



1 HESS Opinions: Improving the evaluation of  
2 groundwater representation in continental to  
3 global scale models

4  
5 **Tom Gleeson**<sup>1,2</sup>, **Thorsten Wagener**<sup>3</sup>, **Petra Döll**<sup>4</sup>, **Samuel C Zipper**<sup>1,5</sup>, **Charles West**<sup>3</sup>, **Yoshihide**  
6 **Wada**<sup>6</sup>, **Richard Taylor**<sup>7</sup>, **Bridget Scanlon**<sup>8</sup>, **Rafael Rosolem**<sup>3</sup>, **Shams Rahman**<sup>3</sup>, **Nurudeen Oshinlaja**<sup>9</sup>,  
7 **Reed Maxwell**<sup>10</sup>, **Min-Hui Lo**<sup>11</sup>, **Hyungjun Kim**<sup>12</sup>, **Mary Hill**<sup>13</sup>, **Andreas Hartmann**<sup>14,3</sup>, **Graham Fogg**<sup>15</sup>,  
8 **James S. Famiglietti**<sup>16</sup>, **Agnès Ducharne**<sup>17</sup>, **Inge de Graaf**<sup>18</sup>, **Mark Cuthbert**<sup>9</sup>, **Laura Condon**<sup>19</sup>, **Etienne**  
9 **Bresciani**<sup>20</sup>, **Marc F.P. Bierkens**<sup>21,22</sup>

10 <sup>1</sup> Department of Civil Engineering, University of Victoria, Canada

11 <sup>2</sup> School of Earth and Ocean Sciences, University of Victoria

12 <sup>3</sup> Department of Civil Engineering, University of Bristol, UK & Cabot Institute, University of Bristol, UK.

13 <sup>4</sup> 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 <sup>5</sup> Kansas Geological Survey, University of Kansas

16 <sup>6</sup> International Institute for Applied Systems Analysis, Laxenburg, Austria

17 <sup>7</sup> Department of Geography, University College London, UK

18 <sup>8</sup> Bureau of Economic Geology, The University of Texas at Austin, USA

19 <sup>9</sup> School of Earth and Ocean Sciences & Water Research Institute, Cardiff University, UK

20 <sup>10</sup> Department of Geology and Geological Engineering, Colorado School of Mines, USA

21 <sup>11</sup> Department of Atmospheric Sciences, National Taiwan University, Taiwan

22 <sup>12</sup> Institute of Industrial Science, The University of Tokyo

23 <sup>13</sup> Department of Geology, University of Kansas, USA

24 <sup>14</sup> Chair of Hydrological Modeling and Water Resources, University of Freiburg, Germany

25 <sup>15</sup> Department of Land, Air and Water Resources - UC Davis

26 <sup>16</sup> School of Environment and Sustainability and Global Institute for Water Security, University of  
27 Saskatchewan, Saskatoon, Canada

28 <sup>17</sup> Sorbonne Université, CNRS, EPHE, IPSL, UMR 7619 METIS, Paris, France

29 <sup>18</sup> Chair or Environmental Hydrological Systems, University of Freiburg, Germany

30 <sup>19</sup> Department of Hydrology & Atmospheric Sciences, University of Arizona, Tucson, Arizona, USA

31 <sup>20</sup> Korea Institute of Science and Technology, Seoul, South Korea

32 <sup>21</sup> Physical Geography, Utrecht University, Utrecht, Netherlands

33 <sup>22</sup> Deltares, Utrecht, Netherlands



34

35

## Abstract

36 Continental- to global-scale hydrologic and land surface models increasingly include representations of  
37 the groundwater system, driven by crucial Earth science and sustainability problems. These models are  
38 essential for examining, communicating, and understanding the dynamic interactions between the Earth  
39 System above and below the land surface as well as the opportunities and limits of groundwater  
40 resources. A key question for this nascent and rapidly developing field is how to evaluate the realism  
41 and performance of such large-scale groundwater models given limitations in data availability and  
42 commensurability. Our objective is to provide clear recommendations for improving the evaluation of  
43 groundwater representation in continental- to global-scale models. We identify three evaluation  
44 approaches, including comparing model outputs with available observations of groundwater levels or  
45 other state or flux variables (observation-based evaluation); comparing several models with each other  
46 with or without reference to actual observations (model-based evaluation); and comparing model  
47 behavior with expert expectations of hydrologic behaviors that we expect to see in particular regions or  
48 at particular times (expert-based evaluation). Based on current and evolving practices in model  
49 evaluation as well as innovations in observations, machine learning and expert elicitation, we argue that  
50 combining observation-, model-, and expert-based model evaluation approaches may significantly  
51 improve the realism of groundwater representation in large-scale models, and thus our quantification,  
52 understanding, and prediction of crucial Earth science and sustainability problems. We encourage  
53 greater community-level communication and cooperation on these challenges, including among global  
54 hydrology and land surface modelers, local to regional hydrogeologists, and hydrologists focused on  
55 model development and evaluation.

