



LJMU Research Online

Brozek, JL, Canelo-Aybar, C, Akl, EA, Bowen, JM, Bucher, J, Chiu, WA, Cronin, MTD, Djulbegovic, B, Falavigna, M, Guyatt, GH, Gordon, AA, Boon, MH, Hutubessy, RCW, Joore, MA, Katikireddi, V, LaKind, J, Langendam, M, Manja, V, Magnuson, K, Mathioudakis, AG, Meerpohl, J, Mertz, D, Mezencev, R, Morgan, R, Morgano, GP, Mustafa, R, O'Flaherty, M, Patlewicz, G, Riva, JJ, Posso, M, Rooney, A, Schlosser, PM, Schwartz, L, Shemilt, I, Tarride, J-E, Thayer, KA, Tsaion, K, Vale, L, Wambaugh, J, Wignall, J, Williams, A, Xie, F, Zhang, Y, Schünemann, HJ and GRADE Working Group,

GRADE Guidelines 30: The GRADE Approach to Assessing the Certainty of Modelled Evidence - an Overview in the Context of Health Decision-making.

<http://researchonline.ljmu.ac.uk/id/eprint/13778/>

Article

Citation (please note it is advisable to refer to the publisher's version if you intend to cite from this work)

Brozek, JL, Canelo-Aybar, C, Akl, EA, Bowen, JM, Bucher, J, Chiu, WA, Cronin, MTD, Djulbegovic, B, Falavigna, M, Guyatt, GH, Gordon, AA, Boon, MH, Hutubessy, RCW, Joore, MA, Katikireddi, V, LaKind, J, Langendam, M, Mania, V, Magnuson, K, Mathioudakis, AG, Meerpohl, J, Mertz, D, Mezencev.

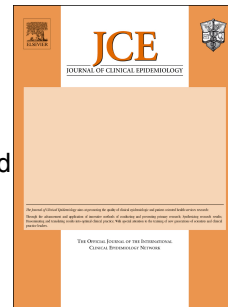
LJMU has developed **LJMU Research Online** for users to access the research output of the University more effectively. Copyright © and Moral Rights for the papers on this site are retained by the individual authors and/or other copyright owners. Users may download and/or print one copy of any article(s) in LJMU Research Online to facilitate their private study or for non-commercial research. You may not engage in further distribution of the material or use it for any profit-making activities or any commercial gain.

<http://researchonline.ljmu.ac.uk/>

The version presented here may differ from the published version or from the version of the record. Please see the repository URL above for details on accessing the published version and note that access may require a subscription.

For more information please contact researchonline@ljmu.ac.uk

Journal Pre-proof



GRADE Guidelines 30: The GRADE Approach to Assessing the Certainty of Modelled Evidence - an Overview in the Context of Health Decision-making

Jan L. Brozek, Carlos Canelo-Aybar, Elie A. Akl, James M. Bowen, John Bucher, Weihsueh A. Chiu, Mark Cronin, Benjamin Djulbegovic, Maicon Falavigna, Gordon H. Guyatt, Ami A. Gordon, Michele Hilton Boon, Raymond C.W. Hutubessy, Manuela A. Joore, Vittal Katikireddi, Judy LaKind, Miranda Langendam, Veena Manja, Kristen Magnuson, Alexander G. Mathioudakis, Joerg Meerpohl, Dominik Mertz, Roman Mezencev, Rebecca Morgan, Gian Paolo Morgano, Reem Mustafa, Martin O'Flaherty, Grace Patlewicz, John J. Riva, Margarita Posso, Andrew Rooney, Paul M. Schlosser, Lisa Schwartz, Ian Shemilt, Jean-Eric Tarride, Kristina A. Thayer, Katya Tsaion, Luke Vale, John Wambough, Jessica Wignall, Ashley Williams, Feng Xie, Yuan Zhang, Holger J. Schünemann, for the GRADE Working Group

PII: S0895-4356(20)31103-3

DOI: <https://doi.org/10.1016/j.jclinepi.2020.09.018>

Reference: JCE 10283

To appear in: *Journal of Clinical Epidemiology*

Received Date: 21 February 2020

Revised Date: 8 September 2020

Accepted Date: 17 September 2020

Please cite this article as: Brozek JL, Canelo-Aybar C, Akl EA, Bowen JM, Bucher J, Chiu WA, Cronin M, Djulbegovic B, Falavigna M, Guyatt GH, Gordon AA, Boon MH, Hutubessy RCW, Joore MA, Katikireddi V, LaKind J, Langendam M, Manja V, Magnuson K, Mathioudakis AG, Meerpohl J, Mertz D, Mezencev R, Morgan R, Morgano GP, Mustafa R, O'Flaherty M, Patlewicz G, Riva JJ, Posso M, Rooney A, Schlosser PM, Schwartz L, Shemilt I, Tarride J-E, Thayer KA, Tsaion K, Vale L, Wambough J, Wignall J, Williams A, Xie F, Zhang Y, Schünemann HJ, for the GRADE Working Group, GRADE Guidelines 30: The GRADE Approach to Assessing the Certainty of Modelled Evidence - an Overview in the Context of Health Decision-making, *Journal of Clinical Epidemiology* (2020), doi: <https://doi.org/10.1016/j.jclinepi.2020.09.018>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier Inc.

Author Statement

All authors analysed and interpreted the data. Jan Brozek and Carlos Canelo-Aybar wrote the first version of the paper. All authors of this paper have read and approved the final version submitted.

Journal Pre-proof

GRADE Guidelines 30: The GRADE Approach to Assessing the Certainty of Modelled Evidence - an Overview in the Context of Health Decision-making

Jan L. Brozek^{# a,b,c}, Carlos Canelo-Aybar^{# d,e}, Elie A. Akl^f, James M. Bowen^{a,g}, John Bucher^h, Weihsueh A. Chiuⁱ, Mark Cronin^j, Benjamin Djulbegovic^k, Maicon Falavigna^l, Gordon H. Guyatt^{a,b,c}, Ami A. Gordon^m, Michele Hilton Boonⁿ, Raymond C. W. Hutubessy^o, Manuela A. Joore^p, Vittal Katikireddiⁿ, Judy LaKind^{q,r}, Miranda Langendam^s, Veena Manja^{a,t,u}, Kristen Magnuson^m, Alexander G. Mathioudakis^v, Joerg Meerpohl^{w,x}, Dominik Mertz^a, Roman Mezencev^y, Rebecca Morgan^a, Gian Paolo Morgano^{a,c}, Reem Mustafa^{a,z}, Martin O'Flaherty^{aa}, Grace Patlewicz^{ab}, John J. Riva^{c,ac}, Margarita Posso^e, Andrew Rooney^h, Paul M. Schlosser^y, Lisa Schwartz^a, Ian Shemilt^{ad}, Jean-Eric Tarride^{a,ae}, Kristina A. Thayer^u, Katya Tsaion^{af}, Luke Vale^{ag}, John Wambough^{ab}, Jessica Wignall^m, Ashley Williams^m, Feng Xie^a, Yuan Zhang^{a,ah}, Holger J. Schünemann^{a,b,c}, for the GRADE Working Group

Co-first author

Affiliations:

^a Department of Health Research Methods, Evidence, and Impact, McMaster University, Hamilton, Ontario, Canada

^b Department of Medicine, McMaster University, Hamilton, Ontario, Canada

^c McMaster GRADE Centre & Michael DeGroot Cochrane Canada Centre, McMaster University, Hamilton, Ontario, Canada

^d Department of Paediatrics, Obstetrics and Gynaecology, Preventive Medicine, and Public Health. PhD Programme in Methodology of Biomedical Research and Public Health. Universitat Autònoma de Barcelona, Bellaterra, Spain.

^e Iberoamerican Cochrane Center, Biomedical Research Institute (IIB Sant Pau-CIBERESP), Barcelona, Spain

^f Department of Internal Medicine, American University of Beirut, Beirut, Lebanon

^g Toronto Health Economics and Technology Assessment (THETA) Collaborative, Toronto, Ontario, Canada

^h National Toxicology Program, National Institute of Environmental Health Sciences, Durham, North Carolina, USA

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

- 35 ⁱ Texas A&M University, College Station, Texas, USA
- 36 ^j Liverpool John Moores University, Liverpool, UK
- 37 ^k Center for Evidence-Based Medicine and Health Outcome Research, Morsani College of
38 Medicine, University of South Florida, Tampa, Florida, USA
- 39 ^l Institute for Education and Research, Hospital Moinhos de Vento, Porto Alegre, Rio Grande do
40 Sul, Brazil
- 41 ^m ICF International, Durham, North Carolina, USA
- 42 ⁿ Institute of Health & Wellbeing, University of Glasgow, Glasgow, UK
- 43 ^o World Health Organization, Geneva, Switzerland
- 44 ^p Clinical Epidemiology and Medical Technology Assessment, Maastricht University Medical
45 Centre+, Maastricht, the Netherlands
- 46 ^q LaKind Associates, LLC, Catonsville, Maryland, USA
- 47 ^r Department of Epidemiology and Public Health, University of Maryland School of Medicine,
48 Baltimore, Maryland, USA
- 49 ^s Department of Clinical Epidemiology, Biostatistics, and Bioinformatics, Academic Medical Center,
50 University of Amsterdam, Amsterdam, the Netherlands
- 51 ^t Department of Surgery, University of California Davis, Sacramento, California, USA
- 52 ^u Department of Medicine, Department of Veterans Affairs, Northern California Health Care
53 System, Mather, California, USA
- 54 ^v Division of Infection, Immunity and Respiratory Medicine, University Hospital of South
55 Manchester, University of Manchester, Manchester, UK
- 56 ^w Institute for Evidence in Medicine, Medical Center, University of Freiburg, Freiburg-am-Breisgau,
57 Germany
- 58 ^x Cochrane Germany, Freiburg-am-Breisgau, Germany
- 59 ^y National Center for Environmental Assessment, U.S. Environmental Protection Agency,
60 Washington, D.C., District of Columbia, USA
- 61 ^z Department of Medicine, University of Kansas Medical Center, Kansas City, Kansas, USA
- 62 ^{aa} Institute of Population Health Sciences, University of Liverpool, Liverpool, UK
- 63 ^{ab} National Center for Computational Toxicology, U.S. Environmental Protection Agency, Durham,
64 North Carolina, USA
- 65 ^{ac} Department of Family Medicine, McMaster University, Hamilton, Ontario, Canada
- 66 ^{ad} EPPI-Centre, Institute of Education, University College London, London, UK
- 67 ^{ae} Programs for Assessment of Technology in Health, McMaster University, Hamilton, Ontario,
68 Canada
- 69 ^{af} Evidence-Based Toxicology Collaboration, Johns Hopkins Bloomberg School of Public Health,
70 Baltimore, Maryland, USA
- 71 ^{ag} Health Economics Group, Institute of Health and Society, Newcastle University, Newcastle upon
72 Tyne, UK

