



<b>Title</b>	Technical assessment and evaluation of environmental models and software : letter to the Editor
<b>Authors(s)</b>	Alexandrov, G. A., Ames, D., Bellocchi, G., Bruen, Michael, Crout, N., Erechtkoukova, M., Hildebrandt, A., Hoffman, F., Jackisch, C., Khaiteer, P., Mannina, G., Matsunaga, T., Purucker, S. T., Rivington, M., Samaniego, L.
<b>Publication date</b>	2011-03
<b>Publication information</b>	Alexandrov, G. A., D. Ames, G. Bellocchi, Michael Bruen, N. Crout, M. Erechtkoukova, A. Hildebrandt, et al. "Technical Assessment and Evaluation of Environmental Models and Software : Letter to the Editor" 26, no. 3 (March, 2011).
<b>Publisher</b>	Elsevier
<b>Item record/more information</b>	<a href="http://hdl.handle.net/10197/3064">http://hdl.handle.net/10197/3064</a>
<b>Publisher's statement</b>	This is the author's version of a work that was accepted for publication in Environmental Modelling & Software. Changes resulting from the publishing process, such as peer review, editing, corrections, structural formatting, and other quality control mechanisms may not be reflected in this document. Changes may have been made to this work since it was submitted for publication. A definitive version was subsequently published in Environmental Modelling & Software Volume 26, Issue 3, March 2011, Pages 328-336 DOI: 10.1016/j.envsoft.2010.08.004
<b>Publisher's version (DOI)</b>	10.1016/j.envsoft.2010.08.004

Downloaded 2023-10-02T04:02:14Z

The UCD community has made this article openly available. Please share how this access benefits you. Your story matters! (@ucd\_oa)



© Some rights reserved. For more information



23 <sup>9</sup> Oak Ridge National Laboratory, Computational Earth Sciences Group, PO Box 2008,  
24 Oak Ridge, TN 37831, USA

25

26 <sup>10</sup> Technical University Munich, Munich, Germany.

27

28 <sup>12</sup> Department of Hydraulic Engineering and Environmental Applications, University of  
29 Palermo, Italy.

30

31 <sup>14</sup> US Environmental Protection Agency, Athens, Georgia, USA.

32

33 <sup>15</sup> Macaulay Land Use Research Institute, Craigiebuckler, Aberdeen. AB15 8QH. United  
34 Kingdom.

35

36

37

38 \*Corresponding author [g.alexandrov@nies.go.jp](mailto:g.alexandrov@nies.go.jp)

39

40

## 41 **Rationale**

42 This letter details the collective views of a number of independent researchers on the  
43 technical assessment and evaluation of environmental models and software. The purpose  
44 is to stimulate debate and initiate action that leads to an improved quality of model  
45 development and evaluation, so increasing the capacity for models to have positive  
46 outcomes from their use. As such, we emphasise the relationship between the model  
47 evaluation process and credibility with stakeholders (including funding agencies) with a  
48 view to ensure continued support for modelling efforts.

49 Many journals, including EM&S, publish the results of environmental modelling studies  
50 and must judge the work and the submitted papers based solely on the material that the  
51 authors have chosen to present and on how they present it. There is considerable variation  
52 in how this is done with the consequent risk of considerable variation in the quality and  
53 usefulness of the resulting publication. Part of the problem is that the review process is  
54 reactive, responding to the submitted manuscript. In this letter, we attempt to be proactive  
55 and give guidelines for researchers, authors and reviewers as to what constitutes best  
56 practice in presenting environmental modelling results. This is a unique contribution to  
57 the organisation and practice of model-based research and the communication of its  
58 results that will benefit the entire environmental modelling community. For a start, our  
59 view is that the community of environmental modellers should have a common vision of  
60 minimum standards that an environmental model must meet. A common vision of what a  
61 good model should be is expressed in various guidelines on Good Modelling Practice.  
62 The guidelines prompt modellers to codify their practice and to be more rigorous in their  
63 model testing. Our statement within this letter deals with another aspect of the issue - it

64 prompts professional journals to codify the peer-review process. Introducing a more  
65 formalized approach to peer-review may discourage reviewers from accepting invitations  
66 to review given the additional time and labour requirements. The burden of proving  
67 model credibility is thus shifted to the authors. Here we discuss how to reduce this burden  
68 by selecting realistic evaluation criteria and conclude by advocating the use of  
69 standardized evaluation tools as this is a key issue that needs to be tackled.

70

71

72

## 73 **Background**

74

75 The use of models for any practical purpose entails the risk of misuse. If a model's  
76 limitations are not completely understood, the model outputs may be easily  
77 misinterpreted (Jakeman et al., 2009). To reduce this risk, every model should be  
78 assessed and evaluated by domain experts - that is, by modellers experienced in model  
79 development and application.

80

81 Such assessment and evaluation is normally undertaken when an article describing a  
82 model (or software) passes through a peer-review system of a professional journal (Fig 1)  
83 such as EM&S, for example. Most experienced modellers are involved in the peer-review  
84 process and periodically evaluate models made by their colleagues. Therefore, the  
85 community of environmental modellers needs to have a common vision of minimum  
86 standards that an environmental model must meet.

87

88 A common vision of what a good model is has been expressed in guidelines on Good  
89 Modelling Practice (STOWA/RIZA, 1999; Murray-Darling Basin Commission, 2000;  
90 Jakeman et al., 2006; Gaber et al., 2008; Robson et al., 2008; Welsch et al., 2008). The  
91 guidelines prompt modellers to codify their practice and to be more rigorous in their  
92 model testing. The purpose of this letter is to deal with another aspect of the issue: that of  
93 prompting professional journals to codify the peer-review process. The objective is to  
94 highlight the obstacles to model evaluation and to provide possible solutions. This is not  
95 however a review on the state of model evaluation, the latter being given in a recent

96 paper by Bellocchi et al. (2010). Rather we seek to promote improvement within the  
97 quality of model evaluation as part of the peer-review process. In doing so, we have also  
98 highlighted issues that potential reviewers need to be aware of.

99

100 Peer review is normally considered as an essential component of research dissemination  
101 and remains the principal mechanism by which the quality of research is judged (Council  
102 of Science Editors, 2006; Müller, 2009). At the same time, there is common  
103 understanding that peer-review cannot be expected to detect fraud and ensure perfection  
104 (Hames, 2007): *“even the most-respected journals have been caught out and, despite  
105 extensive peer review, have ended up publishing fraudulent or seriously flawed material”*  
106 (Wager, 2006). Then, what is the main purpose of peer-review? There is no general  
107 agreement on this issue now. One may suggest that the peer-review system initially  
108 introduced for filtering out unreasonable claims to new research results still serves this  
109 purpose (Walker, 1998; Alexandrov, 2006).

110

111 In the case of complex models, it is likely the process will result in reviews not  
112 evaluating the components of the model (especially if referenced to other sources), on the  
113 defence that journal readers and end users who are specialists on those components will  
114 make their own judgments. In reality, reviewers will not have sufficient time or resources  
115 to conduct detailed evaluation of individual components or the whole model let alone the  
116 software coding. Hence, the emphasis must be on model developers to provide accurate  
117 evidence, covering sufficient complexity interactions, to demonstrate adequate testing to  
118 achieve a stated level of model reliability and utility. We feel this requires clearer

119 statements (supported by evidence) from the model developers on known limitations and  
120 areas of uncertainty. Such an open approach should, if communicated correctly, i.e.  
121 positively through addressing the consequences of any uncertainty, actually increase  
122 credibility with end users rather than diminish it. This is important as perceptions of  
123 uncertainty and how it is handled amongst researchers, policy makers and politicians  
124 have changed recently, especially since the rise of climate change modelling and  
125 planning for adaptation. Previously end users (particularly policy makers and politicians)  
126 were reluctant to deal with the realities of the uncertainty associated with modelling.  
127 Reduction of uncertainties in the models may be pursued by, for instance, model-data  
128 integration techniques (e.g. Wang et al., 2009). However, it is clear that the skill in  
129 handling uncertainty not only lies within statistical and other forms of model testing, but  
130 also in how it is communicated to end users. Hence, our view that the establishment of a  
131 standardized set of criteria and methods for model evaluation is needed to set a minimum  
132 standard for 'proof of testing' that would serve to support uncertainty communication. The  
133 absence of such standardized criteria and methods risk modelling becoming unacceptable  
134 as a form of research for predictive purposes.