### 56 1. WHY AND HOW IS GROUNDWATER MODELED AT CONTINENTAL TO GLOBAL SCALES?

57 Groundwater is the largest human- and ecosystem-accessible freshwater storage component of the  
58 hydrologic cycle (UNESCO, 1978; Margat & Van der Gun, 2013; Gleeson et al., 2016). Therefore, better  
59 understanding of groundwater dynamics is critical at a time when the ‘great acceleration’ (Steffen et al.,  
60 2015) of many human-induced processes is increasing stress on water resources (Wagner et al., 2010;  
61 Montanari et al., 2013; Sivapalan et al., 2014; van Loon et al., 2016), especially in regions with limited  
62 data availability and analytical capacity. Groundwater is often considered to be an inherently regional  
63 rather than global resource or system. This is partially reasonable because local to regional peculiarities  
64 of hydrology, politics and culture are paramount to groundwater resource management (Foster et al.  
65 2013) and groundwater dynamics in different continents are less directly connected and coupled than  
66 atmospheric dynamics. Regional groundwater analysis is a mature, well-established field (Hill &  
67 Tiedeman, 2007; Kresic, 2009; Zhou & Li, 2011; Hiscock & Bense, 2014; Anderson et al. 2015a) and  
68 regional approaches may be preferable for some issues and objectives; yet, important global aspects of  
69 groundwater both as a resource and as part of the Earth System are emerging (Gleeson et al. 2020).  
70 First, our increasingly globalized world trades virtual groundwater and other groundwater-dependent  
71 resources in the food-energy-water nexus, and groundwater often crosses borders in transboundary  
72 aquifers. A solely regional approach can be insufficient to analysing and managing these complex global



73 interlinkages. Second, from an Earth system perspective, groundwater is part of the hydrological cycle  
74 and connected to the atmosphere, oceans and the deeper lithosphere. A solely regional approach is  
75 insufficient to uncover and understand the complex interactions and teleconnections of groundwater  
76 within the Earth System. For example, to assess the impact of groundwater depletion on sea level rise,  
77 groundwater storage loss rate on all continents of the Earth must be aggregated. Thus, we argue that  
78 groundwater is simultaneously a local, regional, and increasingly global resource and system and that  
79 examining groundwater problems, solutions, and interactions at all scales is crucial. As a consequence,  
80 we urgently require predictive understanding about how groundwater, used by humans and connected  
81 with other components of the Earth System, operates at a variety of scales.

82

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

86 (1) To understand and quantify interactions between groundwater and past, present and future  
87 climate. Groundwater systems can have far-reaching effects on climate affecting modulation of  
88 surface energy and water partitioning with a long-term memory (Anyah et al., 2008; Maxwell and  
89 Kollet, 2008; Koirala et al. 2013; Krakauer et al., 2014; Maxwell et al., 2016; Taylor, et al., 2013;  
90 Meixner et et, 2018; Wang et al., 2018; Keune et al., 2018). While there have been significant  
91 advances in understanding the role of lateral groundwater flow on evapotranspiration (Maxwell &  
92 Condon, 2016; Bresciani et al, 2016), the interactions between climate and groundwater over  
93 longer time scales (Cuthbert et al., 2019) as well as between irrigation, groundwater, and climate  
94 (Condon and Maxwell, 2019; Condon et al 2020) remain largely unresolved. Additionally, it is well  
95 established that old groundwater with slow turnover times are common at depth (Befus et al.  
96 2017; Jasechko et al. 2017), but the relationship between groundwater and climate (and how that  
97 potentially impacts current water resources) in the Holocene and Pleistocene is also largely  
98 unresolved.