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

73 ^{ah} Health Quality Ontario, Toronto, Ontario, Canada

74

75

76 Corresponding author:

77 Jan Brozek

78 McMaster University

79 Health Sciences Centre, Area 2C

80 1280 Main Street West

81 Hamilton, ON L8S 4K1, Canada

82

83

84

Journal Pre-proof

Abstract

86

Objectives:

88 To present the Grading of Recommendations Assessment, Development, and Evaluation (GRADE)
89 conceptual approach to the assessment of certainty of evidence from modelling studies (i.e.
90 certainty associated with model outputs).

Study Design and Setting:

92 Expert consultations and, an international multi-disciplinary workshop informed development of a
93 conceptual approach to assessing the certainty of evidence from models within the context of
94 systematic reviews, health technology assessments, and health care decisions. The discussions
95 also clarified selected concepts and terminology used in the GRADE approach and by the
96 modelling community. Feedback from experts in a broad range of modelling and health care
97 disciplines addressed the content validity of the approach.

Results:

99 Workshop participants agreed, that the domains determining the certainty of evidence previously
100 identified in the GRADE approach (risk of bias, indirectness, inconsistency, imprecision, reporting
101 bias, magnitude of an effect, dose-response relation, and the direction of residual confounding)
102 also apply when of assessing the certainty of evidence from models. The assessment depends on
103 the nature of model inputs and the model itself and on whether one is evaluating evidence from a
104 single model or multiple models. We propose a framework for selecting the best available
105 evidence from models: 1) developing *de novo* a model specific to the situation of interest, 2)
106 identifying an existing model the outputs of which provide the highest certainty evidence for the
107 situation of interest, either “off the shelf” or after adaptation, and 3) using outputs from multiple
108 models. We also present a summary of preferred terminology to facilitate communication among
109 modelling and health care disciplines.

Conclusions:

111 This conceptual GRADE approach provides a framework for using evidence from models in health
112 decision making and the assessment of certainty of evidence from a model or models. The GRADE
113 Working Group and the modelling community are currently developing the detailed methods and
114 related guidance for assessing specific domains determining the certainty of evidence from
115 models across health care-related disciplines (e.g. therapeutic decision-making, toxicology,
116 environmental health, health economics).

117

118 Introduction

119

120 When direct evidence to inform health decisions is not available or not feasible to **measure** (e.g.
121 long-term effects of interventions or when studies in certain populations are perceived as
122 unethical), modelling studies may be used to **predict** that “evidence” and inform decision-
123 making.[1, 2] Health decision makers arguably face many more questions than can be reasonably
124 answered with studies that directly measure the outcomes. Modelling studies, therefore, are
125 increasingly used to predict disease dynamics and burden, the likelihood that an exposure
126 represents a health hazard, the impact of interventions on health benefits and harms, or the
127 economic efficiency of health interventions, among others [1]. Irrespective of the modelling
128 discipline, decision makers need to know the best **estimates** of the modelled outcomes and how
129 much **confidence** they may have in each estimate.[3] Knowing to what extent one can trust the
130 outputs of a model is necessary when using them to support health decisions [4].

131

132 Although a number of guidance documents on how to assess the trustworthiness of estimates
133 obtained from models in several health fields have been previously published [5-16], they are
134 limited by failing to distinguish methodological rigor from completeness of reporting, and by
135 failing to clearly distinguish among various components affecting the trustworthiness of model
136 outputs. In particular they lack clarity regarding sources of uncertainty that may arise from **model**
137 **inputs** and from the uncertainty about a **model itself**. Modellers and those using results from
138 models should assess the credibility of both.[4]

139

140 Authors have attempted to develop tools to assess model credibility, but many addressed only
141 selected aspects, such as statistical reproducibility of data, the quality of reporting[17], or a
142 combination of reporting with aspects of good modelling practices[7, 18-21]. Many tools also do
143 not provide sufficiently detailed guidance on how to apply individual domains or criteria. There is
144 therefore a need for further development and validation of such tools in specific disciplines.
145 Sufficiently detailed guidance for making and reporting these assessments is also necessary.

146

147 Models predict outcomes based on model inputs – previous observations, knowledge and
148 assumptions about the situation being modelled. Thus, when developing new models or assessing
149 whether an existing model has been optimally developed, one should specify *a priori* the most
150 appropriate and relevant data sources to inform different parameters required for the model.
151 These may be either (seldom) a single study that provides the most direct information for the
152 situation being modelled or (more commonly) a systematic review of multiple studies that identify
153 all relevant sources of data. The risk of bias, directness and consistency of input data, precision of
154 these estimates, and other domains specified in the Grading of Recommendations Assessment,

155 Development, and Evaluation (GRADE) approach determine the certainty of each of the model
156 inputs.[22-28]

157
158 When assessing the evidence generated, various disciplines in health care and related areas that
159 use modelling face similar challenges may benefit from shared solutions. Table 1 presents
160 examples of selected models used in health-related disciplines in Table 1. Building on the existing
161 GRADE approach, we refined and expand guidance regarding assessment of the certainty of model
162 outputs. We formed a GRADE project group comprised of individuals with expertise in developing
163 models and using model results in health-related disciplines, to create a unified framework for
164 assessing the certainty of model outputs in the context of systematic reviews [29], health
165 technology assessments, health care guidelines, and other health decision-making. In this article,
166 we outline the proposed conceptual approach and clarify key terminology (Table 2). The target
167 audience for this article includes researchers who develop models and those who use models to
168 inform health care-related decisions.

169

170 **What we mean by a model**

171

172 Authors have used the term *model* to describe a variety of different concepts [2] and suggested
173 several broader or narrower definitions [6, 30], so even modellers in the relatively narrow context
174 of health sciences can differ in their views regarding what constitutes a model. Models vary in
175 their structure and degree of complexity. A very simple model might be an equation estimating a
176 variable not directly measured, such as the absolute effect of an intervention estimated as the
177 product of the intervention's relative effect and the assumed baseline risk in a defined population
178 (risk difference equals relative risk reduction multiplied by an assumed baseline risk). On the other
179 end of the spectrum, elaborate mathematical models, such as system dynamics models (e.g.
180 infectious disease transmission) may contain dozens of sophisticated equations that require
181 considerable computing power to solve.

182

183 By their nature, such models only *resemble* the phenomena being modelled – i.e. specific parts of
184 the world that are interesting in the context of a particular decision – with necessary
185 approximations and simplifications, and to the extent that one actually knows and understands
186 the underlying mechanisms.[1] Given the complexity of the world, decision-makers often rely on
187 some sort of a model to answer health-related questions.

188

189 In this article, we focus on quantitative mathematical models defined as “mathematical
190 framework representing variables and their interrelationships to describe observed phenomena or
191 predict future events”[30] used in health-related disciplines for decision-making (Table 1). These
192 may be models of systems representing causal mechanisms (aka mechanistic models), models

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

193 predicting outcomes from input data (aka empirical models), and models combining mechanistic
194 with empirical approaches (aka hybrid models). We do not consider here statistical models used to
195 estimate the associations between measured variables (e.g. proportional hazards models or
196 models used for meta-analysis).

197

198 **The GRADE approach**

199

200 The GRADE working group was established in the year 2000 and continues as a community of
201 people striving to create systematic, and transparent frameworks for assessing and
202 communicating the certainty of the available evidence used in making decisions in healthcare and
203 health-related disciplines.[31] The GRADE Working Group now includes over 600 active members
204 from 40 countries and serves as a think tank for advancing evidence-based decision-making in
205 multiple health-related disciplines (www.gradeworkinggroup.org). GRADE is widely used
206 internationally by over 110 organizations to address topics related to clinical medicine, public
207 health, coverage decisions, health policy, and environmental health.

208

209 The GRADE framework uses concepts familiar to health scientists, grouping specific items to
210 evaluate the certainty of evidence in conceptually coherent domains. Specific approaches to the
211 concepts may differ depending on the nature of the body of evidence (Table 2). GRADE domains
212 include concepts such as risk of bias[28], directness of information [24], precision of an
213 estimate[23], consistency of estimates across studies[25], risk of bias related to selective
214 reporting[26], strength of the association, presence of a dose-response gradient, and the presence
215 of plausible residual confounding that can increase confidence in estimated effects[27].

216

217 The general GRADE approach is applicable irrespective of health discipline. It has been applied to
218 rating the certainty of evidence for management interventions, health care related tests and
219 strategies [32, 33], prognostic information[34], evidence from animal studies[35], use of resources
220 and cost-effectiveness evaluations[36], and values and preferences[37, 38]. Although the GRADE
221 Working Group has begun to address certainty of modelled evidence in the context of test-
222 treatment strategies[39], health care resource use and costs[36], and environmental health[40],
223 more detailed guidance is needed for complex models such as those used in infectious diseases,
224 health economics, public health, and decision analysis.

