

# Biological plausibility in environmental health systematic reviews: a GRADE concept paper

Paul Whaley (1,2), Thomas Piggott (3), Rebecca L. Morgan (3), Sebastian Hoffmann (2), Katya Tsaoun (2), Lukas Schwingshackl (4), Mohammed T. Ansari (5), Kristina A. Thayer (6), Holger Schünemann\* (3,7,8)

## Affiliations

1. Lancaster Environment Centre, Lancaster University, UK

2. Evidence-based Toxicology Collaboration at Johns Hopkins Bloomberg School of Public Health (EBTC)

3. Department of Health Research Methods, Evidence and Impact, McMaster University, 1280 Main St West; Hamilton, ON L8N 3Z5, Canada

4. Institute for Evidence in Medicine, Medical Center - University of Freiburg, Faculty of Medicine, University of Freiburg, Freiburg, Germany

5. School of Epidemiology and Public Health, University of Ottawa. Room 101, 600 Peter Morand Crescent, Ottawa, Ontario K1G 5Z3, Canada

6. U.S. Environmental Protection Agency (US EPA), Office of Research and Development, Center for Public Health and Environmental Assessment (CPHEA), Chemical Pollutant Assessment Division (CPAD), 1200 Pennsylvania Avenue, NW (8623R), Washington, DC 20460, USA

7. Michael G DeGroot Cochrane Canada and McMaster GRADE Centres; McMaster University, HSC-2C, 1280 Main St West; Hamilton, ON L8N 3Z5, Canada

8. Dipartimento di Scienze Biomediche, Humanitas University, Via Rita Levi Montalcini, 4, 20090 Pieve Emanuele, Milan, Italy,

## 24 Abstract

25 **Background:** “Biological plausibility” is a concept frequently referred to in environmental and  
26 public health when researchers are evaluating how confident they are in the results and  
27 inferences of a study or evidence review. Biological plausibility is not, however, a domain of  
28 one of the most widely-used approaches for assessing the certainty of evidence (CoE) which  
29 underpins the findings of a systematic review, the Grading of Recommendations  
30 Assessment, Development and Evaluation (GRADE) CoE Framework. Whether the omission  
31 of biological plausibility is a potential limitation of the GRADE CoE Framework is a topic that  
32 is regularly discussed, especially in the context of environmental health systematic reviews.

33 **Objectives:** We analyse how the concept of “biological plausibility”, as applied in the context  
34 of assessing certainty of the evidence that supports the findings of a systematic review, is  
35 accommodated under the processes of systematic review and the existing GRADE domains.

36 **Results and Discussion:** We argue that “biological plausibility” is a concept which primarily  
37 comes into play when direct evidence about the effects of an exposure on a population of  
38 concern (usually humans) is absent, at high risk of bias, is inconsistent, or limited in other  
39 ways. In such circumstances, researchers look toward evidence from other study designs in  
40 order to draw conclusions. In this respect, we can consider experimental animal and *in vitro*  
41 evidence as “surrogates” for the target populations, exposures, comparators and outcomes  
42 of actual interest. Through discussion of 10 examples of experimental surrogates, we  
43 propose that the concept of biological plausibility consists of two principal aspects: a  
44 “generalisability aspect” and a “mechanistic aspect”. The “generalisability aspect” concerns  
45 the validity of inferences from experimental models to human scenarios, and asks the same  
46 question as does the assessment of external validity or indirectness in systematic reviews.  
47 The “mechanistic aspect” concerns certainty in knowledge of biological mechanisms and  
48 would inform judgements of indirectness under GRADE, and thus the overall CoE. While  
49 both aspects are accommodated under the indirectness domain of the GRADE CoE  
50 Framework, further research is needed to determine how to use knowledge of biological  
51 mechanisms in the assessment of indirectness of the evidence in systematic reviews.

52 **Keywords:** systematic review; biological plausibility; surrogates; environmental health;  
53 toxicology; epidemiology; Bradford Hill;

## 54 Introduction

55 In environmental and public health research, toxicology, and human health chemical risk  
56 assessment (henceforth referred to as “environmental health research”) it is rare to have  
57 direct evidence from studies in humans of the effects that environmental exposures might be  
58 having on people’s health. This elevates the importance in environmental health research of  
59 evidence from experimental animal (*in vivo*) and *in vitro* studies. However, while evidence  
60 from *in vivo* and *in vitro* studies has the advantage that exposure can be controlled, the  
61 laboratory set-up is only indirectly representative of the human situation which it models -  
62 using animals in place of people, artificial cell culture constructs to measure biological  
63 processes, and exposure regimens which are often much higher, shorter and more  
64 regimented than would be seen in human cases (Rhomborg, 2015).

65 There is often, therefore, a need to translate the evidence from laboratory experiments to the  
66 human scenarios they are informing. Our ability to do this correctly is critical in successfully  
67 identifying, quantifying, and limiting health harms from environmental exposures. As  
68 systematic reviews become mainstream in environmental health (Bilotta, Milner and Boyd,  
69 2014; Sheehan and Lam, 2015; Morgan *et al.*, 2016; Whaley *et al.*, 2016; Hoffmann *et al.*,  
70 2017), the need for systematic approaches for translating evidence from the laboratory to the  
71 human context becomes increasingly important (Lewis et al. 2017).

72 One concept which is often applied in assessing causality and translating the findings of  
73 laboratory experiments to human contexts (or, indeed, one epidemiological context to  
74 another) is that of “biological plausibility”. As a concept, biological plausibility was first  
75 formalised in 1965 by Sir Austin Bradford Hill, as one of his considerations for establishing  
76 causality (Hill, 1965). Bradford Hill argued that the presence of biological plausibility can  
77 increase the likelihood that a relationship between an exposure and a health outcome is a  
78 causal one. However, despite the evolution of thinking around the concept and the many  
79 definitions of “biological plausibility” that are available (see Table 1 for some examples),  
80 exactly what constitutes biological plausibility has never been fully or finally characterised.  
81 This is particularly true in the context of conducting environmental health systematic reviews.  
82 Methodologists, including those in the Grading of Recommendations Assessment,  
83 Development and Evaluation (GRADE) Working Group, are frequently challenged by  
84 environmental health practitioners about whether and how the assessment of biological  
85 plausibility is accommodated in the systematic review process (European Food Safety  
86 Authority, 2018).

87

Source	Definition of “biological plausibility”
Bradford Hill (1965)	“It will be helpful if the causation we suspect is biologically plausible. But this is a feature I am convinced we cannot demand. What is biologically plausible depends upon the biological knowledge of the day.”
Wikipedia (Wikipedia contributors, 2014)	“A relationship between a putative cause and an outcome — that is consistent with existing biological and medical knowledge” and “one component of a method of reasoning that can establish a cause-and-effect relationship between a biological factor and a particular disease or adverse event”
European Food Safety Authority (Hardy <i>et al.</i> , 2017)	“Consistency between data and biological theory or mechanism”
Last’s Dictionary of Epidemiology (International Epidemiological Association, 2001)	The “causal consideration that an observed, potentially causal association between an exposure and a health outcome may plausibly be attributed to causation on the basis of existing biomedical and epidemiological knowledge.”
Organisation for Economic Co-operation and Development (OECD, 2016)	Being “consistent with biological knowledge” and “based on extensive previous documentation and broad acceptance”
US Environmental Protection Agency Cancer Guidelines (US Environmental Protection Agency, 2005)	“An inference of causality [which] tends to be strengthened by consistency with data from experimental studies or other sources demonstrating plausible biological mechanisms. A lack of mechanistic data, however, is not a reason to reject causality.”

88

89 **Table 1.** Examples of definitions of “biological plausibility”

90 **GRADE and biological plausibility**

91 The GRADE Framework, originally introduced in 2003, is commonly used in public health  
 92 and healthcare systematic reviews, and increasingly in environmental health (Morgan *et al.*,  
 93 2016; Morgan *et al.* 2019). GRADE contends that assessment of the certainty of evidence  
 94 for answers to research questions can be successfully operationalised (i.e. conducted  
 95 accurately, consistently and transparently by different researchers working in different times  
 96 and places) via systematic consideration of a predefined set of eight “domains” of strengths  
 97 and limitations of the overall evidence base (Guyatt *et al.*, 2008). The domains that reduce  
 98 certainty in a body of evidence summarised in a systematic review are risk of bias,  
 99 inconsistency, indirectness, imprecision and publication bias. The domains which increase  
 100 certainty are large effect size, presence of a dose-response relationship, and residual  
 101 opposing confounding (see Figure 1).

102 These domains are intended to be exhaustive of the concepts necessary for assessing  
 103 certainty in the evidence, operationalised via a structured reasoning process designed to  
 104 produce more consistent and transparent results than is achievable by direct application of  
 105 the considerations of Bradford Hill. Historically, the contention has been that the role played

106 by assessment of biological plausibility in environmental health assessments is already  
 107 accommodated either in the GRADE domains or as part of the systematic review process  
 108 (Schünemann *et al.*, 2011; Hultcrantz *et al.*, 2017). The GRADE Working Group has  
 109 therefore intentionally not included biological plausibility as a domain in rating the certainty of  
 110 the evidence.

-1- Establish initial level of certainty		-2- Consider lowering or raising level of certainty		-3- Final rating for level of certainty
Study Design	Initial level of certainty in an estimate of effect	Reasons for considering lowering or raising certainty		Certainty in an estimate of effect across those considerations
		Lower if	Higher if**	
Randomised trials	High	Risk of bias Unexplained inconsistency Indirectness Imprecision Publication bias	Large effect  Dose response  All plausible confounding and bias would reduce a demonstrated effect or suggest a spurious effect if no effect was observed	High ⊕⊕⊕⊕
	Moderate		⊕⊕⊕⊖	
Observational studies*	Low		⊕⊕⊖⊖	
	Very Low		⊕⊖⊖⊖	

\*Observational studies may start at high certainty if a tool that assesses risk of bias against a target experiment or trial is used (Schünemann *et al.* 2019)  
 \*\*Upgrading criteria are usually applicable to observational studies only, and only applied if the evidence has not already been downgraded

111

112 **Figure 1.** The upgrade and downgrade domains in GRADE and how they are used to determine the overall  
 113 certainty in evidence for a systematic review. Adapted from Morgan *et al.* (2016).