99 (2) To understand and quantify two-way interactions between groundwater, the rest of the  
100 hydrologic cycle, and the broader Earth System. Groundwater connections to the atmosphere are  
101 well documented in modeling studies (e.g. Forrester and Maxwell, 2020). Previous studies have  
102 demonstrated connections between the atmospheric boundary layer and water table depth (e.g.  
103 Maxwell et al 2007; Rahman et al, 2015), under land cover disturbance (e.g. Forrester et al 2018),  
104 under extremes (e.g. Kuene et al 2016) and due to groundwater pumping (Gilbert et al 2017).  
105 While a number of open source platforms have been developed to study these connections (e.g.  
106 Maxwell et al 2011; Shrestha et al 2014; Sulis, 2017) these platforms are regional to continental in  
107 extent. Recent work has shown global impacts of groundwater on atmospheric circulation (Wang  
108 et al 2018), groundwater is still quite simplified. As the main storage component of the  
109 freshwater hydrologic cycle, groundwater systems support baseflow levels in streams and rivers,  
110 and thereby ecosystems and agricultural productivity and other ecosystem services in both  
111 irrigated and rainfed systems (Scanlon et al., 2012; Qiu et al., 2019; Visser, 1959; Zipper et al.,  
112 2015, 2017). When pumped groundwater is transferred to oceans (Konikow 2011; Wada et al.,  
113 2012; Döll et al., 2014a; Wada, 2016; Caceres et al., 2020; Luijendijk et al. 2020), resulting sea-  
114 level rise can impact salinity levels in coastal aquifers, and freshwater and solute inputs to the



- 115 ocean (Moore, 2010; Sawyer et al., 2016). Difficulties are complicated by international trade of  
116 virtual groundwater which causes aquifer stress in disparate regions (Dalin et al., 2017)
- 117 (3) To inform water decisions and policy for large, often transboundary groundwater systems in an  
118 increasingly globalized world (Wada & Heinrich, 2013; Herbert & Döll, 2019). For instance,  
119 groundwater recharge from large-scale models has been used to quantify groundwater resources  
120 in Africa, even though large-scale models do not yet include all recharge processes that are  
121 important in this region (Taylor et al., 2013; Jasechko et al. 2014; Cuthbert et al., 2019; Hartmann  
122 et al., 2017).
- 123 (4) To create visualizations and interactive opportunities that inform citizens and consumers, whose  
124 decisions have global-scale impacts, about the state of groundwater all around the world such as  
125 the World Resources Institute's Aqueduct website (<https://www.wri.org/aqueduct>), a decision-  
126 support tool to identify and evaluate global water risks.

127 The first two purposes are science-focused while the latter two are sustainability-focused. In sum,  
128 continental- to global-scale hydrologic models incorporating groundwater offer a coherent scientific  
129 framework to examine the dynamic interactions between the Earth System above and below the land  
130 surface, and are compelling tools for conveying the opportunities and limits of groundwater resources  
131 to people so that they can better manage the regions they live in, and better understand the world  
132 around them.

133

134 As a result, many global hydrological models and land surface models have incorporated groundwater to  
135 varying levels of complexity depending on the model provenance and purpose. Different from regional-  
136 scale groundwater models that generally focus on subsurface dynamics, the focus of these models is on  
137 estimating either runoff and streamflow (hydrological models) or land-atmosphere water and energy  
138 exchange (land surface models). Simulation of groundwater storages and hydraulic heads mainly serve  
139 to quantify baseflow that affects streamflow during low flow periods or capillary rise that increases  
140 evapotranspiration. Some land-surface models use approaches based on the topographic index to  
141 simulate fast surface and slow subsurface runoff based on the fraction of saturated area in the grid cell  
142 (Clark et al., 2015; Fan et al., 2019); groundwater in these models does not have water storage nor by  
143 hydraulic heads (Famiglietti & Wood, 1994; Koster et al., 2000; Niu et al., 2003; Takata et al., 2003). In  
144 some global hydrological models, groundwater is still represented as a linear reservoir that is fed by  
145 groundwater recharge and drains to a river in each grid cell (Müller Schmied et al., 2014; Gascoin et al.,  
146 2009; Ngo-Duc et al., 2007). Time series of groundwater storage but not hydraulic heads are computed.  
147 This prevents simulation of lateral groundwater flow between grid cells, capillary rise and two-way  
148 exchange flows between surface water bodies and groundwater (Döll et al., 2016). However,  
149 representing groundwater as a water storage compartment that is connected to soil and surface water  
150 bodies by groundwater recharge and baseflow and is affected by groundwater abstractions and returns  
151 enables global-scale assessment of groundwater resources and stress (Herbert and Döll, 2019) and  
152 groundwater depletion (Döll et al., 2014a; Wada et al., 2014; de Graaf et al., 2014). In some land surface  
153 models, the location of the groundwater table with respect to the land surface is simulated within each  
154 grid cell to enable simulation of capillary rise (Niu et al., 2007) but, as in the case of simulating  
155 groundwater as a linear reservoir, lateral groundwater transport or two-way surface water-groundwater  
156 exchange cannot be simulated with this approach.