135

## 136 **The major obstacles to a more formalized approach to model** 137 **evaluation**

138

139 Currently there is no requirement to detail a full set of model specifications. The first task  
140 in standardization therefore would be to have a scheme whereby published models could



141 be fully specified (Fig 2). Depending on the rationale for the research exercise, modelling  
142 for theoretical scientific purposes and modelling for decision-making may follow  
143 separate paths (Haag and Kaupenjohann, 2001) and hence require different specifications  
144 for assessment and evaluation. In general, this would require a definition of the modelling  
145 objective, its formulation, implementation and parameterization, further supplemented by  
146 information on how they have been evaluated and the conclusions of that process.

147

148 A standardized set of criteria with which models should be assessed and evaluated will at  
149 least ensure the minimum of review effort is made. The risk, however, is that a more  
150 formalized approach to peer review requiring the achievement of a minimum standard  
151 may discourage reviewers from accepting invitations to review given the additional time  
152 and labour requirements. Further to this, a limitation of past model development and  
153 application has been that funding organizations have been reluctant to accept the  
154 additional costs of performing appropriate model assessment and evaluation within  
155 proposals from researchers. Hence, model development has often been on tight budgets  
156 causing assessment and evaluation to take a lower priority. Specification of minimum  
157 standards for assessment and evaluation and how it is reported should hence also form the  
158 basis for the minimum level of validation effort written into funding proposals. Similarly,  
159 organizations awarding grants need to include in their calls for proposals more explicit  
160 details on the requirements for evaluation and testing and be prepared to provide the  
161 required funds.

162

163 This implies that funding organizations have to take on board a greater level of  
164 responsibility in supporting modelling work that includes increased assessment and  
165 evaluation efforts. Funders could effectively impose minimum standards for model  
166 development, evaluation and specification, and progress could be made in this area if  
167 'lead' funders introduced such requirements. For example, in the United Kingdom  
168 context if the national research councils were to introduce minimum standards other  
169 funders would follow in time.

170

171 A further obstacle to overcome lies within the community of environmental modellers  
172 itself, which has to take on board a greater level of responsibility in developing standards.  
173 The procedures to perform the evaluation task are not widely accepted (Cheng et al.,  
174 1991) and appear in several forms, depending on data availability, system characteristics  
175 and researchers' opinion (Hsu et al., 1999). Environmental models are made up of  
176 mixtures of rate equations, comprise approaches with different levels of empiricism, aim  
177 at simulating systems which show non-linear behaviour and often require numerical  
178 rather than analytical solutions. Therefore, the computer program, including technical  
179 issues and possible errors, is tested rather than the mathematical model representing the  
180 system (Leffelaar et al., 2003). Hence, given the applied nature of models in representing  
181 a system, their usefulness can be evaluated only in specific case studies. Gardner and  
182 Urban (2003) suggested assessing model usefulness based on its appropriateness and  
183 performance. Model appropriateness describes the extent to which the model meets the  
184 objectives of the study. The appropriateness usually deals with the model structure,  
185 although the necessity to identify model parameters brings observation data into the scene

186 (e.g. Confalonieri et al., 2009b). The availability or unavailability of observational data  
187 largely predetermines the structure of a model. Model performance is evaluated based on  
188 reported testing results in such terms as “goodness of fit” between simulated values of  
189 model variables and observation data and required computational time. Our observation  
190 is that quantification of uncertainty is less often reported.

191

192 Evaluation of model uncertainty is an important part of model assessment, yet application  
193 specific since it depends on model parameterization. On different sets of parameter  
194 values, the same mathematical equations may exhibit substantially different dynamic  
195 features. Thus, in dynamic models, changes in model parameters can trigger a switch  
196 from stable solution to an unstable one, causing a significant increase in model  
197 uncertainty (e.g. van Nes and Scheffer, 2003). Stability analysis of a solution must be a  
198 part of model investigation, but the analysis may become complicated for complex  
199 models, and has therefore not often been undertaken.

200

201 However, we recognise that a complete evaluation of model uncertainty is hardly  
202 possible. Usually, the analysis is confined to quantifiable sources, such as initial values of  
203 state variables and parameters. Indeed as pointed out by Harremoës (2003) not all  
204 uncertainty sources can be ‘quantified’, and that the fraction of uncertainty source terms  
205 being ‘ignored’ might be high in environmental investigations. The investigation of the  
206 model structural uncertainty is uncommon. Even if an estimate of uncertainty is obtained,  
207 its interpretation is not straightforward. The term ‘high uncertainty’ is ambiguous and  
208 was defined rather intuitively. Reichert and Borsuk (2005) considered the uncertainty as

209 'high' when the width of predicted distribution of model solutions is larger than the  
210 difference between expected outcomes of different simulated alternatives. Strictly  
211 speaking, an absolute value of the uncertainty is not important as long as simulations  
212 allow for a clear distinction between considered scenarios and for comparison of  
213 projected outcomes against some known objectives. In other words, the interpretation of  
214 model uncertainty is also application dependent. Codifying the testing process by the  
215 model authors will establish an uncertainty range of at least one application case.

216

217 It is often stated that a clear understanding of the model's purpose is central to its  
218 evaluation. In other words, it should be 'fit for purpose'. In fact the use of this phrase can  
219 be helpful, as it places the purpose 'up front' and emphasizes that generally perfection is  
220 not sought, merely the fitness for the given purpose. It is important to distinguish between  
221 cases where the purpose is prediction to underpin a decision-making process, and those  
222 cases where the model serves as a test bed for scientific hypotheses (even though the  
223 same underlying model may be used in both situations). In the decision support case, the  
224 accuracy of the predictions for the given purpose is important, as are other factors such as  
225 the input data requirements, the safe operating domain of the model and stakeholders'  
226 acceptance of the model. In the scientific method case, the evaluation generally needs to  
227 be more sophisticated. It is not enough to confirm the hypotheses contained in the model  
228 based on the agreement between predictions and observations. A further test is required  
229 to rule out the possibility that alternative model formulations (i.e. different hypotheses)  
230 could also have described the observations available. This relates to the 'equifinality'

231 thesis of Beven and Freer (Beven and Freer, 2001; Beven, 2006) and the issue of choices  
232 in model formulation (e.g. Cox et al., 2006; Crout et al., 2009).

233

## 234 **The solutions: realistic criteria for model evaluation**

235 Since the range of modelling situations is wide, we recognise that generally applicable  
236 standards can be formulated only in a generic form. They form a framework for model  
237 evaluation leaving the details of a particular implementation, such as quantification of  
238 criteria, up to the reviewers. Starting with the technical assessment of a model, we  
239 suggest reviewers first answer the questions below and evaluate the developers' claim on  
240 the usefulness of an environmental model:

- 241 • Do developers delineate the domain of model application?
- 242 • Do they highlight advanced model features against the prior art?
- 243 • Do they provide an example of model application illustrating model performance?

244 Then, they may proceed to assessment of the “proofs” of model usability, which are  
245 expected to show that:

- 246 • The domain of model application is delineated correctly;
- 247 • The model has certain advantages over a prior art;
- 248 • The example of model application shows credibility of the model as a tool for  
249 environmental assessment.

250 In the next sections, we expand our views on the above points.

251 ***How to delineate the domain of model application***

252 An environmental model is normally developed using a four-tier approach: conceptual  
253 scheme, model formulation, computer code, and specific parameterization. Consequently,  
254 delineating the domain of its applicability one should clearly make a distinction between  
255 applicability of each tier. A conceptual scheme may be applied over a large range of  
256 environmental states, whereas its specific parameterization may be intended for use under  
257 very restrictive conditions. In addition, the model code may be suitable for use only  
258 within a certain range of model parameters and inputs.

