



1 GMD Perspective: the quest to improve the
2 evaluation of groundwater representation in
3 continental to global scale models

4
5 Tom Gleeson^{1,2}, Thorsten Wagener³, Petra Döll⁴, Samuel C Zipper^{1,5}, Charles West³, Yoshihide
6 Wada⁶, Richard Taylor⁷, Bridget Scanlon⁸, Rafael Rosolem³, Shams Rahman³, Nurudeen Oshinlaja⁹,
7 Reed Maxwell¹⁰, Min-Hui Lo¹¹, Hyungjun Kim¹², Mary Hill¹³, Andreas Hartmann^{14,3}, Graham Fogg¹⁵,
8 James S. Famiglietti¹⁶, Agnès Ducharne¹⁷, Inge de Graaf^{18,19}, Mark Cuthbert^{9,20}, Laura Condon²¹,
9 Etienne Bresciani²², Marc F.P. Bierkens^{23,24}

10 ¹ Department of Civil Engineering, University of Victoria, Canada

11 ² School of Earth and Ocean Sciences, University of Victoria

12 ³ Department of Civil Engineering, University of Bristol, UK & Cabot Institute, University of Bristol, UK.

13 ⁴ Institut für Physische Geographie, Goethe-Universität Frankfurt am Main and Senckenberg Leibniz
14 Biodiversity and Climate Research Centre Frankfurt (SBIK-F), Frankfurt am Main, Germany

15 ⁵ Kansas Geological Survey, University of Kansas

16 ⁶ International Institute for Applied Systems Analysis, Laxenburg, Austria

17 ⁷ Department of Geography, University College London, UK

18 ⁸ Bureau of Economic Geology, The University of Texas at Austin, USA

19 ⁹ School of Earth and Environmental Sciences & Water Research Institute, Cardiff University, UK

20 ¹⁰ Department of Geology and Geological Engineering, Colorado School of Mines, USA

21 ¹¹ Department of Atmospheric Sciences, National Taiwan University, Taiwan

22 ¹² Institute of Industrial Science, The University of Tokyo

23 ¹³ Department of Geology, University of Kansas, USA

24 ¹⁴ Chair of Hydrological Modeling and Water Resources, University of Freiburg, Germany

25 ¹⁵ Department of Land, Air and Water Resources and Earth and Planetary Sciences, University of
26 California, Davis, USA

27 ¹⁶ School of Environment and Sustainability and Global Institute for Water Security, University of
28 Saskatchewan, Saskatoon, Canada

29 ¹⁷ Sorbonne Université, CNRS, EPHE, IPSL, UMR 7619 METIS, Paris, France

30 ¹⁸ Chair or Environmental Hydrological Systems, University of Freiburg, Germany

31 ¹⁹ Water Systems and Global Change Group, Wageningen University, Wageningen, Netherlands

32 ²⁰ School of Civil and Environmental Engineering, The University of New South Wales, Sydney, Australia

33 ²¹ Department of Hydrology & Atmospheric Sciences, University of Arizona, Tucson, Arizona, USA

34 ²² Center for Advanced Studies in Arid Zones (CEAZA), La Serena, Chile

35 ²³ Physical Geography, Utrecht University, Utrecht, Netherlands

36 ²⁴ Deltares, Utrecht, Netherlands



37

Abstract

38 Continental- to global-scale hydrologic and land surface models increasingly include representations of
39 the groundwater system. Such large-scale models are essential for examining, communicating, and
40 understanding the dynamic interactions between the Earth System above and below the land surface as
41 well as the opportunities and limits of groundwater resources. We argue that both large-scale and
42 regional-scale groundwater models have utility, strengths and limitations so continued modeling at both
43 scales is essential and mutually beneficial. A crucial quest is how to evaluate the realism, capabilities and
44 performance of large-scale groundwater models given their modeling purpose of addressing large-scale
45 science or sustainability questions as well as limitations in data availability and commensurability.
46 Evaluation should identify if, when or where large-scale models achieve their purpose or where
47 opportunities for improvements exists so that such models better achieve their purpose. We suggest
48 that reproducing the spatio-temporal details of regional-scale models and matching local data is not a
49 relevant goal. Instead, it is important to decide on reasonable model expectations regarding when a
50 large scale model is performing 'well enough' in the context of its specific purpose. The decision of
51 reasonable expectations is necessarily subjective even if the evaluation criteria is quantitative. Our
52 objective is to provide recommendations for improving the evaluation of groundwater representation in
53 continental- to global-scale models. We describe current modeling strategies and evaluation practices,
54 and subsequently discuss the value of three evaluation strategies: 1) comparing model outputs with
55 available observations of groundwater levels or other state or flux variables (observation-based
56 evaluation); 2) comparing several models with each other with or without reference to actual
57 observations (model-based evaluation); and 3) comparing model behavior with expert expectations of
58 hydrologic behaviors in particular regions or at particular times (expert-based evaluation). Based on
59 evolving practices in model evaluation as well as innovations in observations, machine learning and
60 expert elicitation, we argue that combining observation-, model-, and expert-based model evaluation



61 approaches, while accounting for commensurability issues, may significantly improve the realism of
62 groundwater representation in large-scale models. Thus advancing our ability for quantification,
63 understanding, and prediction of crucial Earth science and sustainability problems. We encourage
64 greater community-level communication and cooperation on this quest, including among global
65 hydrology and land surface modelers, local to regional hydrogeologists, and hydrologists focused on
66 model development and evaluation.

67 **1. INTRODUCTION: why and how is groundwater modeled at continental to global scales?**

68 Groundwater is the largest human- and ecosystem-accessible freshwater storage component of the
69 hydrologic cycle (UNESCO, 1978; Margat & Van der Gun, 2013; Gleeson et al., 2016). Therefore, better
70 understanding of groundwater dynamics is critical at a time when the ‘great acceleration’ (Steffen et al.,
71 2015) of many human-induced processes is increasing stress on water resources (Wagener et al., 2010;
72 Montanari et al., 2013; Sivapalan et al., 2014; van Loon et al., 2016), especially in regions with limited
73 data availability and analytical capacity. Groundwater is often considered to be an inherently regional
74 rather than global resource or system. This is partially reasonable because local to regional peculiarities
75 of hydrology, politics and culture are paramount to groundwater resource management (Foster et al.
76 2013) and groundwater dynamics in different continents are less directly connected and coupled than
77 atmospheric dynamics. Regional-scale analysis and models are essential for addressing local to regional
78 groundwater issues. Generally, regional scale modeling is a mature, well-established field (Hill &
79 Tiedeman, 2007; Kresic, 2009; Zhou & Li, 2011; Hiscock & Bense, 2014; Anderson et al. 2015a) with clear
80 and robust model evaluation guidelines (e.g. ASTM, 2016; Barnett et al., 2012). Regional models have
81 been developed around the world; for example, Rossman & Zlotnik (2014) and Vergnes et al. (2020)
82 synthesize regional-scale groundwater models across the western United States and Europe,
83 respectively.



84

85 Yet, important global aspects of groundwater both as a resource and as part of the Earth System are
86 emerging (Gleeson et al. 2020). First, our increasingly globalized world trades virtual groundwater and
87 other groundwater-dependent resources in the food-energy-water nexus, and groundwater often
88 crosses borders in transboundary aquifers. A solely regional approach can be insufficient to analysing
89 and managing these complex global interlinkages. Second, from an Earth system perspective,
90 groundwater is part of the hydrological cycle and connected to the atmosphere, oceans and the deeper
91 lithosphere. A solely regional approach is insufficient to uncover and understand the complex
92 interactions and teleconnections of groundwater within the Earth System. Regional approaches
93 generally focus on important aquifers which underlie only a portion of the world's land mass or
94 population and do not include many other parts of the land surface that may be important for processes
95 like surface water-groundwater exchange flows and evapotranspiration. A global approach is also
96 essential to assess the impact of groundwater depletion on sea level rise, since groundwater storage loss
97 rate on all continents of the Earth must be aggregated. Thus, we argue that groundwater is
98 simultaneously a local, regional, and increasingly global resource and system and that examining
99 groundwater problems, solutions, and interactions at all scales is crucial. As a consequence, we urgently
100 require predictive understanding about how groundwater, used by humans and connected with other
101 components of the Earth System, operates at a variety of scales.

102

103 Based on the arguments above for considering global perspectives on groundwater, we see four specific
104 purposes of representing groundwater in continental- to global-scale hydrological or land surface
105 models and their climate modeling frameworks:

106 (1) To understand and quantify interactions between groundwater and past, present and future
107 climate. Groundwater systems can have far-reaching effects on climate affecting modulation of



108 surface energy and water partitioning with a long-term memory (Anyah et al., 2008; Maxwell and
109 Kollet, 2008; Koirala et al. 2013; Krakauer et al., 2014; Maxwell et al., 2016; Taylor, et al., 2013;
110 Meixner et et, 2018; Wang et al., 2018; Keune et al., 2018). While there have been significant
111 advances in understanding the role of lateral groundwater flow on evapotranspiration (Maxwell &
112 Condon, 2016; Bresciani et al, 2016), the interactions between climate and groundwater over
113 longer time scales (Cuthbert et al., 2019) as well as between irrigation, groundwater, and climate
114 (Condon and Maxwell, 2019; Condon et al 2020) remain largely unresolved. Additionally, it is well
115 established that old groundwater with slow turnover times are common at depth (Befus et al.
116 2017; Jasechko et al. 2017). Groundwater connections to the atmosphere are well documented in
117 modeling studies (e.g. Forrester and Maxwell, 2020). Previous studies have demonstrated
118 connections between the atmospheric boundary layer and water table depth (e.g. Maxwell et al
119 2007; Rahman et al, 2015), under land cover disturbance (e.g. Forrester et al 2018), under
120 extremes (e.g. Kuene et al 2016) and due to groundwater pumping (Gilbert et al 2017). While a
121 number of open source platforms have been developed to study these connections (e.g. Maxwell
122 et al 2011; Shrestha et al 2014; Sulis, 2017) these platforms are regional to continental in extent.
123 Recent work has shown global impacts of groundwater on atmospheric circulation (Wang et al
124 2018), but groundwater is still quite simplified in this study.

125 (2) To understand and quantify two-way interactions between groundwater, the rest of the
126 hydrologic cycle, and the broader Earth System. As the main storage component of the freshwater
127 hydrologic cycle, groundwater systems support baseflow levels in streams and rivers, and thereby
128 ecosystems and agricultural productivity and other ecosystem services in both irrigated and
129 rainfed systems (Scanlon et al., 2012; Qiu et al., 2019; Visser, 1959; Zipper et al., 2015, 2017).
130 When pumped groundwater is transferred to oceans (Konikow 2011; Wada et al., 2012; Döll et
131 al., 2014a; Wada, 2016; Caceres et al., 2020; Luijendijk et al. 2020), resulting sea-level rise can



132 impact salinity levels in coastal aquifers, and freshwater and solute inputs to the ocean (Moore,
133 2010; Sawyer et al., 2016). Difficulties are complicated by international trade of virtual
134 groundwater which causes aquifer stress in disparate regions (Dalin et al., 2017)

135 (3) To inform water decisions and policy for large, often transboundary groundwater systems in an
136 increasingly globalized world (Wada & Heinrich, 2013; Herbert & Döll, 2019). For instance,
137 groundwater recharge from large-scale models has been used to quantify groundwater resources
138 in Africa, even though large-scale models do not yet include all recharge processes that are
139 important in this region (Taylor et al., 2013; Jasechko et al. 2014; Cuthbert et al., 2019; Hartmann
140 et al., 2017).

141 (4) To create visualizations and interactive opportunities that inform citizens and consumers, whose
142 decisions have global-scale impacts, about the state of groundwater all around the world such as
143 the World Resources Institute’s Aqueduct website (<https://www.wri.org/aqueduct>), a decision-
144 support tool to identify and evaluate global water risks.

145 The first two purposes are science-focused while the latter two are sustainability-focused. In sum,
146 continental- to global-scale hydrologic models incorporating groundwater offer a coherent scientific
147 framework to examine the dynamic interactions between the Earth System above and below the land
148 surface, and are compelling tools for conveying the opportunities and limits of groundwater resources
149 to people so that they can better manage the regions they live in, and better understand the world
150 around them. We consider both large-scale and regional-scale models to be useful practices that should
151 both continue to be conducted rather than one replacing another. Ideally large-scale and regional-scale
152 models should benefit from the other since each has strengths and weaknesses and together the two
153 practices enrich our understanding and support the management of groundwater across scales (Section
154 2).



155 The challenge of incorporating groundwater processes into continental- or global-scale models is
156 formidable and sometimes controversial. Some of the controversy stems from unanswered questions
157 about how best to represent groundwater in the models whereas some comes from skepticism about
158 the feasibility of modelling groundwater at non-traditional scales. We advocate for the representation of
159 groundwater stores and fluxes in continental to global models for the four reasons described above. We
160 do not claim to have all the answers on how best to meet this challenge. We contend, however, that the
161 hydrologic community needs to work deliberately and constructively towards effective representations
162 of groundwater in global models.

163

164 Driven by the increasing recognition of the purpose of representing groundwater in continental- to
165 global-scale models, many global hydrological models and land surface models have incorporated
166 groundwater to varying levels of complexity depending on the model provenance and purpose. Different
167 from regional-scale groundwater models that generally focus on subsurface dynamics, the focus of these
168 models is on estimating either runoff and streamflow (hydrological models) or land-atmosphere water
169 and energy exchange (land surface models). Simulation of groundwater storages and hydraulic heads
170 mainly serve to quantify baseflow that affects streamflow during low flow periods or capillary rise that
171 increases evapotranspiration. Some land-surface models use approaches based on the topographic
172 index to simulate fast surface and slow subsurface runoff based on the fraction of saturated area in the
173 grid cell (Clark et al., 2015; Fan et al., 2019); groundwater in these models does not have water storage
174 or hydraulic heads (Famiglietti & Wood, 1994; Koster et al., 2000; Niu et al., 2003; Takata et al., 2003).

175 In many hydrological models, groundwater is represented as a linear reservoir that is fed by
176 groundwater recharge and drains to a river in the same grid cell (Müller Schmied et al., 2014; Gascoïn et
177 al., 2009; Ngo-Duc et al., 2007). Time series of groundwater storage but not hydraulic heads are
178 computed. This prevents simulation of lateral groundwater flow between grid cells, capillary rise and



179 two-way exchange flows between surface water bodies and groundwater (Döll et al., 2016). However,
180 representing groundwater as a water storage compartment that is connected to soil and surface water
181 bodies by groundwater recharge and baseflow and is affected by groundwater abstractions and returns,
182 enables global-scale assessment of groundwater resources and stress (Herbert and Döll, 2019) and
183 groundwater depletion (Döll et al., 2014a; Wada et al., 2014; de Graaf et al., 2014). In some land surface
184 models, the location of the groundwater table with respect to the land surface is simulated within each
185 grid cell to enable simulation of capillary rise (Niu et al., 2007) but, as in the case of simulating
186 groundwater as a linear reservoir, lateral groundwater transport or two-way surface water-groundwater
187 exchange cannot be simulated with this approach.

188

189 Increasingly, models for simulating groundwater flows between all model grid cells in entire countries or
190 globally have been developed, either as stand-alone models or as part of hydrological models (Vergnes
191 & Decharme, 2012; Fan et al., 2013; Lemieux et al. 2008; de Graaf et al., 2017; Kollet et al., 2017;
192 Maxwell et al., 2015; Reinecke et al., 2018, de Graaf et al 2019). The simulation of groundwater in large-
193 scale models is a nascent and rapidly developing field with significant computational and
194 parameterization challenges which have led to significant and important efforts to develop and evaluate
195 individual models. It is important to note that herein ‘large-scale models’ refer to models that are
196 laterally extensive across multiple regions (hundreds to thousands of kilometers) and generally include
197 the upper tens to hundreds of meters of subsurface and have resolutions sometimes as small as ~1 km.
198 In contrast, ‘regional-scale’ models (tens to hundreds of kilometers) have long been developed for a
199 specific region or aquifer and can include greater depths and resolutions, more complex
200 hydrostratigraphy and are often developed from conceptual models with significant regional knowledge.
201 Regional-scale models include a diverse range of approaches from stand-alone groundwater models
202 (i.e., representing surface water and vadose zone processes using boundary conditions such as recharge)



203 to fully integrated groundwater-surface water models. In the future, large-scale models could be
204 developed in a number of different directions which we only briefly introduce here to maintain our
205 primary focus on model evaluation. One important direction is clearer representation of three-
206 dimensional geology and heterogeneity including karst (Condon et al. in prep) which should be
207 considered as part of conceptual model development prior to numerical model implementation.
208

209 Now that a number of models that represent groundwater at continental to global scales have been
210 developed and will continue evolving, it is equally important that we advance how we evaluate these
211 models. To date, large-scale model evaluation has largely focused on individual models and lacked the
212 rigor of regional-scale model evaluation, with inconsistent practices between models and little
213 community-level discussion or cooperation. Overall, we have only a partial and piecemeal understanding
214 of the capabilities and limitations of different approaches to representing groundwater in large-scale
215 models. Our objective is to provide clear recommendations for evaluating groundwater representation
216 in continental and global models. We focus on model evaluation because this is the heart of model trust
217 and reproducibility (Hutton et al., 2016) and improved model evaluation will guide how and where it is
218 most important to focus future model development. We describe current model evaluation practices
219 (Section 2) and consider diverse and uncertain sources of information, including observations, models
220 and experts to holistically evaluate the simulation of groundwater-related fluxes, stores and hydraulic
221 heads (Section 3). We stress the need for an iterative and open-ended process of model improvement
222 through continuous model evaluation against the different sources of information. We explicitly
223 contrast the terminology used herein of ‘evaluation’ and ‘comparison’ against terminology such as
224 ‘calibration’ or ‘validation’ or ‘benchmarking’, which suggests a modelling process that is at some point
225 complete. We extend previous commentaries advocating improved hydrologic process representation
226 and evaluation in large-scale hydrologic models (Clark et al. 2015; Melsen et al. 2016) by adding expert-



227 elicitation and machine learning for more holistic evaluation. We also consider model objective and
228 model evaluation across the diverse hydrologic landscapes which can both uncover blindspots in model
229 development. It is important to note that we do not consider water quality or contamination, even
230 though water quality or contamination is important for water resources, management and
231 sustainability, since large-scale water quality models are in their infancy (van Vliet et al., 2019)

232

233 We bring together somewhat disparate scientific communities as a step towards greater community-
234 level cooperation on these challenges, including global hydrology and land surface modelers, local to
235 regional hydrogeologists, and hydrologists focused on model development and evaluation. We see three
236 audiences beyond those currently directly involved in large-scale groundwater modeling that we seek to
237 engage to accelerate model evaluation: 1) regional hydrogeologists who could be reticent about global
238 models, and yet have crucial knowledge and data that would improve evaluation; 2) data scientists with
239 expertise in machine learning, artificial intelligence etc. whose methods could be useful in a myriad of
240 ways; and 3) the multiple Earth Science communities that are currently working towards integrating
241 groundwater into a diverse range of models so that improved evaluation approaches are built directly
242 into model development.