157  
158 Increasingly, models for simulating groundwater flows between all model grid cells in entire countries or  
159 globally have been developed, either as stand-alone models or as part of hydrological models (Vergnes  
160 & Decharme, 2012; Fan et al., 2013; Lemieux et al. 2008; de Graaf et al., 2017; Kollet et al., 2017;  
161 Maxwell et al., 2015; Reinecke et al., 2018, de Graaf et al 2019). The simulation of groundwater in large-  
162 scale models is a nascent and rapidly developing field with significant computational and  
163 parameterization challenges which has led to significant and important efforts to develop and evaluate  
164 individual models. It is important to note that herein ‘large-scale models’ refer to models that are  
165 laterally extensive across multiple regions (hundreds to thousands of kilometers) and generally include  
166 the upper tens to hundreds of meters of subsurface and have resolutions sometimes as small as ~1 km.  
167 In contrast, ‘regional-scale’ models (tens to hundreds of kilometers) have long been developed for a  
168 specific region or aquifer and can include greater depths and resolutions, more complex  
169 hydrostratigraphy and are often developed from conceptual models with significant regional knowledge.  
170 Regional-scale models include a diverse range of approaches from stand-alone groundwater models  
171 (i.e., representing surface water and vadose zone processes using boundary conditions such as recharge)  
172 to fully integrated groundwater-surface water models. We consider both large-scale and regional-scale  
173 models to be useful practices that should both continue to be conducted rather than one replacing  
174 another; ideally both should benefit from the other since each has strengths and weaknesses and  
175 together the two practices enrich our understanding and support the management of groundwater  
176 across scales.

177  
178 Now that a number of models that represent groundwater at continental to global scales have been  
179 developed and are being developed, it is equally important that we advance how we evaluate these  
180 models. To date, large-scale model evaluation has largely focused on individual models and lacked the  
181 rigor of regional-scale model evaluation, with inconsistent practices between models and little  
182 community-level discussion or cooperation. Our objective is to provide clear recommendations for  
183 evaluating groundwater representation in large-scale (continental and global) models. We focus on  
184 model evaluation because this is the heart of model trust and reproducibility (Hutton et al., 2016). We  
185 describe current model evaluation practices (Section 2) and consider diverse and uncertain sources of  
186 information, including observations, models and experts to holistically evaluate the simulation of  
187 groundwater-related fluxes, stores and hydraulic heads (Section 3). We stress the need for an iterative  
188 and open-ended process of model improvement through continuous model evaluation against the  
189 different sources of information. We explicitly contrast the terminology used herein of ‘evaluation’ and  
190 ‘comparison’ against terminology such as ‘calibration’ or ‘validation’ or ‘benchmarking’, which suggests  
191 a modelling process that is at some point complete. We extend previous commentaries advocating  
192 improved hydrologic process representation and evaluation in large-scale hydrologic models (Clark et al.  
193 2015; Melsen et al. 2016) by adding expert-elicitation and machine learning for more holistic evaluation.  
194 We also consider model objective and model evaluation across the diverse hydrologic landscapes which  
195 can both uncover blindspots in model development.

196  
197 We bring together somewhat disparate scientific communities as a step towards greater community-  
198 level cooperation on these challenges, including global hydrology and land surface modelers, local to



199 regional hydrogeologists, and hydrologists focused on model development and evaluation. We see three  
200 audiences beyond those currently directly involved in large-scale groundwater modeling that we seek to  
201 engage to accelerate model evaluation: 1) regional hydrogeologists who could be reticent about global  
202 models, and yet have crucial knowledge and data that would improve evaluation; 2) data scientists with  
203 expertise in machine learning, artificial intelligence etc. whose methods could be useful in a myriad of  
204 ways; and 3) the multiple Earth Science communities that are currently working towards integrating  
205 groundwater into a diverse range of models so that improved evaluation approaches are built directly  
206 into model development.

## 207 **2. CURRENT MODEL EVALUATION PRACTICES**

208 Here we provide a brief overview of evaluation of both large-scale hydrological models as well as  
209 regional-scale models to highlight some of the evaluation differences and opportunities at different  
210 scales. It is important to consider how or if large-scale models are fundamentally different to regional-  
211 scale models, especially in ways that could impact evaluation. As defined above, large-scale models  
212 cover larger areas, often including data-poor areas and generally at coarser resolution compared to  
213 regional-scale models. These differences impact evaluations in at least five relevant ways:

- 214 a) commensurability errors (also called ‘representativeness’ errors) occur either when modelled grid  
215 values are interpolated and compared to an observation ‘point’ or when aggregation of observed  
216 ‘point’ values are compared to a modelled grid value (Beven, 2005; Tustison et al., 2001; Beven,  
217 2016; Pappenberger et al., 2009; Rajabi et al., 2018). For groundwater models in particular,  
218 commensurability error will depend on the number and locations of observation points, the  
219 variability structure of the variables being compared such as hydraulic head and the interpolation or  
220 aggregation scheme applied (Tustison et al., 2001; Pappenberger et al., 2009; Reinecke et al., 2020).  
221 Commensurability is a problem for most scales of modelling, but likely more significant the coarser  
222 the model;
- 223 b) specificity to region and objective because regional-scale models are developed specifically for a  
224 certain region and modeling or management objective whereas large-scale models are often more  
225 general and include different regions leading to greater heterogeneity of processes and parameters;
- 226 c) large-scale models have immense computational requirements which leads to challenges with  
227 uncertainty and sensitivity analysis, while it is important to note that some regional-scale models  
228 also face computational demands;
- 229 d) including data-poor areas in large-scale models leads to challenges when only using observations  
230 for model evaluation; and
- 231 e) regional-scale models routinely include heterogeneous and anisotropic parameterizations which  
232 could be improved in future large-scale models. For example, intense vertical anisotropy routinely  
233 induces vertical flow dynamics from vertical head gradients that are tens to thousands of times  
234 greater than horizontal gradients which profoundly alter the meaning of the deep and shallow  
235 groundwater levels, with only the latter remotely resembling the actual water table.

236 Despite differences between model evaluation at different scales, we suggest that well-established  
237 modeling strategies at regional scales, that we describe more below, can be adapted and built upon to  
238 improve large-scale model evaluation.



239  
240 Evaluation of large-scale models has often focused on streamflow or evapotranspiration observations  
241 but joint evaluation together with groundwater-specific variables is appropriate and necessary (e.g.  
242 Maxwell et al. 2015; Maxwell and Condon, 2016). Groundwater-specific variables useful for evaluating  
243 the groundwater component of large-scale models include a) hydraulic head or water table depth; b)  
244 groundwater storage and groundwater storage changes which refer to long-term, negative or positive  
245 trends in groundwater storage where long-term, negative trends are called groundwater depletion; c)  
246 groundwater recharge; d) flows between groundwater and surface water bodies; and e) human  
247 groundwater abstractions and return flows to groundwater. It is important to note that groundwater  
248 and surface water hydrology communities often have slightly different definitions of terms like recharge  
249 and baseflow (Barthel, 2014); we therefore suggest trying to precisely define the meanings of such  
250 words using the actual hydrologic fluxes which we do below. Table 1 shows the availability of  
251 observational data for these variables but does not evaluate the quality and robustness of observations.  
252 Overall there are significant inherent challenges of commensurability and measurability of groundwater  
253 observations in the evaluation of large-scale models. We describe the current model evaluation  
254 practices for each of these variables here:

- 255
- 256 a) Simulated hydraulic heads or water table depth in large scale models are frequently compared  
257 to well observations, which are often considered the crucial data for groundwater model  
258 evaluation. Hydraulic head observations from a large number groundwater wells (>1 million)  
259 have been used to evaluate the spatial distribution of steady-state heads (Fan et al., 2013, de  
260 Graaf et al., 2015; Maxwell et al., 2015; Reinecke et al., 2019a, 2020). Transient hydraulic heads  
261 with seasonal amplitudes (de Graaf et al. 2017), declining heads in aquifers with groundwater  
262 depletion (de Graaf et al. 2019) and daily transient heads (Tran et al 2020) have also been  
263 compared to well observations. All evaluation with well observations is severely hampered by  
264 the incommensurability of point values of observed head with simulated heads that represent  
265 averages over cells of a size of tens to hundreds square kilometers; within such a large cell, land  
266 surface elevation, which strongly governs hydraulic head, may vary a few hundred meters, and  
267 average observed head strongly depends on the number and location of well within the cell  
268 (Reinecke et al., 2020). Additional concerns with head observations are the 1) strong sampling  
269 bias of wells towards accessible locations, low elevations, shallow water tables, and more  
270 transmissive aquifers in wealthy, generally temperate countries (Fan et al., 2019); 2) the impacts  
271 of pumping which may or may not be well known; 3) observational errors and uncertainty (Post  
272 and von Asmuth, 2013; Fan et al., 2019); and 4) that heads can reflect the poro-elastic effects of  
273 mass loading and unloading rather than necessarily aquifer recharge and drainage (Burgess et al,  
274 2017). To date, simulated hydraulic heads have more often been compared to observed heads  
275 (rather than water table depth) which results in lower relative errors (Reinecke et al., 2020)  
276 because the range of heads (10s to 1000s m head) is much larger than the range of water table  
277 depths (<1 m to 100s m).
- 278
- 279 b) Simulated groundwater storage trends or anomalies in large-scale hydrological models have  
280 been evaluated using observations of groundwater well levels combined with estimates of