259

260 The four-tier description of model domain should answer the following questions:

261

- 262 • Which environmental states may fall within the conceptual scope of the model?
- 263 • Which environmental states may be assessed (or explored) using the current  
264 version of a model in question or its computer code?
- 265 • Which environmental states may be assessed (or explored) using a specific  
266 parameterization of the model?

267

268 The conceptual scheme of a model is derived from the model developer's perceptual  
269 model of the real system at hand. The perceptual model is known to be an approximation  
270 (to a greater or lesser extent). Moreover, it is common to have a range of scientific  
271 opinions regarding the best representation of the perceptual model. Nevertheless, it is  
272 always possible to make the distinction between environmental conditions that may fall  
273 within the conceptual scope of the model and those that may not. For example, if the

274 conceptual scheme of a model does not address some environmental factors, the model  
275 may not assess the environmental impact of this factor.

276

277 In general, the correct description of the model domain must guarantee that the model  
278 will not produce results that go beyond empirically (or theoretically) established bounds.

279 A related part of the technical assessment is to find the “regions” of the declared model  
280 domain, where the model produces obviously erroneous results, or confirm that no such  
281 “regions” were found. The latter helps to evaluate model reliability defined by Mankin et  
282 al. (1975).

283

#### 284 ***How to show that a model has an advantage over the prior art***

285 The purpose of developing a new model is to make visible progress in the state-of-the-art  
286 (Jørgensen et al., 2006). This can be done in different ways. The simplest of them is  
287 improving either the conceptual scheme or computer code of a prior model. In this case,  
288 the advantage over the prior art may be highlighted by providing some proof that:

289     • The model addresses environmental situations that do not fall within the scope of  
290       the prior model(s);

291 or that:

292     • The model code is more efficient than that of the prior model(s) (e.g., needs less  
293       initial information) in addressing some environmental situations;

294 or that:

295     • The specific parameterization of the model shows better performance than that of  
296       the prior model(s) in addressing specific environmental conditions.

297

298 In the well-developed fields of environmental modelling, the multi-model approach is  
299 considered to be more reasonable than the best-model approach. Multi-model  
300 combinations outperform best models. In other words, the progress in the state-of-the-art  
301 is achieved through improving performance of a multi-model ensemble.

302

303 The examples of testing in such cases include multi-model analysis (MMA) for  
304 developing multiple plausible models by considering alternative processes, using  
305 alternative modelling codes, or by defining alternative boundary conditions (Pachepsky et  
306 al., 2006). Quantitative MMA methods assign performance scores to each candidate  
307 model (e.g. Burnham and Anderson, 2002; Ye et al., 2008). The scores are utilized to  
308 rank and select the best models or to assign importance weights (e.g., for use in an  
309 ensemble forecasting). Qualitative MMA methods can also rely on expert elicitation,  
310 stakeholder involvement, and quality assurance/quality control procedures to assess  
311 relative merits of alternative models (Funtowicz and Ravetz, 1990; van der Sluijs, 2007).

312

313 Improving the mathematical formulation of a given conceptual scheme is also a way for  
314 improving the state-of-the art. The selection of a suitable formulation relates to model  
315 comparisons that cannot be fully ‘automated’ or formalized due to a confounding effect.  
316 Confounding appears when two or more factors cause a combined measurable effect and  
317 the contribution of each individual factor cannot be estimated separately. Thus, a  
318 particular value of a model parameter depends not only on the corresponding state  
319 variable and processes included in the model, but also on a given formula used to



320 describe each process. The majority of environmental models require a number of  
321 parameters that must be identified for a given case study. In such a case, the comparison  
322 of different models becomes dubious because it is hard to differentiate (in the overall  
323 model uncertainty) the effect created by model structure from the effect generated by the  
324 assigned values of model parameters.

325

326 Moreover, even a small change in a sub-model introduced to correct its functionality may  
327 produce a different interpretation on simulated processes. The reason for these unwanted  
328 changes lies in the lack of independence/wrong dependencies of parts of the code, which  
329 is not completely avoidable. This aspect might go beyond a simple evaluation by once  
330 again comparing against previously acceptable results (Huth and Holzworth, 2005) and  
331 poses the need for formal model evaluation against observed data at each published stage  
332 of model development (van Oijen, 2002). Each version of a model, throughout its  
333 development life cycle, should be subjected to output testing, designed by identifying test  
334 scenarios, test cases, and/or test data.

335

### 336 ***How to show model credibility***

337 The establishment of credibility is a prerequisite for model acceptance and use.  
338 Credibility is in itself a complex issue extending beyond just model testing (e.g  
339 authenticity of problem ownership, skills and motivation of the research team developing  
340 a model, etc.). Model evaluation is, however, the key starting point for establishing  
341 credibility. Hence, a strengthened peer review procedure will have an essential role in the  
342 credibility building process. However, model evaluation must not be seen as a one-off

343 event or a “once-and-for-all” activity (Janssen and Heuberger, 1995), but as an on-going  
344 process to check for model compatibility to current evidence and variations (e.g. in  
345 spatial, climatic and hydrological conditions). Moreover, according to Sinclair and  
346 Seligman (2000), demonstration that a models’ output more or less fits a set of data is a  
347 necessary but not sufficient indication of validity. This is because model validity is rather  
348 the capability to analyze, clarify, and solve empirical and conceptual problems. Empirical  
349 problems in a domain are, in general, about the observable world in need of explanation  
350 because a model does not adequately solve it, rival models solve it in different ways, or it  
351 is solved/unsolved depending on the model. Conceptual problems arise when the  
352 concepts within a model appear to be logically inconsistent, vague and unclear, or  
353 circularly defined, and when the definition of some phenomenon in a model is hard to  
354 harmonize with an ordinary language or definition (e.g. Parker, 2001). This raises the  
355 issue of widening beyond numerical testing by also including stakeholders’ evaluation  
356 and expert interpretation through soft systems approaches (Bellocchi et al., 2002;  
357 Matthews et al., 2008). For example, non-scientific end users may be more persuaded of  
358 model validity by graphical representations than statistical tests or indices, especially  
359 where historical events or familiar phenomena are shown and are recognizable by them.

360

361 Thus, to evaluate a model as a credible one, a reviewer should confirm at least that:

- 362 • Its conceptual scheme is theoretically adequate to the declared domain of  
363 applicability;
- 364 • Its computer code is verifiable;
- 365 • The accuracy of its specific parameterization is consistent with intended usage.

366

367 **Adequacy and prediction**

368 The model adequacy cannot be assessed regardless of the domain of its applicability  
369 (Rykiel, 1996). The context within which models are used affects the required  
370 functionality and/or accuracy (French and Geldermann, 2005). This is particularly  
371 apparent when comparing models developed to represent the same process at different  
372 scales and for which different qualities of input, parameterization and validation data will  
373 be available, for example soil water balances at plot, farm, catchment and region (e.g.  
374 Keating et al., 2002; Vischel et al., 2007). This has led to the development of application  
375 specific testing of models and the idea of model benchmarking, by comparing simulation  
376 outputs with outputs of another simulation that is accepted as a “standard” (e.g. Vanclay,  
377 1994). Such approaches typically use multi-criteria assessment (e.g. Reynolds and Ford,  
378 1999) with performance criteria weighted by users depending on their relative  
379 importance.

380

381 Such indications of adequacy are essential in relation to the use of models for future  
382 predictive purposes. Papers on modelling often state that they aim to produce an  
383 instrument for prediction (van Oijen, 2002). A fundamental issue is to quantify the degree  
384 to which a model captures an underlying reality and predicts future cases (Marcus and  
385 Elias, 1998; Li et al., 2003). Predictions pose special problems for testing, especially if  
386 prediction focuses on events in the far future. Predictive models can be accepted if they  
387 explain past events (*ex-post* validation). The probability of making reasonable projections  
388 decreases with the length of time looked forward. A continuous exchange of validation

389 data among developers and test teams should either ensure a progressive validation of the  
390 models by time, or highlights the need for updated interpretations of the changed system.