114 Our objectives in this paper are as follows: to further elucidate how the systematic review  
 115 process and the GRADE domains operationalise the assessment of biological plausibility; to  
 116 describe how the concept of biological plausibility maps onto the process of systematically  
 117 reviewing environmental health evidence; and answer the question of how “biomedical”,  
 118 “biological”, or “epidemiological” knowledge, as referred to in the various definitions of  
 119 biological plausibility, contributes to rating certainty in a body of evidence summarised in a  
 120 systematic review.

121 Our argument consists of five parts. Firstly, we argue that consideration of biological  
 122 plausibility is not necessary if the body of evidence that is directly reflective of the  
 123 populations, exposures, comparators and outcomes of concern in a systematic review  
 124 question is sufficiently certain. Secondly, we note that this situation is rare in environmental  
 125 health, and that systematic reviews in this field will often need to include indirect evidence  
 126 from surrogate<sup>1</sup> *in vivo* and *in vitro* experimental models. Thirdly, through 10 examples of the  
 127 use of surrogates, we show what sort of “biological knowledge” is typically used when  
 128 researchers are making judgements about biological plausibility.

<sup>1</sup> We define “surrogate” as any property of a study model that is used to estimate the characteristics of a different property. By “property” we mean any controllable or measurable element of study design. This includes population, exposure or intervention, comparator, outcome, and any individual characteristics thereof respectively. For example, rats might be studied in the laboratory as surrogates for human populations, and IQ might be measured as a surrogate for intellectual capacity.

129 Fourthly, our 10 examples show that the concept of biological plausibility consists of two  
130 connected principle aspects, which we call the “generalisability aspect” and the “mechanistic  
131 aspect”. The “generalisability aspect” of biological plausibility concerns the extent to which  
132 findings from an experimental context apply to a target context of concern. The “mechanistic  
133 aspect” concerns certainty in the evidence of biological mechanisms (i.e. the molecular,  
134 cellular, and organismal events leading to an outcome). Judgements of the generalisability of  
135 a surrogate are informed by evidence of biological mechanisms.

136 Fifthly, we argue that since the generalisability aspect of biological plausibility and the  
137 assessment of indirectness in systematic reviews both concern the external validity of  
138 experimental models, it follows that the generalisability aspect of biological plausibility is  
139 accommodated under the GRADE domain of indirectness. Insofar as judgements of certainty  
140 in biological mechanisms support judgements of the generalisability of a surrogate, then the  
141 mechanistic aspect of biological plausibility should also be operationalised under the  
142 indirectness domain of GRADE.

143 We therefore conclude that, while processes and language may be different, the concepts  
144 involved in the assessment of biological plausibility are covered by the established domains  
145 of GRADE. This means GRADE does not need to introduce additional domains to  
146 accommodate biological plausibility. However, we also recognise that GRADE has not yet  
147 been applied to the assessment of certainty in a way which takes detailed account of  
148 biological mechanisms. We therefore recommend research be conducted to advance  
149 understanding of how knowledge about mechanisms should be applied in determining the  
150 indirectness of evidence.

151 We note that we are not providing a complete account of the concept of biological plausibility  
152 in all contexts and uses, restricting our focus to its application in the conduct of systematic  
153 reviews of exposure-outcome relationships. We also note that there is a potential  
154 relationship between biological plausibility and Bradford Hill’s concept of “coherence”. This is  
155 acknowledged in argument elsewhere that coherence is covered in GRADE under the  
156 domains of inconsistency and indirectness (Schünemann et al., 2011). However, as a  
157 different concept to biological plausibility, coherence it is not a focus of this article.

## 158 “Biological plausibility” and the inclusion of 159 surrogates in systematic reviews

160 Systematic review can be defined as the application of methods designed to minimise risk of  
161 systematic and random error, and maximise transparency of decision-making, when using  
162 existing evidence to answer specific research questions. Asking a specific, focused question  
163 is a fundamental step in the systematic review process. Systematic review questions in  
164 environmental health are generally characterised in terms of the population, exposure,  
165 comparator and outcomes of concern - the PECO mnemonic (Morgan *et al.*, 2018).

166 One of the principal reasons for characterising environmental health questions and the  
167 objectives of systematic reviews in terms of a PECO statement is to facilitate unambiguous

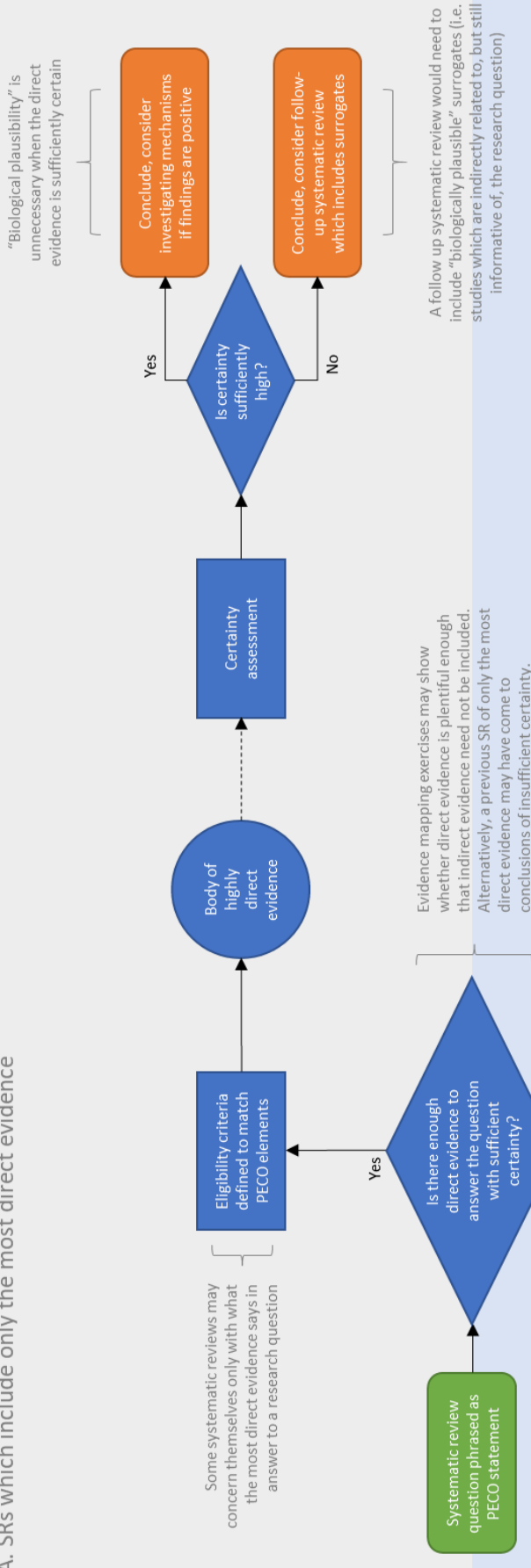
168 characterisation of the types of studies which will be considered by the authors of a  
169 systematic review to be relevant or eligible for answering their question. Studies which are  
170 more directly relevant will be of populations, exposures, comparators and outcomes that  
171 closely match the PECO of the systematic review; those which are less relevant will match  
172 less closely. This concept of fit between a study and the objectives of a systematic review is  
173 “external validity” - the extent to which the findings of a study can be generalised to  
174 populations, exposures and outcomes outside the context of that study (Higgins JPT et al.  
175 (eds), 2019). External validity is part of the indirectness domain in GRADE (Schünemann et  
176 al. 2013).

177 When designing a systematic review, authors need to decide on what their cut-off or  
178 threshold for external validity is going to be, i.e. where they draw the line on a study being  
179 sufficiently generalisable to their target PECO to be worth including in their review. Where  
180 the line is drawn will depend on the review objectives. The way to keep a systematic review  
181 relatively small and simple is to define as eligible only those studies whose designs most  
182 directly match the PECO characterisation of the systematic review question (see Figure 2A).  
183 If the evidence from those studies is sufficiently certain then there is no need to seek out  
184 other evidence in support of the findings of the systematic review - a search for indirect  
185 evidence need not be undertaken.

186 A classic example of this scenario, of direct evidence being of high certainty, is smoking  
187 causing lung cancer. Several observational studies have investigated doctors (P) who  
188 smoke (E), compared them to doctors who do not smoke (C), and assessed the relative risk  
189 of lung cancer (O) between the two groups. The studies are at relatively low risk of bias,  
190 including confounding; multiple studies of similar design give reasonably consistent results;  
191 they are in a representative population; the overall effect size is reasonably precise; there is  
192 no evidence that publication bias exaggerates the observed effect size; there is a dose-  
193 response relationship; and the effect size is large, with smoking increasing lung cancer risk  
194 by a factor of 12-24 (Doll *et al.*, 2005; Pope *et al.*, 2011). These features of the evidence  
195 establish with sufficiently high certainty that a causal relationship has been observed, without  
196 knowledge of the mechanism by which the exposure causes the outcome.

197 In such scenarios, the “biological plausibility” of the exposure-outcome relationship does not  
198 need to be evaluated - it can be assumed that there must be a discoverable biological  
199 mechanism because there is high certainty that the relationship is causal. This is true even  
200 when there is little information about the biological mechanism by which the exposure  
201 causes its outcome. Conversely, that it is not known why or how the exposure causes the  
202 outcome does not undermine certainty that the relationship is causal. This is what we believe  
203 Bradford Hill meant when he stated that establishing biological plausibility is helpful but not  
204 always necessary for a causal claim (Hill, 1965): “It will be helpful if the causation we  
205 suspect is biologically plausible. But this is a feature I am convinced we cannot demand.”