281 storage parameters, such as specific yield; local-scale groundwater modeling; and translation of  
282 regional total water storage trends and anomalies from satellite gravimetry (GRACE: Gravity  
283 Recovery And Climate Experiment) to groundwater storage changes by estimating changes in  
284 other hydrological storages (Döll et al., 2012; 2014a). Groundwater storage changes volumes  
285 and rates have been calculated for numerous aquifers, primarily in the United States, using  
286 calibrated groundwater models, analytical approaches, or volumetric budget analyses (Konikow,  
287 2010). Regional-scale models have also been used to simulate groundwater storage trends  
288 untangling the impacts of water management during drought (Thatch et al. 2020). Satellite  
289 gravimetry (GRACE) is important but has limitations (Alley and Konikow, 2015). First, monthly  
290 time series of very coarse-resolution groundwater storage are indirectly estimated from  
291 observations of total water storage anomalies by satellite gravimetry (GRACE) but only after  
292 model- or observation-based subtraction of water storage changes in glaciers, snow, soil and  
293 surface water bodies (Lo et al., 2016; Rodell et al., 2009; Wada, 2016). As soil moisture, river or  
294 snow dynamics often dominate total water storage dynamics, the derived groundwater storage  
295 dynamics can be so uncertain that severe groundwater drought cannot be detected in this way  
296 (Van Loon et al., 2017). Second, GRACE cannot detect the impact of groundwater abstractions  
297 on groundwater storage unless groundwater depletion occurs (Döll et al., 2014a,b). Third, the  
298 very coarse resolution can lead to incommensurability but in the opposite direction of well  
299 observations. It is important to note that the focus is on storage trends or anomalies since total  
300 groundwater storage to a specific depth (Gleeson et al., 2016) or in an aquifer (Konikow, 2010)  
301 can be estimated but the total groundwater storage in a specific region or cell cannot be  
302 simulated or observed unless the depth of interest is specified (Condon et al., 2020).  
303

304 c) Simulated large-scale groundwater recharge (vertical flux across the water table) has been  
305 evaluated using compilations of point estimates of groundwater recharge, results of regional-  
306 scale models, baseflow indices, and expert opinion (Döll and Fiedler, 2008; Hartmann et al.,  
307 2015) or compared between models (e.g. Wada et al. 2010). In general, groundwater recharge is  
308 not directly measurable except by meter-scale lysimeters (Scanlon et al., 2002), and many  
309 groundwater recharge methods such as water table fluctuations and chloride mass balance also  
310 suffer from similar commensurability issues as water table depth data. Although sometimes an  
311 input or boundary condition to regional-scale models, recharge in many large-scale groundwater  
312 models is simulated and thus can be evaluated.  
313

314 d) The flows between groundwater and surface water bodies (rivers, lakes, wetlands) are  
315 simulated by many models but are generally not evaluated directly against observations of such  
316 flows since they are very rare and challenging. Baseflow (the slowly varying portion of  
317 streamflow originating from groundwater or other delayed sources) or streamflow ‘low flows’  
318 (when groundwater or other delayed sources predominate), generally cannot be used to directly  
319 quantify the flows between groundwater and surface water bodies at large scales. Groundwater  
320 discharge to rivers can be estimated from streamflow observations only in the very dense gauge  
321 network and/or if streamflow during low flow periods is mainly caused by groundwater  
322 discharge and not by water storage in upstream lakes, reservoirs or wetlands. These conditions



323 are rarely met in case of streamflow gauges with large upstream areas that can be used for  
324 comparison to large-scale model output. de Graaf et al. (2019) compared the simulated timing  
325 of changes in groundwater discharge to observations and regional-scale models, but only  
326 compared the fluxes directly between the global- and regional-scale models. Due to the  
327 challenges of directly observing the flows between groundwater and surface water bodies at  
328 large scales, this is not included in the available data in Table 1; instead in Section 3 we highlight  
329 the potential for using baseflow or the spatial distribution of perennial, intermittent and  
330 ephemeral streams in the future.