391

392 In many cases, predictive models are mixed with exploratory models. The distinction  
393 between them can be drawn on the basis of data availability. Predictive models are  
394 normally used in connection with an observing system established for environmental  
395 monitoring. Exploratory models, in contrast, are normally used where observations are  
396 limited.. Therefore, testing methods need to be appropriate for each case, in order to  
397 demonstrate adequacy for each purpose.

398

### 399 **Code verifiability**

400 Computer code is a translation of mathematical clauses from the mathematical language  
401 to a computer language. The one-to-one correspondence is not always achieved. There is  
402 some consensus (after Glasow and Pace, 1999) that component-based development is  
403 indeed an effective and affordable way of creating model applications and conducting  
404 model evaluation.

405

406 In such a case, it is our view that particular emphasis should be placed on designing and  
407 coding object-oriented simulation models to properly transfer simulation control between  
408 entities, resources and system controllers and on techniques for obtaining a  
409 correspondence between simulation code and system behaviour. It is crucial to consider  
410 the issue of model component validity when considering model re-use as it needs to be a  
411 fundamental part of any re-use strategy.

412

413 The distribution of already validated model components (mathematical and coded  
414 algorithms) can substantially decrease the model validation effort when re-used. A key  
415 step in this direction is the coupling between model components and evaluation  
416 techniques, the latter also being implemented into component-based software. Such  
417 evaluation systems should stand at the core of a general framework where the modelling  
418 system (i.e. a set of modelling components) and a data provider supply inputs to an  
419 evaluation tool (e.g. Bellocchi et al., 2006). Such an evaluation tool is also meant as a  
420 component-based system, both communicating with the modelling component and the  
421 data provider via a suitable protocol and allowing the user to interact in some way (e.g.  
422 via a graphical user interface) to choose and parameterize the evaluation tools.

423

424 The output from an evaluation system can be offered to a deliberative process (e.g.  
425 stakeholder review) for interpretation of results. Adjustments in the modelling system or  
426 critical reviewing of data used to evaluate the model can be a next stage, if the results are  
427 assessed as unsatisfactory for the application purpose. A new evaluation-interpretation  
428 cycle can be run any time new versions of the modelling system are developed and  
429 plugged in to the evaluation component. Again, a well-designed, component-based  
430 evaluation system can be easily extended towards including further evaluation  
431 approaches to keep up with evolving methodologies, e.g. statistical, neural networks or  
432 fuzzy-based (e.g. Bellocchi et al., 2008). Hence, further purpose of this letter is to  
433 stimulate debate on the positive and negative aspects of rigid model structures or  
434 component-based ones, and how the review process can best evaluate them.

435

436 **Reliability**

437 Model reliability cannot be assessed regardless of a presumed range of accuracy. A  
438 specific parameterization of a model can be considered as reliable, if it produces results  
439 that fall within a well-defined range of accuracy. In the case of a predictive model, the  
440 range of accuracy can be defined statistically, proceeding from tests against observations.  
441 In the case of an exploratory model, the range of accuracy may be defined through  
442 sensitivity analysis, assuming that inaccuracy results from uncertainty in the values of  
443 model parameters (e.g. Confalonieri et al., 2010).

444

445 Reliability is also a key aspect of credibility, where measures are influenced by the ability  
446 to establish reliability with available past observations. It cannot be assumed, however,  
447 that statistical (or any numerical) analysis is all that is required for model outputs to be  
448 accepted particularly when models are used with and for stakeholders. The numerical  
449 analysis provides credibility within the techno-scientific research community yet, while  
450 necessary; this may be insufficient to achieve credibility with decision makers and other  
451 stakeholders. Possibly a real test of model validity is whether stakeholders have sufficient  
452 confidence in the model to use it as the basis for making management decisions (Vanclay  
453 and Skovsgaard, 1997).

454

455 Reliability can also be interpreted as versatility of the model, that is, how well does the  
456 model perform in situations for which it was not originally designed, or respond to  
457 extreme conditions beyond that which calibration data represent? Sometimes, it is

458 characterized by a ratio of the real world observed data described by the model outputs  
459 (Mankin et al., 1975). The assessment of model versatility is based on the qualitative  
460 analysis of model structure (i.e. mathematical expressions) and potential results the  
461 model in question can generate. In many cases, only qualitative assessment can lead to  
462 subjective conclusions. The quantification of the concept is difficult or hardly possible  
463 due to limited observation data that is insufficient to understand environmental  
464 behaviour, model complexity limiting evaluation of possible model outcomes, and  
465 uncertainty in modelling results.

466

#### 467 ***How to legitimate model usage***

468 For well-developed environmental applications, model evaluation and selection  
469 techniques are heavily influential and can be used to build scientific and perceived  
470 credibility. However, establishing credibility is not straightforward for larger-scale  
471 environmental applications with many sources of uncertainty, decision-makers with  
472 different interests, and plausible future states that can be markedly different from  
473 observed past states. In these cases, credibility can be influenced by subjective measures  
474 and contingencies in the decision-making process (e.g., Aumann, 2008).

475

476 Establishing model credibility with end users / stakeholders can be problematic since they  
477 may have preconceived, and sometimes immovable, conceptions (Carberry et al., 2002).

478 The task then falls on the model developers to show sufficient evidence in a form  
479 understandable by the end user to persuade them to challenge their beliefs and to consider  
480 alternatives.

481

482 Given the number of potential outcomes and stakeholders involved, more inclusive  
483 modelling approaches such as multiple model and ensemble forecasting approaches can  
484 be useful for establishing credibility. In particular, approaches that allow for multiple  
485 model inference where differing models and perspectives are not excluded (e.g., Min and  
486 Hense 2006). Instead, different models are weighted and synthesized using quantitative  
487 criteria such as statistical support. This is important in determining quantitative reliability  
488 and model evaluation (Burnham and Anderson, 2002), but can also assist with qualitative  
489 aspects of credibility when models are used to inform a contentious decision process. The  
490 resulting consideration of multiple models serves as a proxy for including different  
491 scientific and subjective views of how environmental systems function and the resulting  
492 ensemble forecasts are considered to be more broadly representative of the perspectives  
493 of the decision-making participants.

494

495 Our view is that the key to successfully legitimating model usage is making model  
496 outputs be seen by stakeholders as relevant to their decision making process. Legitimacy  
497 of model usage can be seriously compromised when research outputs refer to geographic,  
498 temporal, or organizational scales that do not match those of decision-making. Hence,  
499 though adequate assessment and evaluation of a model in one location may be shown,  
500 acceptance by stakeholders may be limited when applied where testing has not been  
501 conducted.

502



503 Where models are used for decision support or evidence based reasoning, credibility is a  
504 complex mix of social, technological and mathematical aspects that require developers to  
505 include social networking (between developers, researchers and end users / stakeholders)  
506 to determine model rationale, aim, structure etc., and importantly a sense of co-  
507 ownership. Again, evidence of testing and results from a standardised peer-review  
508 procedure aids dialogue with stakeholders, as the researchers applying the models can  
509 demonstrate independent testing.

510 In this respect, a key component to credibility building is that a model should make  
511 available all the key management options that the decision maker considers important and  
512 should to an acceptable degree respond to management interventions in a way that  
513 matches with the decision maker's experience of the real system. In terms of models of  
514 natural processes, management can be substituted with alternatives, such as external  
515 shocks and/or perturbations to the drivers of the system.

516

517 Using the 'see-saw' analogy, where environmental models' estimates are used in  
518 contentious issues, credibility becomes the focal balance point around which opposing  
519 parties construct their arguments. Hence, credible models can serve to unite opposing  
520 parties, rather than serve to allow them to argue at increasing distance from each other's  
521 viewpoints and expertises. For such cases, subjective decisions on the selection and  
522 assessment of evidence may be as important as the accuracy of the measurement or  
523 forecasting of a particular phenomenon (Matthews et al., 2008).