A. SRs which include only the most direct evidence



B. SRs which include evidence from studies of surrogates

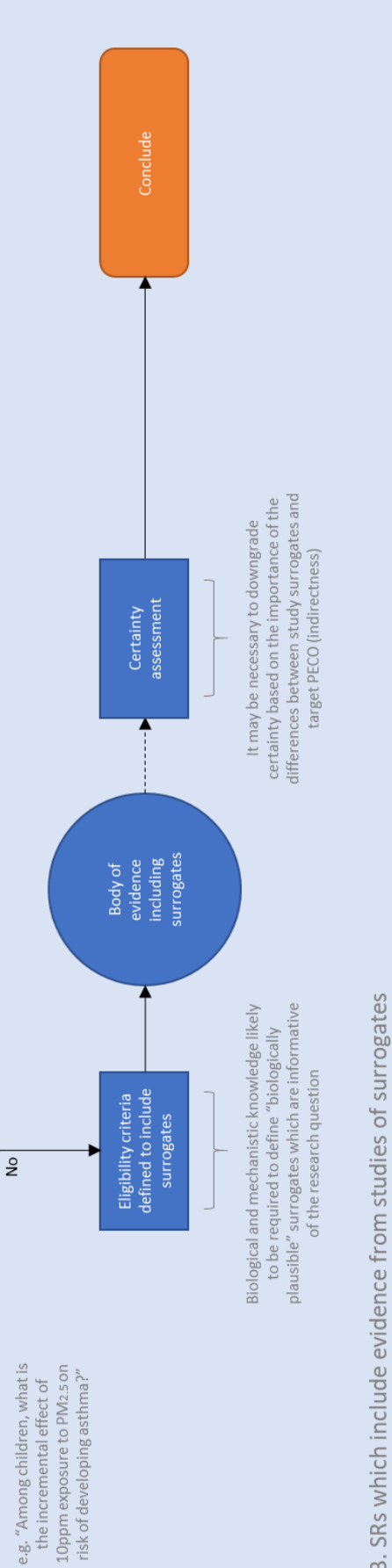


Figure 2. Schematic representation of how it might be decided to include studies of surrogates in a systematic review.



207 The challenge in environmental health research is that high certainty in direct evidence is a  
208 theoretical possibility which is only rarely realised. Usually, environmental health systematic  
209 reviews that focus only on the human evidence for a hypothesised exposure-outcome  
210 relationship would be expected to yield insufficiently conclusive results, due to the human  
211 evidence being highly uncertain or even non-existent. In such circumstances, in order to  
212 further investigate and elucidate potential causal relationships between exposures and  
213 outcomes, it may become necessary to consider indirect evidence in the form of studies of  
214 surrogates (see Figure 2B). This is done in the expectation that including in the systematic  
215 review indirect evidence from studies of surrogates will support an assessment of the  
216 presence of a causal relationship.

217 The use of surrogates is familiar in environmental health contexts, which has long been  
218 reliant on evidence whereby animal models stand in for target human populations,  
219 biomarkers of disease are used in place of observations of clinical health outcomes, and  
220 potential health effects of under-studied chemicals are inferred from their similarity to better-  
221 researched substances. As with any systematic process, decisions on which surrogates to  
222 include in a systematic review should be transparent and well-reasoned, based on evidence  
223 of the validity of the decision, and as far as possible defined in advance of conduct of the  
224 review (Whaley et al. 2020a). Spurious inclusion of surrogate studies is not just a waste of  
225 time and resources: if surrogates are not informative of the question but nonetheless  
226 included in the overall analysis, then the validity of the results of the systematic review may  
227 be compromised; likewise, spurious exclusion of surrogate studies which should have been  
228 included also risks false conclusions.

## 229 The “biological plausibility” of choice of surrogates

230 In conventional environmental health assessments, the consideration of evidence from  
231 surrogates is considered to be justifiable insofar as it provides “biologically plausible” support  
232 for the hypothesised exposure-outcome relationship in the population of concern (European  
233 Food Safety Authority, 2018). In the context of systematic reviews, the GRADE Framework  
234 assesses the importance of the indirectness of the surrogate relative to the question being  
235 asked. Evidence from surrogates which is too indirect would be excluded from a systematic  
236 review; evidence from surrogates which is direct enough to be informative would be included  
237 but might be rated down for indirectness (Guyatt *et al.*, 2011).

238 Here we present 10 examples of the use of surrogates in environmental health  
239 assessments. We frame the examples in terms of biological plausibility and describe how the  
240 indirectness of the surrogates might be interpreted in the GRADE approach. We then use  
241 these examples to show how judgements of biological plausibility map onto the concepts of  
242 systematic review. The ten examples and related analyses are summarised in Table 2 and  
243 Table 3.

244 Note that this is a conceptual article describing how ratings of indirectness may be described  
245 using the GRADE approach. This information should not be used for decision-making. As for  
246 any GRADE concept article, the particular approach described here will require further

247 validation before it may become official GRADE guidance. Thus, the importance of the  
 248 indirectness of a surrogate when assessed as part of a systematic review, and therefore any  
 249 judgements to exclude or downgrade evidence based on its indirectness, may turn out to be  
 250 different to the judgements which have been made in each of the examples we present.

	Surrogates of higher biological plausibility, for which indirectness is less important	Surrogates of lower biological plausibility, for which indirectness is more important
<b>Population</b>	Animal models for human carcinogenicity of 2-nitropropane	Rat models for human bladder carcinogenicity of saccharine
<b>Exposure (dose)</b>	Extrapolating from high doses to low doses of genotoxic substances	Extrapolating from high doses to low doses of endocrine disrupting chemicals
<b>Exposure (route)</b>	Oral administration of bisphenol-A via gavage, or availability of a pharmacokinetic model to translate intravenous dose to oral equivalent	Intravenous administration of bisphenol-A in absence of pharmacokinetic model to translate intravenous dose to oral equivalent
<b>Exposure (substance)</b>	Inferring estrogenic potential of other bisphenols and from studies of bisphenol-A	Inferring neurotoxicity of organophosphate flame retardants from studies of organophosphate pesticides
<b>Outcome</b>	Maternal serum thyroxine (T4) for child neurodevelopmental outcomes	Biomarkers of Alzheimer's Disease progression in place of clinical measures

251

252 **Table 2:** Summary of the 10 examples used in this manuscript to show how discussion of biological plausibility  
 253 maps onto the concepts of systematic review.

## 254 Surrogate populations

255 Toxicology has a long history of use of animal models for investigating potential harm to  
 256 human health from exposure to chemical substances. This is due to the ethical prohibition on  
 257 conducting experiments in humans that are designed to potentially cause harm, combined  
 258 with the need for evidence to inform evaluation of chemical health risks, e.g. for regulatory  
 259 approval and compliance.

260 One example of where surrogate animal and *in vitro* populations are accepted as providing  
 261 evidence for health outcomes in human populations of concern is in the assessment of the  
 262 carcinogenicity of 2-nitropropane. While there is no direct evidence of carcinogenicity in  
 263 humans, animal and *in vitro* evidence is considered to be sufficiently certain to justify  
 264 classifying 2-nitropropane as a human carcinogen (Papameletiou *et al.*, 2017). Although the  
 265 authors did not have a complete account of the mechanism by which 2-nitropropane is a  
 266 genotoxic carcinogen, they judged it sufficiently biologically plausible that observations in  
 267 surrogate experimental populations would also be seen in humans that they felt able to draw  
 268 a conclusion of carcinogenicity.

269 In a contrasting example, the US Food and Drug Administration (FDA) has deemed  
 270 evidence from rat models as not relevant for the assessment of saccharin as a bladder  
 271 carcinogen. This is due to the mechanism by which saccharin causes tumour growth in rats  
 272 not being present in humans (US National Research Council, 2014). The rat model was  
 273 originally considered to be predictive but, once the mechanism by which saccharin causes  
 274 cancer in rats was determined not to be present in humans, the US FDA excluded the rat  
 275 model from assessment. The US FDA judged the hypothesis that the mechanism by which  
 276 saccharin causes cancer in rats also occurs in humans as not biologically plausible.

277 Expressing the reasoning around 2-nitropropane in the concepts of GRADE, it would be to  
278 say that surrogate evidence from animal studies is sufficiently direct to be included in a  
279 systematic review of the human carcinogenicity of 2-nitropropane. A conclusion about  
280 carcinogenicity in humans can be made on the basis of this animal evidence, in spite of its  
281 indirectness and a lack of a complete account of the mechanism of carcinogenicity. For  
282 saccharin it would be to say that, based on the evidence about the underlying mechanism,  
283 the indirectness of the surrogate animal model is unacceptable, i.e. the rat model is too  
284 indirect, does not generalise to humans, and therefore is not eligible for inclusion in a  
285 systematic review of whether saccharin is a bladder carcinogen.

286 We again emphasise that the purpose of these examples is to survey how judgements of  
287 biological plausibility map onto the concepts and processes of GRADE and systematic  
288 reviews. We are not validating any of the judgements that have been made by others in the  
289 selected examples. The purpose of the examples is to understand what is involved when  
290 researchers are making judgements of biological plausibility, not to determine whether those  
291 judgements are valid.

## 292 Surrogate outcomes

293 Surrogate outcomes are used in environmental health research because it is often easier or  
294 more ethical in experimental and observational studies to measure biomarkers of disease  
295 than clinical outcomes of interest. This is the case when health outcomes may have long  
296 latency periods in the population of concern (such as for many cancer types), for particular  
297 study designs (e.g. the use of *in vivo* models for allergic contact dermatitis that focus on the  
298 induction phase only), or when the observed population may not manifest the apical  
299 outcome of interest (e.g. when non-animal test methods address downstream key events  
300 relating skin sensitisation).

301 One example of the use of a surrogate outcome is in a systematic review of the  
302 developmental and reproductive toxicity of the biocide triclosan by Johnson et al. (2016). In  
303 this case, serum thyroxine concentrations in pregnant women were chosen as a surrogate  
304 for the neurodevelopmental health of children. The authors' reasoning was that maternal  
305 thyroid hormone levels during pregnancy are predictive of the subsequent  
306 neurodevelopmental health of the child - an association described in another systematic  
307 review as being "biologically plausible" (Thompson et al. 2018). This can be taken as a  
308 judgment by the authors that there is a sufficiently "biologically plausible" relationship  
309 between maternal serum thyroxine and neurodevelopment that the former can be treated as  
310 a surrogate outcome for the latter.

