and Translation.  
555 Available: <http://ntp.niehs.nih.gov/go/38673> [accessed 9 April 2013].
- 556 NTP (National Toxicology Program). 2013g. OHAT Implementation of Systematic Review. Available:  
557 <http://ntp.niehs.nih.gov/go/38673> [accessed 16 June 2013].

- 558 Rhomberg LR, Goodman JE, Bailey LA, Prueitt RL, Beck NB, Bevan C, et al. 2013. A survey of frameworks for  
559 best practices in weight-of-evidence analyses. *Crit Rev Toxicol* 43:753-784.
- 560 Schünemann H, Hill S, Guyatt G, Akl EA, Ahmed F. 2011. The GRADE approach and Bradford Hill's criteria for  
561 causation. *J Epidemiol Community Health* 65:392-395.
- 562 Schünemann HJ, Oxman AD, Vist GE, Higgins JPT, Deeks JJ, P. G, et al. 2012. Chapter 12: Interpreting results  
563 and drawing conclusions. In: *Cochrane Handbook for Systematic Reviews of Interventions* Version 5.10 [updated  
564 March 2011], (Higgins JPT, Green S, eds):The Cochrane Collaboration.
- 565 Shelby MD. 2005. National Toxicology Program Center for the Evaluation of Risks to Human Reproduction:  
566 guidelines for CERHR expert panel members. *Birth Defects Res B Dev Reprod Toxicol* 74:9-16.
- 567 Stern JM, Simes RJ. 1997. Publication bias: evidence of delayed publication in a cohort study of clinical research  
568 projects. *Br Med J* 315:640-645.
- 569 Twombly R. 1998. New NTP centers meet the need to know. *Environ Health Perspect* 106:A480-483.
- 570 Viswanathan M, Ansari M, Berkman ND, Chang S, Hartling L, McPheeters LM, et al. 2012. Assessing the risk of  
571 bias of individual studies when comparing medical interventions. Agency for Healthcare Research and Quality  
572 (AHRQ). Available: [http://www.effectivehealthcare.ahrq.gov/index.cfm/search-for-guides-reviews-and-](http://www.effectivehealthcare.ahrq.gov/index.cfm/search-for-guides-reviews-and-reports/?pageaction=displayproduct&productid=998)  
573 [reports/?pageaction=displayproduct&productid=998](http://www.effectivehealthcare.ahrq.gov/index.cfm/search-for-guides-reviews-and-reports/?pageaction=displayproduct&productid=998) [accessed 3 January 2013].
- 574 Whitlock EP, Lopez SA, Chang S, Helfand M, Eder M, Floyd N. 2010. AHRQ series paper 3: identifying, selecting,  
575 and refining topics for comparative effectiveness systematic reviews: AHRQ and the effective health-care program.  
576 *J Clin Epidemiol* 63:491-501.
- 577 Woodruff TJ, Sutton P. 2011. An evidence-based medicine methodology to bridge the gap between clinical and  
578 environmental health sciences. *Health Aff* 30:931-937.
- 579
- 580
- 581

582

583 **Tables**

584

585

586

587 **Table 1.** OHAT Risk-of-bias Questions

588 The OHAT risk-of-bias questions are applied to evaluate the risk of bias of studies on an outcome basis. The study design determines  
589 which questions are applicable as indicated in the table by an “X” for each question that applies to a given study design. Risk-of-bias  
590 ratings are developed by answering each applicable question with one of four options (definitely low, probably low, probably high, or  
591 definitely high risk of bias).

592

593

594

595

596

597

**Table 1:** OHAT Risk-of-bias Questions.

The OHAT risk-of-bias questions are applied to evaluate the risk of bias of studies on an outcome basis. The study design determines which questions are applicable as indicated in the table by an “X” for each question that applies to a given study design. Risk-of-bias ratings are developed by answering each applicable question with one of four options (definitely low, probably low, probably high, or definitely high risk of bias). Answering “Yes” indicates lower risk of bias, while “No” indicates higher risk of bias for that question)

	Experimental Animal <sup>1</sup> Human Controlled Trials <sup>2</sup>	Cohort	Case-control	Cross-sectional <sup>3</sup>	Case Series
<b>Selection BIAS</b>					
<p><b>Was administered dose or exposure level adequately randomized?</b> Randomization requires that each human subject or animal had an equal chance of being assigned to any study group including controls (e.g., use of random number table or computer generated randomization).</p>	X	X			
<p><b>Was allocation to study groups adequately concealed?</b> Allocation concealment requires that research personnel do not know which administered dose or exposure level is assigned at the start of a study. Human studies also require that allocation be concealed from human subjects prior to entering the study. <i>Note: 1) a question under performance bias addresses blinding of personnel and human subjects to treatment during the study; 2) a question under detection bias addresses blinding of outcome assessors.</i></p>	X	X			
<p><b>Were the comparison groups appropriate?</b> Comparison group appropriateness refers to having similar baseline characteristics between the groups aside from the exposures and outcomes under study.</p>			X	X	X
<b>Confounding BIAS</b>					
<p><b>Did the study design or analysis account for important confounding and modifying variables?</b> <i>Note: a parallel question under detection bias addresses reliability of the measurement of confounding variables.</i></p>	X	X	X	X	X
<p><b>Did researchers adjust or control for other exposures that are anticipated to bias results?</b></p>	X	X	X	X	X
<b>Performance BIAS</b>					
<p><b>Were experimental conditions identical across study groups?</b></p>	X				
<p><b>Did researchers adhere to the study protocol?</b></p>	X	X	X	X	X
<p><b>Were the research personnel and human subjects blinded to the study group during the study?</b> Blinding requires that study scientists do not know which administered dose or exposure level the human subject or animal is being given (i.e., study group). Human studies require blinding of the human subjects when possible.</p>	X	X			