524

525 Lack of transparency is frequently cited as the reason for the failure of model based  
526 approaches. It is important to challenge some of the assumptions and conclusions that are  
527 drawn on how to respond to the issue of transparency. One response is to make models  
528 simpler and hence the argument becomes easier to understand. Yet while simplicity is in  
529 itself desirable (Raupach and Finnigan, 1988) and the operation of simpler models may  
530 indeed be easier to understand, it may well be that the interpretation of their outputs is no  
531 simpler and indeed their simplicity may mean that they lack the capability to provide  
532 secondary data which can ease the process of interpretation. There is also a trade-off  
533 between simplicity and flexibility and this flexibility may be a crucial factor in allowing  
534 the tools to be relevant for counter-factual analyses. The current best practice for  
535 balancing simplicity and flexibility within the model development process seems to be  
536 the reusable component approach combined with a flexible model integration / evaluation  
537 environment. A set of standards applied within the peer review procedure therefore needs  
538 to address the issues of simple versus complex models and so look beyond the numerical  
539 testing and consider the flexibility of the model, ability of it to shed new light on an  
540 environmental issue, and aid the process of interpretation.

541

542 A constraint to both scientific credibility and transparency of models is the necessarily  
543 inherent inter-dependency of the modelled processes. A 'fault' in a model may be  
544 difficult to locate as many other related modelled processes confound it. Similarly, an  
545 effective modelled description of a specific sub-process may not be readily identified due  
546 to its dependence on less than satisfactory descriptions of other system features. Ranges  
547 of sensitivity and uncertainty analyses have been deployed to address this issue, although

548 the results are not always easy to interpret in terms of the original model formulation.  
549 Comparison of alternative model formulations can provide useful information in this  
550 context (e.g. Confalonieri et al., 2009a) but still suffers from the difficulty of  
551 disentangling inter-dependencies in the model. Crout et al. (2009) have proposed simple  
552 model reduction methods based on the approach of Cox et al (2006) which systematically  
553 explored the role of individual model variables on the models' predictive performance.  
554 This procedure frequently locates variables whose formulation has a detrimental effect on  
555 model performance.

556

## 557 **The way forward: standardized evaluation tools**

558 Based on the above views and highlighting of issues, we now explore options for future  
559 model evaluation. Turning back to the fact that model developers are normally lacking  
560 resources for adequate model evaluation, we conclude that introducing a more formalized  
561 set of evaluation criteria demands a standardised set of evaluation tools.

### 562 ***Identifying the prior art***

563 It is common to believe that the total amount of environmental models is huge and that  
564 they cover almost all environmental situations. Nevertheless, the recent review of models  
565 used by the European Environment Agency in its recent environmental assessments and  
566 reports identified gaps in the availability, accessibility and applicability of current  
567 modelling tools (EEA, 2008). Indeed, journal articles reporting modelling efforts  
568 normally focus on the scientific interpretation of the findings, not on model  
569 documentation. There is no guarantee that a model, on which a published article several

570 years ago, is still readily available or even existing in any form that makes it possible to  
571 reproduce reported results. Can we therefore consider models that are not readily and  
572 completely available as the prior art? Scientific etiquette suggests that model  
573 documentations must be conveniently accessible, complete and mutually comparable  
574 (Benz and Knorrenschild, 1997; Voinov et al., 2009). We therefore suggest a register of  
575 models that can be considered as the prior art is needed for the technical assessment and  
576 evaluation of newly developed models.

577

### 578 ***Developing a comprehensive numerical library***

579 A disciplined approach, effective management, and well-educated personnel are some of  
580 the key factors affecting the success of a software development project. Professionals in  
581 environmental modelling can learn a lot from software engineering, commercial product  
582 testing (especially in aircraft design and other areas where there is a very high safety  
583 standard required), stakeholders' deliberation and scientific developments from other  
584 disciplines. In so doing, we can expand our horizons to include the necessary knowledge  
585 to conduct successful model evaluation. Whilst some research has been undertaken  
586 focusing on establishing a baseline for evaluation practice, rather less work has been done  
587 to develop a basic, scientifically rigorous approach to be able to meet the technical  
588 challenges we currently face. We believe model evaluation software tools can valuably  
589 support this activity, allowing consolidated experience in evaluating models to be formed  
590 and shared. Whether model evaluation is a scheduled action in modelling projects, little  
591 work is published in the open literature (e.g., conference proceedings and journals)  
592 describing the evaluation experience accumulated by modelling teams (including

593 interactions with the stakeholders). Failing to disseminate the evaluation experience may  
594 result in the repetition of the same mistakes in future modelling projects. Based on past  
595 experience, establishing a better quality assurance program for a new modelling project  
596 may certainly increase the probability of success for that project. Learning from the  
597 experience of others is an excellent and cost-effective educational tool. The return on  
598 such an investment can easily be realized by preventing the failures of modelling projects  
599 and by avoiding wrong simulation-based decisions. Where complex models are to be  
600 evaluated, options are available to combine detailed numeric and statistical tests of  
601 components and sub-processes with a deliberative approach for overall model  
602 acceptance. Future model development should aim to incorporate automated evaluation  
603 checks using embedded software tools, with the aim of achieving greater cost and time  
604 efficiency and to achieve a higher level of credibility. Information from evaluation tools  
605 employed by the model developers needs to be made available to the peer review process.  
606 Beyond this, providing third parties with the capability of extending methodologies  
607 without re-compiling the component will ensure greater transparency and ease of  
608 maintenance, also providing functionalities such as the test of input data versus their  
609 definition prior to computing any simple or integrated evaluation metric. Making it in  
610 agreement with the most modern developments in software engineering, components for  
611 model evaluation will better serve as a convenient means to support collaborative model  
612 testing among the network of scientists involved in creating component-oriented models  
613 in the environmental field.  
614

615 ***Moving from software to webware***

616 Modern information and communication technologies offer the opportunity for a  
617 revolution in the area of technical assessments of environmental models (Alexandrov and  
618 Matsunaga, 2008; Hoffman et al., 2008). Moving from software to webware makes  
619 models (and data used to test them) available through a web-browser. It seems a time has  
620 come to think seriously about an Environmental Modelling Server (EMS) - a  
621 supercomputer (or a computing grid) for deploying environmental models and running  
622 them through web-browsers. The EMS may also do the routine work on technical  
623 assessment of models, providing the necessary resource currently lacking.

624 **Our position**

625 Concluding this letter, we emphasize that having standardized evaluation tools is the  
626 issue that needs to be tackled. Standardised model evaluation can consist of evaluation  
627 tools for use during and after the model development process, which can feed into a  
628 codified procedure during peer review of articles based on the model. The articles  
629 published in this thematic volume of EM&S are suggesting “*the evaluation of models*  
630 *should be a central part of the model development process, not an afterthought*” (Crout et  
631 al., 2009). This implies a clear demand for relevant software tools and acceptance by  
632 journals to adopt a minimum standard for peer review. Evaluation tools, in contrast to  
633 models, are generic by their nature, based on shared information and on re-using data  
634 from previous research exercises. The burden of developing evaluation tools is too hard  
635 for every single modeller. This and improving the peer review process are tasks that need

636 a communal effort based on International Environmental modelling and Software  
637 Society's (iEMSs) leadership.

638

## 639 **References**

640

641 Alexandrov, G.A., 2006. The purpose of peer review in the case of an open-access  
642 publication. *Carbon Balance and Management*, 1:10

643 [<http://www.cbmjournal.com/content/1/1/10>]

644

645 Alexandrov, G.A., Matsunaga, T., 2008. Evaluating consistency of biosphere models:  
646 software tools for a web-based service. In: Sànchez-Marrè, M., Béjar, Comas, J.J.,  
647 Rizzoli, A.E., Guariso, E., (Eds.), *Proceedings of the iEMSs Fourth Biennial Meeting:*  
648 *International Congress on Environmental Modelling and Software (iEMSs 2008).*

649 *International Environmental Modelling and Software Society, Barcelona, Catalonia.*

650 [[http://www.iemss.org/iemss2008/uploads/Main/S12-02-Alexandrov\\_et\\_al-](http://www.iemss.org/iemss2008/uploads/Main/S12-02-Alexandrov_et_al-)

651 [IEMSS2008.pdf](http://www.iemss.org/iemss2008/uploads/Main/S12-02-Alexandrov_et_al-IEMSS2008.pdf)]

652

653 Aumann, C., 2008. A methodology for building Credible Models for Policy Evaluation.

654 In: Sànchez-Marrè, M., Béjar, Comas, J.J., Rizzoli, A.E., Guariso, E., (Eds.),

655 *Proceedings of the iEMSs Fourth Biennial Meeting: International Congress on*

656 *Environmental Modelling and Software (iEMSs 2008).* *International Environmental*

657 *Modelling and Software Society, Barcelona, Catalonia.*

658 [[http://www.iemss.org/iemss2008/uploads/Main/S12-01-Aumann\\_et\\_al-IEMSS2008.pdf](http://www.iemss.org/iemss2008/uploads/Main/S12-01-Aumann_et_al-IEMSS2008.pdf)]

659

660 Bellocchi, G., Confalonieri, R., Donatelli, M., 2006. Crop modelling and validation:  
661 integration of IRENE\_DLL in the WARM environment. Italian Journal of  
662 Agrometeorology 3, 35-39.