311 In contrast, a systematic review of biomarkers for Alzheimer's Disease found insufficient  
312 evidence to be able to recommend any biomarker for use as a surrogate outcome for  
313 disease progression (McGhee *et al.*, 2014). While it might appear to be "biologically  
314 plausible" that Alzheimer's Disease results in specific changes to physical brain structure  
315 detectable in an MRI scan (Downey et al. 2017), there seems to be a lack of empirical  
316 evidence that directly connects the biomarker to the outcome of concern.

317 Expressing the triclosan example in the concepts of the GRADE Framework, it would be to  
318 say that, in spite of their indirectness, studies which investigate the surrogate outcome can  
319 be considered eligible for inclusion in a systematic review of neurodevelopmental toxicity  
320 and contribute towards its findings (note that Johnson et al. (2016) used a modified version  
321 of the GRADE Framework and did not downgrade for indirectness). For Alzheimer's disease,  
322 it would be to say that due to uncertainty around how the surrogate biomarker predicts the  
323 ultimate outcome of concern, studies of brain structure biomarkers should either be  
324 downgraded more than once for indirectness or excluded from a systematic review if the  
325 indirectness is judged to be unacceptable.

## 326 *Surrogate exposures*

327 Selecting and attributing appropriate importance to surrogate exposures is a complex issue  
328 in environmental health systematic reviews. We briefly discuss three aspects of surrogate  
329 exposure: route of exposure; administered dose; and active substance. These should  
330 provide sufficient illustration of principle, although we note that other aspects of exposure  
331 such as measurement of metabolites vs. parent compound, timing of exposure, and other  
332 issues, will need consideration in environmental health systematic reviews (Cohen Hubal et  
333 al., 2020).

### 334 *Route*

335 Extrapolating from experimental routes of exposure to the actual routes of exposure likely to  
336 be encountered by target populations is a major preoccupation of toxicological risk  
337 assessment. For example, toxicology studies which administer bisphenol-A (BPA) to animal  
338 test subjects via oral gavage are considered to be of direct relevance to assessing outcomes  
339 from dietary exposure. In contrast, intravenous (IV) administration of BPA is typically  
340 considered not to be relevant to such assessment, due to the avoidance of first-pass  
341 metabolism in the liver (European Food Safety Authority, 2015). However, the relevance of  
342 studies using IV administration can increase if knowledge of how BPA is metabolised allows  
343 equivalent oral doses to be calculated from IV doses, as this provides what can be  
344 interpreted as a "biologically plausible" account of how the two doses are related (Taylor,  
345 Welshons and Vom Saal, 2008).

346 Expressing this in the conceptual framework of GRADE, we would say the indirectness of  
347 the route of exposure becomes less important when the exposures of concern can be  
348 determined from surrogate exposure routes. Physiologically-based pharmacokinetic (PBPK)  
349 models to aid in route-to-route extrapolation are encouraged in chemical assessments  
350 (Meek et al., 2013; US Environmental Protection Agency, 2002). The availability of such  
351 models may lead to indirect evidence from studies using IV exposure routes being included  
352 in a systematic review and potentially rated down fewer levels for indirectness than for  
353 scenarios in which such models are unavailable.

## 354 *Dose*

355 In toxicological research, experiments are often conducted using high doses that are not  
356 considered environmentally or occupationally relevant. Many bioassays also merely aim at  
357 identifying a maximum tolerated dose of a chemical substance in order to provide a  
358 benchmark of toxicity. High dose regimens can raise critical concerns about the indirectness  
359 of a study, if the toxicokinetic and toxicodynamic factors by which the administered dose  
360 causes an outcome are different from those operating at the dose level of concern (Slikker et  
361 al. 2004).

362 This is a key point of debate about the potential health effects of exposure to endocrine  
363 disrupting chemicals: if the administered high dose overwhelms the biological pathway that  
364 is involved in the endocrine activity of the active substance, triggering nonspecific pathways  
365 that are responsible for the observed outcomes, then there may be critical concerns about  
366 the indirectness of the surrogate dose for determining whether the chemical of concern is an  
367 endocrine disruptor (Lagarde *et al.*, 2015). An example of this is the causing of endocrine  
368 effects via direct damage to the liver (Marty et al. 2018). This would lead to concern that  
369 disease induction at high doses via endocrine disruption is not a biologically plausible  
370 mechanism, due to differences between the mechanism by which the dose of concern  
371 causes the outcome of interest as compared to the mechanism by which the surrogate dose  
372 causes the outcome.

373 In contrast, chemicals which cause cancer by a genotoxic mechanism are considered to  
374 operate according to the same mechanism of action at high and low doses (Crump, 1996). In  
375 this case, extrapolation from across the dose range is taken to be unproblematic.

376 Expressing the example of endocrine disruption and genotoxicity in the concepts of the  
377 GRADE Framework, the indirectness of a surrogate dose becomes more important when  
378 there is evidence of different toxicokinetic and toxicodynamic processes operating at  
379 different dose levels. The presence of such differences may result in a decision to exclude  
380 evidence from studies using surrogate doses in a systematic review because of  
381 unacceptable levels of indirectness. If, on the other hand, it is decided to include the studies  
382 that use surrogate doses, rating down for indirectness by two or more levels would be more  
383 likely in the example of endocrine disruptors than for genotoxic carcinogens.

## 384 *Substance*

385 There are many chemicals to which people are potentially exposed which have very few  
386 associated toxicology studies. One means for anticipating the potential toxicity of under-  
387 studied substances is by extrapolation from evidence of the toxicity of suitably similar  
388 chemicals. Often this is based on the demonstration of a common mode of action of toxicity,  
389 or sufficient likeness of the surrogate chemical in terms of physical properties that a common  
390 mode of action can reasonably be inferred.

391 For example, the UK Committee on Toxicity (COT) recently evaluated evidence of the  
392 neurotoxicity of organophosphate flame retardants (OPFRs) (UK Committee on Toxicity,

393 2019). Part of their assessment concerned whether the neurotoxicity of OPFRs could be  
394 extrapolated from studies of the neurotoxicity of organophosphate pesticides (OPPs). COT  
395 determined that OPPs are not a good surrogate exposure for OPFRs, because OPFRs do  
396 not inhibit acetylcholinesterase to the same degree as OPPs. COT concluded that there is  
397 no “biologically plausible” explanation for how OPPs and OPFRs can cause the same effect,  
398 and therefore determined that conclusions about the neurotoxicity of OPFRs should not be  
399 derived from evidence of the neurotoxicity of OPPs.

400 In contrast, since the phase-out of consumer uses of the plastic additive bisphenol-A due to  
401 concerns about its potential to act as an oestrogen, considerable research has been  
402 conducted into whether replacements such as bisphenol-AF and bisphenol-C may have  
403 similar estrogenic potential. Enough similarities in biological effects have been observed for  
404 some researchers to suggest that, at least as a group, exposure to some bisphenols may be  
405 predictive of the effects of exposure to others (Pelch *et al.*, 2019). Similar suggestions have  
406 been made for polyfluorinated compounds (Cousins *et al.*, 2020). When it is more  
407 “biologically plausible” that different chemical substances share the same mechanisms by  
408 which they exert health effects, then it might be acceptable to use one as a surrogate (also  
409 referred to in the environmental health field as an “analogue”) for the other.

410 Expressing the example of OPFRs in the concepts of the GRADE Framework, the absence  
411 of explanation for a shared mechanism by which OPPs and OPFRs would exert a neurotoxic  
412 effect increases the indirectness of OPFRs as a surrogate for OPPs. If the level of  
413 indirectness is unacceptable, it would lead to studies of OPFRs being excluded from a  
414 systematic review of their neurotoxicity; if very high, evidence from the surrogate exposure  
415 might be included but would be rated down for indirectness, potentially two or three times.  
416 For bisphenols and polyfluorinated compounds, if indirectness of surrogates is deemed less  
417 important, they may be included in a systematic review and rated down only once for  
418 indirectness, or possibly not at all.

## 419 Discussion

### 420 Biological plausibility as a dual-aspect concept

421 Our examples show that “biological plausibility” is a concept that can be deployed in multiple  
422 scenarios in environmental health assessments. In general, judgements of biological  
423 plausibility seem to support judgements of causality insofar as studies of causal relationships  
424 in surrogates can be generalised to the target populations, exposures and outcomes of  
425 actual concern. These uses extend beyond the definitions of biological plausibility as  
426 provided by Bradford Hill and Last’s Dictionary of Epidemiology, which define biological  
427 plausibility exclusively in terms of biological explanations of a causal relationship between  
428 exposure and outcome (see Table 1). When translated into the conceptual underpinnings of  
429 GRADE, the uses centre on judging the indirectness of surrogates and describing the impact  
430 on certainty in the evidence for the effect an exposure has on a health outcome in a  
431 population of concern. These judgements not only govern decisions about the eligibility of

432 surrogates for a systematic review but also the extent to which a body of evidence based on  
433 those surrogates should be downgraded for indirectness.

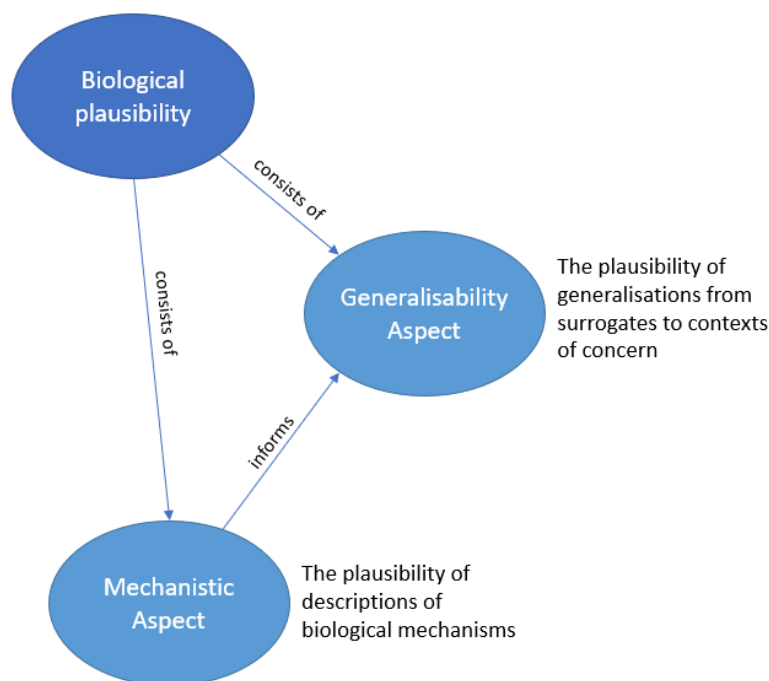
434 Our examples also demonstrate that the judgements being made when assessing  
435 indirectness are complex. Not only do judgements need to be made about the  
436 generalisability of a surrogate, there are also judgements that need to be made about  
437 certainty in descriptions of biological mechanism. These judgements are intrinsically  
438 connected, as absence of mechanistic explanation limits the ability to generalise from a  
439 study surrogate to a target context of concern. This is perhaps most clearly illustrated in the  
440 examples of considering whether studies of OPPs are relevant to characterising the potential  
441 neurotoxicity of OPFRs, and whether studies of rats are relevant to characterising saccharin  
442 as a bladder carcinogen in humans.