243 **2. CURRENT MODEL EVALUATION PRACTICES**

244 Here we provide a brief overview of the synergies and differences between regional-scale and large-
245 scale model evaluation and development as well as the imitations of current evaluation practices for
246 large-scale models.

247

248 **2.1 Synergies between regional-scale and large-scales**



249 Regional-scale and large-scale groundwater models are both governed by the same physical equations
250 and share many of the same challenges. Like large-scale models, some regional-scale models have
251 challenges with representing important regional hydrologic processes such as mountain block recharge
252 (Markovich et al. 2019), and data availability challenges (such as the lack of reliable subsurface
253 parameterization and hydrologic monitoring data) are common. We propose there are largely untapped
254 potential synergies between regional-scale and large-scale models based on these commonalities and
255 the inherent strengths and limitations of each scale (Section 1).

256

257 Much can be learned from regional-scale models to inform the development and evaluation of large-
258 scale groundwater models. Regional-scale models are evaluated using a variety of data types, some of
259 which are available and already used at the global scale and some of which are not. In general, the most
260 common data types used for regional-scale groundwater model evaluation match global-scale
261 groundwater models: hydraulic head and either total streamflow or baseflow estimated using
262 hydrograph separation approaches (eg. RRCA, 2003; Woolfenden and Nishikawa, 2014; Tolley et al.,
263 2019). However, numerous data sources unavailable or not currently used at the global scale have also
264 been applied in regional-scale models, such as elevation of surface water features (Hay et al., 2018),
265 existing maps of the potentiometric surface (Meriano and Eyles, 2003), and dendrochronology (Schilling
266 et al., 2014) - these and other 'non-classical' observations (Schilling et al. 2019) could be the inspiration
267 for model evaluation of large-scale models in the future but are beyond our scope to discuss. Further,
268 given the smaller domain size of regional-scale models, expert knowledge and local ancillary data
269 sources can be more directly integrated and automated parameter estimation approaches such as PEST
270 are tractable (Leaf et al., 2015; Hunt et al., 2013). We directly build upon this practice of integration of
271 expert knowledge below in Section 3.3.

272



273 We propose that there may also be potential benefits of large-scale models for the development of
274 regional-scale models. For instance, the boundary conditions of some regional-scale models could be
275 improved with large-scale model results. The boundary conditions of regional-scale models are often
276 assumed, calibrated or derived from other models or data. In a regional-scale model, increasing the
277 model domain (moving the boundary conditions away from region of interests) or incorporating more
278 hydrologic processes (for example, moving the boundary condition from recharge to the land surface
279 incorporating evapotranspiration and infiltration) both can reduce the impact of boundary conditions on
280 the region and problem of interest. Another potential benefit of large-scale models for regional-scale
281 models is the more fulsome inclusion of large-scale hydrologic and human processes that could further
282 enhance the ability of regional-scale models to address both the science-focused and sustainability-
283 focused purposes described in Section 1. For example, the stronger representation of large-scale
284 atmospheric processes means that the downwind impact of groundwater irrigation on
285 evapotranspiration on precipitation and streamflow can be assessed (DeAngelis et al., 2010; Kustu et al.,
286 2011). Or, the effects of climate change and increased water use that affect the inflow of rivers into the
287 regional modelling domain can be taken from global scale analyses (Wada and Bierkens, 2014). Also,
288 regional groundwater depletion might be largely driven by virtual water trade which can be better
289 represented in global analysis and models than regional-scale models (Dalin et al. 2017). Therefore the
290 processes and results of large-scale models could be used to make regional-scale models even more
291 robust and better address key science and sustainability questions.

292

293 Given the strengths of regional models, a potential alternative to development of large-scale
294 groundwater models would be combining or aggregating multiple regional models in a patchwork
295 approach (as in Zell and Sanford, 2020) to provide global coverage. This would have the advantage of
296 better respecting regional differences but potentially create additional challenges because the regional



297 models would have different conceptual models, governing equations, boundary conditions etc. in
298 different regions. Some challenges of this patchwork approach include 1) the required collaboration of a
299 large number of experts from all over the world over a long period of time; 2) regional groundwater flow
300 models alone are not sufficient, they need to be integrated into a hydrological model so that
301 groundwater-soil water and the surface water-groundwater interactions can be simulated; 3) the extent
302 of regional aquifers does not necessarily coincide with the extent of river basins; and 4) the bias of
303 regional groundwater models towards important aquifers which as described above, underlie only a
304 portion of the world's land mass or population and may bias estimates of fluxes such as surface water-
305 groundwater exchange or evapotranspiration. Given these challenges, we argue that a patchwork
306 approach of integrating multiple regional models is a compelling idea but likely insufficient to achieve
307 the purposes of large-scale groundwater modeling described in Section 1. Although this nascent idea of
308 aggregating regional models is beyond the scope of this manuscript, we consider this an important
309 future research avenue, and encourage further exploration and improvement of regional-scale model
310 integration from the groundwater modeling community.

311

312 **2.2 Differences between regional-scale and large-scales**

313 Although there are important similarities and potential synergies across scales, it is important to
314 consider how or if large-scale models are fundamentally different to regional-scale models, especially in
315 ways that could impact evaluation. The primary differences between large-scale and regional-scale
316 models are that large-scale models (by definition) cover larger areas and, as a result, typically include
317 more data-poor areas and are generally built at coarser resolution. These differences impact evaluations
318 in at least five relevant ways:

319 1) Commensurability errors (also called 'representativeness' errors) occur either when modelled grid
320 values are interpolated and compared to an observation 'point' or when aggregation of observed



321 'point' values are compared to a modelled grid value (Beven, 2005; Tustison et al., 2001; Beven,
322 2016; Pappenberger et al., 2009; Rajabi et al., 2018). For groundwater models in particular,
323 commensurability error will depend on the number and locations of observation points, the
324 variability structure of the variables being compared such as hydraulic head and the interpolation or
325 aggregation scheme applied (Tustison et al., 2001; Pappenberger et al., 2009; Reinecke et al., 2020).
326 Commensurability is a problem for most scales of modelling, but likely more significant the coarser
327 the model. Regional-scale groundwater models typically have fewer (though not insignificant)
328 commensurability issues due to smaller grid cell sizes compared to large-scale models.

329 2) Specificity to region, objective and model evaluation criteria because regional-scale models are
330 developed specifically for a certain region and modeling or management objective whereas large-
331 scale models are often more general and include different regions. As a result, large-scale models
332 often have greater heterogeneity of processes and parameters, may not adopt the same calibration
333 targets and variables, and are not subject to the policy or litigation that sometimes drives model
334 evaluation of regional-scale models.

335 3) Computational requirements can be immense for large-scale models which leads to challenges with
336 uncertainty and sensitivity analysis. While some regional-scale models also have large
337 computational demands, large-scale models cover larger domains and are therefore more
338 vulnerable to this potential constraint.

339 4) Data availability for large-scale models can be limited because they typically include data-poor
340 areas, which leads to challenges when only using observations for model evaluation. While data
341 availability also affects regional-scale models, they are often developed for regions with known
342 hydrological challenges based on existing data and/or modeling efforts are preceded by significant
343 regional data collection from detailed sources (such as local geological reports) that are not often
344 included in continental to global datasets used for large-scale model parameterization.



345 5) Subsurface detail in regional-scale models routinely include heterogeneous and anisotropic
346 parameterizations which could be improved in future large-scale models. For example, intense
347 vertical anisotropy routinely induces vertical flow dynamics from vertical head gradients that are
348 tens to thousands of times greater than horizontal gradients which profoundly alter the meaning of
349 the deep and shallow groundwater levels, with only the latter remotely resembling the actual water
350 table. In contrast, currently most large-scale models use a single vertically homogeneous value for
351 each grid cell, or at best have two layers (de Graaf et al., 2017)

352

353 **2.3 Limitations of current evaluation practices for large-scale models**

354 Evaluation of large-scale models has often focused on streamflow or evapotranspiration observations
355 but joint evaluation together with groundwater-specific variables is appropriate and necessary (e.g.
356 Maxwell et al. 2015; Maxwell and Condon, 2016). Groundwater-specific variables useful for evaluating
357 the groundwater component of large-scale models include a) hydraulic head or water table depth; b)
358 groundwater storage and groundwater storage changes which refer to long-term, negative or positive
359 trends in groundwater storage where long-term, negative trends are called groundwater depletion; c)
360 groundwater recharge; d) flows between groundwater and surface water bodies; and e) human
361 groundwater abstractions and return flows to groundwater. It is important to note that groundwater
362 and surface water hydrology communities often have slightly different definitions of terms like recharge
363 and baseflow (Barthel, 2014); we therefore suggest trying to precisely define the meanings of such
364 words using the actual hydrologic fluxes which we do below. Table 1 shows the availability of
365 observational data for these variables but does not evaluate the quality and robustness of observations.
366 Overall there are significant inherent challenges of commensurability and measurability of groundwater
367 observations in the evaluation of large-scale models. We describe the current model evaluation
368 practices for each of these variables here:



369

370 a) Simulated hydraulic heads or water table depth in large scale models are frequently compared
371 to well observations, which are often considered the crucial data for groundwater model
372 evaluation. Hydraulic head observations from a large number groundwater wells (>1 million)
373 have been used to evaluate the spatial distribution of steady-state heads (Fan et al., 2013, de
374 Graaf et al., 2015; Maxwell et al., 2015; Reinecke et al., 2019a, 2020). Transient hydraulic heads
375 with seasonal amplitudes (de Graaf et al. 2017), declining heads in aquifers with groundwater
376 depletion (de Graaf et al. 2019) and daily transient heads (Tran et al 2020) have also been
377 compared to well observations. All evaluation with well observations is severely hampered by
378 the incommensurability of point values of observed head with simulated heads that represent
379 averages over cells of a size of tens to hundreds square kilometers; within such a large cell, land
380 surface elevation, which strongly governs hydraulic head, may vary a few hundred meters, and
381 average observed head strongly depends on the number and location of well within the cell
382 (Reinecke et al., 2020). Additional concerns with head observations are the 1) strong sampling
383 bias of wells towards accessible locations, low elevations, shallow water tables, and more
384 transmissive aquifers in wealthy, generally temperate countries (Fan et al., 2019); 2) the impacts
385 of pumping which may or may not be well known; 3) observational errors and uncertainty (Post
386 and von Asmuth, 2013; Fan et al., 2019); and 4) that heads can reflect the poro-elastic effects of
387 mass loading and unloading rather than necessarily aquifer recharge and drainage (Burgess et al,
388 2017). To date, simulated hydraulic heads have more often been compared to observed heads
389 (rather than water table depth) which results in lower relative errors (Reinecke et al., 2020)
390 because the range of heads (10s to 1000s m head) is much larger than the range of water table
391 depths (<1 m to 100s m).

392



393 b) Simulated groundwater storage trends or anomalies in large-scale hydrological models have
394 been evaluated using observations of groundwater well levels combined with estimates of
395 storage parameters, such as specific yield; local-scale groundwater modeling; and translation of
396 regional total water storage trends and anomalies from satellite gravimetry (GRACE: Gravity
397 Recovery And Climate Experiment) to groundwater storage changes by estimating changes in
398 other hydrological storages (Döll et al., 2012; 2014a). Groundwater storage changes volumes
399 and rates have been calculated for numerous aquifers, primarily in the United States, using
400 calibrated groundwater models, analytical approaches, or volumetric budget analyses (Konikow,
401 2010). Regional-scale models have also been used to simulate groundwater storage trends
402 untangling the impacts of water management during drought (Thatch et al. 2020). Satellite
403 gravimetry (GRACE) is important but has limitations (Alley and Konikow, 2015). First, monthly
404 time series of very coarse-resolution groundwater storage are indirectly estimated from
405 observations of total water storage anomalies by satellite gravimetry (GRACE) but only after
406 model- or observation-based subtraction of water storage changes in glaciers, snow, soil and
407 surface water bodies (Lo et al., 2016; Rodell et al., 2009; Wada, 2016). As soil moisture, river or
408 snow dynamics often dominate total water storage dynamics, the derived groundwater storage
409 dynamics can be so uncertain that severe groundwater drought cannot be detected in this way
410 (Van Loon et al., 2017). Second, GRACE cannot detect the impact of groundwater abstractions
411 on groundwater storage unless groundwater depletion occurs (Döll et al., 2014a,b). Third, the
412 very coarse resolution can lead to incommensurability but in the opposite direction of well
413 observations. It is important to note that the focus is on storage trends or anomalies since total
414 groundwater storage to a specific depth (Gleeson et al., 2016) or in an aquifer (Konikow, 2010)
415 can be estimated but the total groundwater storage in a specific region or cell cannot be
416 simulated or observed unless the depth of interest is specified (Condon et al., 2020).



417

418 c) Simulated large-scale groundwater recharge (vertical flux across the water table) has been
419 evaluated using compilations of point estimates of groundwater recharge, results of regional-
420 scale models, baseflow indices, and expert opinion (Döll and Fiedler, 2008; Hartmann et al.,
421 2015) or compared between models (e.g. Wada et al. 2010). In general, groundwater recharge is
422 not directly measurable except by meter-scale lysimeters (Scanlon et al., 2002), and many
423 groundwater recharge methods such as water table fluctuations and chloride mass balance also
424 suffer from similar commensurability issues as water table depth data. Although sometimes an
425 input or boundary condition to regional-scale models, recharge in many large-scale groundwater
426 models is simulated and thus can be evaluated.

427

428 d) The flows between groundwater and surface water bodies (rivers, lakes, wetlands) are
429 simulated by many models but are generally not evaluated directly against observations of such
430 flows since they are very rare and challenging. Baseflow (the slowly varying portion of
431 streamflow originating from groundwater or other delayed sources) or streamflow ‘low flows’
432 (when groundwater or other delayed sources predominate), generally cannot be used to directly
433 quantify the flows between groundwater and surface water bodies at large scales. Groundwater
434 discharge to rivers can be estimated from streamflow observations only in the very dense gauge
435 network and/or if streamflow during low flow periods is mainly caused by groundwater
436 discharge and not by water storage in upstream lakes, reservoirs or wetlands. These conditions
437 are rarely met in case of streamflow gauges with large upstream areas that can be used for
438 comparison to large-scale model output. de Graaf et al. (2019) compared the simulated timing
439 of changes in groundwater discharge to observations and regional-scale models, but only
440 compared the fluxes directly between the global- and regional-scale models. Due to the



441 challenges of directly observing the flows between groundwater and surface water bodies at
442 large scales, this is not included in the available data in Table 1; instead in Section 3 we highlight
443 the potential for using baseflow or the spatial distribution of perennial, intermittent and
444 ephemeral streams in the future.

445

446 e) Groundwater abstractions have been evaluated by comparison to national, state and county
447 scale statistics in the U.S. (Wada et al. 2010, Döll et al., 2012, 2014a, de Graaf et al. 2014).
448 Irrigation is the dominant groundwater use sector in many regions; however, irrigation pumpage
449 is generally estimated from crop water demand and rarely metered although GRACE and other
450 remote sensing data have been used to estimate the irrigation water demand (Anderson et al.
451 2015b). The lack of records or observations of abstraction introduces significant uncertainties
452 into large-scale models and is simulated and thus can be evaluated. Human groundwater
453 abstractions and return flows as well as groundwater recharge and the flows between
454 groundwater and surface water bodies are necessary to simulate storage trends (described
455 above). But each of these are considered separate observations since they each have different
456 data sources and assumptions. Groundwater abstraction data at the well scale are severely
457 hampered by the incommensurability like hydraulic head and recharge described above.

458 3. HOW TO IMPROVE THE EVALUATION OF LARGE-SCALE GROUNDWATER MODELS

459 Based on Section 2, we argue that the current model evaluation practices are insufficient to robustly
460 evaluate large-scale models. We therefore propose evaluating large-scale models using at least three
461 strategies (pie-shapes in Figure 1): observation-, model-, and expert-driven evaluation which are
462 potentially mutually beneficial because each strategy has its strengths and weaknesses. We are not
463 proposing a brand new evaluation method here but rather separating strategies to consider the problem



464 of large-scale model evaluation from different but highly interconnected perspectives. All three
465 strategies work together for the common goal of ‘improved model large-scale model evaluation’ which
466 is what is the centre of Figure 1.