225

226 **Methods**

227

228 On May 15 and 16, 2017, health scientists participated in a GRADE modelling project group
229 workshop in Hamilton, Ontario, Canada, to initiate a collaboration in developing common

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

230 principles for the application of the GRADE assessment of certainty of evidence to modelled
231 outputs. The National Toxicology Program of the Department of Health and Human Services in the
232 USA and the MacGRADE Center in the Department of Health Research Methods, Evidence, and
233 Impact at McMaster University sponsored the workshop which was co-organized by MacGRADE
234 Center and ICF International.

235
236 Workshop participants were selected to ensure a broad representation of all modelling related
237 fields (Appendix). Participants had expertise in modelling in the context of clinical practice
238 guidelines, public health, environmental health, dose-response modelling, physiologically based
239 pharmacokinetic (PBPK) modelling, environmental chemistry, physical/chemical property
240 prediction, evidence integration, infectious disease, computational toxicology, exposure
241 modelling, prognostic modelling, diagnostic modelling, cost effectiveness modelling, biostatistics,
242 and health ethics.

243
244 Leading up to the workshop, we held three webinars to introduce participants to the GRADE
245 approach. Several workshop participants (VM, KT, JB, AR, JW, JLB, HJS) collected and summarized
246 findings from literature and the survey of experts as background material that provided a starting
247 point for discussion. The materials included collected terminology representing common concepts
248 across multiple disciplines that relate to evaluating modelled evidence, and a draft framework for
249 evaluating modelled evidence. Participants addressed specific tasks in small groups and large
250 group discussion sessions and agreed on key principles both during the workshop and through
251 written documents.

252

253 **Results**

254

255 **Terminology**

256

257 Workshop participants agreed on the importance of clarifying terminology to facilitate
258 communication among modellers, researchers, and users of model outputs from different
259 disciplines. Modelling approaches evolved somewhat independently, resulting in different terms
260 being used to describe the same or very similar concepts or the same term being used to describe
261 different concepts. For instance, the concept of extrapolating from the available data to the
262 context of interest has been referred to as directness, applicability, generalizability, relevance, or
263 external validity. The lack of standardized terminology leads to confusion and hinders effective
264 communication and collaboration among modellers and users of models.

265

266 Overcoming these obstacles would require clarifying the definitions of concepts and agreeing on
267 terminology across disciplines. Realizing that this involves changing established customary use of
268 terms in several disciplines, workshop participants suggested accepting the use of alternative
269 terminology while always being clear about the preferred terms to be used and the underlying
270 concept to which it refers (Table 2). Experts attending a World Health Organization's consultation
271 have very recently suggested a more extensive set of terms [41]. To facilitate future
272 communication, participants of this workshop will further collaborate to build a comprehensive
273 glossary of terminology related to modelling.

274

275 **Outline of an approach to using model outputs for decision making**

276

277 Workshop participants suggested an approach to incorporate model outputs in health-related
278 decision making (Figure 1). In this article we describe only the general outline of the suggested
279 approach – in subsequent articles we will discuss the details of the approach and provide more
280 specific guidance on its application to different disciplines and contexts.

281

282 Researchers should start by **conceptualizing the problem and the ideal target model** that would
283 best represent the actual phenomenon or decision problem they are considering [13]. This
284 conceptualization would either guide the development of a new model or serve as a reference
285 against which existing models could be compared. The ideal target model should reflect: 1) the
286 relevant population (e.g., patients receiving some diagnostic procedure or exposed to some
287 hazardous substance), 2) the exposures or health interventions being considered, 3) the outcomes
288 of interest in that context, and 4) their relationships. [42]. Conceptualizing the model will also
289 reduce the risk of intentional or unintentional development of data-driven models, in which inputs
290 and structure would be determined only by what is feasible to develop given the available data at
291 hand.

292

293 Participants identified 3 options in which users may incorporate model outputs in health decision-
294 making (Figure 1):

295 **1. Develop a model de novo designed specifically to answer the very question at hand.**

296 Workshop participants agreed that in an ideal situation such an approach would almost always
297 be the most appropriate. Following this approach, however, requires suitable skills, ample
298 resources, and time being available. It also requires enough knowledge about the
299 phenomenon being modelled to be able to tell whether or not the new model would have any
300 advantage over already existing models.

301 **2. Search for an existing model describing the same or a very similar problem and use it “off-
302 the-shelf”** or adapt it appropriately in order to answer the current question. In practice many
303 researchers initially use this approach because of the above limitations of developing a new

304 model. However, it is often not possible to find an existing model that would be directly
305 relevant to the problem at hand and/or it is not feasible to adapt an existing model when
306 found. Any adaptation of a model requires availability of input data relevant for current
307 problem, appropriate expertise and resources, and access to the original model. The latter is
308 often not available (e.g. proprietary model or no longer maintained) or the structure of the
309 original model is not being transparent enough to allow adaptation (“black-box”).

310 **3. Use the results from multiple existing models** found in the literature [43]. This approach may
311 be useful when a limited knowledge about the phenomenon being modelled makes it
312 impossible to decide which of the available models is more relevant, or when many alternative
313 models are relevant but use different input parameters. In such situations, one may be
314 compelled to rely on the results of several models, because selection of the single, seemingly
315 “best” model may provide incorrect estimates of outputs and lead to incorrect decisions.

316 Identifying existing models that are similar to the ideal target model often requires performing a
317 scoping of the literature or a complete **systematic review** of potentially relevant models – a
318 structured process following a standardized set of methods with a goal to identify and assess all
319 available models that are accessible, transparently reported, and fulfil the pre-specified eligibility
320 criteria based on the conceptual ideal target model. Some prefer the term **systematic survey** that
321 differs from a systematic review in the initial intention to use the results: in systematic reviews the
322 initial intention is to combine the results across studies either statistically through a meta-analysis
323 or narratively summarizing their results when appropriate, whereas in a systematic survey the
324 initial intention is to examine the various ways that an intervention or exposure has been
325 modelled, to review the input evidence that has been used, and ultimately to identify a single
326 model that fits the conceptual ideal target model the best or requires the least adaptation; only
327 when one cannot identify a single such model will it be necessary to use the results of multiple
328 existing models.

329
330 If a systematic search revealed one or more models meeting the eligibility criteria, then
331 researchers would assess the certainty of outputs from each model. Depending on this
332 assessment, researchers may be able to use the results of a single most direct and lowest risk of
333 bias model “off-the-shelf” or proceed to adapt that model. If researchers failed to find an existing
334 model that would be sufficiently direct and low risk of bias, then they would ideally develop their
335 own model de novo.

336
337 Assessing the certainty of outputs from a single model

338
339 When researchers develop their own model or when they identify a single model that is
340 considered sufficiently direct to the problem at hand, they should assess the certainty of its

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

341 outputs (i.e. evidence generated from that model). Note, that if a model estimates multiple
342 outputs, researchers needs to assess the certainty of each output separately [23-28]. Workshop
343 participants agreed that all GRADE domains are applicable to assess the certainty of model
344 outputs, but further work is needed to identify examples and develop specific criteria to be
345 assessed, which may differ depending on the model being used and/or situation being modelled.

346

347 *Risk of bias in a single model*

348

349 The risk of bias of model outputs (i.e. model outputs being systematically overestimated or
350 underestimated) is determined by the credibility of a model itself and the certainty of evidence for
351 each of model inputs.

352

353 The **credibility of a model**, also referred to as the quality of a model (Table 2) is influenced by its
354 conceptualization, structure, calibration, validation, and other factors. Determinants of model
355 credibility are likely to be specific to a modelling discipline (e.g., health economic models have
356 different determinants of their credibility than PBPK models). There are some discipline-specific
357 guidelines or checklists developed for the assessment of credibility of a model and other factors
358 affecting the certainty of model outputs such as the framework to assess adherence to good
359 practice guidelines in decision-analytic modelling [18], the questionnaire to assess relevance and
360 credibility of modelling studies [18, 44, 45], good research practices for modelling in health
361 technology assessment [5, 6, 8, 9, 12-14], the approaches to assessing uncertainty in read-across
362 [46], and the quantitative structure-activity relationships [47] in predictive toxicology. Workshop
363 participants agreed that there is a need for comprehensive tools developed specifically to assess
364 credibility of various types of models in different modelling disciplines.

365

366 The **certainty of evidence in each of the model inputs** is another critical determinant of the risk of
367 bias in a model. A model has several types of input data – bodies of evidence used to populate a
368 model (Table 2). When researchers develop their model *de novo*, in order to minimize the risk of
369 bias they need to specify those input parameters to which the model outputs are the most
370 sensitive. For instance, in economic models these key parameters may include health effects,
371 resource use, utility values, and baseline risks of outcomes. Model inputs should reflect the entire
372 body of relevant evidence satisfying clear pre-specified criteria rather than an arbitrarily selected
373 evidence that is based on convenience (“any available evidence”) or picked in any other non-
374 systematic way (e.g., “first evidence found” – single studies that researchers happen to know
375 about or are the first hits in a database search).

376

377 The appropriate approach will depend on the type of data and may require performing a
378 systematic review of evidence on each important or crucial input variable [48-50]. Some inputs

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

379 may have a very narrow inclusion criteria and therefore evidence from single epidemiological
380 survey or population surveillance may provide all relevant data for the population of interest (e.g.
381 baseline population incidence or prevalence).