443 Based on these observations, we posit that the concept of “biological plausibility” in fact  
444 consists of two principle aspects. We call these the “generalisability aspect” and the  
445 “mechanistic aspect”.

446 We define the *generalisability aspect* of biological plausibility as concerning the validity of  
447 generalisations from a surrogate population, exposure, comparator or outcome to a target  
448 population, exposure, comparator or outcome of concern, respectively. The generalisability  
449 aspect is not about the plausibility of causal claims about the effect of exposures on  
450 outcomes, but instead about the extent to which an observation in a surrogate population  
451 plausibly generalises to a target population, a surrogate exposure generalises to a target  
452 exposure, etc. The generalisability aspect supports judgements of causality insofar as  
453 observations are made in studies of surrogates, and the surrogates then generalise to the  
454 target contexts of concern.

455 We define the *mechanistic aspect* as concerning certainty in biological mechanism. Our  
456 examples show that judgements of whether a surrogate plausibly generalises to a target  
457 context are informed by knowledge of relevant biological mechanisms, i.e. how an exposure  
458 causes an outcome in a given biological or experimental system. While this knowledge is not  
459 often available, when it is, it has a significant impact on judgements about the  
460 generalisability of observations in a surrogate: the higher is the certainty in the knowledge of  
461 relevant biological mechanisms (e.g. that similar mechanisms are present in humans and  
462 surrogate animal species), the higher is the certainty that a generalisation from a given  
463 surrogate to a target context is valid or not. The mechanistic aspect informs the  
464 generalisability aspect, as knowledge of mechanism helps determine the validity of  
465 generalising from surrogates to target contexts of concern.

466 These two aspects are different but fundamentally linked: judgements of the plausibility of  
467 generalisations are informed by judgements of the plausibility of mechanisms. This  
468 connection is illustrated in Figure 3.



469

470 **Figure 3.** The relationship between the generalisability and mechanistic aspects of biological plausibility.

## 471 Biological plausibility and GRADE

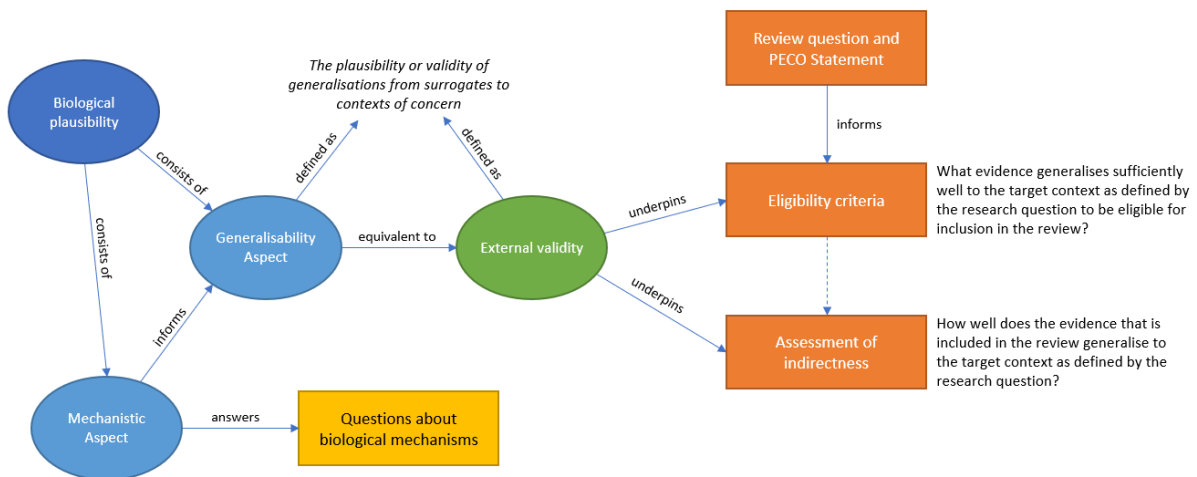
### 472 *How biological plausibility is accommodated within GRADE*

473 Both the generalisability and mechanistic aspects of biological plausibility can be  
 474 accommodated in the indirectness domain of GRADE. This is because the concept of  
 475 “generalisability” is the same as the concepts of external validity and indirectness of  
 476 evidence already familiar in systematic reviews, i.e. the extent to which the results of an  
 477 experimental or observational study apply to a target context outside of that study (Higgins et  
 478 al. 2019; Schünemann et al. 2013). The difference is in vocabulary, whereby systematic  
 479 reviewers talk about the “validity” rather than “plausibility” of a generalisation. Since the  
 480 generalisability aspect of biological plausibility is asking the same question as the  
 481 assessment of external validity in systematic reviews, and external validity is subsumed  
 482 under the GRADE domain of indirectness, it follows that there is no need to extend GRADE  
 483 to accommodate the generalisability aspect of biological plausibility.

484 The mechanistic aspect of biological plausibility, because it informs judgements of  
 485 indirectness or generalisability, is also logically positioned under the indirectness domain of  
 486 GRADE. This relationship is shown in Figure 4. This means GRADE does not need an  
 487 additional domain to accommodate assessment of certainty in biological mechanisms.  
 488 However, this is an interesting category of question for which systematic methods have only  
 489 recently begun to be explored (Whaley et al. 2020) and we recommend further research on  
 490 this issue.



491 Conceptualising matters in this way suggests an operational definition of biological  
 492 plausibility that maps the concept onto the GRADE framework, as follows: “Biological  
 493 plausibility is a dual-aspect concept operationalised in systematic reviews as (1) the validity  
 494 of generalisations from studies of surrogates to target contexts of concern, and (2) certainty  
 495 in biological mechanisms. Certainty in biological mechanisms informs judgement of the  
 496 validity of generalisations. When knowledge of biological mechanisms is available, it can  
 497 have significant impact on judgements of the validity of generalisations.”



498

499 **Figure 4.** How biological plausibility maps onto the processes of systematic review via the shared concept of  
 500 external validity, included in GRADE’s indirectness domain. While questions about biological mechanisms (e.g.  
 501 how an exposure causes an outcome) are independent of a given systematic review, answers to those questions  
 502 can be highly informative in judging the external validity or indirectness of evidence.

503 *How biological plausibility is accommodated within the processes of systematic review*

504 We have established that biological plausibility maps onto judgements of the generalisability  
 505 or indirectness of evidence in a systematic review. These judgements are informed by  
 506 certainty in biological mechanisms. Our task now is to clarify when these judgements are  
 507 made in the systematic review process. This will show how biological plausibility, in the form  
 508 of judgements of generalisability informed by knowledge of biological mechanisms, is  
 509 accounted for when conducting a systematic review. The general principles are articulated  
 510 below and the specific steps described in Box 1.

511 Judgements of generalisability or indirectness occur at two stages in systematic reviews, as  
 512 illustrated in Figure 2. The first stage is in the formulation of eligibility criteria for the inclusion  
 513 of evidence in a systematic review. Here, the authors decide what study designs are  
 514 sufficiently generalisable to their question to be worth including in their systematic review.  
 515 These criteria may be narrowly defined around the most direct evidence, if the authors are  
 516 attempting to keep their review focused and/or they are confident that looking only at the  
 517 most direct evidence will provide sufficiently conclusive results. Otherwise, the eligibility  
 518 criteria may be quite broadly defined, if the authors consider indirect evidence to be of value  
 519 for their review objectives. Even then, there will be limits to eligibility, as many studies will be  
 520 so irrelevant to the objective that it would be a waste of time and resources to include them.

521 Setting these limits, i.e. judging what sort of study designs are informative enough of the  
522 research question to be worth including in the systematic review, is a judgement of  
523 acceptable indirectness. Mechanistic data, if available, may be of high value in making these  
524 judgements.

525 The second stage is in the judgement of indirectness of the evidence that has been included  
526 in the review, when the authors are determining certainty in the evidence on which their  
527 results are based. As we show with our example of smoking and lung cancer, in systematic  
528 reviews of very direct evidence indirectness is trivial and assessment of biological plausibility  
529 unnecessary - the generalisability of findings is a given and mechanistic information is not  
530 needed to support judgements of certainty when certainty is already high.

531 The situation is different for systematic reviews with broadly-defined eligibility criteria.  
532 Indirectness in a systematic review with broadly-defined eligibility criteria rapidly becomes  
533 very important (due to potentially significant differences between target and surrogate) and  
534 complex (due to there being numerous potential differences between the characteristics of  
535 the question as formulated in the PECO vs. the included studies). The generalisability of  
536 surrogates to the target context of concern is a non-trivial issue and needs to be carefully  
537 evaluated. Information about mechanisms is of high value in making these judgements.