331  
332 e) Groundwater abstractions have been evaluated by comparison to national, state and county  
333 scale statistics in the U.S. (Wada et al. 2010, Döll et al., 2012, 2014a, de Graaf et al. 2014).  
334 Irrigation is the dominant groundwater use sector in many regions; however, irrigation pumpage  
335 is generally estimated from crop water demand and rarely metered although GRACE and other  
336 remote sensing data have been used to estimate the irrigation water demand (Anderson et al.  
337 2015b). Groundwater abstraction uncertainties introduce significant uncertainties into large-  
338 scale models and is simulated and thus can be evaluated. Human groundwater abstractions and  
339 return flows as well as groundwater recharge and the flows between groundwater and surface  
340 water bodies are necessary to simulate storage trends (described above). But each of these are  
341 considered separate observations since they each have different data sources and assumptions.  
342 Groundwater abstraction data at the well scale are severely hampered by the  
343 incommensurability like hydraulic head and recharge described above.

344  
345 Regional-scale groundwater models typically have fewer (though not insignificant) commensurability  
346 issues due to smaller grid cell sizes compared to global-scale models, and may have different challenges  
347 related to data availability, such as the lack of reliable hydrologic monitoring data in many regions.  
348 Regional-scale models are evaluated using a variety of data types, some of which are available and  
349 already used at the global scale and some of which are not. In general, the most common data types  
350 used for regional-scale groundwater model evaluation match global-scale groundwater models:  
351 hydraulic head and either total streamflow or baseflow estimated using hydrograph separation  
352 approaches (eg. RRCA, 2003; Woolfenden and Nishikawa, 2014; Tolley et al., 2019). However, numerous  
353 data sources unavailable or not currently used at the global scale have also been applied in regional-  
354 scale models, such as elevation of surface water features (Hay et al., 2018), existing maps of the  
355 potentiometric surface (Meriano and Eyles, 2003), and dendrochronology (Schilling et al., 2014) - these  
356 and other 'non-classical' observations (Schilling et al. 2019) could be inspiration for model evaluation of  
357 large-scale models in the future but are beyond our scope to discuss. Further, given the smaller domain  
358 size of regional-scale models, expert knowledge and local ancillary data sources can be more directly  
359 integrated and automated parameter estimation approaches such as PEST are tractable (Leaf et al.,  
360 2015; Hunt et al., 2013). We directly build upon this practice of integration of expert knowledge below  
361 in Section 3.3.



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

363

364 Based on Section 2, we argue that the current model evaluation practices are insufficient to robustly  
365 evaluate large-scale models. We therefore propose evaluating large-scale models using at least three  
366 strategies (pillars of Figure 1): observation-, model-, and expert-driven evaluation which are potentially  
367 mutually beneficial because each strategy has its strengths and weaknesses. Across all three model  
368 evaluation strategies, we advocate three principles underpinning model evaluation (base of Figure 1),  
369 none of which we are the first to suggest but we highlight here as a reminder: 1) model objectives, such  
370 as the groundwater science or groundwater sustainability objective summarised in Section 1, are  
371 important to model evaluation because they provide the context through which relevance of the  
372 evaluation outcome is set; 2) all sources of information (observations, models and experts) are  
373 uncertain and this uncertainty needs to be quantified for robust evaluation; and 3) regional differences  
374 are likely important for large-scale model evaluation - understanding these differences is crucial for the  
375 transferability of evaluation outcomes to other places or times. For example, in assessing climate change  
376 impacts on groundwater the objective is relatively clear, uncertainty is an integral part of the evaluation,  
377 and regional differences are common.

378

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

394

395 We advocate for considering regional differences more explicitly in model evaluation since likely no  
396 single model will perform consistently across the diverse hydrologic landscapes of the world (Van  
397 Werkhoven et al., 2008). Considering regional differences in large-scale model evaluation is motivated  
398 by recent model evaluation results and is already starting to be practiced. Two recent sensitivity  
399 analyses of large-scale models reveal how sensitivities to input parameters vary in different regions for  
400 both hydraulic heads and flows between groundwater and surface water (de Graaf et al. 2019; Reinecke  
401 et al., 2020). In mountain regions, large-scale models tend to underestimate steady-state hydraulic  
402 head, possibly due to over-estimated hydraulic conductivity in these regions, which highlights that



403 model performance varies in different hydrologic landscapes. (de Graaf et al., 2015; Reinecke et al.  
404 2019b). Additionally, there are significant regional differences in performance with low flows for a  
405 number of large-scale models (Zaherpour et al. 2018) likely because of diverse implementations of  
406 groundwater and baseflow schemes. Large-scale model evaluation practice is starting to shift towards  
407 highlighting regional differences as exemplified by two different studies that explicitly mapped  
408 hydrologic landscapes to enable clearer understanding of regional differences. Reinecke et al. (2019b)  
409 identified global hydrological response units which highlighted the spatially distributed parameter  
410 sensitivities in a computationally expensive model, whereas Hartmann et al. (2017) developed and  
411 evaluated models for karst aquifers in different hydrologic landscapes based on different a priori system  
412 conceptualizations. Considering regional differences in model evaluation suggests that global models  
413 could in the future consider a patchwork approach of different conceptual models, governing equations,  
414 boundary conditions etc. in different regions. Although beyond the scope of this manuscript, we  
415 consider this an important future research avenue.  
416