382
383 The certainty of evidence for each input needs to be assessed following the established GRADE
384 approach specific to that type of evidence (e.g. estimates of intervention effects or baseline risk of
385 outcomes)[22, 32, 34, 37]. Following the logic of the GRADE approach that the overall certainty of
386 evidence cannot be higher than the lowest certainty for any body of evidence that is critical for a
387 decision [51], the overall rating of certainty of evidence across model inputs should be limited by
388 the lowest certainty rating for any body of evidence (in this case input data) to which the model
389 output(s) was proved sensitive.

390
391 Application of this approach requires a priori consideration of likely critical and/or important
392 inputs when specifying the *conceptual ideal target model* and the examination of the results of
393 *back-end* sensitivity analyses. It further requires deciding how to judge whether results are or are
394 not sensitive to alternative input parameters. Authors have described several methods to identify
395 the most influential parameters including global sensitivity analysis to obtain “parameter
396 importance measures” (i.e. information based measures) [52]; or alternatively by varying one
397 parameter at a time and assessing their influence in “base case” outputs [52] For example, in a
398 model-based economic evaluation one might be looking for the influence of sensitivity analysis on
399 cost-effectiveness ratios at a specified willingness-to-pay threshold.

400
401 *Indirectness in a single model*

402 By directness or relevance, we mean the extent to which model outputs directly represent the
403 phenomenon being modelled. To evaluate the relevance of a model, one needs to compare it
404 against the conceptual ideal target model. When there are concerns about the directness of the
405 model or there is limited understanding of the system being modelled making it difficult to assess
406 directness, then one may have lower confidence in model outputs.

407
408 Determining the directness of model outputs includes assessing to what extent the modelled
409 population, the assumed interventions and comparators, the time horizon, the analytic
410 perspective, as well as the outcomes being modelled reflect those that are current interest. For
411 instance, if the question is about the risk of birth defects in children of mothers chronically
412 exposed to a certain substance, there may be concerns about the directness of the evidence if the
413 model assumed short-term exposure, the route of exposure was different, or the effects of
414 exposure to a similar but not the same substance were measured.

415
416 Assessing indirectness in a single model also requires evaluating two separate sources of
417 indirectness:

418 1. indirectness of input data with respect to the ideal target model’s inputs.

419 2. indirectness of model outputs with respect to the decision problem at hand.

420

421 This conceptual distinction is important because, although they are interrelated, one needs to
422 address each type of indirectness separately. Even if the outputs might be direct to the problem
423 of interest, the final assessment should consider if the inputs used were also direct for the target
424 model.

425

426 Using an existing model has potential limitations: its inputs might have been direct for the decision
427 problem addressed by its developers but are not direct with respect to the problem currently at
428 hand. In this context, sensitivity analysis can help to assess to what extent model outputs are
429 robust to the changes in input data or assumptions used in model development.

430

431 *Inconsistency in a single model*

432

433 A single model may yield inconsistent outputs owing to unexplained variability in the results of
434 individual studies informing the pooled estimates of input variables. For instance, when
435 developing a health economic model, a systematic review may yield several credible, but
436 discrepant, utility estimates in the population of interest. If there is no plausible explanation for
437 that difference in utility estimates, outputs of a model based on those inputs may also be
438 qualitatively inconsistent. Again, sensitivity analysis may help to make a judgment to what extent
439 such inconsistency of model inputs would translate into a meaningful inconsistency in model
440 outputs with respect to the decision problem at hand.

441

442 *Imprecision in a single model*

443

444 Sensitivity analysis characterizes the response of model outputs to parameter variation, and helps
445 to determine the robustness of model's qualitative conclusions [52, 53]. The overall certainty of
446 model outputs may also be lower when the outputs are estimated imprecisely. For quantitative
447 outputs one should examine not only the point estimate (e.g., average predicted event) but also
448 the variability of that estimate (e.g., results of the probabilistic sensitivity analysis based in the
449 distribution of the input parameters). It is essential that a report from a modelling study always
450 includes information about output variability. Further guidance on how to assess imprecision in
451 model outputs will need to take into account if the conclusions change according to that specific
452 parameter. In some disciplines, for instance in environmental health, model inputs are frequently
453 qualitative. Users of such models may assess "adequacy" of the data, i.e. the degree of "richness"
454 and quantity of data supporting particular outputs of a model.

455

456 *Risk of publication bias in the context of a single model*

457
458 The risk of publication bias, also known as “reporting bias”, “non-reporting bias”, or “bias owing to
459 missing results”, as it is currently called in the Cochrane Handbook [54], is the likelihood that
460 relevant models have been constructed but were not published or otherwise made publicly
461 available. Risk of publication bias may not be relevant when assessing the certainty of outputs of a
462 single model constructed de novo. However, when one intends to reuse an existing model but is
463 aware or strongly suspects that similar models had been developed but are not available, then
464 one may be inclined to think that their outputs might have systematically differed from the model
465 that is available. In such a case, one may have lower confidence in the outputs of the identified
466 model if there is no reasonable explanation for the inability to obtain those other models.

467
468 *Domains that increase the certainty of outputs from a single model*

469
470 The GRADE approach to rating the certainty of evidence recognized three situations when the
471 certainty of evidence can increase: large magnitude of an estimated effect, presence of a dose-
472 response gradient in an estimated effect, and an opposite direction of plausible residual
473 confounding.[27] Workshop participants agreed that presence of a **dose-response gradient** in
474 model outputs may be applicable in some modelling disciplines (e.g., environmental health).
475 Similarly, whether or not a **large magnitude of an effect** in model outputs increases the certainty
476 of the evidence may depend on the modelling discipline. The **effect of an opposite direction of**
477 **plausible residual confounding** seems theoretically also applicable in assessing the certainty of
478 model outputs (i.e. a conservative model not incorporating input data parameter in favour of an
479 intervention but still finding favorable outputs) but an actual example of this phenomenon in
480 modelling studies is still under discussion.

481
482 *Assessing the certainty of outputs across multiple models*

483
484 Not infrequently, particularly in disciplines relying on mechanistic models, the current knowledge
485 about the real system being modelled is very limited precluding the ability to determine which of
486 the available existing models generates higher certainty outputs. Therefore, it may be necessary
487 to rely on the results across multiple models. Other examples include using multiple models when
488 no model was developed for the population directly of interest (e.g. the European Breast Cancer
489 Guideline for Screening and Diagnosis relied on a systematic review of modelling studies that
490 compared different mammography screening intervals [55]) or when multiple models of the same
491 situation exist but vary in structure, complexity, and parameter choices (e.g. HIV Modelling
492 Consortium compared several different mathematical models simulating the same antiretroviral

493 therapy program and found that all models predicted that the program has the potential to
494 reduce new HIV infections in the population [56]).

495
496 When researchers choose or are compelled to include outputs from several existing models, they
497 should assess the certainty of outputs across all included models. This assessment may be more
498 complex than for single models and single bodies of evidence. The feasibility of GRADE's guidance
499 to judge the certainty of evidence lies in the availability of accepted methods for assessing most
500 bodies of evidence from experimental to observational studies. However, the methods for
501 systematic reviews of modelling studies are less well-established, some stages of the process are
502 more complex, the number of highly skilled individuals with experience in such systematic reviews
503 is far lower, and there is larger variability in the results [57]. Additionally, researchers must be
504 careful to avoid "double counting" the same model as if it were multiple models. For instance, the
505 same model (i.e. same structure and assumptions) may have been used in several modelling
506 studies, in which investigators relied on different inputs. When facing this scenario, researchers
507 may need to decide which of the inputs are the most direct to their particular question and
508 include in only this model in the review.

509

510 *Risk of bias across multiple models*

511

512 The assessment of risk of bias across models involves an assessment of the risk of bias in each
513 individual model (see above discussion of risk of bias in single model) and subsequently making a
514 judgement about the overall risk of bias across all included models. Specific methods for
515 operationalizing this integration remain to be developed.

516

517 *Indirectness across multiple models*

518

519 As for the risk of bias, researchers need to assess indirectness of outputs initially for each of
520 included models and then integrate the judgements across models. Likewise, specific methods for
521 operationalizing this integration still remain to be developed. During this assessment researchers
522 may find some models too indirect to be informative for their current question and decide to
523 exclude them from further consideration. However, the criteria to determine which models are
524 too indirect should be developed a priori, before the search for the models is performed and their
525 results are known.

526

527 *Imprecision across multiple models*

528

529 The overall certainty of model outputs may also be lower when model outputs are not estimated
530 precisely. If researchers attempt a quantitative synthesis of outputs across models, they will

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

531 report the range of estimates and variability of that estimates. When researchers choose to
532 perform only a qualitative summary of the results across models, it is desirable that they report
533 some estimate of variability in the outputs of individual models and an assessment of how severe
534 the variability is (e.g. range of estimated effects).

535

536 *Inconsistency of outputs across multiple models*

537

538 The assessment of inconsistency should focus on unexplained differences across model outputs
539 for a given outcome. If multiple existing models addressing the same issue produce considerably
540 different outputs or reach contrasting conclusions, then careful comparison of the models may
541 lead to a deeper understanding of the factors that drive outputs and conclusions. Ideally, the
542 different modelling groups that developed relevant models would come together to explore the
543 importance of differences in the type and structure of their models, and of the data used as model
544 inputs.

545

546 Invariably there will be some differences among the estimates from different models. Researchers
547 will need to assess whether or not these differences are important, i.e. whether they would lead
548 to different conclusions. If the differences are important but can be explained by model structure,
549 model inputs, the certainty of the evidence of the input parameters or other relevant reasons, one
550 may present the evidence separately for the relevant subgroups. If differences are important, but
551 cannot be clearly explained, the certainty of model outputs may be lower.