**How the assessment of biological plausibility is operationalised in systematic reviews of the health effects of environmental exposures which use the GRADE approach for assessing certainty in the evidence**

1. Define the systematic review question as a PECO statement: "In population P, what effect does exposure E have on outcome O in comparison to comparator C?"
2. Define as ineligible study models that do not sufficiently generalise to the scenario described in the research question (the generalisability aspect of biological plausibility). These judgements may be informed by knowledge of biological mechanisms (the mechanistic aspect of biological plausibility).
3. Determine the effect of the exposure on the outcome in the studies included in the systematic review. This may require studies to be grouped by design characteristics.
4. Evaluate how well the included evidence generalises to the situation described in the research question for each element of the PECO statement (generalisability aspect). These judgements may be informed by knowledge of biological mechanisms (mechanistic aspect).
5. If it is not certain that the evidence generalises to the research question, rate down the evidence one or more times for indirectness depending on the level of this uncertainty.

539 **Box 1.** Explanation of how the concept of biological plausibility is operationalised in the conduct of a systematic  
540 review and assessment of certainty in the evidence using the GRADE Framework. Step 2 is the first place where  
541 judgements of indirectness may be made. Here, high certainty that an indirect study model does not generalise to  
542 the research question may lead to such models being excluded. The US FDA exclusion of rat models from  
543 assessments of the bladder carcinogenicity of saccharin is an example of this. Otherwise, systematic reviews of  
544 health effects of environmental exposures will likely include indirect study models. (The exception is for  
545 deliberately narrowly-focused reviews, as illustrated in Figure 2.) Indirectness judgements are next made in Steps  
546 4 and 5, where the included evidence is assessed under the GRADE domain of indirectness. The UK COT  
547 analysis of neurotoxicity of OPFRs based on neurotoxicity of OPPs is an example of when a judgement of lack of  
548 certainty in shared biological mechanism results in evidence effectively being rated down for certainty due to  
549 indirectness. We note that regulatory frameworks tend to assume a high level of generalisability of a surrogate  
550 model unless there is a high level of evidence to the contrary.

### 551 *Research requirements: judging indirectness at the level of individual studies*

552 The 10 examples in this manuscript show that judgements of external validity are complex  
553 and potentially need to be made across multiple related domains of population, exposure,  
554 comparator, outcome, and subdomains thereof. Instruments which would facilitate  
555 transparent, consistent, and accurate judgements across these domains are not yet  
556 available for study-level judgements of indirectness in environmental health and should be  
557 developed.

558 Answering questions about biological mechanisms draws on a wide variety of information  
559 about the absorption, distribution, metabolism and excretion (“ADME”) of chemical  
560 substances, knowledge of mechanisms by which chemicals cause outcomes in both target  
561 and observed populations, information about interaction between chemicals and target sites,  
562 and the extent to which biomarkers of disease are predictive of clinical outcomes, among  
563 many other issues. If mechanistic knowledge is informative of judgements of external  
564 validity, and therefore of the indirectness domain in GRADE, it follows that we need to  
565 develop methods for assessing certainty in mechanistic knowledge.

566 Assessing certainty in biological mechanisms would be an important and interesting  
567 extension of the GRADE indirectness domain. Unlike questions about associations which  
568 are of the form “is X associated with Y?”, questions about mechanisms are of the form “how  
569 does X cause Y?”. Answering this form of question involves describing sequences of  
570 biological events, one of which is associated with the next. In principle, event-event  
571 associations should be approachable in the same way as exposure-outcome associations,  
572 and therefore be amenable to the GRADE approach. A particular challenge we can foresee  
573 is in handling the sheer volume of data involved in systematically assessing multiple  
574 associated biological events, if very indirect evidence is permitted to enter into the  
575 assessment (Whaley et al. 2020b).

576 We note that developments in the Adverse Outcome Pathway framework may be informative  
577 for operationalising the assessment of certainty in biological mechanisms and interpreting  
578 indirectness of evidence in systematic reviews (de Vries et al. 2021). Alternatively, the Key  
579 Characteristics framework may also provide a structured approach to assessing indirectness  
580 via similarity of biological mechanisms (Smith et al. 2016; Guyton et al. 2018).

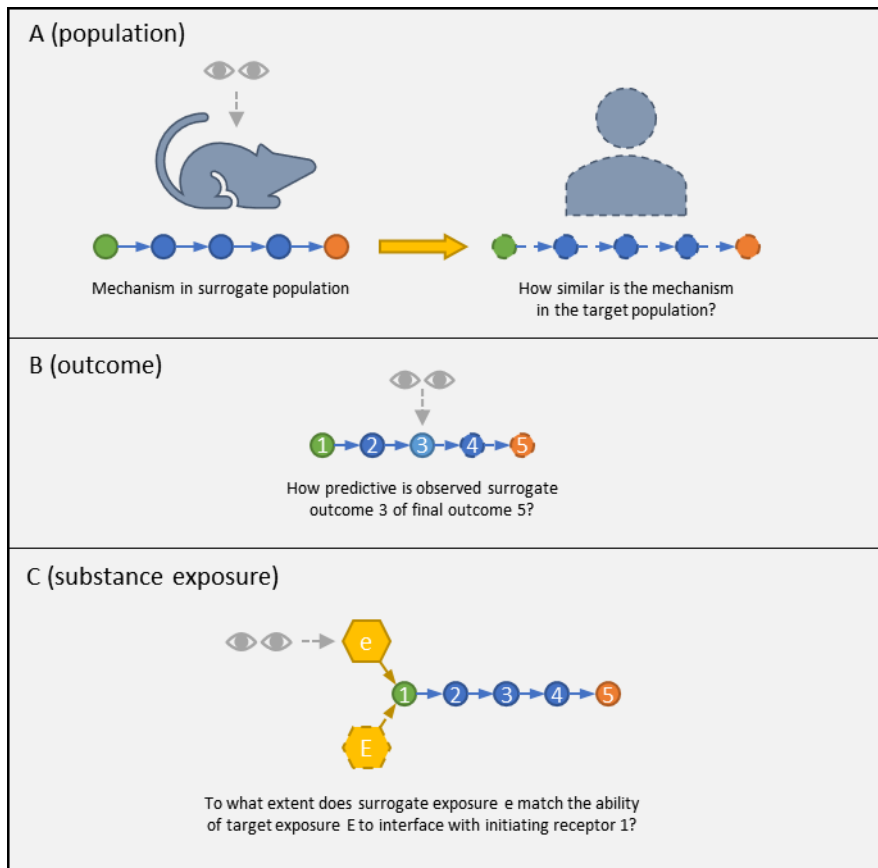
581 The 10 examples discussed above give us indications of some the considerations which  
 582 may reduce concerns about the indirectness of a study included in a systematic review.  
 583 These are outlined in Table 3 and illustrated, where feasible, in Figure 5. While these are  
 584 only suggestive selections from the examples we have used in this manuscript, they do  
 585 illustrate how much of this discussion is already familiar in toxicology and environmental  
 586 health. This experience should provide a robust platform for further research and should  
 587 draw on the experience of GRADE and the environmental health communities.

Potential influencing factors in judging the biological plausibility or external validity of study surrogates	
<b>Population</b>	The extent to which the biological pathway connecting exposure to outcome is operating in both the surrogate population and the target population (Figure 5A)
<b>Exposure – dose</b>	The similarity of the toxicodynamic and toxicokinetic processes by which the surrogate dose acts in comparison to that of the dose range of interest
<b>Exposure – route</b>	The similarity by which an organism absorbs and metabolises the substance of concern via the surrogate route as opposed to the target route; or the reliability with which exposure from the surrogate route can be transformed to values which match exposure from the route of interest
<b>Exposure – substance</b>	The extent to which the surrogate molecule influences the biological processes by which the target molecule is thought to elicit its biological effects (Figure 5C)
<b>Outcome</b>	The extent to which a surrogate outcome is predictive of the target outcome of concern (Figure 5B)

588

589 **Table 3:** Summary of potential influencing factors in judging biological plausibility or external validity of study  
 590 surrogates, as suggested by the examples in this manuscript

591 Finally, we note that sufficient biological knowledge to permit high-certainty judgements of  
 592 mechanism and external validity, and thus avoiding rating down for indirectness (and,  
 593 conversely, being certain that evidence from a surrogate is not relevant), is rare. Absent  
 594 explanations of mechanism, evidence would either (a) end up being excluded from a  
 595 systematic review because there is no theoretical route (apart from presumption of  
 596 relevance) to considering it as eligible, or (b) evidence would be included but its external  
 597 validity would be unclear, indirectness higher as a consequence, and certainty lower overall.  
 598 In these cases of low certainty due to unclear external validity of the included studies,  
 599 significant mechanistic research may be required before it is possible to determine whether  
 600 one experimental model is more externally valid than another. Currently, in regulatory  
 601 circumstances where making decisions in the face of uncertainty is important, and  
 602 mechanistic evidence to support judgements of external validity in a health assessment is  
 603 limited, the external validity of a choice of surrogate is often assumed unless there is  
 604 compelling evidence to the contrary (US Environmental Protection Agency, 2005).



605

606 **Figure 5:** Illustrations of the potential influencing factors in judging biological plausibility or external validity of  
 607 study surrogates, as suggested by the examples in this manuscript. Circles represent individual events in a  
 608 biological pathway by which activation of a receptor (green) results in a health outcome (orange). Eyes indicate  
 609 what surrogate is being observed in an experimental model in lieu of the target of concern (dashed outlines).

610 **Limitations**

611 The attentive reader of our source material will notice that the concept of biological  
 612 plausibility is rarely clearly applied, even when it appears that it is being discussed. This  
 613 phenomenon has been observed by other researchers (Dailey, Rosman and Silbergeld,  
 614 2018). We have therefore had to impute the concept of biological plausibility to some of our  
 615 examples - particularly for Alzheimer's Disease, bisphenols, and neurodevelopment - based  
 616 on surrounding literature and our general understanding of how discussion of biological  
 617 plausibility is conducted. We believe our imputation to be consistent with use of the concept  
 618 and intent of the source material, and it is anyway not necessary for the specific term  
 619 "biological plausibility" to have been used for the concept to have been applied. While a  
 620 greater number of direct examples could be gathered from a systematic survey of the use of  
 621 the concept of biological plausibility in the literature, we do not expect that they would  
 622 invalidate our argument.