663

664 Bellocchi, G., Acutis, M., Fila, G., Donatelli, M., 2002. An indicator of solar radiation  
665 model performance based on a fuzzy expert system. Agron. J. 94, 1222-1233.  
666 [<http://agron.scijournals.org/cgi/reprint/94/6/1222.pdf>]

667

668 Bellocchi G., E. Habyarimana, M. Donatelli, M. Acutis, 2008. A software component for  
669 model output evaluation. In: Sánchez-Marrè, M., Béjar, Comas, J.J., Rizzoli, A.E.,  
670 Guariso, E., (Eds.), Proceedings of the iEMSs Fourth Biennial Meeting: International  
671 Congress on Environmental Modelling and Software (iEMSs 2008). International  
672 Environmental Modelling and Software Society, Barcelona, Catalonia.  
673 [[http://www.iemss.org/iemss2008/uploads/Main/S12-06-Donatelli\\_et\\_al-](http://www.iemss.org/iemss2008/uploads/Main/S12-06-Donatelli_et_al-IEMSS2008.pdf)  
674 [IEMSS2008.pdf](http://www.iemss.org/iemss2008/uploads/Main/S12-06-Donatelli_et_al-IEMSS2008.pdf)]

675

676 Bellocchi G., Rivington M., Donatelli M., Matthews K.B., 2010. Validation of  
677 biophysical models: issues and methodologies. A review. Agron. Sustain. Dev., 30, 109-  
678 130..

679

680 Benz, J., Knorrnschild, M., 1997. Call for a common model documentation etiquette.  
681 Ecological Modelling 97, 141-143.



682

683 Beven, K., 2006. A manifesto for the equifinality thesis. *Journal of Hydrology* 320, 18–  
684 36.

685

686 Beven K.J., Freer J., 2001. Equifinality, data assimilation, and uncertainty estimation  
687 in mechanistic modelling of complex environmental systems using the GLUE  
688 methodology. *J. Hydrol.* 249, 11-29.

689

690 Burnham, K.P., Anderson, D.R., 2002. *Model selection and multimodel inference: a  
691 practical information-theoretic approach*, 2nd ed., Springer-Verlag: New York (NY).

692

693 Carberry, P.S., Hochman, Z., McCown, R.L., Dalglish, N.P. and others, 2002. The  
694 FARMSCAPE approach to decision support: farmers', advisers, researchers' monitoring,  
695 simulation, communication and performance evaluation. *Agricultural Systems* 74, 141–  
696 177.

697

698 Cheng, R.T., Burau, J.R., Gartner, J.W., 1991. Interfacing data analysis and numerical  
699 modelling for tidal hydrodynamic phenomena. In: Parker, B.B. (Ed.) *Tidal  
700 hydrodynamics*, John Wiley & Sons: New York, pp. 201-219.

701

702 Confalonieri, R., Acutis, M., Bellocchi, G., Donatelli, M., 2009a. Multi-metric evaluation  
703 of the models WARM, CropSyst, and WOFOST for rice. *Ecological Modelling* 220,  
704 1395-1410.

705

706 Confalonieri, R., Bellocchi, G., Boschetti, M., Acutis, M., 2009b. Evaluation of  
707 parameterization strategies for rice modelling. *Spanish Journal of Agricultural Research*  
708 7, 680-686.

709

710 Confalonieri, R., Bellocchi, G., Tarantola, S., Acutis, M., Donatelli, M., Genovese, G.,  
711 2010. Sensitivity analysis of the rice model WARM in Europe: Exploring the effects of  
712 different locations, climates and methods of analysis on model sensitivity to crop  
713 parameters. *Environmental Modelling & Software* 25, 479-488.

714

715 Council of Science Editors, 2006. *CSE's White Paper on Promoting Integrity in Scientific*  
716 *Journal Publications*, CSE, Reston, Va.

717 [[http://www.councilscienceeditors.org/editorial\\_policies/whitepaper/entire\\_whitepaper.p](http://www.councilscienceeditors.org/editorial_policies/whitepaper/entire_whitepaper.pdf)  
718 [df](http://www.councilscienceeditors.org/editorial_policies/whitepaper/entire_whitepaper.pdf)]

719

720 Cox, G.M., Gibbons, J.M., Wood, A.T.A., Craigon, J., Crout, N.M.J., 2006. Towards the  
721 systematic simplification of mechanistic models. *Ecological Modelling* 198, 240–246.

722

723 Crout, N., Kokkonen, T., Jakeman, A.J., Norton, J.P., Newham, L.T.H., Anderson, R.,  
724 Assaf, H., Croke, B.F.W., Gaber, N., Gibbons, J., Holzworth, D., Mysiak, J., Reichl, J.,  
725 Seppelt, R., Wagener, T., Whitfield, P., 2009. Good modelling practice. In: Jakeman,  
726 A.J., Voinov, A.A., Rizzoli, A.E., Chen, S.H. (Eds.) *Environmental modelling, software*  
727 *and decision support*. Elsevier, pp.15-31.

728

729 EEA, 2008. Modelling environmental change in Europe: towards a model inventory  
730 (SEIS/Forward). EEA Technical Report 11/2008, European Environment Agency,  
731 Copenhagen.

732

733 French, S., Geldermann, J., 2005. The varied contexts of environmental decision  
734 problems and their implications for decision support. *Environ. Sci. Policy* 8, 378-391.

735

736 Funtowicz, S.O., Ravetz, J.R., 1990. *Uncertainty and Quality in Science for Policy*.  
737 Kluwer Academic Press: Dordrecht (The Netherlands).

738

739 Gaber, N., Pascual, P., Stiber, N., Sunderland, E., Cope, B., Nold, A., 2008. *Guidance on*  
740 *the Development, Evaluation and Application of Environmental Models*, Council for  
741 *Regulatory Environmental Modeling*, U.S. Environmental Protection Agency,  
742 Washington.

743

744 Gardner, R.H., Urban, D.L., 2003. Model validation and testing: Past lessons, present  
745 concerns, future prospects. In: Canham, C.D., J.J. Cole, Lauenroth, W.K., (Eds.), *Models*  
746 *in Ecosystem Science*. Princeton (NJ): Princeton University Press, pp. 184-203.

747

748 Glasow, P.A., Pace, D.K., 1999. SIMVAL '99: making VV&A effective and affordable  
749 workshop, The Simulation Validation Workshop 1999, January 26-29, Laurel, MD, USA.