467

468 When evaluating large-scale models, it is necessary to first consider reasonable expectations or how to
469 know a model is ‘well enough’. Reasonable expectations should be based on the modeling purpose,
470 hydrologic process understanding and the plausibly achievable degree of model realism. First, model
471 evaluation should be clearly linked to the four science- or sustainability-focused purposes of
472 representing groundwater in large-scale models (Section 1) and second, to our understanding of
473 relevant hydrologic processes. The objective of large-scale models cannot be to reproduce the spatio-
474 temporal details that regional-scale models can reproduce. Determining the reasonable expectations is
475 necessarily subjective, but can be approached using observation-, model-, and expert-driven evaluation.
476 As a simple first step in setting realistic expectations, we propose that three physical variables can be
477 used to form more convincing arguments that a large-scale model is well enough: change in
478 groundwater storage, water table depth, and regional fluxes between groundwater and surface water.
479 Below we explore in more detail additional variables and approaches that can support this simple
480 approach.

481

482 Across all three model evaluation strategies of observation-, model-, and expert-driven evaluation, we
483 advocate three principles underpinning model evaluation (base of Figure 1), none of which we are the
484 first to suggest but we highlight here as a reminder: 1) model objectives, such as the groundwater
485 science or groundwater sustainability objective summarised in Section 1, are important to model
486 evaluation because they provide the context through which relevance of the evaluation outcome is set;
487 2) all sources of information (observations, models and experts) are uncertain and this uncertainty



488 needs to be quantified for robust evaluation; and 3) regional differences are likely important for large-
489 scale model evaluation - understanding these differences is crucial for the transferability of evaluation
490 outcomes to other places or times.

491

492 We stress that we see the consideration and quantification of uncertainty as an essential need across all
493 three types of model evaluation we describe below, so we discuss it here rather than with model-driven
494 model evaluation (Section 3.2) where uncertainty analysis more narrowly defined would often be
495 discussed. We further note that large-scale models have only been assessed to a very limited degree
496 with respect to understanding, quantifying, and attributing relevant uncertainties. Expanding computing
497 power, developing computationally frugal methods for sensitivity and uncertainty analysis, and
498 potentially employing surrogate models can enable more robust sensitivity and uncertainty analysis
499 such as used in regional-scale models (Habets et al., 2013; Hill, 2006; Hill & Tiedeman, 2007; Reinecke et
500 al., 2019b). For now, we suggest applying computationally frugal methods such as the elementary effect
501 test or local sensitivity analysis (Hill, 2006; Morris, 1991; Saltelli et al., 2000). Such sensitivity and
502 uncertainty analyses should be applied not only to model parameters and forcings but also to model
503 structural properties (e.g. boundary conditions, grid resolution, process simplification, etc.) (Wagener
504 and Pianosi, 2019). This implies that the (independent) quantification of uncertainty in all model
505 elements (observations, parameters, states, etc.) needs to be improved and better captured in available
506 metadata.

507

508 We advocate for considering regional differences more explicitly in model evaluation since likely no
509 single model will perform consistently across the diverse hydrologic landscapes of the world (Van
510 Werkhoven et al., 2008). Considering regional differences in large-scale model evaluation is motivated
511 by recent model evaluation results and is already starting to be practiced. Two recent sensitivity



512 analyses of large-scale models reveal how sensitivities to input parameters vary in different regions for
513 both hydraulic heads and flows between groundwater and surface water (de Graaf et al. 2019; Reinecke
514 et al., 2020). In mountain regions, large-scale models tend to underestimate steady-state hydraulic
515 head, possibly due to over-estimated hydraulic conductivity in these regions, which highlights that
516 model performance varies in different hydrologic landscapes. (de Graaf et al., 2015; Reinecke et al.
517 2019b). Additionally, there are significant regional differences in performance with low flows for a
518 number of large-scale models (Zaherpour et al. 2018) likely because of diverse implementations of
519 groundwater and baseflow schemes. Large-scale model evaluation practice is starting to shift towards
520 highlighting regional differences as exemplified by two different studies that explicitly mapped
521 hydrologic landscapes to enable clearer understanding of regional differences. Reinecke et al. (2019b)
522 identified global hydrological response units which highlighted the spatially distributed parameter
523 sensitivities in a computationally expensive model, whereas Hartmann et al. (2017) developed and
524 evaluated models for karst aquifers in different hydrologic landscapes based on different a priori system
525 conceptualizations. Considering regional differences in model evaluation suggests that global models
526 could in the future consider a patchwork approach of different conceptual models, governing equations,
527 boundary conditions etc. in different regions. Although beyond the scope of this manuscript, we
528 consider this an important future research avenue.

529 **3.1 Observation-based model evaluation**

530 Observation-based model evaluation is the focus of most current efforts and is important because we
531 want models to be consistent with real-world observations. Section 2 and Table 1 highlight both the
532 strengths and limitations of current practices using observations. Despite existing challenges, we foresee
533 significant opportunities for observation-based model evaluation and do not see data scarcity as a
534 reason to exclude groundwater in large-scale models or to avoid evaluating these models. It is important



535 to note that most so-called ‘observations’ are modeled or derived quantities, and often at the wrong
536 scale for evaluating large-scale models (Table 1; Beven, 2019). Given the inherent challenges of direct
537 measurement of groundwater fluxes and stores especially at large scales, herein we consider the word
538 ‘observation’ loosely as any measurements of physical stores or fluxes that are combined with or filtered
539 through models for an output. For example, GRACE gravity measurements are combined with model-
540 based estimates of water storage changes in glaciers, snow, soil and surface water for ‘groundwater
541 storage change observations’ or streamflow measurements are filtered through baseflow separation
542 algorithms for ‘baseflow observations’. The strengths and limitations as well as the data availability and
543 spatial and temporal attributes of different observations are summarized in Table 1 which we hope will
544 spur more systematic and comprehensive use of observations.

545

546 Here we highlight nine important future priorities for improving evaluation using available observations.
547 The first five priorities focus on current observations (Table 1) whereas the latter four focus on new
548 methods or approaches:

549 1) Focus on transient observations of the water table depth rather than hydraulic head
550 observations that are long-term averages or individual times (often following well
551 drilling). Water table depth are likely more robust evaluation metrics than hydraulic
552 head because water table depth reveals great discrepancies and is a complex function of
553 the relationship between hydraulic head and topography that is crucial to predicting
554 system fluxes (including evapotranspiration and baseflow). Comparing transient
555 observations and simulations instead of long-term averages or individual times
556 incorporates more system dynamics of storage and boundary conditions as temporal
557 patterns are more important than absolute values (Heudorfer et al. 2019). For regions
558 with significant groundwater depletion, comparing to declining water tables is a useful



559 strategy (de Graaf et al. 2019), whereas in aquifers without groundwater depletion,
560 seasonally varying water table depths are likely more useful observations (de Graaf et
561 al. 2017).

562 2) Use baseflow, the slowly varying portion of streamflow originating from groundwater or
563 other delayed sources. Döll and Fiedler (2008) included the baseflow index in evaluating
564 recharge and baseflow has been used to calibrate the groundwater component of a land
565 surface model (Lo et al. 2008, 2010). But the baseflow index (BFI), linear and nonlinear
566 baseflow recession behavior or baseflow fraction (Gnann et al., 2019) have not been
567 used to evaluate any large-scale model that simulates groundwater flows between all
568 model grid cells. There are limitations of using BFI and baseflow recession characteristics
569 to evaluate large-scale models (Table 1). Using baseflow only makes sense when the
570 baseflow separation algorithm is better than the large-scale model itself, which may not
571 be the case for some large-scale models and only in time periods that can be assumed
572 to be dominated by groundwater discharge. Similarly, using recession characteristics is
573 dependent on an appropriate choice of recession extraction methods. But this remains
574 available and obvious data derived from streamflow or spring flow observations that has
575 been under-used to date.

576 3) Use the spatial distribution of perennial, intermittent, and ephemeral streams as an
577 observation, which to our best knowledge has not been done by any large-scale model
578 evaluation. The transition between perennial and ephemeral streams is an important
579 system characteristic in groundwater-surface water interactions (Winter et al. 1998), so
580 we suggest that this might be a revealing evaluation criteria although there are similar
581 limitations to using baseflow. The results of both quantifying baseflow and mapping
582 perennial streams depend on the methods applied, they are not useful for quantifying



583 groundwater-surface water interactions when there is upstream surface water storage,
584 and they do not directly provide information about fluxes between groundwater and
585 surface water.

586 4) Use data on land subsidence to infer head declines or aquifer properties for regions
587 where groundwater depletion is the main cause of compaction (Bierkens and Wada,
588 2019). Lately, remote sensing methods such as GPS, airborne and space borne radar and
589 lidar are frequently used to infer land subsidence rates (Erban et al., 2014). Also, a
590 number of studies combine geomechanical modelling (Ortega-Guerrero et al 1999;
591 Minderhoud et al 2017) and geodetic data to explain the main drivers of land
592 subsidence. A few papers (e.g. Zhang and Burbey 2016) use a geomechanical model
593 together with a withdrawal data and geodetic observations to estimate hydraulic and
594 geomechanical subsoil properties.

595 5) Consider using socio-economic data for improving model input. For example, reported
596 crop yields in areas with predominant groundwater irrigation could be used to evaluate
597 groundwater abstraction rates. Or using well depth data (Perrone and Jasechko, 2019)
598 to assess minimum aquifer depths or in coastal regions and deltas, the presence of
599 deeper fresh groundwater under semi-confining layers.

600 6) Derive additional new datasets using meta-analysis and/or geospatial analysis such as
601 gaining or losing stream reaches (e.g., from interpolated head measurements close to
602 the streams), springs and groundwater-dependent surface water bodies, or tracers.
603 Each of these new data sources could in principle be developed from available data
604 using methods already applied at regional scales but do not currently have an 'off the
605 shelf' global dataset. For example, some large-scale models have been explicitly
606 compared with residence time and tracer data (Maxwell et al., 2016) which have also



607 been recently compiled globally (Gleeson et al., 2016; Jasechko et al., 2017). This could
608 be an important evaluation tool for large-scale models that are capable of simulating
609 flow paths, or can be modified to do although a challenge of this approach is the
610 conservativity of tracers. Future meta-analyses data compilations should report on the
611 quality of the data and include possible uncertainty ranges as well as the mean
612 estimates.

613 7) Use machine learning to identify process representations (e.g. Beven, 2020) or
614 spatiotemporal patterns, for example of perennial streams, water table depths or
615 baseflow fluxes, which might not be obvious in multi-dimensional datasets and could be
616 useful in evaluation. For example, Yang et al. (2019) predicted the state of losing and
617 gaining streams in New Zealand using random forests. A staggering variety of machine
618 learning tools are available and their use is nascent yet rapidly expanding in geoscience
619 and hydrology (Reichstein et al., 2019; Shen, 2018; Shen et al., 2018; Wagener et al.,
620 2020). While large-scale groundwater models are often considered ‘data-poor’, it may
621 seem strange to propose using data-intensive machine learning methods to improve
622 model evaluation. But some of the data sources are large (e.g over 2 million water level
623 measurements in Fan et al. 2013 although biased in distribution) whereas other
624 observations such as evapotranspiration (Jung et al., 2011) and baseflow (Beck et al.
625 2013) are already interpolated and extrapolated using machine learning. Moving
626 forwards, it is important to consider commensurability while applying machine learning
627 in this context.

628 8) Consider comparing models against hydrologic signatures - indices that provide insight
629 into the functional behavior of the system under study (Wagener et al., 2007; McMillan,
630 2020). The direct comparison of simulated and observed variables through statistical



631 error metrics has at least two downsides. One, the above mentioned unresolved
632 problem of commensurability, and two, the issue that such error metrics are rather
633 uninformative in a diagnostic sense - simply knowing the size of an error does not tell
634 the modeller how the model needs to be improved, only that it does (Yilmaz et al.,
635 2009). One way to overcome these issues, is to derive hydrologically meaningful
636 signatures from the original data, such as the signatures derived from transient
637 groundwater levels by Heudorfer et al. (2019). For example, recharge ratio (defined as
638 the ratio of groundwater recharge to precipitation) might be hydrologically more
639 informative than recharge alone (Jasechko et al., 2014) or the water table ratio and
640 groundwater response time (Cuthbert et al. 2019; Opie et al., 2020) which are spatially-
641 distributed signatures of groundwater systems dynamics. Such signatures might be used
642 to assess model consistency (Wagener & Gupta, 2005; Hrachowitz et al.2014) by looking
643 at the similarity of patterns or spatial trends rather than the size of the aggregated
644 error, thus reducing the commensurability problem.

645 9) Understand and quantify commensurability error issues better so that a fairer
646 comparison can be made across scales using existing data. As described above,
647 commensurability errors will depend on the number and locations of observation
648 points, the variability structure of the variables being compared such as hydraulic head
649 and the interpolation or aggregation scheme applied. While to some extent we may
650 appreciate how each of these factors affect commensurability error in theory, in
651 practice their combined effects are poorly understood and methods to quantify and
652 reduce commensurability errors for groundwater model purposes remain largely
653 undeveloped. As such, quantification of commensurability error in (large-scale)
654 groundwater studies is regularly overlooked as a source of uncertainty because it cannot



655 be satisfactorily evaluated (Tregoning et al., 2012). Currently, evaluation of simulated
656 groundwater heads is plagued by, as yet, poorly quantified uncertainties stemming from
657 commensurability errors and we therefore recommend future studies focus on
658 developing solutions to this problem. An additional, subtle but important and
659 unresolved commensurability issue can stem from conceptual models. Different
660 hydrogeologists examining different scales, data or interpreting geology differently can
661 produce quite different conceptual models of the same region (Troldborg et al. 2007).

662 We recommend evaluating models with a broader range of currently available data sources (with
663 explicit consideration of data uncertainty and regional differences) while also simultaneously working to
664 derive new data sets. Using data (such as baseflow, land subsidence, or the spatial distribution of
665 perennial, intermittent, and ephemeral streams) that is more consistent with the scale modelled grid
666 resolution will hopefully reduce the commensurability challenges. However, data distribution and
667 commensurability issues will likely still be present, which underscores the importance of the two
668 following strategies.

669 3.2. Model-based model evaluation

670 Model-based model evaluation, which includes model intercomparison projects (MIP) and model
671 sensitivity and uncertainty analysis, can be done with or without explicitly using observations. We
672 describe both inter-model and inter-scale comparisons which could be leveraged to maximize the
673 strengths of each of these approaches.

674

675 The original MIP concept offers a framework to consistently evaluate and compare models, and
676 associated model input, structural, and parameter uncertainty under different objectives (e.g., climate
677 change, model performance, human impacts and developments). Early model intercomparisons of



678 groundwater models focused on nuclear waste disposal (SKI, 1984). Since the Project for the
679 Intercomparison of Land-Surface Parameterization Schemes (PILPS; Sellers et al., 1993), the first large-
680 scale MIP, the land surface modeling community has used MIPs to deepen understanding of land
681 physical processes and to improve their numerical implementations at various scales from regional (e.g.,
682 Rhône-aggregation project; Boone et al., 2004) to global (e.g., Global Soil Wetness Project; Dirmeyer,
683 2011). Two examples of recent model intercomparison efforts illustrate the general MIP objectives and
684 practice. First, ISIMIP (Schewe et al., 2014; Warszawski et al., 2014) assessed water scarcity at different
685 levels of global warming. Second, IH-MIP2 (Kollet et al., 2017) used both synthetic domains and an
686 actual watershed to assess fully-integrated hydrologic models because these cannot be validated easily
687 by comparison with analytical solutions and uncertainty remains in the attribution of hydrologic
688 responses to model structural errors. Model comparisons have revealed differences, but it is often
689 unclear whether these stem from differences in the model structures, differences in how the
690 parameters were estimated, or from other modelling choices (Duan et al., 2006). Attempts for modular
691 modelling frameworks to enable comparisons (Wagener et al., 2001; Leavesley et al., 2002; Clark et al.,
692 2008; Fenicia et al., 2011; Clark et al., 2015) or at least shared explicit modelling protocols and boundary
693 conditions (Refsgaard et al., 2007; Ceola et al., 2015; Warszawski et al., 2014) have been proposed to
694 reduce these problems.

695

696 Inter-scale model comparison - for example, comparing a global model to a regional-scale model - is a
697 potentially useful approach which is emerging for surface hydrology models (Hattermann et al., 2017;
698 Huang et al., 2017) and could be applied to large-scale models with groundwater representation. For
699 example, declining heads and decreasing groundwater discharge have been compared between a
700 calibrated regional-scale model (RRCA, 2003) and a global model (de Graaf et al., 2019). A challenge to
701 inter-scale comparisons is that regional-scale models often have more spatially complex subsurface



702 parameterizations because they have access to local data which can complicate model inter-
703 comparison. Another approach which may be useful is running large-scale models over smaller
704 (regional) domains at a higher spatial resolution (same as a regional-scale model) so that model
705 structure influences the comparison less. In the future, various variables that are hard to directly
706 observe at large scales but routinely simulated in regional-scale models such as baseflow or recharge
707 could be used to evaluate large-scale models. In this way, the output fluxes and intermediate spatial
708 scale of regional models provide a bridge across the “river of incommensurability” between highly
709 location-specific data such as well observations and the coarse resolution of large-scale models. It is
710 important to consider that regional-scale models are not necessarily or inherently more accurate than
711 large-scale models since problems may arise from conceptualization, groundwater-surface water
712 interactions, scaling issues, parameterization etc.