623 **Conclusion**

624 We asked what sort of "biomedical", "biological", or "epidemiological" knowledge may  
 625 influence certainty in the evidence of a systematic review of an exposure-outcome

626 relationship. Our answer is knowledge of biological mechanisms, which informs judgements  
627 of the indirectness of a body of evidence constituted from studies of surrogates. We also set  
628 out to determine whether biological plausibility, when applied as a concept relating to  
629 certainty in the evidence for the findings of a systematic review, is accommodated by the  
630 GRADE domain of indirectness. We have argued that it is, although its full operationalisation  
631 will require additional study.

632 In answering these questions, we have elucidated Bradford Hill's proposition that  
633 establishing biological plausibility is helpful but not always necessary for a causal claim (Hill,  
634 1965): "It will be helpful if the causation we suspect is biologically plausible. But this is a  
635 feature I am convinced we cannot demand." We have shown that biological plausibility is  
636 indeed not necessary for determining that an exposure causes an outcome, so long as the  
637 direct evidence for the exposure-outcome relationship is sufficiently certain. The presence of  
638 "biological plausibility" can nonetheless be "helpful" to establishing causation. This happens  
639 when sufficient information about mechanisms is available to characterise the  
640 generalisability of a surrogate, thereby supporting judgements about indirectness of the  
641 evidence and potentially permitting a more certain answer to the review question than would  
642 be yielded by inclusion of the most direct evidence alone.

643 Our analysis also broadens the scope of discussion in GRADE of study surrogates.  
644 Currently, GRADE guidance only explicitly addresses surrogate outcomes (Guyatt *et al.*,  
645 2011): "Guideline developers should consider surrogate outcomes only when high-quality  
646 evidence regarding important outcomes is lacking. When such evidence is lacking [...] they  
647 should specify the important outcomes and the associated surrogates they must use as  
648 substitutes. [...] the necessity to substitute the surrogate may ultimately lead to rating down  
649 the quality of the evidence because of indirectness." Here, we have extended discussion of  
650 eligibility and potential grading of surrogate outcomes to also cover surrogate populations  
651 and surrogate exposures.

652 We have argued that judgements of biological plausibility, at least in their application to  
653 determining the relevance of evidence to answering a focused research question, are  
654 accommodated under the operational procedures of systematic review and the GRADE  
655 domain of indirectness. While vocabulary and processes may differ, we feel confident that  
656 there is nothing in biological plausibility that, for this context, is "missing" from GRADE. What  
657 is needed, however, are means to operationalise the assessment of the indirectness of  
658 included studies and certainty in evidence for biological mechanisms, the outputs of which  
659 can be used in determining the extent to which evidence should be rated down for  
660 indirectness. Such methods would help bring shape to the amorphous nature of mechanistic  
661 evidence and aid in its exploitation in environmental health systematic reviews.

662 As a final point, we observe a clear parallel between the clinical and public health contexts in  
663 which GRADE was developed and the environmental health context in which it is here being  
664 applied. The difference is that in clinical contexts, GRADE is nearly always used to evaluate  
665 human evidence where treatments are being trialled in people, far downstream from the pre-  
666 clinical *in vitro* and animal research that is used to justify conducting a human trial. While

667 treatments are advanced to human trials based on evidence from preclinical studies, this  
668 evidence is often many years old by the time a systematic review is conducted - and  
669 therefore preclinical evidence is not needed. In contrast, *in vitro* and *in vivo* research  
670 constitutes in many environmental health contexts most of the evidence being dealt with.  
671 The fundamental principles for systematically reviewing this evidence are no different to  
672 systematic reviews of human evidence, it is just the availability of human evidence that is  
673 more limited and mechanisms are often not known. In the context of environmental health,  
674 GRADE is, therefore, being applied to a more indirect evidence base which is often focused  
675 on events that are further upstream than those dealt with by most healthcare systematic  
676 reviews.

## 677 Acknowledgements

678 The authors would like to thank the GRADE Environmental Health Project Group and  
679 GRADE Working Group for their contributions to this manuscript, and the Evidence-based  
680 Toxicology Collaboration (EBTC) at Johns Hopkins Bloomberg School of Public Health for  
681 providing funding to cover the time of PW, KT and SH in working on this manuscript. We  
682 would like to thank Dr Andrew Kraft and Dr Michelle Angrish for their technical review of the  
683 manuscript, and Dr Daniele Wikoff (ToxStrategies) for detailed discussion of the ideas and  
684 concepts we present here. The authors would also thank the European Food Safety  
685 Authority and EBTC for organising the Scientific Colloquium, and the participants who  
686 contributed to discussions therein, which gave genesis to the concept of this manuscript  
687 (European Food Safety Authority, 2018).

688

---

## 689 Bibliography

- 690 Bilotta, G. S., Milner, A. M. and Boyd, I. (2014) 'On the use of systematic reviews to inform  
691 environmental policies', *Environmental science & policy*, 42, pp. 67–77. doi:  
692 10.1016/j.envsci.2014.05.010.
- 693 Braun, J. M. and Gray, K. (2017) 'Challenges to studying the health effects of early life  
694 environmental chemical exposures on children's health', *PLoS biology*, 15(12), p. e2002800.  
695 doi: 10.1371/journal.pbio.2002800.
- 696 Burns, J. *et al.* (2020) 'Interventions to reduce ambient air pollution and their effects on  
697 health: An abridged Cochrane systematic review', *Environment international*, 135, p.  
698 105400. doi: 10.1016/j.envint.2019.105400.
- 699 Cohen Hubal, E. A. *et al.* (2020) 'Advancing systematic-review methodology in exposure  
700 science for environmental health decision making', *Journal of exposure science &  
701 environmental epidemiology*. doi: 10.1038/s41370-020-0236-0.
- 702 Cousins, I. T. *et al.* (2020) 'Strategies for grouping per- and polyfluoroalkyl substances  
703 (PFAS) to protect human and environmental health', *Environmental science. Processes &  
704 impacts*. doi: 10.1039/d0em00147c.
- 705 Crump, K. S. (1996) 'The linearized multistage model and the future of quantitative risk

- 706 assessment', *Human & experimental toxicology*, 15(10), pp. 787–798. doi:  
707 10.1177/096032719601501001.
- 708 Dailey, J., Rosman, L. and Silbergeld, E. K. (2018) "Evaluating biological plausibility in  
709 supporting evidence for action through systematic reviews in public health," *Public health*.  
710 Elsevier BV, 165, pp. 48–57. doi: 10.1016/j.puhe.2018.08.015.
- 711 de Vries et al. (2021) "Applying evidence-based methods to the development and use of  
712 adverse outcome pathways", *ALTEX - Alternatives to animal experimentation*, 38(2), pp.  
713 336-347. doi: 10.14573/altex.2101211
- 714 Doll, R. *et al.* (2005) 'Mortality from cancer in relation to smoking: 50 years observations on  
715 British doctors', *British journal of cancer*, 92(3), pp. 426–429. doi: 10.1038/sj.bjc.6602359.
- 716 Downey, A. *et al.* (2017) 'communicating with the public about interventions to prevent  
717 cognitive decline and dementia', in *Preventing Cognitive Decline and Dementia: A Way  
718 Forward*. National Academies Press (US).
- 719 European Food Safety Authority (2015) 'Scientific Opinion on the risks to public health  
720 related to the presence of bisphenol A (BPA) in foodstuffs: PART II - Toxicological  
721 assessment and risk characterisation', 13(1). doi: 10.2903/j.efsa.2015.3978.
- 722 European Food Safety Authority (2018) 'EFSA Scientific Colloquium 23 – Joint European  
723 Food Safety Authority and Evidence-Based Toxicology Collaboration Colloquium Evidence  
724 integration in risk assessment: the science of combining apples and oranges 25–26 October  
725 2017 Lisbon, Portugal', *EFSA Supporting Publications*, 15(3). doi: 10.2903/sp.efsa.2018.EN-  
726 1396.
- 727 Gauderat, G. *et al.* (2017) 'Prediction of human prenatal exposure to bisphenol A and  
728 bisphenol A glucuronide from an ovine semi-physiological toxicokinetic model', *Scientific  
729 reports*, 7(1), p. 15330. doi: 10.1038/s41598-017-15646-5.
- 730 Guyatt, G. H. *et al.* (2008) 'GRADE: an emerging consensus on rating quality of evidence  
731 and strength of recommendations', *BMJ*, 336(7650), pp. 924–926. doi:  
732 10.1136/bmj.39489.470347.AD.
- 733 Guyatt, G. H. *et al.* (2011) 'GRADE guidelines: 8. Rating the quality of evidence--  
734 indirectness', *Journal of clinical epidemiology*, 64(12), pp. 1303–1310.
- 735 Guyton et al. (2018) Application of the key characteristics of carcinogens in cancer hazard  
736 identification, *Carcinogenesis*, Volume 39, Issue 4, April 2018, Pages 614–622,  
737 <https://doi.org/10.1093/carcin/bgy031>
- 738 Hardy, A. *et al.* (2017) 'Guidance on the use of the weight of evidence approach in scientific  
739 assessments', *EFSA Journal*, 15(8). doi: 10.2903/j.efsa.2017.4971.
- 740 Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, Welch VA (ed.) (2019)  
741 *Cochrane Handbook for Systematic Reviews of Interventions version 6.0 (updated July  
742 2019)*. Cochrane. Available at: [www.training.cochrane.org/handbook](http://www.training.cochrane.org/handbook).
- 743 Hill, A. B. (1965) 'The Environment and Disease: Association or Causation?', *Proceedings of  
744 the Royal Society of Medicine*, 58, pp. 295–300. Available at:  
745 <https://www.ncbi.nlm.nih.gov/pubmed/14283879>.
- 746 Hoffmann, S. *et al.* (2017) 'A primer on systematic reviews in toxicology', *Archives of  
747 toxicology*, 91(7), pp. 2551–2575. doi: 10.1007/s00204-017-1980-3.



748 Hultcrantz, M. *et al.* (2017) 'The GRADE Working Group clarifies the construct of certainty of  
749 evidence', *Journal of clinical epidemiology*, 87, pp. 4–13. doi: 10.1016/j.jclinepi.2017.05.006.

750 International Epidemiological Association (2001) *A Dictionary of Epidemiology*. Oxford  
751 University Press. Available at:  
752 <https://play.google.com/store/books/details?id=nQmhQgAACAAJ>.