750

751 Haag, D., Kaupenjohann, M., 2001. Parameters, prediction, post-normal science and the  
752 precautionary principle - a roadmap for modelling for decision-making. *Ecological*  
753 *Modelling* 144, 45-60.  
754

755 Hames, I., 2007. Peer-review and manuscript management in scientific journals:  
756 guidelines for good practice. Blackwell Publishing; Oxford.  
757

758 Harremoës, P., 2003. The role of uncertainty in application of integrated urban water  
759 modeling. In: Proceedings of the International IMUG Conference, Tilburg, 23–25 April  
760 2003.  
761

762 Hoffman, F., G. Bonan, C. Covey, I. Fung, Y.-H. Lee, J. Randerson, S. Running, P.  
763 Thornton, The Carbon-Land Model Intercomparison Project (C-LAMP): A Protocol and  
764 Evaluation Metrics for Global Terrestrial Biogeochemistry Models. In: Sánchez-Marrè,  
765 M., Béjar, Comas, J.J., Rizzoli, A.E., Guariso, E., (Eds.), Proceedings of the iEMSs  
766 Fourth Biennial Meeting: International Congress on Environmental Modelling and  
767 Software (iEMSs 2008). International Environmental Modelling and Software Society,  
768 Barcelona, Catalonia. [[http://www.iemss.org/iemss2008/uploads/Main/S12-03-  
769 Hoffmann\\_et\\_al-IEMSS2008.pdf](http://www.iemss.org/iemss2008/uploads/Main/S12-03-Hoffmann_et_al-IEMSS2008.pdf)]  
770

771 Hsu, M.H., Kuo, A.Y., Kuo, J.T., Liu, W.C., 1999. Procedure to calibrate and verify  
772 numerical models of estuarine hydrodynamics. *J. Hydr. Engrg.* 125, 166-182.  
773

774 Huth, N., Holzworth, D., 2005. Common sense in model testing. In: Zerger, A., Argent,  
775 R.M. (Eds.), Proceedings of MODSIM 2005 International Congress on Modelling and  
776 Simulation: Advances and applications for management and decision making, 12-15  
777 December, Melbourne, Australia, p.2804-2809.

778

779 Jakeman, A., Chen, S., Rizzoli, A., Voinov, A.A., 2009. Modelling and software as  
780 instruments for advancing sustainability. In: Jakeman, A.J., Voinov, A.A., Rizzoli, A.E.,  
781 Chen, S.H. (Eds.) Environmental Modelling, Software and Decision Support, Elsevier,  
782 pp. 345-366.

783

784 Jakeman, A.J., Letcher, R.A., Norton, J.P., 2006. Ten iterative steps in development and  
785 evaluation of environmental models. Environmental Modelling and Software 21 (5), 602–  
786 614.

787

788 Janssen, P.H.M., Heuberger, P.S.C., 1995. Calibration of process-oriented models. Ecol.  
789 Model. 83, 55-66.

790

791 Jørgensen, S.E., Fath, B.D., Grant, W., Nielsen, S.N., 2006. The editorial policy of  
792 Ecological Modelling. Ecological Modelling 199, 1-3.

793

794 Keating, B.A., Gaydon, D., Huth, N.I., Probert, M.E., Verburg, K., Smith, C.J., Bond,  
795 W., 2002. Use of modelling to explore the water balance of dryland farming systems in  
796 the Murray-Darling Basin, Australia. Eur. J. Agron. 18, 159-169.

797

798 Leffelaar, P.A., Meike, H., Smith, P., Wallach, D., 2003. Modelling cropping systems -  
799 highlights of the symposium and preface to the special issues. 3. Session B. Model  
800 parameterisation and testing. *Eur. J. Agron.* 18, 189-191.

801

802 Li, W., Arena, V.C, Sussman, N.B., Mazumdar, S. 2003. Model validation software for  
803 classification models using repeated partitioning: MVREP. *Comput. Meth Prog. Bio.* 72,  
804 81-87.

805

806 Mankin, J.B., O'Neill, R.V., Shugart, H.H. and Rust B.W. 1975. The importance of  
807 validation in ecosystem analysis. In *New directions in the analysis of ecological systems,*  
808 *Part 1.* La Jolla, CA: Society for computer simulation, v. 5: 63-71.

809

810 Marcus, A.H., Elias, R.W., 1998. Some useful statistical methods for model validation.  
811 *Environ. Health Persp.* 106, 1541-1550.

812

813 Matthews, K.B., Rivington, M., Blackstock, K., Buchan, K., Miller, D.G., 2008. Raising  
814 the bar - Is evaluating the outcomes of decision and information support tools a bridge  
815 too far? In: Sánchez-Marrè, M., Béjar, Comas, J.J., Rizzoli, A.E., Guariso, E., (Eds.),  
816 *Integrating sciences and information technology for environmental assessment and*  
817 *decision making. Proceedings of the 4th Biennial Meeting of the International*  
818 *Environmental Modelling and Software Society, 7-10 July, Barcelona, Spain, Vol. 1,*  
819 *p.948-955.*

820

821 Min, S.-K., Hense, A., 2006. A Bayesian Assessment of Climate Change Using  
822 Multimodel Ensembles. Part I: Global Mean Surface Temperature. *Journal of Climate*,  
823 19:3237-3256.

824

825 Müller, U. T., 2009. Peer-Review-Verfahren zur Qualitätssicherung von Open-Access-  
826 Zeitschriften – systematische Klassifikation und empirische Untersuchung, PhD thesis,  
827 Humboldt-Universität zu Berlin, 2009; available at [http://edoc.hu-](http://edoc.hu-berlin.de/docviews/abstract.php?id=29636)  
828 [berlin.de/docviews/abstract.php?id=29636](http://edoc.hu-berlin.de/docviews/abstract.php?id=29636)

829

830 Murray-Darling Basin Commission, 2000. Murray-Darling Basin Commission,  
831 Groundwater Flow Modelling Guideline, Murray-Darling Basin Commission, Project no.  
832 125. Canberra.

833

834 Pachepsky, Y.A., Guber, A.K., Van Genuchten, M.T., Nicholson, T.J., Cady, R.E.,  
835 Simunek, J., Schaap, M.G., 2006. Model abstraction techniques for soil water flow and  
836 transport, 175 pp, Nuclear Regulatory Commission, Washington, DC; NUREG/CR-6884.

837

838 Parker, V.T., 2001. Conceptual problems and scale limitations of defining ecological  
839 communities: a critique of the CI concept (Community of Individuals). *Perspect. Plant*  
840 *Ecol. Evol. Syst.* 4, 80-96.

841

842 Raupach, M.R., Finnigan, J.J., 1988. Single layer models of evaporation from plant  
843 canopies are incorrect, but useful, whereas multilayer models are correct, but useless:  
844 discussion. *Aust. J. Plant Physiol.* 15, 705-716.  
845

846 Reichert, P., Borsuk, M.E., 2005. Does high forecast uncertainty preclude effective  
847 decision support? *Environmental Modelling and Software*,20: 991-1001.  
848

849 Reynolds, J.F., Ford, E.D., 1999. Multi-criteria assessment of ecological process models.  
850 *Ecology* 80, 538-553.  
851

852 Robson, B., Hamilton, D., Webster, I., Chan, T., 2008. Ten steps applied to development  
853 and evaluation of process-based biogeochemical models of estuaries. *Environmental*  
854 *Modelling & Software* 23 (4), 369–384.  
855

856 Rykiel, E.J. Jr., 1996. Testing ecological models: The meaning of validation. *Ecological*  
857 *Modelling* 90 229–244.  
858

859 Sinclair, T.R., Seligman, N., 2000. Criteria for publishing papers on crop modelling.  
860 *Field Crops Res.* 68, 165-172.  
861

862 STOWA/RIZA, 1999. STOWA/RIZA, Smooth Modelling in Water Management, Good  
863 Modelling Practice Handbook STOWA Report 99–05, Dutch Department of Public



864 Works, Institute for Inland Water Management and Waste Water Treatment (1999) ISBN  
865 90-5773-056-1 Report 99.036.  
866  
867 Vanclay, J.K., 1994. Modelling forest growth and yield, CAB International, Wallingford,  
868 United Kingdom.  
869  
870 Vanclay, J.K., Skovsgaard, J.P., 1997. Evaluating forest growth models. *Ecol. Modell.*  
871 98, 1-12.  
872  
873 Van der Sluijs, J. P., 2007. Uncertainty and Precaution in Environmental Management:  
874 insights from the UPEM conference. *Environ Modelling and Software*, 22(5), 590-598.  
875  
876 Van Nes, E.H. and Scheffer, M., 2003. Alternative attractors may boost uncertainty and  
877 sensitivity in ecological models. *Ecological Modelling*, 159, 117–124.  
878  
879 Van Oijen, M., 2002. On the use of specific publication criteria for papers on process-  
880 based modelling in plant science. *Field Crops Res.* 74, 197-205.  
881  
882 Vischel, T., Pegram, G., Sinclair, S., Wagner, W., Bartsch, A., 2007. Comparison of soil  
883 moisture fields estimated by catchment modelling and remote sensing: a case study in  
884 South Africa. *Hydrol. Earth Syst. Sci. Discuss.* 4, 2273–2306.  
885

886 Voinov, A., Hood, R.R., Daves, J.D., Assaf, H., Stewart, R., 2009. Building a community  
887 modelling and information sharing culture. In: Jakeman, A.J., Voinov, A.A., Rizzoli,  
888 A.E., Chen, S.H. (Eds.) Environmental Modelling, Software and Decision Support,  
889 Elsevier, pp. 345-366.