713

714 In order for a regional-scale model to provide a useful evaluation of a large-scale model, there are
715 several important documentation and quality characteristics it should meet. At a bare minimum, the
716 regional-scale model must be accessible and therefore meet basic replicability requirements including
717 open and transparent input and output data and model code to allow large-scale modelers to run the
718 model and interpret its output. Documentation through peer review, either through a scientific journal
719 or agency such as the US Geological Survey, would be ideal. It is particularly important that the
720 documentation discusses limitations, assumptions and uncertainties in the regional-scale model so that
721 a large-scale modeler can be aware of potential weaknesses and guide their comparison accordingly.

722 Second, the boundary conditions and/or parameters being evaluated need to be reasonably comparable
723 between the regional- and large-scale models. For example, if the regional-scale model includes human
724 impacts through groundwater pumping while the large-scale model does not, a comparison of baseflow
725 between the two models may not be appropriate. Similarly, there needs to be consistency in the time



726 period simulated between the two models. Finally, as with data-driven model evaluation, the purpose of
727 the large-scale model needs to be consistent with the model-based evaluation; matching the hydraulic
728 head of a regional-scale model, for instance, does not indicate that estimates of stream-aquifer
729 exchange are valid. Ideally, we recommend developing a community database of regional-scale models
730 that meet this criteria. It is important to note that Rossman & Zlotnik (2014) review 88 regional-scale
731 models while a good example of such a repository is the California Groundwater Model Archive
732 ([https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-](https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-modeling.html)
733 [modeling.html](https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-modeling.html)).

734

735 In addition to evaluating whether models are similar in terms of their outputs, e.g. whether they
736 simulate similar groundwater head dynamics, it is also relevant to understand whether the influence of
737 controlling parameters are similar across models. This type of analysis provides insights into process
738 controls as well as dominant uncertainties. Sensitivity analysis provides the mathematical tools to
739 perform this type of model evaluation (Saltelli et al., 2008; Pianosi et al., 2016; Borgonovo et al., 2017).
740 Recent applications of sensitivity analysis to understand modelled controls on groundwater related
741 processes include the study by Reinecke et al. (2019b) trying to understand parametric controls on
742 groundwater heads and flows within a global groundwater model. Maples et al. (2020) demonstrated
743 that parametric controls on groundwater recharge can be assessed for complex models, though over a
744 smaller domain. As highlighted by both of these studies, more work is needed to understand how to
745 best use sensitivity analysis methods to assess computationally expensive, spatially distributed and
746 complex groundwater models across large domains (Hill et al., 2016). In the future, it would be useful to
747 go beyond parameter uncertainty analysis (e.g. Reinecke et al. 2019b) to begin to look at all of the
748 modelling decisions holistically such as the forcing data (Weiland et al., 2015) and digital elevation
749 models (Hawker et al., 2018). Addressing this problem requires advancements in statistics (more



750 efficient sensitivity analysis methods), computing (more effective model execution), and access to large-
751 scale models codes (Hutton et al. 2016), but also better utilization of process understanding, for
752 example to create process-based groups of parameters which reduces the complexity of the sensitivity
753 analysis study (e.g. Hartmann et al., 2015; Reinecke et al., 2019b).

754 3.3 Expert-based model evaluation

755 A path much less traveled is expert-based model evaluation which would develop hypotheses of
756 phenomena (and related behaviors, patterns or signatures) we expect to emerge from large-scale
757 groundwater systems based on expert knowledge, intuition, or experience. In essence, this model
758 evaluation approach flips the traditional scientific method around by using hypotheses to test the
759 simulation of emergent processes from large-scale models, rather than using large-scale models to test
760 our hypotheses about environmental phenomena. This might be an important path forward for regions
761 where available data is very sparse or unreliable. The recent discussion by Fan et al. (2019) shows how
762 hypotheses about large-scale behavior might be derived from expert knowledge gained through the
763 study of smaller scale systems such as critical zone observatories. While there has been much effort to
764 improve our ability to make hydrologic predictions in ungauged locations through the regionalization of
765 hydrologic variables or of model parameters (Bloeschl et al., 2013), there has been much less effort to
766 directly derive expectations of hydrologic behavior based on our perception of the systems under study.
767

768 Large-scale models could then be evaluated against such hypotheses, thus providing a general
769 opportunity to advance how we connect hydrologic understanding with large-scale modeling - a strategy
770 that could also potentially reduce epistemic uncertainty (Beven et al., 2019), and which may be
771 especially useful for groundwater systems given the data limitations described above. Developing
772 appropriate and effective hypotheses is crucial and should likely focus on large-scale controlling factors



773 or relationships between controlling factors and output in different parts of the model domain;
774 hypotheses that are too specific may only be able to be tested by certain model complexities or in
775 certain regions. To illustrate the type of hypotheses we are suggesting, we list some examples of
776 hypotheses drawn from current literature:

- 777 • water table depth and lateral flow strongly affect transpiration partitioning (Famiglietti and
778 Wood, 1994; Salvucci and Entekhabi, 1995; Maxwell & Condon, 2016);
- 779 • the percentage of inter-basinal regional groundwater flow increases with aridity or decreases
780 with frequency of perennial streams (Gleeson & Manning, 2008; Goderniaux et al, 2013; Schaller
781 and Fan, 2008); or
- 782 • human water use systematically redistributes water resources at the continental scale via non-
783 local atmospheric feedbacks (Al-Yaari et al., 2019; Keune et al., 2018).

784 Alternatively, it might be helpful to also include hypotheses that have been shown to be incorrect since
785 models should also not show relationships that have been shown to not exist in nature. For example of
786 a hypotheses that has recently been shown to be incorrect is that the baseflow fraction (baseflow
787 volume/precipitation volume) follows the Budyko curve (Gnann et al. 2019) . As yet another alternative,
788 hydrologic intuition could form the basis of model experiments, potentially including extreme model
789 experiments (far from the natural conditions). For example, an experiment that artificially lowers the
790 water table by decreasing precipitation (or recharge directly) could hypothesize the spatial variability
791 across a domain regarding how 'the drainage flux will increase and evaporation flux will decrease as the
792 water table is lowered'. These hypotheses are meant only for illustrative purposes and we hope future
793 community debate will clarify the most appropriate and effective hypotheses. We believe that the
794 debate around these hypotheses alone will lead to advance our understanding, or, at least highlight
795 differences in opinion.

796



797 Formal approaches are available to gather the opinions of experts and to integrate them into a joint
798 result, often called expert elicitation (Aspinall, 2010; Cooke, 1991; O’Hagan, 2019). Expert elicitation
799 strategies have been used widely to describe the expected behavior of environmental or man-made
800 systems for which we have insufficient data or knowledge to build models directly. Examples include
801 aspects of future sea-level rise (Bamber and Aspinall, 2013), tipping points in the Earth system (Lenton
802 et al., 2018), or the vulnerability of bridges to scour due to flooding (Lamb et al., 2017). In the
803 groundwater community, expert opinion is already widely used to develop system conceptualizations
804 and related model structures (Krueger et al., 2012; Rajabi et al., 2018; Refsgaard et al., 2007), or to
805 define parameter priors (Ross et al., 2009; Doherty and Christensen, 2011; Brunner et al., 2012;
806 Knowling and Werner, 2016; Rajabi and Ataie-Ashtiani, 2016). The term expert opinion may be
807 preferable to the term expert knowledge because it emphasizes a preliminary state of knowledge
808 (Krueger et al., 2012).

809

810 A critical benefit of expert elicitation is the opportunity to bring together researchers who have
811 experienced very different groundwater systems around the world. It is infeasible to expect that a single
812 person could have gained in-depth experience in modelling groundwater in semi-arid regions, in cold
813 regions, in tropical regions etc. Being able to bring together different experts who have studied one or a
814 few of these systems to form a group would certainly create a whole that is bigger than the sum of its
815 parts. If captured, it would be a tremendous source of knowledge for the evaluation of large-scale
816 groundwater models. Expert elicitation also has a number of challenges including: 1) formalizing this
817 knowledge in such a way that it is still usable by third parties that did not attend the expert workshop
818 itself; and 2) perceived or real differences in perspectives, priorities and backgrounds between regional-
819 scale and large-scale modelers.

820



821 So, while expert opinion and judgment play a role in any scientific investigation (O'Hagan, 2019),
822 including that of groundwater systems, we rarely use formal strategies to elicit this opinion. It is also less
823 common to use expert opinion to develop hypotheses about the dynamic behavior of groundwater
824 systems, rather than just priors on its physical characteristics. Yet, it is intuitive that information about
825 system behavior can help in evaluating the plausibility of model outputs (and thus of the model itself).
826 This is what we call expert-based evaluation herein. Expert elicitation is typically done in workshops with
827 groups of a dozen or so experts (e.g. Lamb et al., 2018). Upscaling such expert elicitation in support of
828 global modeling would require some web-based strategy and a formalized protocol to engage a
829 sufficiently large number of people. Contributors could potentially be incentivized to contribute to the
830 web platform by publishing a data paper with all contributors as co-authors and a secondary analysis
831 paper with just the core team as coauthors. We recommend the community develop expert elicitation
832 strategies to identify effective hypotheses that directly link to the relevant large-scale hydrologic
833 processes of interest.

834 **4. CONCLUSIONS: towards a holistic evaluation of groundwater representation in large-scale models**

835 Ideally, all three strategies (observation-based, model-based, expert-based) should be pursued
836 simultaneously because the strengths of one strategy might further improve others. For example,
837 expert- or model-based evaluation may highlight and motivate the need for new observations in certain
838 regions or at new resolutions. Or observation-based model evaluation could highlight and motivate
839 further model development or lead to refined or additional hypotheses. We thus recommend the
840 community significantly strengthens efforts to evaluate large-scale models using all three strategies.
841 Implementing these three model evaluation strategies may require a significant effort from the scientific
842 community, so we therefore conclude with two tangible community-level initiatives that would be



843 excellent first steps that can be pursued simultaneously with efforts by individual research groups or
844 collaborations of multiple research groups.

845

846 First, we need to develop a ‘Groundwater Modeling Data Portal’ that would both facilitate and
847 accelerate the evaluation of groundwater representation in continental to global scale models (Bierkens,
848 2015). Existing initiatives such as IGRAC’s Global Groundwater Monitoring Network ([https://www.un-
849 igrac.org/special-project/ggm-global-groundwater-monitoring-network](https://www.un-
849 igrac.org/special-project/ggm-global-groundwater-monitoring-network)) and HydroFrame
850 (www.hydroframe.org), are an important first step but were not designed to improve the evaluation of
851 large-scale models and the synthesized data remains very heterogeneous - unfortunately, even
852 groundwater level time series data often remains either hidden or inaccessible for various reasons. This
853 open and well documented data portal should include:

- 854 a) observations for evaluation (Table 1) as well as derived signatures (Section 3.1);
- 855 b) regional-scale models that meet the standards described above and could facilitate inter-scale
856 comparison (Section 3.2) and be a first step towards linking regional models (Section 2.1);
- 857 c) Schematizations, conceptual or perceptual models of large-scale models since these are the
858 basis of computational models; and
- 859 d) Hypothesis and other results derived from expert elicitation (Section 3.3).

860 Meta-data documentation, data tagging, aggregation and services as well as consistent data structures
861 using well-known formats (netCDF, .csv, .txt) will be critical to developing a useful, dynamic and evolving
862 community resource. The data portal should be directly linked to harmonized input data such as forcings
863 (climate, land and water use etc.) and parameters (topography, subsurface parameters etc.), model
864 codes, and harmonized output data. Where possible, the portal should follow established protocols,
865 such as the Dublin Core Standards for metadata (<https://dublincore.org>) and ISIMIP protocols for
866 harmonizing data and modeling approach, and would ideally be linked to or contained within an existing



867 disciplinary repository such as HydroShare (<https://www.hydroshare.org/>) to facilitate discovery,
868 maintenance, and long-term support. Additionally, an emphasis on model objective, uncertainty and
869 regional differences as highlighted (Section 3) will be important in developing the data portal. Like
870 expert-elicitation, contribution to the data portal could be incentivized through co-authorship in data
871 papers and by providing digital object identifiers (DOIs) to submitted data and models so that they are
872 citable. By synthesizing and sharing groundwater observations, models, and hypotheses, this portal
873 would be broadly useful to the hydrogeological community beyond just improving global model
874 evaluation.

875

876 Second, we suggest ISIMIP, or a similar model intercomparison project, could be harnessed as a
877 platform to improve the evaluation of groundwater representation in continental to global scale models.
878 For example, in ISIMIP (Warszawski et al., 2014), modelling protocols have been developed with an
879 international network of climate-impact modellers across different sectors (e.g. water, agriculture,
880 energy, forestry, marine ecosystems) and spatial scales. Originally, ISIMIP started with multi-model
881 comparison (model-based model evaluation), with a focus on understanding how model projections
882 vary across different sectors and different climate change scenarios (ISIMIP Fast Track). However, more
883 rigorous model evaluation came to attention more recently with ISIMIP2a, and various observation data,
884 such as river discharge (Global Runoff Data Center), terrestrial water storage (GRACE), and water use
885 (national statistics), have been used to evaluate historical model simulation (observation-based model
886 evaluation). To better understand model differences and to quantify the associated uncertainty sources,
887 ISIMIP2b includes evaluating scenarios (land use, groundwater use, human impacts, etc) and key
888 assumptions (no explicit groundwater representation, groundwater availability for the future, water
889 allocation between surface water and groundwater), highlighting that different types of hypothesis
890 derived as part of the expert-based model evaluation could possibly be simulated as part of the ISIMIP



891 process in the future. While there has been a significant amount of research and publications on MIPs
892 including surface water availability, limited multi-model assessments for large-scale groundwater
893 studies exist. Important aspects of MIPs in general could facilitate all three model evaluation strategies:
894 community-building and cooperation with various scientific communities and research groups, and
895 making the model input and output publicly available in a standardized format.

896

897 Large-scale hydrologic and land surface models increasingly represent groundwater, which we envision
898 will lead to a better understanding of large-scale water systems and to more sustainable water resource
899 use. We call on various scientific communities to join us in this effort to improve the evaluation of
900 groundwater in continental to global models. As described by examples above, we have already started
901 this journey and we hope this will lead to better outcomes especially for the goals of including
902 groundwater in large-scale models that we started with above: improving our understanding of Earth
903 system processes; and informing water decisions and policy. Along with the community currently
904 directly involved in large-scale groundwater modeling, above we have made pointers to other
905 communities who we hope will engage to accelerate model evaluation: 1) regional hydrogeologists, who
906 would be useful especially in expert-based model evaluation (Section 3.3); 2) data scientists with
907 expertise in machine learning, artificial intelligence etc. whose methods could be useful especially for
908 observation- and model-based model evaluation (Sections 3.1 and 3.2); and 3) the multiple Earth
909 Science communities that are currently working towards integrating groundwater into a diverse range of
910 models so that improved evaluation approaches are built directly into model development. Together we
911 can better understand what has always been beneath our feet, but often forgotten or neglected.

912

913

914



915 **Competing interests:** The authors declare that they have no conflict of interest.

916

917 **Acknowledgements:**

918 The commentary is based on a workshop at the University of Bristol and significant debate and

919 discussion before and after. This community project was directly supported by a Benjamin Meaker

920 Visiting Professorship at the Bristol University to TG and a Royal Society Wolfson Award to TW

921 (WM170042). We thank many members of the community who contributed to the discussions,

922 especially at the IGEM (Impact of Groundwater in Earth System Models) workshop in Taiwan.

923

924 **Author Contributions:** (using the [CRediT taxonomy](#) which offers standardized descriptions of author

925 contributions) conceptualization and writing original draft: TG, TW and PD; writing - review and

926 editing:all co-authors. Authors are ordered by contribution for the first three coauthors (TG, TW and PD)

927 and then ordered in reverse alphabetical order for all remaining coauthors.

928

929 **Code and data availability:** This Perspective paper does not present any computational results. There is

930 therefore no code or data associated with this paper.