Table 1 OHAT Risk-of-bias Questions continued

	Experimental Animal <sup>1</sup>	Human Controlled Trials <sup>2</sup>	Cohort	Case-control	Cross-sectional <sup>3</sup>	Case Series
<b>Attrition/Exclusion BIAS</b>						
<b>Were outcome data complete without attrition or exclusion from analysis?</b> Attrition rates are required to be similar and uniformly low across groups with respect to withdrawal or exclusion from analysis.	X	X	X	X	X	
<b>Detection BIAS</b>						
<b>Were the outcome assessors blinded to study group or exposure level?</b> Blinding requires that outcome assessors do not know the study group or exposure level of the human subject or animal when the outcome was assessed.	X	X	X	X	X	X
<b>Were confounding variables assessed consistently across groups using valid and reliable measures?</b> Consistent application of valid, reliable, and sensitive methods of assessing important confounding or modifying variables is required across study groups. <i>Note, a parallel question under selection bias addresses whether design or analysis account for confounding.</i>	X	X	X	X	X	X
<b>Can we be confident in the exposure characterization?</b> Confidence requires valid, reliable, and sensitive methods to measure exposure applied consistently across groups.	X	X	X	X	X	X
<b>Can we be confident in the outcome assessment?</b> Confidence requires valid, reliable, and sensitive methods to assess the outcome and the methods should be applied consistently across groups.	X	X	X	X	X	X
<b>Selective Reporting BIAS</b>						
<b>Were all measured outcomes reported?</b>	X	X	X	X	X	X
<b>Other</b>						
<b>Were there no other potential threats to internal validity (e.g., statistical methods were appropriate)?</b> On a project specific basis, additional questions for other potential threats to internal validity can be added and applied to study designs as appropriate.						

599  
600  
601  
602  
603

<sup>1</sup> Experimental animal studies are controlled exposure studies. Non-human animal observational studies could be evaluated using the design features of observational human studies such as cross-sectional study design.

<sup>2</sup> Human Controlled Trials (HCTs): studies in humans with a controlled exposure, including Randomized Controlled Trials (RCTs) and non-randomized experimental studies.

<sup>3</sup> Cross-sectional studies include population surveys with individual data (e.g., National Health and Nutrition Examination Survey or NHANES) and population surveys with aggregate data (i.e., air pollution exposure estimated by zip code).

604

605 **Table 2.** Confidence Ratings in the Bodies of Evidence

<b>Confidence Rating</b>	<b>Definition</b>
<b>High Confidence (++++)</b>	High confidence in the association between exposure to the substance and the outcome. The true effect is <u>highly likely to be</u> reflected in the apparent relationship.
<b>Moderate Confidence (+++)</b>	Moderate confidence in the association between exposure to the substance and the outcome. The true effect <u>may be</u> reflected in the apparent relationship.
<b>Low Confidence (++)</b>	Low confidence in the association between exposure to the substance and the outcome. The true effect <u>may be different</u> than the apparent relationship.
<b>Very Low Confidence (+)</b>	Very low confidence in the association between exposure to the substance and the outcome. The true effect <u>is highly likely to be</u> different than the apparent relationship.

606

607

608

**Table 3.** Aspects of the Hill considerations on causality within the OHAT Approach

<b>Hill Consideration</b>	<b>Relationship to the OHAT Approach</b>
Strength	Considered in upgrading the confidence rating for the body of evidence for <b>large magnitude of effect</b> and downgrading the confidence rating for <b>imprecision</b> .
Consistency	Considered in upgrading confidence rating for the body of evidence for <b>consistency across study types, across dissimilar populations, or across animal species</b> ; and in integrating the body of evidence among human, animal, and other relevant data; also in downgrading confidence rating for the body of evidence for <b>unexplained inconsistency</b> .
Temporality	Considered in <b>initial confidence ratings</b> by key features of study design, for example experimental studies have an initial rating of “High Confidence” because of the increased confidence that the controlled exposure preceded outcome.
Biological gradient	Considered in upgrading the confidence rating for the body of evidence for evidence of a <b>dose–response</b> relationship.
Biological plausibility	Considered in examining non monotonic <b>dose–response</b> relationships and developing confidence rating conclusions across biologically related outcomes, particularly outcomes along a pathway to disease. Other relevant data that inform plausibility such as physiologically based pharmacokinetic and mechanistic studies are considered in integrating the body of evidence. Also considered in downgrading the confidence rating for the body of evidence for <b>indirectness</b> .
Experimental evidence	Considered in setting <b>initial confidence ratings</b> by key features of study design and downgrading the confidence rating for <b>risk of bias</b> .

609

610

## 611 **Figure Legend**

612 **Figure 1.** The OHAT Approach for Systematic Review and Evidence Integration for Literature-  
613 Based Environmental Health Science Assessments

614

615