552

553 *Risk of publication bias across multiple models*

554

555 The assessment is similar to that of the risk of publication bias in the context of a single model.

556

557 *Domains that increase the certainty of outputs across multiple models*

558

559 All considerations are the same to those in the context of a single model.

560

561 **Discussion**

562

563 The goal of the GRADE project group on modelling is to provide concepts and operationalization of
564 how to rate the certainty of evidence in model outputs. This article provides an overview of the
565 conclusions of the project group. This work is important because there is a growing need and
566 availability of modelled information resulting from a steadily increasing knowledge of the
567 complexity of the structure and interactions in our environment, and computational power to

568 construct and run models. Users of evidence obtained from modelling studies need to know how
569 much trust they may have in model outputs. There is a need to improve the methods of
570 constructing models and to develop methods for assessing the certainty in model outputs. In this
571 article we have attempted to clarify the most important concepts related to developing and using
572 model outputs to inform health-related decision-making. Our preliminary work identified
573 confusion about terminology, lack of clarity of what is a model, and need for methods to assess
574 certainty in model outputs as priorities to be addressed in order to improve the use of evidence
575 from modelling studies.

576
577 In some situations, decision-makers might be better off developing a new model specifically
578 designed to answer their current question. However, we suggest that it is not always feasible to
579 develop a new model or that developing a new model might not be any better than using already
580 existing models, when the knowledge of the real life system to be modelled is limited precluding
581 the ability to choose one model that would be better than any other. Thus, sometimes it may be
582 necessary or more appropriate to use one or multiple existing models depending on their
583 availability, credibility, and relevance to the decision-making context. The assessment of the
584 certainty of model outputs will be conceptually similar when a new model is constructed, or one
585 existing model is used. The main difference between the latter two approaches is the availability
586 of information to perform a detailed assessment. That is, information for one's own model may be
587 easily accessible, but information required to assess someone else's model will often be more
588 difficult to obtain. Assessment of the certainty evidence across models can build on existing
589 GRADE domains but requires different operationalization.

590
591 Because it builds on an existing, widely used framework that includes a systematic and
592 transparent evaluation process, modelling disciplines' adoption of the GRADE approach and
593 further development of methods to assess the certainty of model outputs may be beneficial for
594 health decision making. Systematic approaches improve rigor of research, reducing the risk of
595 error and its potential consequences; transparency of the approach increases its trustworthiness.
596 There may be additional benefits related to other aspects of the broader GRADE approach, for
597 instance a potential to reduce unnecessary complexity and workload in modelling by careful
598 consideration of the most direct evidence as model inputs. This may allow, for instance,
599 optimization of the use of different streams of evidence as model inputs. Frequently, authors
600 introduce unnecessary complexity by considering multiple measures of the same outcome when
601 focus could be on the most direct outcome measure.

602
603 The GRADE working group will continue developing methods and guidance for using model
604 outputs in health-related decision-making. In subsequent articles we will provide more detailed
605 guidance about choosing the "best" model when multiple models are found, using multiple

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

606 models, integrating the certainty of evidence from various bodies of evidence with credibility of
607 the model and arriving at the overall certainty in model outputs, how to assess the credibility of
608 various types of models themselves, and further clarification of terminology. In the future we aim
609 to develop and publish the detailed guidance for assessing certainty of evidence from models, the
610 specific guidance for the use of modelling across health care-related disciplines (e.g. toxicology,
611 environmental health or health economics), validation of the approach, and accompanying
612 training materials and examples.
613

614 Acknowledgments

615
616 AR was supported by the National Institutes of Health, National Institute of Environmental Health
617 Sciences.

618 **Table 1.** Examples of modelling methods in health-related disciplines (not comprehensive)*
 619

Decision analysis models	Structured model representing health care pathways examining effects of an intervention on outcomes of interest.
	<p>Types</p> <ul style="list-style-type: none"> ▪ Decision tree models ▪ State transition models <ul style="list-style-type: none"> ○ Markov cohort simulation ○ Individual based microsimulation (first-order Monte Carlo) ▪ Discrete event simulation ▪ Dynamic transmission models ▪ Agent based models
	<p>Examples</p> <ul style="list-style-type: none"> ▪ Estimation of long-term benefits and harms outcomes from complex intervention, e.g. minimum unit pricing of alcohol ▪ Estimation of benefits and harms of population mammography screening based in microsimulation model, e.g. Wisconsin model from CISNET collaboration[58] ▪ Susceptible-Infectious-Recovery transmission dynamic model to assess effectiveness of lockdown during the SARS-CoV-2 pandemic[59]
Pharmacology and toxicology models	Computational models developed to organize, analyse, simulate, visualize or predict toxicological and ecotoxicological effects of chemicals. In some cases, these models are used to estimate the toxicity of a substance even before it has been synthesized.
	<p>Types</p> <ul style="list-style-type: none"> ▪ Structural alerts and rule-based models ▪ Read-Across ▪ Dose response and Time response ▪ Toxicokinetic (TK) and toxicodynamic(TD) ▪ Uncertainty factors ▪ Quantitative structure activity relationship (QSAR) ▪ Biomarker-based toxicity models
	<p>Examples</p> <ul style="list-style-type: none"> • Structural alerts for mutagenicity and skin sensitisation • Read-across for complex endpoints such as chronic toxicity • Pharmacokinetic (PK) models to calculate concentrations of substances in organs, following a variety of exposures QSAR models for carcinogenicity • TGx-DDI biomarker to detect DNA damage-inducing agents
Environmental models	The EPA defined these models as: ‘A simplification of reality that is constructed to gain insights into select attributes of a physical, biological, economic, or social system.’ It involves the application of multidisciplinary knowledge to explain, explore and predict the Earth’s response to environmental change, and the interactions between human activities and natural processes.

	<p>Classification (based on the CREM guidance document):</p> <ul style="list-style-type: none"> • Human activity models • Natural systems process • Emission models • Fate and transport models • Exposure models • Human health effects models • Ecological effects models • Economic impact models • Noneconomic impact models
	<p>Examples</p> <ul style="list-style-type: none"> • Land use regression models • IH SkinPerm [60] • ConsExpo [61] • other exposure models [62]
Other	<ul style="list-style-type: none"> • HopScore: An Electronic Outcomes-Based Emergency Triage System [63] • Computational general equilibrium (CGE) models [64]
<p>*Although not described in this classification simple calculations incorporating two or more pieces of evidence as for example the multiplication of a RR by the baseline risk to obtain the absolute risk difference of an intervention is a model, although pragmatic, with their respective assumptions.</p>	

620
621
622 **Table 2.** Selected commonly used and potentially confusing terms used in the context of modelling
623 and the GRADE approach
624

Term	General definition
Sources of evidence (may come from in vitro or in vivo experiment or a mathematical model)	
Streams of evidence	Parallel information about the same outcome that may have been obtained using different methods of estimating that outcome. For instance, evidence of the increased risk for developing lung cancer in humans after an exposure to certain chemical compound may come from several streams of evidence: 1) mechanistic evidence – models of physiological mechanisms, 2) studies in animals – observations and experiments in animals from different phyla, classes, orders, families, genera, and species (e.g., bacteria, nematodes, insects, fish, mice, rats), and 3) studies in humans.
Bodies of evidence	Information about multiple different aspects around a decision about the best course of action. For instance, in order to decide whether or not a given diagnostic test should be used in some people, one needs to integrate the bodies of evidence about: the accuracy of the test, the prevalence of the conditions being suspected, the natural history of these conditions, the effects of potential treatments, values and preferences of affected individuals, cost, feasibility, etc.
Quality	

(may refer to many concepts, thus alternative terms are preferred to reduce confusion)	
<p>Certainty of model outputs</p> <p>Alternative terms:</p> <ul style="list-style-type: none"> ▪ certainty of modelled evidence ▪ quality of evidence ▪ quality of model output ▪ strength of evidence ▪ confidence in model outputs 	<p>In the context of health decision-making, the certainty of evidence (term preferred over “quality” in order to avoid confusion with the risk of bias in an individual study) reflects the extent to which one’s confidence in an estimate of an effect is adequate to make a decision or a recommendation. Decisions are influenced not only by the best estimates of the expected desirable and undesirable consequences but also by one’s confidence in these estimates. In the context of evidence syntheses of separate bodies of evidence (e.g., systematic reviews), the certainty of evidence reflects the extent of confidence that an estimate of effect is correct. For instance, the attributable national risk of cardiovascular mortality resulting from exposure to air pollution measured in selected cities.</p> <p>The GRADE Working Group published several articles explaining the concept in detail.[22-28, 65] Note that the phrase “confidence in an estimate of an effect” does not refer to statistical confidence intervals. Certainty of evidence is always assessed for the whole body of evidence rather than on a single study level (single studies are assessed for risk of bias and indirectness).</p>
<p>Certainty of model inputs</p> <p>Alternative term:</p> <ul style="list-style-type: none"> ▪ quality of model inputs 	<p>Characteristics of data that are used to develop, train, or run the model, e.g., source of input values, their manipulation prior to input into a model, quality control, risk of bias in data, etc.</p>
<p>Credibility of a model</p> <p>Alternative terms:</p> <ul style="list-style-type: none"> ▪ quality of a model ▪ risk of bias in a model ▪ validity of a model 	<p>To avoid confusion and keep with terminology used by modelling community[7] we suggest using the term <i>credibility</i> rather than <i>quality</i> of a model. The concept refers to the characteristics of a model itself – its design or execution – that affect the risk that the results may overestimate or underestimate the true effect. Various factors influence the overall credibility of a model, such as its structure, the analysis and the validation of the assumptions made during modelling.</p>
<p>Quality of reporting</p>	<p>Refers to how comprehensively and clearly model inputs, a model itself, and model outputs have been documented and described such that they can be critically evaluated and used for decision-making. Quality of reporting and quality of a model are separate concepts: a model with a low quality of reporting is not necessarily a low-quality model and vice versa.</p>
<p>Directness</p> <p>Directness of a model</p> <p>Alternative terms:</p> <ul style="list-style-type: none"> ▪ relevance ▪ external validity ▪ applicability ▪ generalizability ▪ transferability ▪ translatability 	<p>By directness of a model we mean the extent to which the model represents the real-life situation being modelled which is dependent on how well the input data and the model structure reflect the scenario of interest.</p> <p>Directness is the term used in the GRADE approach, because each of the alternatives has been used usually in a narrower meaning.</p>