931

932

933

934

935



936 **Table 1. Available observations for evaluating the groundwater component of large-scale models**
 937

Data type	Strengths	Limitations	Data availability and spatial resolution
Available observations already used to evaluate large-scale models			
Hydraulic heads or water table depth (averages or single times)	Direct observation of groundwater levels and storage	observations biased towards North America and Europe; non- commensurable with large-scale models; mixture of observation times	IGRAC Global Groundwater Monitoring Network ; Fan et al., 2013; USGS Point measurements at existing wells
Hydraulic heads or water table depth (transient)	Direct observation of changing groundwater levels and storage	As above	time-series available in a few regions, especially through USGS and European Groundwater Drought Initiative Point measurements at existing wells
Total water storage anomalies (GRACE)	Globally available and regionally integrated signal of water storage trends and anomalies	Groundwater changes are uncertain model remainder; very coarse spatial resolution and limited period	Various mascons gridded with resolution of ~100,000 km ² (Scanlon et al. 2016) which are then processed as groundwater storage change
Storage change (regional aquifers)	Regionally integrated response of aquifer	Bias towards North America and Europe	Konikow 2011 Döll et al., 2014a Regional aquifers (10,000s to 100,000s km ²)
Recharge	Direct inflow of groundwater system	Challenging to measure and upscale	Döll and Fiedler, 2008; Hartmann et al. 2017; Mohan et al. 2018; Moeck et al. 2020 Point to small basin
Abstractions	Crucial for groundwater depletion and sustainability studies	National scale data highly variable in quality; downscaling uncertain	de Graaf et al. 2014 Döll et al. 2014 National-scale data down-scaled to grid
Streamflow or spring flow observations	Widely available at various scales; low flows can be related to groundwater	Challenging to quantify the flows between groundwater and surface water from streamflow	Global Runoff Data Centre (GRDC) or other data sources ; large to small basin; Olariño et al. 2020 point measurements of spring flow

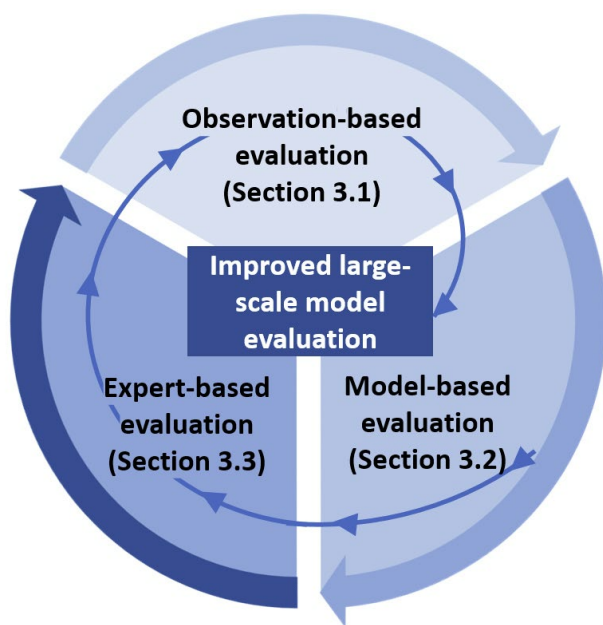


Evapotranspiration	Widely available; related to groundwater recharge or discharge (for shallow water tables)	Not a direct groundwater observations	Various datasets (Miralles et al., 2016); gridded
Available observations not being used to evaluate large-scale models			
Baseflow index (BFI) or (non-)linear baseflow recession behavior	Possible integrator of groundwater contribution to streamflow over a basin	BFI and k values vary with method; baseflow may be dominated by upstream surface water storage rather than groundwater inflow; can not identify losing river conditions	Beck et al. (2013) Point observations extrapolated by machine learning
Perennial stream map	Ephemeral streams are losing streams, whereas perennial streams could be gaining (or impacted by upstream surface water storage)	Mapping perennial streams requires arbitrary streamflow and duration cutoffs; not all perennial streams reaches are groundwater-influenced; does not provide information about magnitude of inflows/outflows.	Schneider et al. (2017) Cuthbert et al. (2019); Spatially continuous along stream networks
Gaining or losing stream reaches	Multiple techniques for measurement (interpolated head measurements, streamflow data, water chemistry). Constrains direction of fluxes at groundwater system boundaries	Relevant processes occur at sub-grid-cell resolution.	Not globally available but see Bresciani et al. (2018) for a regional example; Spatially continuous along stream networks
Springs and groundwater-dependent surface water bodies	Constrains direction of fluxes at groundwater system boundaries	Relevant processes occur at sub-grid-cell resolution.	Springs available for various regions (e.g. Springer, & Stevens, 2009) but not globally; Point measurements at water feature locations
Tracers (heat, isotopes or other geochemical)	Provides information about temporal aspects of groundwater systems (e.g. residence time)	No large-scale models simulate transport processes (Table S1)	Isotopic data compiled (Gleeson et al., 2016; Jasechko et al., 2017) but no global data for heat or other chemistry; Point measurements at existing wells or surface water features



Surface elevation data (leveling, GPS, radar/lidar) an in particular land subsidence observations	Provides information about changes in surface elevation that are related to groundwater head variations or groundwater head decline	Provides indirect information and needs a geomechanical model to translate to head. Introduces additional uncertainty of geomechanical properties.	Leveling data, GPS data and lidar observations mostly limited to areas of active subsidence (e.g. Minderhoud et al., 2019,2020) and not always open. Global data on elevation change are available from the Sentinel 1 mission.
---	---	--	---

938
 939
 940



Improved model evaluation rests of three core principles:

- 1) Modelling purpose or objective are paramount
- 2) All sources of information are uncertain
- 3) Regional differences are important

941
 942
 943
 944
 945
 946
 947
 948
 949
 950

Figure 1: Improved large-scale model evaluation rests on three pillars: observation-, model-, and expert-based model evaluation. We argue that each pillar is an essential strategy so that all three should be simultaneously pursued by the scientific community. The three pillars of model evaluation all rest on three core principles related to 1) model objectives, 2) uncertainty and 3) regional differences.



951 **References**

- 952 Addor, N., & Melsen, L. A. (2018). Legacy, Rather Than Adequacy, Drives the Selection of Hydrological Models.
953 *Water Resources Research*, 0(0). <https://doi.org/10.1029/2018WR022958>
954
- 955 Al-Yaari, A., Ducharne, A., Cheruy, F., Crow, W.T. & Wigneron, J.P. (2019). Satellite-based soil moisture provides
956 missing link between summertime precipitation and surface temperature biases in CMIP5 simulations over
957 conterminous United States. *Scientific Reports*, 9, article number 1657, doi:10.1038/s41598-018-38309-5
958
- 959 Anderson, M. P., Woessner, W. W. & Hunt, R. (2015a). *Applied groundwater modeling- 2nd Edition*. San Diego:
960 Academic Press.
- 961 Anderson, R. G., Min-Hui Lo, Swenson, S., Famiglietti, J. S., Tang, Q., Skaggs, T. H., Lin, Y.-H., and Wu, R.-J. (2015b),
962 Using satellite-based estimates of evapotranspiration and groundwater changes to determine anthropogenic
963 water fluxes in land surface models, *Geosci. Model Dev.*, 8, 3021-3031, doi:10.5194/gmd-8-3021-2015. Alley, W.M.
964 and LF Konikow (2015) Bringing GRACE down to earth. *Groundwater* 53 (6): 826–829
- 965 Anyah, R. O., Weaver, C. P., Miguez-Macho, G., Fan, Y., & Robock, A. (2008). Incorporating water table dynamics in
966 climate modeling: 3. Simulated groundwater influence on coupled land-atmosphere variability. *J. Geophys. Res.*,
967 113. Retrieved from <http://dx.doi.org/10.1029/2007JD009087>
- 968 Archfield, S. A., Clark, M., Arheimer, B., Hay, L. E., McMillan, H., Kiang, J. E., et al. (2015). Accelerating advances in
969 continental domain hydrologic modeling. *Water Resources Research*, 51(12), 10078–10091.
970 <https://doi.org/10.1002/2015WR017498>
- 971 Aspinall, W. (2010). A route to more tractable expert advice. *Nature*, 463, 294–295.
972 <https://doi.org/10.1038/463294a>
- 973 ASTM (2016), Standard Guide for Conducting a Sensitivity Analysis for a Groundwater Flow Model Application,
974 ASTM International D5611-94, West Conshohocken, PA, 2016, www.astm.org
- 975 Bamber, J.L. and Aspinall, W.P. (2013). An expert judgement assessment of future sea level rise from the ice
976 sheets. *Nature Climate Change*. 3(4), 424-427.
- 977 Barnett, B., Townley, L.R., Post, V.E.A., Evans, R.E., Hunt, R.J., Peeters, L., Richardson, S., Werner, A.D., Knapton, A.,
978 Boronkay, A. (2012). Australian groundwater modelling guidelines, National Water Commission, Canberra, 203
979 pages
- 980 Barthel, R. (2014). HESS Opinions “Integration of groundwater and surface water research: an interdisciplinary
981 problem?” *Hydrology and Earth System Sciences*, 18(7), 2615–2628.
- 982 Beck, H. et al (2013). Global patterns in base flow index and recession based on streamflow observations from
983 3394 catchments. *Water Resources Research*.
- 984 Befus, K., Jasechko, S., Luijendijk, E., Gleeson, T., Cardenas, M.B. (2017) The rapid yet uneven turnover of Earth's
985 groundwater. (2017) *Geophysical Research Letters* 11: 5511-5520 doi: 10.1002/2017GL073322
- 986 Best, M. J., Pryor, M., Clark, D. B., Rooney, G. G., Essery, R. L. H., Ménard, C. B., Edwards, J. M., Hendry, M. A.,
987 Porson, A., Gedney, N., Mercado, L. M., Sitch, S., Blyth, E., Boucher, O., Cox, P. M., Grimmond, C. S. B., & Harding,



- 988 R. J. (2011). The Joint UK Land Environment Simulator (JULES), model description – Part 1: Energy and water fluxes,
989 *Geosci. Model Dev.*, 4, 677-699. <https://doi.org/10.5194/gmd-4-677-2011>
- 990 Beven, K. (2000). Uniqueness of place and process representations in hydrological modelling. *Hydrology and Earth*
991 *System Sciences*, 4(2), 203–213.
- 992 Beven, K. (2005). On the concept of model structural error. *Water Science & Technology*, 52(6), 167–175.
- 993 Beven, K. (2016). Facets of uncertainty: epistemic uncertainty, nonstationarity, likelihood, hypothesis testing, and
994 communication. *Hydrological Sciences Journal*, 61(9), 1652-1665, DOI: 10.1080/02626667.2015.1031761
- 995 Beven, K. (2019) How to make advances in hydrological modelling. In: *Hydrology Research*. 50, 6, p. 1481-1494. 14
996 p.
- 997 Beven, K. (2020). Deep learning, hydrological processes and the uniqueness of place. *Hydrological Processes*,
998 34(16), 3608–3613. <https://doi.org/10.1002/hyp.13805>
- 999 Beven, K. J., and H. L. Cloke (2012), Comment on “Hyperresolution global land surface modeling: Meeting a grand
1000 challenge for monitoring Earth’s terrestrial water” by Eric F. Wood et al., *Water Resour.Res.*, 48, W01801,
1001 doi:10.1029/2011WR010982.
- 1002 Beven, K.J., Aspinall, W.P., Bates, P.D., Borgomeo, E., Goda, K., Hall, J.W., Page, T., Phillips, J.C., Simpson, M., Smith,
1003 P.J., Wagener, T. and Watson, M. 2018. Epistemic uncertainties and natural hazard risk assessment – Part 2: What
1004 should constitute good practice? *Natural Hazards and Earth System Sciences*, 18, 10.5194/nhess-18-1-2018
- 1005 Bierkens, M. F. P. (2015). Global hydrology 2015: State, trends, and directions. *Water Resources Research*, 51(7),
1006 4923–4947. <https://doi.org/10.1002/2015WR017173>
- 1007 Bierkens, M. F. P. & Wada, Y. (2019). Non-renewable groundwater use and groundwater depletion: A review.
1008 *Environmental Research Letters*, 14(6), 063002
- 1009 Boone, A. A., Habets, F., Noilhan, J., Clark, D., Dirmeyer, P., Fox, S., Gusev, Y., Haddeland, I., Koster, R., Lohmann,
1010 D., Mahanama, S., Mitchell, K., Nasonova, O., Niu, G. Y., Pitman, A., Polcher, J., Shmakin, A. B., Tanaka, K., Van Den
1011 Hurk, B., Vérant, S., Verseghy, D., Viterbo, P. and Yang, Z. L.: The Rhône-aggregation land surface scheme
1012 intercomparison project: An overview, *J. Clim.*, 17(1), 187–208, doi:10.1175/1520-
1013 0442(2004)017<0187:TRLSSI>2.0.CO;2, 2004.
- 1014 Borgonovo, E. Lu, X. Plischke, E. Rakovec, O. and Hill, M. C. (2017). Making the most out of a hydrological model
1015 data set: Sensitivity analyses to open the model black-box. *Water Resources Research*.
1016 DOI:10.1002/2017WR020767
- 1017 Bresciani, E., P. Goderniaux, and O. Batelaan (2016), Hydrogeological controls of water table-land surface
1018 interactions, *Geophysical Research Letters*, 43, 9653-9661.
- 1019 Bresciani, E., Cranswick, R. H., Banks, E. W., Batlle-Aguilar, J., et al. (2018). Using hydraulic head, chloride and
1020 electrical conductivity data to distinguish between mountain-front and mountain-block recharge to basin aquifers.
1021 *Hydrology and Earth System Sciences*, 22(2), 1629–1648.



- 1022 Brunner, P., J. Doherty, and C. T. Simmons (2012), Uncertainty assessment and implications for data acquisition in
1023 support of integrated hydrologic models, *Water Resources Research*, 48.
1024
- 1025 Burgess, W. G., Shamsudduha, M., Taylor, R. G., Zahid, A., Ahmed, K. M., Mukherjee, A., et al. (2017). Terrestrial
1026 water load and groundwater fluctuation in the Bengal Basin. *Scientific Reports*, 7(1), 3872.
- 1027 Caceres, D., Marzeion, B., Malles, J.H., Gutknecht, B., Müller Schmied, H., Döll, P. (2020): Assessing global water
1028 mass transfers from continents to oceans over the period 1948–2016. *Hydrol. Earth Syst. Sci. Discuss.*
1029 doi:10.5194/hess-2019-664
- 1030 Ceola, S., Arheimer, B., Baratti, E., Blöschl, G., Capell, R., Castellarin, A., et al. (2015). Virtual laboratories: new
1031 opportunities for collaborative water science. *Hydrology and Earth System Sciences*, 19(4), 2101–2117.
- 1032 Clark, M. P., A. G. Slater, D. E. Rupp, R. A. Woods, J. A. Vrugt, H. V. Gupta, T. Wagener, and L. E. Hay (2008)
1033 Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between
1034 hydrological models, *Water Resour. Res.*, 44, W00B02, doi:10.1029/2007WR006735.
- 1035 Clark, M. P., et al. (2015), A unified approach for process-based hydrologic modeling: 1. Modeling concept, *Water*
1036 *Resources Research*, 51, 2498–2514, doi:10.1002/2015WR017198
- 1037 Condon, L. E., & Maxwell, R. M. (2019). Simulating the sensitivity of evapotranspiration and streamflow to large-
1038 scale groundwater depletion. *Science Advances*, 5(6), eaav4574. <https://doi.org/10.1126/sciadv.aav4574>
- 1039 Condon, LE et al Evapotranspiration depletes groundwater under warming over the contiguous United States
1040 *Nature Comm*, 2020, <https://doi.org/10.1038/s41467-020-14688-0>
- 1041 Condon, L. E., Markovich, K. H., Kelleher, C. A., McDonnell, J. J., Ferguson, G., & McIntosh, J. C. (2020). Where Is the
1042 Bottom of a Watershed? *Water Resources Research*, 56(3). <https://doi.org/10.1029/2019wr026010>
- 1043 Condon, L.E., Stefan Kollet, Marc F.P. Bierkens, Reed M. Maxwell, Mary C. Hill, Anne Verhoef, Anne F. Van Loon,
1044 Graham E. Fogg, Mauro Sulis, Harrie-Jan Hendricks Franssen; Corinna Abesser. Global groundwater modeling and
1045 monitoring?: Opportunities and challenges (in preparation)
- 1046 Cooke, R. (1991). *Experts in uncertainty: opinion and subjective probability in science*. Oxford University Press on
1047 Demand.
- 1048 Cuthbert, M. O., Gleeson, T., Moosdorf, N., Befus, K. M., Schneider, A., Hartmann, J., & Lehner, B. (2019). Global
1049 patterns and dynamics of climate–groundwater interactions. *Nature Climate Change*, 9, 137–141
1050 <https://doi.org/10.1038/s41558-018-0386-4>
- 1051 Cuthbert, M. O., et al. (2019) Observed controls on resilience of groundwater to climate variability in sub-Saharan
1052 Africa. *Nature* 572: 230–234
- 1053 Dalin, C., Wada, Y., Kastner, T., & Puma, M. J. (2017). Groundwater depletion embedded in international food
1054 trade. *Nature*, 543(7647), 700–704. <https://doi.org/10.1038/nature21403>
- 1055 DeAngelis, A., Dominguez, F., Fan, Y., Robock, A., Kustu, M. D., & Robinson, D. (2010). Evidence of enhanced
1056 precipitation due to irrigation over the Great Plains of the United States. *Journal of Geophysical Research*:
1057 *Atmospheres*, 115(D15).