753 Johnson, P. I. *et al.* (2016) 'Application of the Navigation Guide systematic review  
754 methodology to the evidence for developmental and reproductive toxicity of triclosan',  
755 *Environment international*. The Authors, 92-93, pp. 716–728. doi:  
756 10.1016/j.envint.2016.03.009.

757 Lagarde, F. *et al.* (2015) 'Non-monotonic dose-response relationships and endocrine  
758 disruptors: a qualitative method of assessment', *Environmental health: a global access  
759 science source*, 14, p. 13. doi: 10.1186/1476-069X-14-13.

760 Lewis, S. J. *et al.* (2017) 'Developing the WCRF International/University of Bristol  
761 Methodology for Identifying and Carrying Out Systematic Reviews of Mechanisms of  
762 Exposure–Cancer Associations', *Cancer epidemiology, biomarkers & prevention: a  
763 publication of the American Association for Cancer Research, cosponsored by the American  
764 Society of Preventive Oncology*, 26(11), pp. 1667–1675.

765 Lin, J. H. (2008) 'CSF as a surrogate for assessing CNS exposure: an industrial  
766 perspective', *Current drug metabolism*, 9(1), pp. 46–59. doi: 10.2174/138920008783331077.

767 McGhee, D. J. M. *et al.* (2014) 'A systematic review of biomarkers for disease progression in  
768 Alzheimer's disease', *PloS one*, 9(2), p. e88854. doi: 10.1371/journal.pone.0088854.

769 Meek, M. E. B. *et al.* (2013) 'Case study illustrating the WHO IPCS guidance on  
770 characterization and application of physiologically based pharmacokinetic models in risk  
771 assessment', *Regulatory toxicology and pharmacology: RTP*, 66(1), pp. 116–129. doi:  
772 10.1016/j.yrtph.2013.03.005.

773 Morgan, R. L. *et al.* (2016) 'GRADE: Assessing the quality of evidence in environmental and  
774 occupational health', *Environment international*. Elsevier Ltd, 92-93, pp. 1–6. doi:  
775 10.1016/j.envint.2016.01.004. Morgan, R. L. *et al.* (2018) 'Identifying the PECO: A framework  
776 for formulating good questions to explore the association of environmental and other  
777 exposures with health outcomes', *Environment international*. Elsevier, (July), pp. 1–5. doi:  
778 10.1016/j.envint.2018.07.015.

779 Morgan, R. L. *et al.* (2019) 'GRADE guidelines for environmental and occupational health: A  
780 new series of articles in Environment International', *Environment international*, 128, pp. 11–  
781 12.

782 OECD (2016) 'Users' Handbook supplement to the Guidance Document for developing and  
783 assessing Adverse Outcome Pathways', *Env/Jm/Mono(2016) 12*, (OECD Series on Adverse  
784 Outcome Pathways1), p. 63. doi: 10.1787/5jlvl1m9d1g32-en.

785 Papameletiou, D. *et al.* (2017) 'SCOEL/REC/300 2-Nitropropane - Recommendation from  
786 the Scientific Committee on Occupational Exposure Limits'. Directorate-General for  
787 Employment, Social Affairs and Inclusion (European Commission) , Scientific Committee on  
788 Occupational Exposure Limits. doi: 10.2767/841951.

789 Pelch, K. E. *et al.* (2019) 'Characterization of Estrogenic and Androgenic Activities for  
790 Bisphenol A-like Chemicals (BPs): In Vitro Estrogen and Androgen Receptors  
791 Transcriptional Activation, Gene Regulation, and Binding Profiles', *Toxicological sciences*:

- 792 an official journal of the Society of Toxicology. doi: 10.1093/toxsci/kfz173.
- 793 Pope, C. A., 3rd *et al.* (2011) 'Lung cancer and cardiovascular disease mortality associated  
794 with ambient air pollution and cigarette smoke: shape of the exposure-response  
795 relationships', *Environmental health perspectives*, 119(11), pp. 1616–1621. doi:  
796 10.1289/ehp.1103639.
- 797 Prozialeck, W. C. (2013) 'Biomarkers for Cadmium', in Kretsinger, R. H., Uversky, V. N., and  
798 Permyakov, E. A. (eds) *Encyclopedia of Metalloproteins*. New York, NY: Springer New York,  
799 pp. 272–277. doi: 10.1007/978-1-4614-1533-6\_33.
- 800 Prozialeck, W. C. and Edwards, J. R. (2010) 'Early biomarkers of cadmium exposure and  
801 nephrotoxicity', *Biometals: an international journal on the role of metal ions in biology,*  
802 *biochemistry, and medicine*, 23(5), pp. 793–809. doi: 10.1007/s10534-010-9288-2.
- 803 Radke, E. G. *et al.* (accepted) 'Application of US EPA IRIS systematic review methods to the  
804 health effects of phthalates: lessons learned and path forward', *Environment international*.
- 805 Rhomberg, L. (2015) 'Hypothesis-Based Weight of Evidence: An Approach to Assessing  
806 Causation and its Application to Regulatory Toxicology', *Risk analysis: an official publication*  
807 *of the Society for Risk Analysis*, 35(6), pp. 1114–1124. doi: 10.1111/risa.12206.
- 808 Schünemann, H. *et al.* (2011) 'The GRADE approach and Bradford Hill's criteria for  
809 causation', *Journal of Epidemiology & Community Health*, 65(5), pp. 392–395. doi:  
810 10.1136/jech.2010.119933.
- 811 Schünemann, H. J. *et al.* (2013) 'Non-randomized studies as a source of complementary,  
812 sequential or replacement evidence for randomized controlled trials in systematic reviews on  
813 the effects of interventions', *Research synthesis methods*, 4(1), pp. 49–62.
- 814 Schünemann, H. J. *et al.* (2019) 'GRADE guidelines: 18. How ROBINS-I and other tools to  
815 assess risk of bias in nonrandomized studies should be used to rate the certainty of a body  
816 of evidence', *Journal of clinical epidemiology*, 111, pp. 105–114.
- 817 Sheehan, M. C. and Lam, J. (2015) 'Use of Systematic Review and Meta-Analysis in  
818 Environmental Health Epidemiology: a Systematic Review and Comparison with Guidelines',  
819 *Current environmental health reports*, 2(3), pp. 272–283. doi: 10.1007/s40572-015-0062-z.
- 820 Slikker, W., Jr *et al.* (2004) 'Dose-dependent transitions in mechanisms of toxicity: case  
821 studies', *Toxicology and applied pharmacology*, 201(3), pp. 226–294.
- 822 Smith *et al.* (2016) 'Key characteristics of carcinogens as a basis for organizing data on  
823 mechanisms of carcinogenesis. *Environ Health Perspect* 124:713–721;  
824 <http://dx.doi.org/10.1289/ehp.1509912>
- 825 Marty, M. S. *et al.* (2018) 'Distinguishing between endocrine disruption and non-specific  
826 effects on endocrine systems', *Regulatory toxicology and pharmacology: RTP*, 99, pp. 142–  
827 158.
- 828 Taylor, J. A., Welshons, W. V. and Vom Saal, F. S. (2008) 'No effect of route of exposure  
829 (oral; subcutaneous injection) on plasma bisphenol A throughout 24h after administration in  
830 neonatal female mice', *Reproductive toxicology*, 25(2), pp. 169–176. doi:  
831 10.1016/j.reprotox.2008.01.001.
- 832 Thompson, W. *et al.* (2018) 'Maternal thyroid hormone insufficiency during pregnancy and  
833 risk of neurodevelopmental disorders in offspring: A systematic review and meta-analysis',

834 *Clinical endocrinology*, 88(4), pp. 575–584.

835 UK Committee on Toxicity (COT) (2019) ‘Statement on phosphate-based flame retardants  
836 and the potential for neurodevelopmental toxicity’. Available at:  
837 [https://cot.food.gov.uk/cotstatements/cotstatementsyrs/cot-statements-2019/cot-phosphate-](https://cot.food.gov.uk/cotstatements/cotstatementsyrs/cot-statements-2019/cot-phosphate-based-flame-retardants-statement)  
838 [based-flame-retardants-statement](https://cot.food.gov.uk/cotstatements/cotstatementsyrs/cot-statements-2019/cot-phosphate-based-flame-retardants-statement).

839 US Environmental Protection Agency (2002) ‘A review of the reference dose and reference  
840 concentration processes’. Report Number EPA/630/P-02/002F. Washington, DC. Available  
841 at: <https://www.epa.gov/sites/production/files/2014-12/documents/rfd-final.pdf>

842 US Environmental Protection Agency (2005) ‘Guidelines for Carcinogen Risk Assessment’.  
843 Available at: <https://www.epa.gov/risk/guidelines-carcinogen-risk-assessment>.

844 US National Research Council (2014) *Review of EPA’s Integrated Risk Information System*  
845 *(IRIS) Process*. The National Academies Press. Available at:  
846 [http://www.nap.edu/openbook.php?record\\_id=18764](http://www.nap.edu/openbook.php?record_id=18764).

847 Whaley, P. *et al.* (2016) ‘Implementing systematic review techniques in chemical risk  
848 assessment: Challenges, opportunities and recommendations’, *Environment international*.  
849 The Authors, 92-93, pp. 556–564. doi: 10.1016/j.envint.2015.11.002.

850 Whaley, P. *et al.* (2020a) ‘Recommendations for the conduct of systematic reviews in  
851 toxicology and environmental health research (COSTER)’, *Environment international*, 143, p.  
852 105926.

853 Whaley, P. *et al.* (2020b) ‘Knowledge Organization Systems for Systematic Chemical  
854 Assessments’, *Environmental health perspectives*, 128(12), p. 125001.

855 Wikipedia contributors (2014) *Biological plausibility*, *Wikipedia, The Free Encyclopedia*.  
856 Available at:  
857 [https://en.wikipedia.org/w/index.php?title=Biological\\_plausibility&oldid=614374435](https://en.wikipedia.org/w/index.php?title=Biological_plausibility&oldid=614374435)  
858 (Accessed: 16 October 2019).

859

860