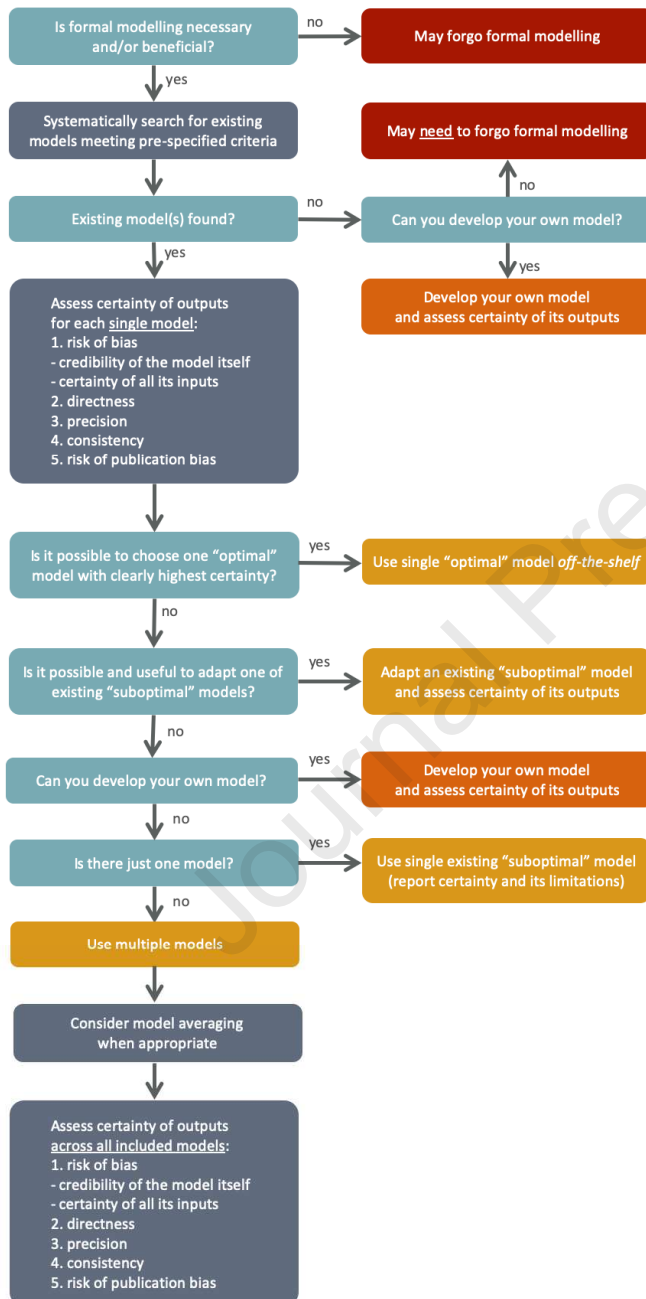
GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

626 * There may be either subtle or fundamental differences among some disciplines in how these
627 terms are being used; for the purposes of this article, these terms are generalized rather than
628 discipline specific.
629
630

Journal Pre-proof

631 **Figure 1.** The general approach to using modelled evidence and assessing its certainty in health-
 632 related disciplines.

633



634

635

636 Appendix. List of workshop participants

- 637
- 638 Elie Akl (EA)– American University of Beirut, Lebanon
- 639 Jim Bowen (JMB)– McMaster University, Canada
- 640 Chris Brinkerhoff (CB)– US Environmental Protection Agency, USA
- 641 Jan Brozek (JLB)– McMaster University, Canada
- 642 John Bucher (JB)– US National Toxicology Program, USA
- 643 Carlos Canelo-Aybar (CCA)– Iberoamerican Cochrane Centre, Spain
- 644 Marcy Card (MC)– US Environmental Protection Agency, USA
- 645 Weihsueh A. Chiu (WCh)– Texas A&M University, USA
- 646 Mark Cronin (MC)– Liverpool John Moores University, UK
- 647 Tahira Devji (TD)– McMaster University, Canada
- 648 Ben Djulbegovic (BD)– University of South Florida, USA
- 649 Ken Eng (KE)– Public Health Agency of Canada
- 650 Gerald Gartlehner (GG)– Donau-Universität Krems, Austria
- 651 Gordon Guyatt (GGu)– McMaster University, Canada
- 652 Raymond Hutubessy (RH)– World Health Organization Initiative for Vaccine Research, Switzerland
- 653 Manuela Joore (MJ)– Maastricht University, the Netherlands
- 654 Richard Judson (RJ)– US Environmental Protection Agency, USA
- 655 S. Vittal Katikireddi (SK)– University of Glasgow, UK
- 656 Nicole Kleinstreuer (NK)– US National Toxicology Program, USA
- 657 Judy LaKind (JL)– University of Maryland, USA
- 658 Miranda Langendam (ML)– University of Amsterdam, the Netherlands
- 659 Zbyszek Leś (ZL)– Evidence Prime Inc., Canada
- 660 Veena Manja (VM)– McMaster University, Canada
- 661 Joerg Meerpohl (JM)– GRADE Center Freiburg, Cochrane Germany, University Medical Center
662 Freiburg
- 663 Dominik Mertz (DM)– McMaster University, Canada
- 664 Roman Mezencev (RM)– US Environmental Protection Agency, USA
- 665 Rebecca Morgan (RMo)– McMaster University, Canada
- 666 Gian Paolo Morgano (GPM)– McMaster University, Canada
- 667 Reem Mustafa (RMu)– University of Kansas, USA
- 668 Bhash Naidoo (BN)– National Institute for Health and Clinical Excellence, UK
- 669 Martin O'Flaherty (MO)– Public Health and Policy, University of Liverpool, UK
- 670 Grace Patlewicz (GP)– US Environmental Protection Agency, USA
- 671 John Riva (JR)– McMaster University, Canada
- 672 Alan Sasso (AS)– US Environmental Protection Agency, USA
- 673 Paul Schlosser (PS)– US Environmental Protection Agency, USA

GRADE approach to modelled data • DRAFT: DO NOT DISTRIBUTE

- 674 Holger Schünemann (HJS)– McMaster University, Canada
675 Lisa Schwartz (LS)– McMaster University, Canada
676 Ian Shemilt (IS)– University College London, UK
677 Marek Smieja (MS)– McMaster University, Canada
678 Ravi Subramaniam (RS)– US Environmental Protection Agency, USA
679 Jean-Eric Tarride (JT)– McMaster University, Canada
680 Kris Thayer (KAT)– US Environmental Protection Agency, USA
681 Katya Tsaïoun (KT)– John Hopkins University, USA
682 Bernhard Ultsch (BU)– Robert Koch Institute, Germany
683 John Wambaugh (JW)– US Environmental Protection Agency, USA
684 Jessica Wignall (JWi)– ICF, USA
685 Ashley Williams (AW)– ICF, USA
686 Feng Xie (FX)– McMaster University, Canada
687