- 1058 Dirmeyer, P. A.: A History and Review of the Global Soil Wetness Project (GSWP), *J. Hydrometeorol.*, 12(5),
1059 110404091221083, doi:10.1175/jhm-d-10-05010, 2011
- 1060 Doherty, J., and S. Christensen (2011), Use of paired simple and complex models to reduce predictive bias and
1061 quantify uncertainty, *Water Resources Research*, 47(12),
- 1062 Döll, P., Fiedler, K. (2008): Global-scale modeling of groundwater recharge. *Hydrol. Earth Syst. Sci.*, 12, 863-885,
1063 doi: 10.5194/hess-12-863-2008
- 1064 Döll, P., Douville, H., Güntner, A., Müller Schmied, H., Wada, Y. (2016): Modelling freshwater resources at the
1065 global scale: Challenges and prospects. *Surveys in Geophysics*, 37(2), 195-221. doi: 10.1007/s10712-015-9343-1
- 1066 Döll, P., Müller Schmied, H., Schuh, C., Portmann, F. T., & Eicker, A. (2014a). Global-scale assessment of
1067 groundwater depletion and related groundwater abstractions: Combining hydrological modeling with information
1068 from well observations and GRACE satellites. *Water Resources Research*, 50(7), 5698–5720.
1069 <https://doi.org/10.1002/2014WR015595>
- 1070 Döll, P., Fritsche, M., Eicker, A., Müller Schmied, H. (2014b): Seasonal water storage variations as impacted by
1071 water abstractions: Comparing the output of a global hydrological model with GRACE and GPS observations.
1072 *Surveys in Geophysics*, 35(6), 1311-1331, doi: 10.1007/s10712-014-9282-2.
- 1073 Döll, P., Hoffmann-Dobrev, H., Portmann, F.T., Siebert, S., Eicker, A., Rodell, M., Strassberg, G., Scanlon, B. (2012):
1074 Impact of water withdrawals from groundwater and surface water on continental water storage variations. *J.*
1075 *Geodyn.* 59-60, 143-156, doi:10.1016/j.jog.2011.05.001.
- 1076 Duan Q., Schaake, J., Andreassian, V., Franks, S., Gupta, H.V., Gusev, Y.M., Habets, F., Hall, A., Hay, L., Hogue, T.S.,
1077 Huang, M., Leavesley, G., Liang, X., Nasonova, O.N., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T. and Wood,
1078 E.F. (2006). Model Parameter Estimation Experiment (MOPEX): Overview and Summary of the Second and Third
1079 Workshop Results. *Journal of Hydrology*, 320(1-2), 3-17.
- 1080 Enemark, T., Peeters, L. J. M., Mallants, D., & Batelaan, O. (2019). Hydrogeological conceptual model building and
1081 testing: A review. *Journal of Hydrology*, 569, 310–329. <https://doi.org/10.1016/j.jhydrol.2018.12.007>
- 1082 Erban L E, Gorelick S M and Zebker H A 2014 Groundwater extraction, land subsidence, and sea-level rise in the
1083 Mekong Delta, Vietnam *Environ. Res. Lett.* 9 084010
- 1084 Famiglietti, J. S., & E. F. Wood (1994). Multiscale modeling of spatially variable water and energy balance
1085 processes, *Water Resour. Res.*, 30(11), 3061–3078, <https://doi.org/10.1029/94WR01498>
- 1086 Fan, Y. et al., (2019) Hillslope hydrology in global change research and Earth System modeling. *Water Resources*
1087 *Research*, doi.org/10.1029/2018WR023903
- 1088 Fan, Y. (2015). Groundwater in the Earth's critical zone: Relevance to large-scale patterns and processes. *Water*
1089 *Resources Research*, 51(5), 3052–3069. <https://doi.org/10.1002/2015WR017037>
- 1090 Fan, Y., & Miguez-Macho, G. (2011). A simple hydrologic framework for simulating wetlands in climate and earth
1091 system models. *Climate Dynamics*, 37(1–2), 253–278.



- 1092 Fan, Y., Li, H., & Miguez-Macho, G. (2013). Global patterns of groundwater table depth. *Science*, 339(6122), 940–
1093 943.
- 1094 Fenicia, F., D. Kavetski, and H. H. G. Savenije (2011), Elements of a flexible approach for conceptual hydrological
1095 modeling: 1. Motivation and theoretical development, *Water Resources Research*, 47(11), W11510,
1096 10.1029/2010wr010174.
- 1097 Forrester, M.M. and Maxwell, R.M. Impact of lateral groundwater flow and subsurface lower boundary conditions
1098 on atmospheric boundary layer development over complex terrain. *Journal of Hydrometeorology*,
1099 doi:10.1175/JHM-D-19-0029.1, 2020.
- 1100 Forrester, M.M., Maxwell, R.M., Bearup, L.A., and Gochis, D.J. Forest Disturbance Feedbacks from Bedrock to
1101 Atmosphere Using Coupled Hydro-Meteorological Simulations Over the Rocky Mountain Headwaters. *Journal of*
1102 *Geophysical Research-Atmospheres*, 123:9026-9046, doi:10.1029/2018JD028380 2018.
- 1103 Freeze, R. A., & Witherspoon, P. A. (1966). Theoretical analysis of regional groundwater flow, 1. Analytical and
1104 numerical solutions to a mathematical model. *Water Resources Research*, 2, 641–656.
- 1105 Foster, S., Chilton, J., Nijsten, G.-J., & Richts, A. (2013). Groundwater — a global focus on the ‘local resource.’
1106 *Current Opinion in Environmental Sustainability*, 5(6), 685–695. doi.org/10.1016/j.cosust.2013.10.010
- 1107 Garven, G. (1995). Continental-scale groundwater flow and geologic processes. *Annual Review of Earth and*
1108 *Planetary Sciences*, 23, 89–117.
- 1109 Gascoïn, S., Ducharne, A., Ribstein, P., Carli, M., Habets, F. (2009). Adaptation of a catchment-based land surface
1110 model to the hydrogeological setting of the Somme River basin (France). *Journal of Hydrology*, 368(1-4), 105-116.
1111 <https://doi.org/10.1016/j.jhydrol.2009.01.039>
- 1112 Genereux, D. (1998). Quantifying uncertainty in tracer-based hydrograph separations. *Water Resources Research*,
1113 34(4), 915–919.
- 1114 Gilbert, J.M., Maxwell, R.M. and Gochis, D.J. Effects of water table configuration on the planetary boundary layer
1115 over the San Joaquin River watershed, California. *Journal of Hydrometeorology*, 18:1471-1488, doi:10.1175/JHM-
1116 D-16-0134.1, 2017.
- 1117 Gleeson, T. et al. (2020) HESS Opinions: Improving the evaluation of groundwater representation in continental to
1118 global scale models. <https://hess.copernicus.org/preprints/hess-2020-378/>
- 1119 Gleeson, T., & Manning, A. H. (2008). Regional groundwater flow in mountainous terrain: Three-dimensional
1120 simulations of topographic and hydrogeologic controls. *Water Resources Research*, 44. Retrieved from
1121 <http://dx.doi.org/10.1029/2008WR006848>
- 1122 Gleeson, T., Befus, K. M., Jasechko, S., Luijendijk, E., & Cardenas, M. B. (2016). The global volume and distribution
1123 of modern groundwater. *Nature Geosci*, 9(2), 161–167.
- 1124 de Graaf, I. E. M., van Beek, L. P. H., Wada, Y., & Bierkens, M. F. P. (2014). Dynamic attribution of global water
1125 demand to surface water and groundwater resources: Effects of abstractions and return flows on river discharges.
1126 *Advances in Water Resources*, 64(0), 21–33. <https://doi.org/10.1016/j.advwatres.2013.12.002>



- 1127 de Graaf, I. E. M., Sutanudjaja, E. H., Van Beek, L. P. H., & Bierkens, M. F. P. (2015). A high-resolution global-scale
1128 groundwater model. *Hydrology and Earth System Sciences*, 19(2), 823–837.
- 1129 de Graaf, I. E. M., van Beek, L. P. H., Gleeson, T., Moosdorf, N., Schmitz, O., Sutanudjaja, E. H., & Bierkens, M. F. P.
1130 (2017). A global-scale two-layer transient groundwater model: Development and application to groundwater
1131 depletion. *Advances in Water Resources*, 102, 53–67. <https://doi.org/10.1016/j.advwatres.2017.01.011>
- 1132 de Graaf, I. E. M., Gleeson, T., Beek, L. P. H. (Rens) van, Sutanudjaja, E. H., & Bierkens, M. F. P. (2019).
1133 Environmental flow limits to global groundwater pumping. *Nature*, 574(7776), 90–94.
1134 <https://doi.org/10.1038/s41586-019-1594-4>
- 1135 Gnann, S. J., Woods, R. A., & Howden, N. J. (2019). Is there a baseflow Budyko curve? *Water Resources Research*,
1136 55(4), 2838–2855.
- 1137 Goderniaux, P., P. Davy, E. Bresciani, J.-R. de Dreuzy, and T. Le Borgne (2013), Partitioning a regional groundwater
1138 flow system into shallow local and deep regional flow compartments, *Water Resources Research*, 49(4), 2274-
1139 2286.
- 1140 Gosling, S. N., Zaherpour, J., Mount, N. J., Hattermann, F. F., Dankers, R., Arheimer, B., et al. (2017). A comparison
1141 of changes in river runoff from multiple global and catchment-scale hydrological models under global warming
1142 scenarios of 1 °C, 2 °C and 3 °C. *Climatic Change*, 141(3), 577–595. <https://doi.org/10.1007/s10584-016-1773-3>
- 1143 Guimberteau, M., Ducharne, A., Ciais, P., Boisier, J. P., Peng, S., De Weirtdt, M., & Verbeeck, H. (2014). Testing
1144 conceptual and physically based soil hydrology schemes against observations for the Amazon Basin, *Geosci. Model*
1145 *Dev.*, 7, 1115–1136. <https://doi.org/10.5194/gmd-7-1115-2014>
- 1146 Habets, F., Boé, J., Déqué, M., Ducharne, A., Gascoïn, S., Hachour, A., Martin, E., Pagé, C., Sauquet, E., Terray, L.,
1147 Thiéry, D., Oudin, L. & Viennot, P. (2013). Impact of climate change on surface water and ground water of two
1148 basins in Northern France: analysis of the uncertainties associated with climate and hydrological models, emission
1149 scenarios and downscaling methods. *Climatic Change*, 121, 771–785. <https://doi.org/10.1007/s10584-013-0934-x>
- 1150 Hartmann, A., Gleeson, T., Rosolem, R., Pianosi, F., Wada, Y., & Wagener, T. (2015). A large-scale simulation model
1151 to assess karstic groundwater recharge over Europe and the Mediterranean. *Geosci. Model Dev.*, 8(6), 1729–1746.
1152 <https://doi.org/10.5194/gmd-8-1729-2015>
- 1153 Hartmann, Andreas, Gleeson, T., Wada, Y., & Wagener, T. (2017). Enhanced groundwater recharge rates and
1154 altered recharge sensitivity to climate variability through subsurface heterogeneity. *Proceedings of the National*
1155 *Academy of Sciences*, 114(11), 2842–2847. <https://doi.org/10.1073/pnas.1614941114>
- 1156 Hattermann, F. F., Krysanova, V., Gosling, S. N., Dankers, R., Daggupati, P., Donnelly, C., et al. (2017). Cross-scale
1157 intercomparison of climate change impacts simulated by regional and global hydrological models in eleven large
1158 river basins. *Climatic Change*, 141(3), 561–576. <https://doi.org/10.1007/s10584-016-1829-4>
- 1159 Hay, L., Norton, P., Viger, R., Markstrom, S., Regan, R. S., & Vanderhoof, M. (2018). Modelling surface-water
1160 depression storage in a Prairie Pothole Region. *Hydrological Processes*, 32(4), 462–479.
1161 <https://doi.org/10.1002/hyp.11416>
- 1162 Henderson-Sellers, A., Z. L. Yang, and R. E. Dickinson: The Project for Intercomparison of Land-Surface Schemes
1163 (PILPS). *Bull. Amer. Meteor. Soc.*, 74, 1335–1349, 1993



- 1164 Herbert, C., & Döll, P. (2019). Global assessment of current and future groundwater stress with a focus on
1165 transboundary aquifers. *Water Resources Research*, 55, 4760–4784. <https://doi.org/10.1029/2018WR023321>
- 1166 Heudorfer, B., Haaf, E., Stahl, K., & Barthel, R. (2019). Index-based characterization and quantification of
1167 groundwater dynamics. *Water Resources Research*, 55, 5575–5592. <https://doi.org/10.1029/2018WR024418>
- 1168 Hill, M. C. (2006). The practical use of simplicity in developing ground water models. *Ground Water*, 44(6), 775–
1169 781. <https://doi.org/10.1111/j.1745-6584.2006.00227.x>
- 1170 Hill, M. C., & Tiedeman, C. R. (2007). *Effective groundwater model calibration*. Wiley.
- 1171 Hill, M. C., Kavetski, D. Clark, M. Ye, M. Arabi, M. Lu, D. Foglia, L. & Mehl, S. (2016). Practical use of computationally
1172 frugal model analysis methods. *Groundwater*. DOI:10.1111/gwat.12330
1173
- 1174 Hiscock, K. M., & Bense, V. F. (2014). *Hydrogeology—principles and practice* (2nd edition). Blackwell.
- 1175 Huang, S., Kumar, R., Flörke, M., Yang, T., Hundecha, Y., Kraft, P., et al. (2017). Evaluation of an ensemble of
1176 regional hydrological models in 12 large-scale river basins worldwide. *Climatic Change*, 141(3), 381–397.
1177 <https://doi.org/10.1007/s10584-016-1841-8>
- 1178 Hrachowitz, M., Fovet, O., Ruiz, L., Euser, T., Gharari, S., Nijzink, R., Freer, J., Savenije, H.H.G. and Gascuel-Odoux, C.
1179 (2014). Process Consistency in Models: the Importance of System Signatures, Expert Knowledge and Process
1180 Complexity. *Water Resources Research* 50:7445-7469.
- 1181 Hunt, R. J., Walker, J. F., Selbig, W. R., Westenbroek, S. M., & Regan, R. S. (2013). Simulation of climate-change
1182 effects on streamflow, lake water budgets, and stream temperature using GSFLOW and SNTMP, Trout Lake
1183 Watershed, Wisconsin. USGS Scientific Investigations Report No. 2013–5159. Reston, VA: U.S. Geological Survey.
- 1184 Hutton, C., Wagener, T., Freer, J., Han, D., Duffy, C., & Arheimer, B. (2016). Most computational hydrology is not
1185 reproducible, so is it really science? *Water Resources Research*, 52(10), 7548–7555.
1186 <https://doi.org/10.1002/2016WR019285>
- 1187 Jasechko, S., Birks, S.J., Gleeson, T., Wada, Y., Sharp, Z.D., Fawcett, P.J., McDonnell, J.J., Welker, J.M. (2014)
1188 Pronounced seasonality in the global groundwater recharge. *Water Resources Research*. 50, 8845–8867 doi:
1189 10.1002/2014WR015809
- 1190 Jasechko, S., Perrone, D., Befus, K. M., Bayani Cardenas, M., Ferguson, G., Gleeson, T., et al. (2017). Global aquifers
1191 dominated by fossil groundwaters but wells vulnerable to modern contamination. *Nature Geoscience*, 10(6), 425–
1192 429. <https://doi.org/10.1038/ngeo2943>
- 1193 Jung, M., et al. (2011). Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible
1194 heat derived from eddy covariance, satellite, and meteorological observations. *J. Geophys. Res.*, 116,
1195 G00J07, doi:10.1029/2010JG001566.
- 1196 Keune, J., Sulis, M., Kollet, S., Siebert, S., & Wada, Y. (n.d.). Human Water Use Impacts on the Strength of the
1197 Continental Sink for Atmospheric Water. *Geophysical Research Letters*, 45(9), 4068–4076.
1198 <https://doi.org/10.1029/2018GL077621>



- 1199 Keune, J., F. Gasper, K. Goergen, A. Hense, P. Shrestha, M. Sulis, and S. Kollet, 2016, Studying the influence of
1200 groundwater representations on land surface-atmosphere feedbacks during the European heat wave in 2003, J.
1201 Geophys. Res. Atmos., 121, 13, 301–13,325, doi:10.1002/2016JD025426. doi:10.1002/2016JD025426.
- 1202 Knowling, M. J., and A. D. Werner (2016), Estimability of recharge through groundwater model calibration: Insights
1203 from a field-scale steady-state example, *Journal of Hydrology*, 540, 973–987.
- 1204 Koirala et al. (2013) Global-scale land surface hydrologic modeling with the representation of water table
1205 dynamics, *JGR Atmospheres* <https://doi.org/10.1002/2013JD020398>
- 1206 Koirala, S., Kim, H., Hirabayashi, Y., Kanae, S. and Oki, T. (2019) Sensitivity of Global Hydrological Simulations to
1207 Groundwater Capillary Flux Parameterizations, *Water Resour. Res.*, 55(1), 402–425, doi:10.1029/2018WR023434,
- 1208 Kollet, S. J., & Maxwell, R. M. (2008). Capturing the influence of groundwater dynamics on land surface processes
1209 using an integrated, distributed watershed model. *Water Resources Research*, 44(2).
- 1210 Kollet, S., Sulis, M., Maxwell, R. M., Paniconi, C., Putti, M., Bertoldi, G., et al. (2017). The integrated hydrologic
1211 model intercomparison project, IH-MIP2: A second set of benchmark results to diagnose integrated hydrology and
1212 feedbacks. *Water Resources Research*, 53(1), 867–890.
- 1213 Konikow, L. F. (2011), Contribution of global groundwater depletion since 1900 to sea-level rise, *Geophys. Res.*
1214 *Lett.*, 38, L17401, doi: 10.1029/2011GL048604.
- 1215 Koster, R.D., Suarez, M.J., Ducharne, A., Praveen, K., & Stieglitz, M. (2000). A catchment-based approach to
1216 modeling land surface processes in a GCM - Part 1: Model structure. *Journal of Geophysical Research*, 105 (D20),
1217 24809–24822.
- 1218 Konikow, L.F. (2011) Contribution of global groundwater depletion since 1900 to sea-level rise. *Geophysical*
1219 *Research Letters* <https://doi.org/10.1029/2011GL048604>
- 1220 Krakauer, N. Y., Li, H., & Fan, Y. (2014). Groundwater flow across spatial scales: importance for climate modeling.
1221 *Environmental Research Letters*, 9(3), 034003.
- 1222 Kresic, N. (2009). *Groundwater resources: sustainability, management and restoration*. McGraw-Hill.
- 1223 Krueger, T., T. Page, K. Hubacek, L. Smith, and K. Hiscock (2012), The role of expert opinion in environmental
1224 modelling, *Environmental Modelling & Software*, 36, 4–18.
1225
- 1226 Kustu, M. D., Fan, Y., & Rodell, M. (2011). Possible link between irrigation in the US High Plains and increased
1227 summer streamflow in the Midwest. *Water Resources Research*, 47(3).
- 1228 Lamb, R., Aspinall, W., Odbert, H. and Wagener, T. (2017). Vulnerability of bridges to scour: Insights from an
1229 international expert elicitation workshop. *Natural Hazards and Earth System Sciences*. 17(8), 1393–1409.
- 1230 Leaf, A. T., Fienen, M. N., Hunt, R. J., & Buchwald, C. A. (2015). Groundwater/surface-water interactions in the Bad
1231 River Watershed, Wisconsin. USGS Numbered Series No. 2015–5162. Reston, VA: U.S. Geological Survey.
- 1232 Leavesley, G. H., S. L. Markstrom, P. J. Restrepo, and R. J. Viger (2002), A modular approach for addressing model
1233 design, scale, and parameter estimation issues in distributed hydrological modeling, *Hydrol. Processes*, 16, 173–
1234 187, doi:10.1002/hyp.344.