890

891 Wager, E., 2006. Ethics: What is it for? Nature: Web Debate – Peer-review.  
892 [<http://www.nature.com/nature/peerreview/debate/nature04990.html>]

893

894 Walker, T.J., 1998. Free Internet Access to Traditional Journals. American Scientist,  
895 86:463

896

897 Wang, Y.-P., Trudinger, C.M., Enting, I.G., 2009. A review of applications of model-data  
898 fusion to studies of terrestrial carbon fluxes at different scales. Agricultural and Forest  
899 Meteorology, 149, 1829-1842.

900

901 Welsh, W., 2008. Water balance modelling in Bowen, Queensland, and the ten iterative  
902 steps in model development and evaluation. Environmental Modelling & Software 23 (2),  
903 195–205.

904

905 Ye, M., Meyer, P.D., Neuman, S.P. ,2008. On model selection criteria in multimodel  
906 analysis, Water Resources Research 44, W03428.

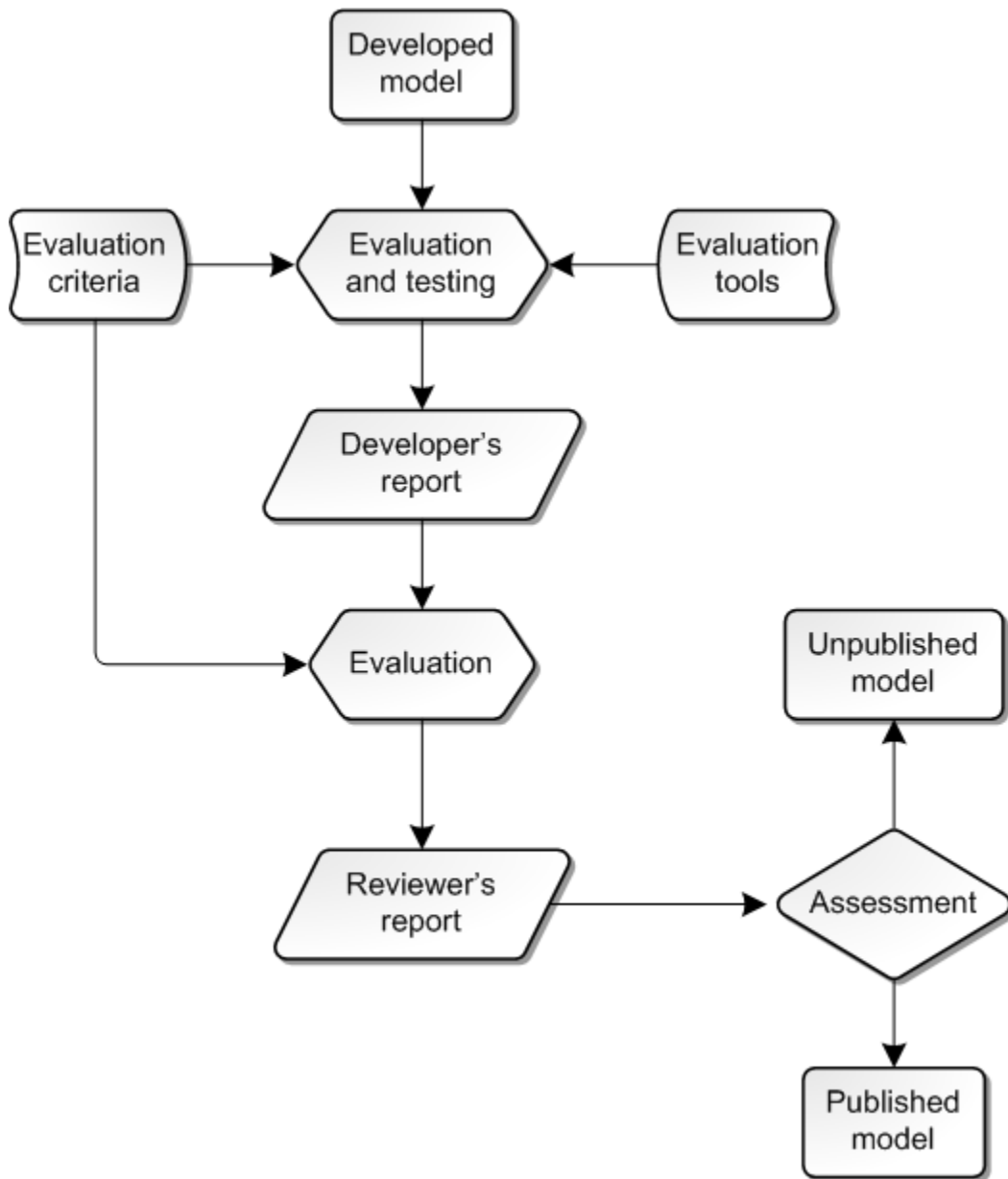
907

908

909

910

911



912

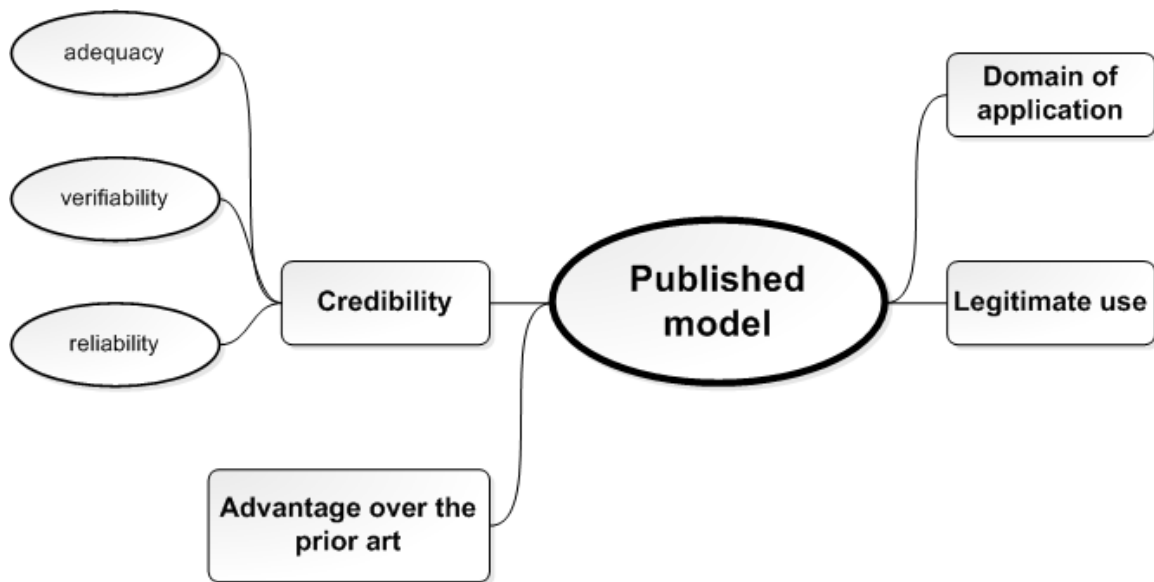
913

914 Figure 1. A flow-chart of peer-review process

915

916

917  
918  
919  
920  
921  
922



923  
924  
925  
926  
927  
928

929 Figure 2. A scheme for specifying a published model.

930