688 **References**

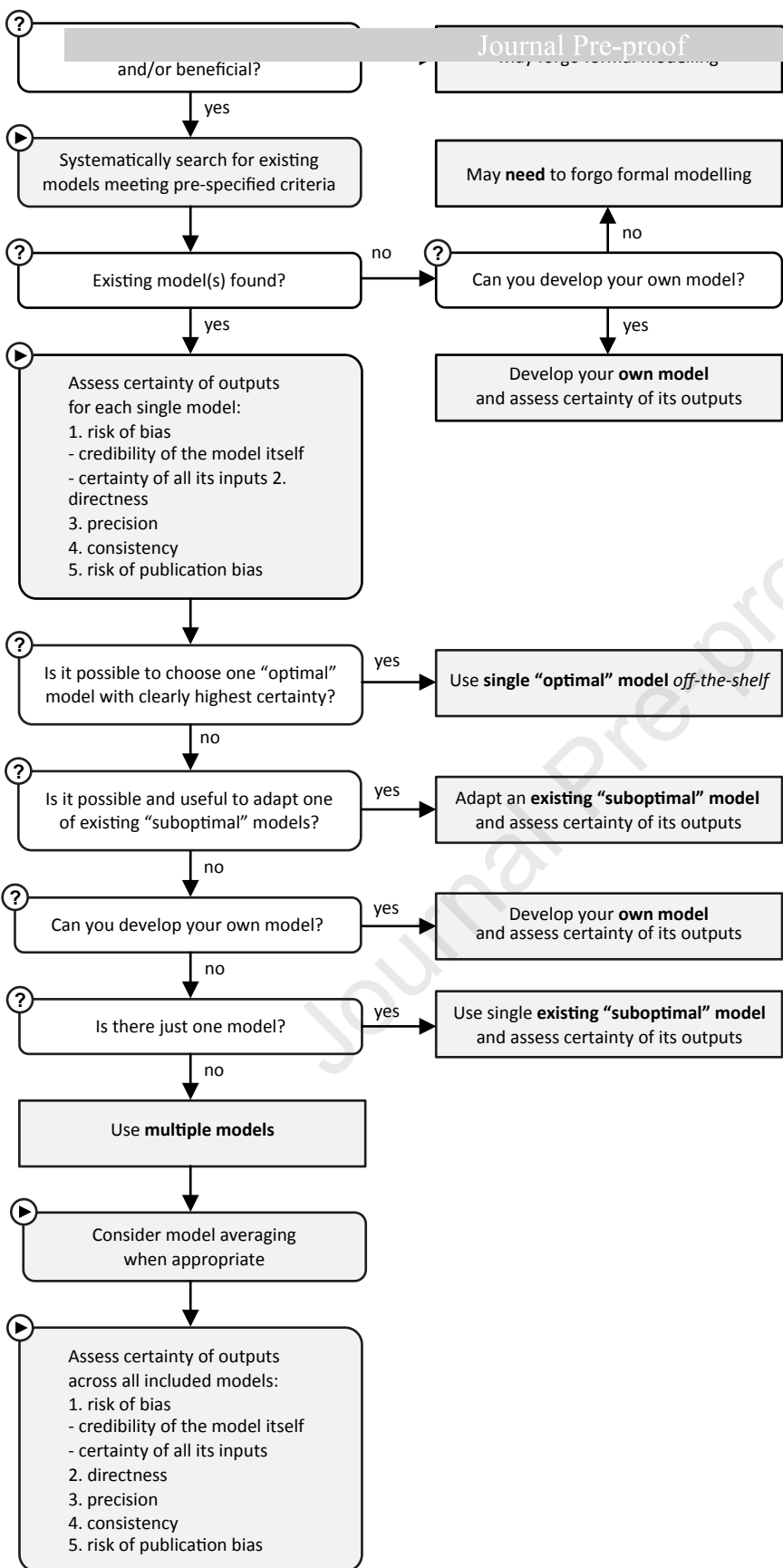
- 689
- 690 [1] Oreskes N. The role of quantitative models in science. In: Canham CD, Cole JJ, Lauenroth
691 WK, editors. *Models in ecosystem science*: Princeton University Press; 2003. p. 13–31.
- 692 [2] Frigg R, Hartmann S. Models in Science. In: Zalta EN, editor. *The Stanford Encyclopedia of*
693 *Philosophy* (Spring 2017 Edition)2017.
- 694 [3] Guyatt GH, Oxman AD, Kunz R, Vist GE, Falck-Ytter Y, Schunemann HJ, et al. What is
695 "quality of evidence" and why is it important to clinicians? *BMJ*. 2008;336:995-8.
- 696 [4] Oreskes N. Evaluation (not validation) of quantitative models. *Environ Health Perspect*.
697 1998;106 Suppl 6:1453-60.
- 698 [5] Briggs AH, Weinstein MC, Fenwick EA, Karnon J, Sculpher MJ, Paltiel AD, et al. Model
699 parameter estimation and uncertainty: a report of the ISPOR-SMDM Modeling Good Research
700 Practices Task Force--6. *Value Health*. 2012;15:835-42.
- 701 [6] Caro JJ, Briggs AH, Siebert U, Kuntz KM, Force I-SMGRPT. Modeling good research
702 practices--overview: a report of the ISPOR-SMDM Modeling Good Research Practices Task
703 Force-1. *Med Decis Making*. 2012;32:667-77.
- 704 [7] Caro JJ, Eddy DM, Kan H, Kaltz C, Patel B, Eldessouki R, et al. Questionnaire to assess
705 relevance and credibility of modeling studies for informing health care decision making: an
706 ISPOR-AMCP-NPC Good Practice Task Force report. *Value Health*. 2014;17:174-82.
- 707 [8] Eddy DM, Hollingworth W, Caro JJ, Tsevat J, McDonald KM, Wong JB, et al. Model
708 transparency and validation: a report of the ISPOR-SMDM Modeling Good Research Practices
709 Task Force-7. *Med Decis Making*. 2012;32:733-43.
- 710 [9] Karnon J, Stahl J, Brennan A, Caro JJ, Mar J, Moller J. Modeling using discrete event
711 simulation: a report of the ISPOR-SMDM Modeling Good Research Practices Task Force-4. *Med*
712 *Decis Making*. 2012;32:701-11.
- 713 [10] Marshall DA, Burgos-Liz L, MJ IJ, Crown W, Padula WV, Wong PK, et al. Selecting a
714 dynamic simulation modeling method for health care delivery research-part 2: report of the
715 ISPOR Dynamic Simulation Modeling Emerging Good Practices Task Force. *Value Health*.
716 2015;18:147-60.
- 717 [11] Marshall DA, Burgos-Liz L, MJ IJ, Osgood ND, Padula WV, Higashi MK, et al. Applying
718 dynamic simulation modeling methods in health care delivery research-the SIMULATE
719 checklist: report of the ISPOR simulation modeling emerging good practices task force. *Value*
720 *Health*. 2015;18:5-16.
- 721 [12] Pitman R, Fisman D, Zaric GS, Postma M, Kretzschmar M, Edmunds J, et al. Dynamic
722 transmission modeling: a report of the ISPOR-SMDM Modeling Good Research Practices Task
723 Force Working Group-5. *Med Decis Making*. 2012;32:712-21.
- 724 [13] Roberts M, Russell LB, Paltiel AD, Chambers M, McEwan P, Krahn M, et al.
725 Conceptualizing a model: a report of the ISPOR-SMDM Modeling Good Research Practices
726 Task Force-2. *Med Decis Making*. 2012;32:678-89.
- 727 [14] Siebert U, Alagoz O, Bayoumi AM, Jahn B, Owens DK, Cohen DJ, et al. State-transition
728 modeling: a report of the ISPOR-SMDM Modeling Good Research Practices Task Force-3. *Med*
729 *Decis Making*. 2012;32:690-700.
- 730 [15] Vemer P, van Voom GA, Ramos IC, Krabbe PF, Al MJ, Feenstra TL. Improving model
731 validation in health technology assessment: comments on guidelines of the ISPOR-SMDM
732 modeling good research practices task force. *Value Health*. 2013;16:1106-7.

- 733 [16] Weinstein MC, O'Brien B, Hornberger J, Jackson J, Johannesson M, McCabe C, et al.
734 Principles of good practice for decision analytic modeling in health-care evaluation: report of
735 the ISPOR Task Force on Good Research Practices--Modeling Studies. *Value Health*. 2003;6:9-
736 17.
- 737 [17] Bennett C, Manuel DG. Reporting guidelines for modelling studies. *BMC Med Res*
738 *Methodol*. 2012;12:168.
- 739 [18] Peñaloza Ramos MC, Barton P, Jowett S, Sutton AJ. A Systematic Review of Research
740 Guidelines in Decision-Analytic Modeling. *Value Health*. 2015;18:512-29.
- 741 [19] Philips Z, Bojke L, Sculpher M, Claxton K, Golder S. Good practice guidelines for decision-
742 analytic modelling in health technology assessment: a review and consolidation of quality
743 assessment. *Pharmacoeconomics*. 2006;24:355-71.
- 744 [20] LaKind JS, O'Mahony C, Armstrong T, Tibaldi R, Blount BC, Naiman DQ. ExpoQual:
745 Evaluating measured and modeled human exposure data. *Environ Res*. 2019;171:302-12.
- 746 [21] Husereau D, Drummond M, Petrou S, Carswell C, Moher D, Greenberg D, et al.
747 Consolidated Health Economic Evaluation Reporting Standards (CHEERS) statement. *BMJ*.
748 2013;346:f1049.
- 749 [22] Balshem H, Helfand M, Schunemann HJ, Oxman AD, Kunz R, Brozek J, et al. GRADE
750 guidelines: 3. Rating the quality of evidence. *J Clin Epidemiol*. 2011;64:401-6.
- 751 [23] Guyatt GH, Oxman AD, Kunz R, Brozek J, Alonso-Coello P, Rind D, et al. GRADE guidelines:
752 6. Rating the quality of evidence--imprecision. *J Clin Epidemiol*. 2011;64:1283-93.
- 753 [24] Guyatt GH, Oxman AD, Kunz R, Woodcock J, Brozek J, Helfand M, et al. GRADE guidelines:
754 8. Rating the quality of evidence--indirectness. *J Clin Epidemiol*. 2011;64:1303-10.
- 755 [25] Guyatt GH, Oxman AD, Kunz R, Woodcock J, Brozek J, Helfand M, et al. GRADE guidelines:
756 7. Rating the quality of evidence--inconsistency. *J Clin Epidemiol*. 2011;64:1294-302.
- 757 [26] Guyatt GH, Oxman AD, Montori V, Vist G, Kunz R, Brozek J, et al. GRADE guidelines: 5.
758 Rating the quality of evidence--publication bias. *J Clin Epidemiol*. 2011;64:1277-82.
- 759 [27] Guyatt GH, Oxman AD, Sultan S, Glasziou P, Akl EA, Alonso-Coello P, et al. GRADE
760 guidelines: 9. Rating up the quality of evidence. *J Clin Epidemiol*. 2011;64:1311-6.
- 761 [28] Guyatt GH, Oxman AD, Vist G, Kunz R, Brozek J, Alonso-Coello P, et al. GRADE guidelines:
762 4. Rating the quality of evidence--study limitations (risk of bias). *J Clin Epidemiol*.
763 2011;64:407-15.
- 764 [29] Lasserson TJ, Thomas J, Higgins JPT. Chapter 1: Starting a review. In: Higgins JPT, Thomas
765 J, Chandler J, Cumpston M, Li T, Page MJ, et al., editors. *Cochrane Handbook for Systematic*
766 *Reviews of Interventions* version 60 (updated July 2019): Cochrane; 2019.
- 767 [30] Eykhoff P. System identification: parameter and state estimation: Wiley-Interscience;
768 1974.
- 769 [31] Schunemann HJ, Best D, Vist G, Oxman AD, Group GW. Letters, numbers, symbols and
770 words: how to communicate grades of evidence and recommendations. *CMAJ*. 2003;169:677-
771 80.
- 772 [32] Schunemann HJ, Mustafa R, Brozek J, Santesso N, Alonso-Coello P, Guyatt G, et al. GRADE
773 Guidelines: 16. GRADE evidence to decision frameworks for tests in clinical practice and
774 public health. *J Clin Epidemiol*. 2016;76:89-98.
- 775 [33] Schunemann HJ, Oxman AD, Brozek J, Glasziou P, Jaeschke R, Vist GE, et al. Grading
776 quality of evidence and strength of recommendations for diagnostic tests and strategies. *BMJ*.
777 2008;336:1106-10.