- 1235 Lemieux, J. M., Sudicky, E. A., Peltier, W. R., & Tarasov, L. (2008). Dynamics of groundwater recharge and seepage
1236 over the Canadian landscape during the Wisconsinian glaciation. *J. Geophys. Res.*, 113. Retrieved from
1237 <http://dx.doi.org/10.1029/2007JF000838>
- 1238 Lenton, T.M. et al. (2008). Tipping elements in the Earth's climate system. *Proceedings of the National Academy of*
1239 *Sciences* 105 (6), 1786-1793.
- 1240 Liang, X., Z. Xie, and M. Huang (2003). A new parameterization for surface and groundwater interactions and its
1241 impact on water budgets with the variable infiltration capacity (VIC) land surface model, *J. Geophys. Res.*, 108,
1242 8613, D16. <https://doi.org/10.1029/2002JD003090>
- 1243 Lo, M.-H., Famiglietti, J. S., Reager, J. T., Rodell, M., Swenson, S., & Wu, W.-Y. (2016). GRACE-Based Estimates of
1244 Global Groundwater Depletion. In Q. Tang & T. Oki (Eds.), *Terrestrial Water Cycle and Climate Change* (pp. 135–
1245 146). John Wiley & Sons, Inc. <https://doi.org/10.1002/9781118971772.ch7>
- 1246 Lo, M.-H., Yeh, P. J.-F., & Famiglietti, J. S. (2008). Constraining water table depth simulations in a land surface
1247 model using estimated baseflow. *Advances in Water Resources*, 31(12), 1552–1564.
- 1248 Lo, M. and J. S. Famiglietti, (2010) Effect of water table dynamics on land surface hydrologic memory, *J. Geophys.*
1249 *Res.*, 115, D22118, doi:10.1029/2010JD014191
- 1250 Lo, M.-H., J. S. Famiglietti, P. J.-F. Yeh, and T. H. Syed (2010), Improving Parameter Estimation and Water Table
1251 Depth Simulation in a Land Surface Model Using GRACE Water Storage and Estimated Baseflow Data, *Water*
1252 *Resour. Res.*, 46, W05517, doi:10.1029/2009WR007855.
- 1253 Loheide, S. P., Butler Jr, J. J., & Gorelick, S. M. (2005). Estimation of groundwater consumption by phreatophytes
1254 using diurnal water table fluctuations: A saturated-unsaturated flow assessment. *Water Resources Research*, 41(7).
- 1255 Luijendijk, E., Gleeson, T. and Moosdorf, N. (2020) Fresh groundwater discharge insignificant for the world's oceans
1256 but important for coastal ecosystems *Nature Communications*, 11, 1260 (2020). doi: 10.1038/s41467-020-15064-8
1257
- 1258 Maples, S., Foglia, L., Fogg, G.E. and Maxwell, R.M. (2020). Sensitivity of Hydrologic and Geologic Parameters on
1259 Recharge Processes in a Highly-Heterogeneous, Semi-Confined Aquifer System. *Hydrology and Earth Systems*
1260 *Sciences*, in press.
1261
- 1262 Margat, J., & Van der Gun, J. (2013). *Groundwater around the world: a geographic synopsis*. London: CRC Press
- 1263 Markovich, KH, AH Manning, LE Condon, JC McIntosh (2019). Mountain-block Recharge: A Review of Current
1264 Understanding. *Water Resources Research*, 55, <https://doi.org/10.1029/2019WR025676>
- 1265 Maxwell, R. M., Condon, L. E., and Kollet, S. J. (2015) A high-resolution simulation of groundwater and surface
1266 water over most of the continental US with the integrated hydrologic model ParFlow v3, *Geosci. Model Dev.*, 8,
1267 923–937, <https://doi.org/10.5194/gmd-8-923-2015>.
- 1268 Maxwell, R.M., Chow, F.K. and Kollet, S.J., The groundwater-land-surface-atmosphere connection: soil moisture
1269 effects on the atmospheric boundary layer in fully-coupled simulations. *Advances in Water Resources* 30(12),
1270 doi:10.1016/j.advwatres.2007.05.018, 2007.



- 1271 Maxwell, R. M., & Condon, L. E. (2016). Connections between groundwater flow and transpiration partitioning.
1272 *Science*, 353(6297), 377–380.
- 1273 Maxwell, R. M., Condon, L. E., Kollet, S. J., Maher, K., Haggerty, R., & Forrester, M. M. (2016). The imprint of
1274 climate and geology on the residence times of groundwater. *Geophysical Research Letters*, 43(2), 701–708.
1275 <https://doi.org/10.1002/2015GL066916>
- 1276 McMillan, H. (2020). Linking hydrologic signatures to hydrologic processes: A review. *Hydrological Processes*. 34,
1277 1393– 1409.
- 1278 Meixner, T., Manning, A. H., Stonestrom, D. A., Allen, D. M., Ajami, H., Blasch, K. W., et al. (2016). Implications of
1279 projected climate change for groundwater recharge in the western United States. *Journal of Hydrology*, 534, 124–
1280 138.
- 1281 Melsen, L. A., A. J. Teuling, P. J. J. F. Torfs, R. Uijlenhoet, N. Mizukami, and M. P. Clark, 2016a: HESS Opinions: The
1282 need for process-based evaluation of large-domain hyper-resolution models. *Hydrology and Earth System
1283 Sciences*, doi:10.5194/hess-20-1069-2016.
- 1284 Meriano, M., & Eyles, N. (2003). Groundwater flow through Pleistocene glacial deposits in the rapidly urbanizing
1285 Rouge River-Highland Creek watershed, City of Scarborough, southern Ontario, Canada. *Hydrogeology Journal*,
1286 11(2), 288–303. <https://doi.org/10.1007/s10040-002-0226-4>
- 1287 Milly, P.C., S.L. Malyshev, E. Shevliakova, K.A. Dunne, K.L. Findell, T. Gleeson, Z. Liang, P. Philipps, R.J. Stouffer, & S.
1288 Swenson (2014). An Enhanced Model of Land Water and Energy for Global Hydrologic and Earth-System Studies. *J.
1289 Hydrometeor.*, 15, 1739–1761. <https://doi.org/10.1175/JHM-D-13-0162.1>
- 1290 Minderhoud P S J, Erkens G, Pham Van H, Bui Tran V, Erban L E, Kooi, H and Stouthamer E (2017) Impacts of 25
1291 years of groundwater extraction on subsidence in the Mekong delta, Vietnam *Environ. Res. Lett.* 12 064006
- 1292 Minderhoud, P.S.J., Coumou, L., Erkens, G., Middelkoop, H. & Stouthamer, E. (2019). Mekong delta much lower
1293 than previously assumed in sea-level rise impact assessments. *Nature Communications* 10, 3847.
- 1294 Minderhoud, P.S.J., Middelkoop, H., Erkens, G. and Stouthamer, E. Groundwater (2020). extraction may drown
1295 mega-delta: projections of extraction-induced subsidence and elevation of the Mekong delta for the 21st century.
1296 *Environ. Res. Commun.* 2, 011005.
- 1297 Miralles, D. G., Jimenez, C., Jung, M., Michel, D., Ershadi, A., McCabe, M. F., et al. (2016). The WACMOS-ET project -
1298 Part 2: Evaluation of global terrestrial evaporation data sets. *Hydrology and Earth System Sciences*, 20(2), 823-842.
1299 doi:10.5194/hess-20-823-2016.
- 1300 Moeck, C. Nicolas Grech-Cumbo, Joel Podgorski, Anja Bretzler, Jason J. Gurdak ,Michael Berg, Mario Schirmer
1301 (2020) A global-scale dataset of direct natural groundwater recharge rates: A review of variables, processes and
1302 relationships. *Science of The Total Environment* <https://doi.org/10.1016/j.scitotenv.2020.137042>
- 1303 Mohan, C., Wei, Y., & Saft, M. (2018). Predicting groundwater recharge for varying land cover and climate
1304 conditions—a global meta-study. *Hydrology and Earth System Sciences*, 22(5), 2689–2703.



- 1305 Montanari, A., Young, G., Savenije, H.H.G., Hughes, D., Wagener, T., Ren, L.L., Koutsoyiannis, D., Cudennec, C.,
1306 Toth, E., Grimaldi, S., et al. (2013). “Panta Rhei—Everything Flows”: Change in hydrology and society—The IAHS
1307 Scientific Decade 2013–2022. *Hydrological Sciences Journal* 58, 1256–1275.
- 1308 Moore, W. S. (2010). The effect of submarine groundwater discharge on the ocean. *Annual Review of Marine*
1309 *Science*, 2, 59–88.
- 1310 Morris, M. D. (1991). Factorial sampling plans for preliminary computational experiments. *Technometrics*, 33(2),
1311 161–174.
- 1312 Müller Schmied, H., Eisner, S., Franz, D., Wattenbach, M., Portmann, F.T., Flörke, M., Döll, P. (2014):
1313 Sensitivity of simulated global-scale freshwater fluxes and storages to input data, hydrological model
1314 structure, human water use and calibration. *Hydrol. Earth Syst. Sci.*, 18, 3511–3538, doi: 10.5194/hess-
1315 18-3511-2014.
- 1316 Niu, G.-Y., Z.-L. Yang, R. E. Dickinson, and L. E. Gulden (2005), A simple TOPMODEL-based runoff parameterization
1317 (SIMTOP) for use in global climate models. *J. Geophys. Res.*, 110, D21106, doi:10.1029/2005JD006111
- 1318 Niu GY, Yang ZL, Dickinson RE, Gulden LE, Su H (2007) Development of a simple groundwater model for use in
1319 climate models and evaluation with Gravity Recovery and Climate Experiment data. *J Geophys Res* 112:D07103.
1320 doi:10.1029/2006JD007522
- 1321 Ngo-Duc, T., Laval, K. Ramillien, G., Polcher, J. & Cazenave, A. (2007). Validation of the land water storage
1322 simulated by Organising Carbon and Hydrology in Dynamic Ecosystems (ORCHIDEE) with Gravity Recovery and
1323 Climate Experiment (GRACE) data. *Water Resour. Res.*, 43, W04427. <https://doi.org/10.1029/2006WR004941>
- 1324 O’Hagan, A. (2019). Expert Knowledge Elicitation: Subjective but Scientific. *The American Statistician*, 73,
1325 doi.org/10.1080/00031305.2018.1518265
1326
- 1327 Olarinoye, T., et al. (2020): Global karst springs hydrograph dataset for research and management of the world’s
1328 fastest-flowing groundwater, *Sci. Data*, 7(1), doi:10.1038/s41597-019-0346-5.
1329
- 1330 Opie, S., Taylor, R. G., Brierley, C. M., Shamsudduha, M., & Cuthbert, M. O. (2020). Climate–groundwater dynamics
1331 inferred from GRACE and the role of hydraulic memory. *Earth System Dynamics*, 11(3), 775–791.
1332 <https://doi.org/10.5194/esd-11-775-2020>
- 1333 Ortega-Guerrero A, Rudolph D L and Cherry J A 1999 Analysis of long-term land subsidence near Mexico City: field
1334 investigations and predictive modeling *Water Resour. Res.* 353327–41
- 1335 Pan, M., Sahoo, A. K., Troy, T. J., Vinukollu, R. K., Sheffield, J., & Wood, F. E. (2012). Multisource estimation of long-
1336 term terrestrial water budget for major global river basins. *J. Climate*, 25, 3191–3206.
1337 <https://doi.org/10.1175/JCLI-D-11-00300.1>
1338
- 1339 Pappenberger, F., Ghelli, A., Buizza, R. and Bodis, K. (2009). The Skill of Probabilistic Precipitation Forecasts under
1340 Observational Uncertainties within the Generalized Likelihood Uncertainty Estimation Framework for Hydrological
1341 Applications. *Journal of Hydrometeorology*, DOI: 10.1175/2008JHM956.1



- 1342 Pellet, V., Aires, F., Munier, S., Fernández Prieto, D., Jordá, G., Dorigo, W. A., Polcher, J., & Brocca, L. (2019).
1343 Integrating multiple satellite observations into a coherent dataset to monitor the full water cycle – application to
1344 the Mediterranean region. *Hydrol. Earth Syst. Sci.*, 23, 465-491. <https://doi.org/10.5194/hess-23-465-2019>
- 1345 Perrone, D. and Jasechko (2019). Deeper well drilling an unsustainable stopgap to groundwater depletion. *Nature*
1346 *Sustain.* 2, 773-782.
- 1347 Person, M. A., Raffensperger, J. P., Ge, S., & Garven, G. (1996). Basin-scale hydrogeologic modeling. *Reviews of*
1348 *Geophysics*, 34(1), 61–87.
- 1349 Pianosi, F., Beven, K., Freer, J., Hall, J. W., Rougier, J., Stephenson, D. B., & Wagener, T. (2016). Sensitivity analysis
1350 of environmental models: A systematic review with practical workflow. *Environmental Modelling & Software*, 79,
1351 214–232.
- 1352 Post, V. E., & von Asmuth, J. R. (2013). Hydraulic head measurements—new technologies, classic pitfalls.
1353 *Hydrogeology Journal*, 21(4), 737–750.
- 1354 Qiu J. Q., Zipper, S.C., Motew M., Booth, E.G., Kucharik, C.J., & Loheide, S.P. (2019). Nonlinear groundwater
1355 influence on biophysical indicators of ecosystem services. *Nature Sustainability*, in press, doi: 10.1038/s41893-019-
1356 0278-2
- 1357
- 1358 Rajabi, M. M., and B. Ataie-Ashtiani (2016), Efficient fuzzy Bayesian inference algorithms for incorporating expert
1359 knowledge in parameter estimation, *Journal of Hydrology*, 536, 255-272.
- 1360
- 1361 Rajabi, M. M., B. Ataie-Ashtiani, and C. T. Simmons (2018), Model-data interaction in groundwater studies: Review
1362 of methods, applications and future directions, *Journal of Hydrology*, 567, 457-477.
- 1363
- 1364 Rashid, M., Chien, R.Y., Ducharne, A., Kim, H., Yeh, P.J.F., Peugeot, C., Boone, A., He, X., Séguis, L., Yabu, Y., Boukari,
1365 M. & Lo, M.H. (2019). Evaluation of groundwater simulations in Benin from the ALMIP2 project. *J. Hydromet.*,
1366 accepted.
- 1367 Refsgaard, J.C., van der Sluijs, J.P., Højberg, A.L., and Vanrolleghem, P.A. (2007). Uncertainty in the environmental
1368 modelling process—a framework and guidance. *Environmental Modelling & Software*, 22(11), 1543-1556
- 1369 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat. (2019). Deep learning
1370 and process understanding for data-driven Earth system science. *Nature*, 566(7743), 195–204.
- 1371 Reinecke, R., Foglia, L., Mehl, S., Trautmann, T., Cáceres, D., & Döll, P. (2019a). Challenges in developing a global
1372 gradient-based groundwater model. (G³M v1.0) for the integration into a global hydrological model. *Geosci. Model*
1373 *Dev.*, 12, 2401-2418. doi: 10.5194/gmd-12-2401-2019
- 1374 Reinecke, R., Foglia, L., Mehl, S., Herman, J., Wachholz, A., Trautmann, T., and Döll, P. (2019b) Spatially distributed
1375 sensitivity of simulated global groundwater heads and flows to hydraulic conductivity, groundwater recharge and
1376 surface water body parameterization, *Hydrology and Earth System Sciences*, (23) 4561–4582. 2019.
- 1377 Reinecke, R., Wachholz, A., Mehl, S., Foglia, L., Niemann, C., Döll, P. (2020). Importance of spatial resolution in
1378 global groundwater modeling. *Groundwater*. doi: 10.1111/gwat.12996