- 778 [34] Iorio A, Spencer FA, Falavigna M, Alba C, Lang E, Burnand B, et al. Use of GRADE for
779 assessment of evidence about prognosis: rating confidence in estimates of event rates in
780 broad categories of patients. *BMJ*. 2015;350:h870.
- 781 [35] Hooijmans CR, de Vries RBM, Ritskes-Hoitinga M, Rovers MM, Leeflang MM, IntHout J, et
782 al. Facilitating healthcare decisions by assessing the certainty in the evidence from preclinical
783 animal studies. *PLoS One*. 2018;13:e0187271.
- 784 [36] Brunetti M, Shemilt I, Pregno S, Vale L, Oxman AD, Lord J, et al. GRADE guidelines: 10.
785 Considering resource use and rating the quality of economic evidence. *J Clin Epidemiol*.
786 2013;66:140-50.
- 787 [37] Zhang Y, Alonso-Coello P, Guyatt GH, Yepes-Nunez JJ, Akl EA, Hazlewood G, et al. GRADE
788 Guidelines: 19. Assessing the certainty of evidence in the importance of outcomes or values
789 and preferences-Risk of bias and indirectness. *J Clin Epidemiol*. 2018.
- 790 [38] Zhang Y, Coello PA, Guyatt GH, Yepes-Nunez JJ, Akl EA, Hazlewood G, et al. GRADE
791 guidelines: 20. Assessing the certainty of evidence in the importance of outcomes or values
792 and preferences-inconsistency, imprecision, and other domains. *J Clin Epidemiol*. 2018.
- 793 [39] World Health Organization. WHO guidelines for screening and treatment of precancerous
794 lesions for cervical cancer prevention. Geneva, Switzerland: World Health Organization; 2013.
- 795 [40] Thayer KA, Schunemann HJ. Using GRADE to respond to health questions with different
796 levels of urgency. *Environ Int*. 2016;92-93:585-9.
- 797 [41] Porgo TV, Norris SL, Salanti G, Johnson LF, Simpson JA, Low N, et al. The use of
798 mathematical modeling studies for evidence synthesis and guideline development: A glossary.
799 *Res Synth Methods*. 2019;10:125-33.
- 800 [42] (NICE) NifHaCE. The reference case. Guide to the methods of technology appraisal 2013:
801 NICE; 2013.
- 802 [43] Eyles H, Ni Mhurchu C, Nghiem N, Blakely T. Food pricing strategies, population diets,
803 and non-communicable disease: a systematic review of simulation studies. *PLoS Med*.
804 2012;9:e1001353.
- 805 [44] Jaime Caro J, Eddy DM, Kan H, Kaltz C, Patel B, Eldessouki R, et al. Questionnaire to assess
806 relevance and credibility of modeling studies for informing health care decision making: an
807 ISPOR-AMCP-NPC Good Practice Task Force report. *Value Health*. 2014;17:174-82.
- 808 [45] (NICE) NifHaCE. Appendix G: Methodology checklist: economic evaluations. The
809 guidelines manual: NICE; 2012.
- 810 [46] Schultz TW, Richarz A-N, Cronin MTD. Assessing uncertainty in read-across: Questions to
811 evaluate toxicity predictions based on knowledge gained from case studies. *Computational*
812 *Toxicology*. 2019;9:1-11.
- 813 [47] Cronin MTD, Richarz AN, Schultz TW. Identification and description of the uncertainty,
814 variability, bias and influence in quantitative structure-activity relationships (QSARs) for
815 toxicity prediction. *Regul Toxicol Pharmacol*. 2019;106:90-104.
- 816 [48] Brazier J, Ara R, Azzabi I, Busschbach J, Chevrou-Severac H, Crawford B, et al.
817 Identification, Review, and Use of Health State Utilities in Cost-Effectiveness Models: An
818 ISPOR Good Practices for Outcomes Research Task Force Report. *Value Health*. 2019;22:267-
819 75.
- 820 [49] Kaltenthaler E, Tappenden P, Paisley S, Squires H. NICE DSU Technical Support
821 Document 13: Identifying and Reviewing Evidence to Inform the Conceptualisation and
822 Population of Cost-Effectiveness Models. London 2011.

- 823 [50] Paisley S. Identification of Evidence for Key Parameters in Decision-Analytic Models of
824 Cost Effectiveness: A Description of Sources and a Recommended Minimum Search
825 Requirement. *Pharmacoeconomics*. 2016;34:597-608.
- 826 [51] Guyatt G, Oxman AD, Sultan S, Brozek J, Glasziou P, Alonso-Coello P, et al. GRADE
827 guidelines: 11. Making an overall rating of confidence in effect estimates for a single outcome
828 and for all outcomes. *J Clin Epidemiol*. 2013;66:151-7.
- 829 [52] Bilcke J, Beutels P, Brisson M, Jit M. Accounting for methodological, structural, and
830 parameter uncertainty in decision-analytic models: a practical guide. *Med Decis Making*.
831 2011;31:675-92.
- 832 [53] Saltelli A, Ratto M, Andres T, Campolongo F, Cariboni J, Gatelli D, Saisana M & Tarantola S.
833 2008. *Global sensitivity analysis. The primer*. Chichester, UK: John Wiley & Sons.
- 834 [54] Page MJ, Higgins JPT, Sterne JAC. Chapter 13: Assessing risk of bias due to missing results
835 in a synthesis. In: Higgins JPT, Thomas J, Chandler J, Cumpston M, Li T, Page MJ, et al., editors.
836 *Cochrane Handbook for Systematic Reviews of Interventions version 60 (updated July 2019)*:
837 Cochrane; 2019.
- 838 [55] Schünemann HJ, Lerda D, Quinn C, Follmann M, Alonso-Coello P, Rossi PG, et al. Breast
839 Cancer Screening and Diagnosis: A Synopsis of the European Breast Guidelines. *Annals of*
840 *Internal Medicine*. 2020;172:46-56.
- 841 [56] Eaton JW, Johnson LF, Salomon JA, Barnighausen T, Bendavid E, Bershteyn A, et al. HIV
842 treatment as prevention: systematic comparison of mathematical models of the potential
843 impact of antiretroviral therapy on HIV incidence in South Africa. *PLoS Med*.
844 2012;9:e1001245.
- 845 [57] Gomersall JS, Jadotte YT, Xue Y, Lockwood S, Riddle D, Preda A. Conducting systematic
846 reviews of economic evaluations. *Int J Evid Based Healthc*. 2015;13:170-8.
- 847 [58] Mandelblatt JS, Stout NK, Schechter CB, van den Broek JJ, Miglioretti DL, Krapcho M, et al.
848 Collaborative Modeling of the Benefits and Harms Associated With Different U.S. Breast
849 Cancer Screening Strategies. *Ann Intern Med*. 2016;164:215-25.
- 850 [59] Davies NG, Kucharski AJ, Eggo RM, Gimma A, Edmunds WJ, Centre for the Mathematical
851 Modelling of Infectious Diseases C-wg. Effects of non-pharmaceutical interventions on COVID-
852 19 cases, deaths, and demand for hospital services in the UK: a modelling study. *Lancet Public*
853 *Health*. 2020;5:e375-e85.
- 854 [60] Tibaldi R, ten Berge W, Drolet D. Dermal absorption of chemicals: estimation by IH
855 SkinPerm. *J Occup Environ Hyg*. 2014;11:19-31.
- 856 [61] Young BM, Tulse NS, Egeghy PP, Driver JH, Zartarian VG, Johnston JE, et al. Comparison of
857 four probabilistic models (CARES((R)), Calendex, ConsExpo, and SHEDS) to estimate
858 aggregate residential exposures to pesticides. *J Expo Sci Environ Epidemiol*. 2012;22:522-32.
- 859 [62] United States Environmental Protection Agency. Human Exposure Modeling - Overview.
860 In: United States Environmental Protection Agency, editor.
- 861 [63] Levin S, Dugas A, Gurses A, Kirsch T, Kelen G, Hinson J, et al. HOPSCORE: AN ELECTRONIC
862 OUTCOMES-BASED EMERGENCY TRIAGE SYSTEM. Agency for Healthcare Research and
863 Quality; 2018.
- 864 [64] Smith RD, Keogh-Brown MR, Barnett T, Tait J. The economy-wide impact of pandemic
865 influenza on the UK: a computable general equilibrium modelling experiment. *BMJ*.
866 2009;339:b4571.
- 867 [65] Hultcrantz M, Rind D, Akl EA, Treweek S, Mustafa RA, Iorio A, et al. The GRADE Working
868 Group clarifies the construct of certainty of evidence. *J Clin Epidemiol*. 2017;87:4-13.

Journal Pre-proof



What is new

1. General concepts determining the certainty of evidence in the GRADE approach (risk of bias, indirectness, inconsistency, imprecision, reporting bias, magnitude of an effect, dose-response relation, and the direction of residual confounding) also apply in the context of assessing the certainty of evidence from models (model outputs).
2. Detailed assessment of the certainty of evidence from models differs for the assessment of outputs from a single model compared to the assessment of outputs across multiple models.
3. We propose a framework for selecting the best available evidence from models to inform health care decisions: to develop a model de novo, to identify an existing model the outputs of which provide the highest certainty evidence, or to use outputs from multiple models.
4. We suggest that the modelling and health care decision making communities collaborate further to clarify terminology used in the context of modelling and make it consistent across the disciplines to facilitate communication.