- 1379 Rodell, M., Velicogna, I., & Famiglietti, J. S. (2009). Satellite-based estimates of groundwater depletion in India.
1380 *Nature*, 460(7258), 999–1002.
- 1381 Rodell, M., Famiglietti, J. S., Wiese, D. N., Reager, J. T., Beaudoin, H. K., Landerer, F. W., & Lo, M.-H. (2018).
1382 Emerging trends in global freshwater availability. *Nature*, 557(7707), 651.
- 1383 Rosolem, R., Hoar, T., Arellano, A., Anderson, J. L., Shuttleworth, W. J., Zeng, X., and Franz, T. E.: Translating
1384 aboveground cosmic-ray neutron intensity to high-frequency soil moisture profiles at sub-kilometer scale, *Hydrol.*
1385 *Earth Syst. Sci.*, 18, 4363-4379
- 1386 Ross, J. L., M. M. Ozbek, and G. F. Pinder (2009), Aleatoric and epistemic uncertainty in groundwater flow and
1387 transport simulation, *Water Resources Research*, 45(12).
1388
- 1389 Rossmann, N., & Zlotnik, V. (2013). Review: Regional groundwater flow modeling in heavily irrigated basins of
1390 selected states in the western United States. *Hydrogeology Journal*, 21(6), 1173–1192.
1391 <https://doi.org/10.1007/s10040-013-1010-3>
- 1392 RRCA. (2003). Republican River Compact Administration Ground Water Model. Retrieved from
1393 <http://www.republicanrivercompact.org/>
- 1394 Saltelli, A., Chan, K., & Scott, E. M. (Eds.). (2000). *Sensitivity analysis*. Wiley.
- 1395 Salvucci, G. D., & Entekhabi, D. (1995). Hillslope and climatic controls on hydrologic fluxes. *Water Resources*
1396 *Research*, 31(7), 1725–1739.
- 1397 Sawyer, A. H., David, C. H., & Famiglietti, J. S. (2016). Continental patterns of submarine groundwater discharge
1398 reveal coastal vulnerabilities. *Science*, 353(6300), 705–707.
- 1399 Scanlon, B., Healy, R., & Cook, P. (2002). Choosing appropriate techniques for quantifying groundwater recharge.
1400 *Hydrogeology Journal*, 10(1), 18–39.
- 1401 Scanlon, B. R., Keese, K. E., Flint, A. L., Flint, L. E., Gaye, C. B., Edmunds, W. M., & Simmers, I. (2006). Global
1402 synthesis of groundwater recharge in semiarid and arid regions. *Hydrological Processes*, 20, 3335–3370.
- 1403 Scanlon, B. R., Faunt, C. C., Longuevergne, L., Reedy, R. C., Alley, W. M., McGuire, V. L., & McMahon, P. B. (2012).
1404 Groundwater depletion and sustainability of irrigation in the US High Plains and Central Valley. *Proceedings of the*
1405 *National Academy of Sciences*, 109(24), 9320–9325. <https://doi.org/10.1073/pnas.1200311109>
- 1406 Scanlon, B. R., Zhang, Z., Save, H., Wiese, D. N., Landerer, F. W., Long, D., et al. (2016). Global evaluation of new
1407 GRACE mascon products for hydrologic applications. *Water Resources Research*, 52(12), 9412–9429.
- 1408 Scanlon, B. R., Zhang, Z., Save, H., Sun, A. Y., Müller Schmied, H., van Beek, L. P., et al. (2018). Global models
1409 underestimate large decadal declining and rising water storage trends relative to GRACE satellite data. *Proceedings*
1410 *of the National Academy of Sciences*, 201704665.
- 1411 Schaller, M., and Y. Fan (2009) River basins as groundwater exporters and importers: Implications for water cycle
1412 and climate modeling. *Journal of Geophysical Research-Atm*, 114, D04103, doi: 10.1029/2008 JD010636



- 1413 Schewe, J., Heinke, J., Gerten, D., Haddeland, I., Arnell, N. W., Clark, D. B., et al. (2014). Multimodel assessment of
1414 water scarcity under climate change. *Proceedings of the National Academy of Sciences*, 111(9), 3245–3250.
1415 <https://doi.org/10.1073/pnas.1222460110>
- 1416 Schilling, O. S., Doherty, J., Kinzelbach, W., Wang, H., Yang, P. N., & Brunner, P. (2014). Using tree ring data as a
1417 proxy for transpiration to reduce predictive uncertainty of a model simulating groundwater–surface water–
1418 vegetation interactions. *Journal of Hydrology*, 519, Part B, 2258–2271.
1419 <https://doi.org/10.1016/j.jhydrol.2014.08.063>
- 1420 Schilling, O.S., Cook, P.G., Brunner, P., 2019. Beyond classical observations in hydrogeology: The advantages of
1421 including exchange flux, temperature, tracer concentration, residence time, and soil moisture observations in
1422 groundwater model calibration. *Reviews of Geophysics*, 57(1): 146-182.
- 1423 Schneider, A.S., Jost, A., Coulon, C., Silvestre, M., Théry, S., & Ducharne, A. (2017). Global scale river network
1424 extraction based on high-resolution topography, constrained by lithology, climate, slope, and observed drainage
1425 density. *Geophysical Research Letters*, 44, 2773–2781. <https://doi.org/10.1002/2016GL071844>
- 1426 Shen, C. (2018). A transdisciplinary review of deep learning research and its relevance for water resources
1427 scientists. *Water Resources Research*, 54(11), 8558–8593.
- 1428 Shen, C., Laloy, E., Elshorbagy, A., Albert, A., Bales, J., Chang, F.-J., et al. (2018). HESS Opinions: Incubating deep-
1429 learning-powered hydrologic science advances as a community. *Hydrology and Earth System Sciences*, 22(11).
- 1430 SKI (1984). Intracoin - International Nuclide Transport Code Intercomparison Study (No. SKI--84-3). Swedish
1431 Nuclear Power Inspectorate. Retrieved from http://inis.iaea.org/Search/search.aspx?orig_q=RN:16046803
- 1432 Springer, A., & Stevens, L. (2009). Spheres of discharge of springs. *Hydrogeology Journal*, 17(1), 83–93.
1433 <https://doi.org/10.1007/s10040-008-0341-y>
- 1434 Steffen, W., Broadgate, W., Deutsch, L., Gaffney, O., & Ludwig, C. (2015). The trajectory of the Anthropocene: the
1435 great acceleration. *The Anthropocene Review*, 2(1), 81–98.
- 1436 Sutanudjaja, E. H., Beek, R. van, Wanders, N., Wada, Y., Bosmans, J. H., Drost, N., et al. (2018). PCR-GLOBWB 2: a 5
1437 arcmin global hydrological and water resources model. *Geoscientific Model Development*, 11(6), 2429–2453.
- 1438 Takata, K., Emori, S. and Watanabe, T.: Development of the minimal advanced treatments of surface interaction
1439 and runoff, *Glob. Planet. Change*, 38(1–2), 209–222, doi:10.1016/S0921-8181(03)00030-4, 2003.
- 1440 Tallaksen, L. M. (1995). A review of baseflow recession analysis. *Journal of Hydrology*, 165(1–4), 349–370.
1441 [https://doi.org/10.1016/0022-1694\(94\)02540-R](https://doi.org/10.1016/0022-1694(94)02540-R)
- 1442 Taylor, R. G., Todd, M. C., Kongola, L., Maurice, L., Nahozya, E., Sanga, H., & MacDonald, A. M. (2013). Evidence of
1443 the dependence of groundwater resources on extreme rainfall in East Africa. *Nature Clim. Change*, 3(4), 374–378.
1444 <https://doi.org/10.1038/nclimate1731>
- 1445 Taylor, R. G., Scanlon, B., Doll, P., Rodell, M., van Beek, R., Wada, Y., et al. (2013). Groundwater and climate
1446 change. *Nature Clim. Change*, 3(4), 322–329. <https://doi.org/10.1038/nclimate1744>



- 1447 Thatch, L. M., Gilbert, J. M., & Maxwell, R. M. (2020). Integrated hydrologic modeling to untangle the impacts of
1448 water management during drought. *Groundwater*, 58(3), 377–391.
- 1449 Thomas, Z., Rousseau-Gueutin, P., Kolbe, T., Abbott, B.W., Marçais, J., Peiffer, S., Frei, S., Bishop, K., Pichelin, P.,
1450 Pinay, G., de Dreuzy, J.R. (2016). Constitution of a catchment virtual observatory for sharing flow and transport
1451 models outputs. *Journal of Hydrology*, 543, Pages 59-66. <https://doi.org/10.1016/j.jhydrol.2016.04.067>
- 1452 Tolley, D., Foglia, L., & Harter, T. (2019). Sensitivity Analysis and Calibration of an Integrated Hydrologic Model in
1453 an Irrigated Agricultural Basin with a Groundwater-Dependent Ecosystem. *Water Resources Research*.
1454 <https://doi.org/10.1029/2018WR024209>
- 1455 Tóth, J. (1963). A theoretical analysis of groundwater flow in small drainage basins. *Journal of Geophysical*
1456 *Research*, 68(16), 4795–4812.
- 1457 Tran, H., Jun Zhang, Jean-Martial Cohard, Laura E. Condon, Reed M. Maxwell (2020) Simulating groundwater-
1458 Streamflow Connections in the Upper Colorado River Basin Groundwater, 2020
1459 <https://doi.org/10.1111/gwat.13000>
- 1460 Tregoning, P., McClusky, S., van Dijk, A.I.J.M. and Crosbie, R.S. (2012). Assessment of GRACE satellites for
1461 groundwater estimation in Australia. *Waterlines Report Series No 71*, National Water Commission, Canberra
- 1462 Troldborg, L., Refsgaard, J. C., Jensen, K. H., & Engesgaard, P. (2007). The importance of alternative
1463 conceptual models for simulation of concentrations in a multi-aquifer system. *Hydrogeology Journal*,
1464 15(5), 843–860.
- 1465 Tustison, B., Harris, D. and Foufoula-Georgiou, E. (2001). Scale issues in verification of precipitation
1466 forecasts. *Journal of geophysical Research*, 106(D11), 11775-11784.
- 1467 UNESCO. (1978). *World water balance and water resources of the earth* (Vol. USSR committee for the international
1468 hydrologic decade). Paris: UNESCO.
- 1469 van Vliet, M. T., Flörke, M., Harrison, J. A., Hofstra, N., Keller, V., Ludwig, F., et al. (2019). Model inter-comparison
1470 design for large-scale water quality models. *Current Opinion in Environmental Sustainability*, 36, 59–67.
1471 <https://doi.org/10.1016/j.cosust.2018.10.013>
- 1472 Van Werkhoven, K., Wagener, T., Tang, Y., and Reed, P. 2008. Understanding watershed model behavior across
1473 hydro-climatic gradients using global sensitivity analysis. *Water Resources Research*, 44, W01429,
1474 doi:10.1029/2007WR006271.
- 1475 Van Loon, A.F. et al. (2016) [Drought in the Anthropocene](#). *Nature Geoscience* 9: 89-91 doi: 10.1038/ngeo2646.
- 1476 van Loon, Anne F.; Kumar, Rohini; Mishra, Vimal (2017): Testing the use of standardised indices and GRACE
1477 satellite data to estimate the European 2015 groundwater drought in near-real time. In *Hydrol. Earth Syst. Sci.* 21
1478 (4), pp. 1947–1971. DOI: 10.5194/hess-21-1947-2017.
- 1479 Vergnes, J.-P., & Decharme, B. (2012). A simple groundwater scheme in the TRIP river routing model: global off-line
1480 evaluation against GRACE terrestrial water storage estimates and observed river discharges. *Hydrol. Earth Syst.*
1481 *Sci.*, 16, 3889-3908. <https://doi.org/10.5194/hess-16-3889-2012>



- 1482 Vergnes, J.-P., B. Decharme, & F. Habets (2014). Introduction of groundwater capillary rises using subgrid spatial
1483 variability of topography into the ISBA land surface model, *J. Geophys. Res. Atmos.*, 119, 11,065–11,086.
1484 <https://doi.org/10.1002/2014JD021573>
- 1485 Vergnes, J.-P., Roux, N., Habets, F., Ackerer, P., Amraoui, N., Besson, F., et al. (2020). The AquifR
1486 hydrometeorological modelling platform as a tool for improving groundwater resource monitoring over France:
1487 evaluation over a 60-year period. *Hydrology and Earth System Sciences*, 24(2), 633–654.
1488 <https://doi.org/10.5194/hess-24-633-2020>
- 1489 Visser, W. C. (1959). Crop growth and availability of moisture. *Journal of the Science of Food and Agriculture*, 10(1),
1490 1–11.
- 1491 Wada, Y., L. P. H. van Beek, C. M. van Kempen, J. W. T. M. Reckman, S. Vasak, M. F. P. Bierkens, (2010) Global
1492 depletion of groundwater resources. *Geophys. Res. Lett.* 37, L20402.
- 1493 Wada, Y.; Wisser, D.; Bierkens, M. F. P. (2014). Global modeling of withdrawal, allocation and consumptive use of
1494 surface water and groundwater resources. *Earth System Dynamics Discussions*, volume 5, issue 1, pp. 15 - 40
- 1495 Wada, Y. (2016). Modeling Groundwater Depletion at Regional and Global Scales: Present State and Future
1496 Prospects. *Surveys in Geophysics*, 37(2), 419–451. <https://doi.org/10.1007/s10712-015-9347-x>
- 1497 Wada, Y., & Bierkens, M. F. P. (2014). Sustainability of global water use: past reconstruction and future projections.
1498 *Environmental Research Letters*, 9(10), 104003. <https://doi.org/10.1088/1748-9326/9/10/104003>
- 1499 Wada, Y., & Heinrich, L. (2013). Assessment of transboundary aquifers of the world—vulnerability arising from
1500 human water use. *Environmental Research Letters*, 8(2), 024003.
- 1501 Wagener, T. 2003. Evaluation of catchment models. *Hydrological Processes*, 17, 3375–3378.
- 1502 Wagener, T., & Gupta, H. V. (2005). Model identification for hydrological forecasting under uncertainty. *Stochastic
1503 Environmental Research and Risk Assessment*, 19(6), 378–387.
- 1504 Wagener, T., Sivapalan, M., Troch, P. and Woods, R. (2007). Catchment classification and hydrologic similarity.
1505 *Geography Compass*, 1(4), 901, doi:10.1111/j.1749-8198.2007.00039.x
- 1506 Wagener, T. and Pianosi, F. (2019) What has Global Sensitivity Analysis ever done for us? A systematic review to
1507 support scientific advancement and to inform policy-making in earth system modelling. *Earth-Science Reviews*,
1508 194, 1-18. doi.org/10.1016/j.earscirev.2019.04.006
- 1509 Wagener, T., Boyle, D.P., Lees, M.J., Wheeler, H.S., Gupta, H.V. and Sorooshian, S. (2001). A framework for
1510 development and application of hydrological models. *Hydrology and Earth System Sciences*, 5(1), 13-26.
- 1511 Wagener, T., Sivapalan, M., Troch, P. A., McGlynn, B. L., Harman, C. J., Gupta, H. V., et al. (2010). The future of
1512 hydrology: An evolving science for a changing world. *Water Resources Research*, 46(5).
- 1513 Wagener, T., Gleeson, T., et al. On doing large-scale hydrology with lions: perceptual models and knowledge
1514 accumulation. submitted to *Water Wires and preprint*: <https://eartharxiv.org/zdy5n/>



- 1515 Wang, F., Ducharne, A., Cheruy, F., Lo, M.H., & Grandpeix, J.L. (2018). Impact of a shallow groundwater table on
1516 the global water cycle in the IPSL land-atmosphere coupled model, *Climate Dynamics*, 50, 3505-3522,
1517 <https://doi.org/10.1007/s00382-017-3820-9>
- 1518 Warszawski, L., Frieler, K., Huber, V., Piontek, F., Serdeczny, O., & Schewe, J. (2014). The Inter-Sectoral Impact
1519 Model Intercomparison Project (ISI-MIP): Project framework. *Proceedings of the National Academy of Sciences*,
1520 111(9), 3228–3232. <https://doi.org/10.1073/pnas.1312330110>
- 1521 Winter, T. C., Harvey, J. W., Franke, O. L., & Alley, W. M. (1998). *Ground water and surface water: a single resource*
1522 (p. 79). U.S. Geological Survey circular 1139
- 1523 Woolfenden, L. R., & Nishikawa, T. (2014). Simulation of groundwater and surface-water resources of the Santa
1524 Rosa Plain watershed, Sonoma County, California. USGS Scientific Investigations Report 2014–5052). Reston, VA:
1525 U.S. Geological Survey.
- 1526 Yang, J., Griffiths, J., & Zammit, C. (2019). National classification of surface–groundwater interaction using random
1527 forest machine learning technique. *River Research and Applications*, 35(7), 932–943.
1528 <https://doi.org/10.1002/rra.3449>
- 1529 Yeh, P. J.-F. and J. Famiglietti, Regional groundwater evapotranspiration in Illinois, *J. Hydrometeorology*, 10(2),
1530 464–478, 2010
- 1531 Yilmaz, K., Gupta, H.V. and Wagener, T. 2009. Towards improved distributed modeling of watersheds: A process
1532 based diagnostic approach to model evaluation. *Water Resources Research*, 44, W09417,
1533 doi:10.1029/2007WR006716.
- 1534 Young, P., Parkinson, S. and Lees, M. (1996). Simplicity out of complexity in environmental modelling: Occam's
1535 razor revisited. *Journal of Applied Statistics*, 23(2-3), 165-210. <https://doi.org/10.1080/02664769624206>
- 1536 Zell, W. O., & Sanford, W. E. (2020). Calibrated Simulation of the Long-Term Average Surficial Groundwater System
1537 and Derived Spatial Distributions of its Characteristics for the Contiguous United States. *Water Resources*
1538 *Research*, 56(8), e2019WR026724. <https://doi.org/10.1029/2019WR026724>
- 1539 Zipper, S. C., Soylyu, M. E., Booth, E. G., & Loheide, S. P. (2015). Untangling the effects of shallow groundwater and
1540 soil texture as drivers of subfield-scale yield variability. *Water Resources Research*, 51(8), 6338–6358.
- 1541 Zipper, S. C., Soylyu, M. E., Kucharik, C. J., & Loheide, S. P. (2017). Quantifying indirect groundwater-mediated
1542 effects of urbanization on agroecosystem productivity using MODFLOW-AgroIBIS (MAGI), a complete critical zone
1543 model. *Ecological Modelling*, 359, 201-219
- 1544 Zhang, M and Burbey T J 2016 Inverse modelling using PS-InSAR data for improved land subsidence simulation in
1545 Las Vegas Valley, Nevada *Hydrol. Process.* 30 4494–516
- 1546 Zhou, Y., Li, W., 2011. A review of regional groundwater flow modeling. *Geoscience Frontiers*, 2(2): 205-214.