### 417 **3.1 Observation-based model evaluation**

418 Observation-based model evaluation is the focus of most current efforts and is important because we  
419 want models to be consistent with real-world observations. Section 2 and Table 1 highlight both the  
420 strengths and limitations of current practices using observations. Despite existing challenges, we foresee  
421 significant opportunities for observation-based model evaluation and do not see data scarcity as a  
422 reason to exclude groundwater in large-scale models or to avoid evaluating these models. It is important  
423 to note that most so-called ‘observations’ are modeled or derived quantities, and often at the wrong  
424 scale for evaluating large-scale models (Table 1; Beven, 2019). Given the inherent challenges of direct  
425 measurement of groundwater fluxes and stores especially at large scales, herein we consider the word  
426 ‘observation’ loosely as any measurements of physical stores or fluxes that are combined with or filtered  
427 through models for an output. For example, GRACE gravity measurements are combined with model-  
428 based estimates of water storage changes in glaciers, snow, soil and surface water for ‘groundwater  
429 storage change observations’ or streamflow measurements are filtered through baseflow separation  
430 algorithms for ‘baseflow observations’. The strengths and limitations as well as the data availability and  
431 spatial and temporal attributes of different observations are summarized in Table 1 which we hope will  
432 spur more systematic and comprehensive use of observations.  
433

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

- 437 1) Focus on transient observations of the water table depth rather than hydraulic head  
438 observations that are long-term averages or individual times (often following well  
439 drilling). Water table depth are likely more robust evaluation metrics than hydraulic  
440 head because water table depth reveals great discrepancies and is a complex function of  
441 the relationship between hydraulic head and topography that is crucial to predicting  
442 system fluxes (including evapotranspiration and baseflow). Comparing transient



443 observations and simulations instead of long-term averages or individual times  
444 incorporates more system dynamics of storage and boundary conditions as temporal  
445 patterns are more important than absolute values (Heudorfer et al. 2019). For regions  
446 with significant groundwater depletion, comparing to declining water tables is a useful  
447 strategy (de Graaf et al. 2019), whereas in aquifers without groundwater depletion,  
448 seasonally varying water table depths are likely more useful observations (de Graaf et  
449 al. 2017).

450 2) Use baseflow, the slowly varying portion of streamflow originating from groundwater or  
451 other delayed sources. Döll and Fiedler (2008) included the baseflow index in evaluating  
452 recharge and baseflow has been used to calibrate the groundwater component of a land  
453 surface model (Lo et al. 2008, 2010). But the baseflow index (BFI), baseflow recession (k)  
454 or baseflow fraction (Gnann et al., 2019) have not been used to evaluate any large-scale  
455 model that simulates groundwater flows between all model grid cells. There are  
456 limitations of using BFI and baseflow recession (k) to evaluate large-scale models (Table  
457 1) and this only makes sense when the baseflow separation algorithm is better than the  
458 large-scale model itself, which may not be the case for some large-scale models. But this  
459 remains available and obvious data derived from streamflow observations that has been  
460 under-used to date.

461 3) Use the spatial distribution of perennial, intermittent, and ephemeral streams as an  
462 observation, which to our best knowledge has not been done by any large-scale model  
463 evaluation. The transition between perennial and ephemeral streams is an important  
464 system characteristic in groundwater-surface water interactions (Winter et al. 1998), so  
465 we suggest that this might be a revealing evaluation criteria although there are similar  
466 limitations to using baseflow. The results of both quantifying baseflow and mapping  
467 perennial streams depend on the methods applied, they are not useful for quantifying  
468 groundwater-surface water interactions when there is upstream surface water storage,  
469 and they do not directly provide information about fluxes between groundwater and  
470 surface water.

471 4) Use data on land subsidence to infer head declines or aquifer properties for regions  
472 where groundwater depletion is the main cause of compaction (Bierkens and Wada,  
473 2019). Lately, remote sensing methods such as GPS, airborne and space borne radar and  
474 lidar are frequently used to infer land subsidence rates (Erban et al., 2014). Also, a  
475 number of studies combine geomechanical modelling (Ortega-Guerrero et al 1999;  
476 Minderhoud et al 2017) and geodetic data to explain the main drivers of land  
477 subsidence. A few papers (e.g. Zhang and Burbey 2016) use a geomechanical model  
478 together with a withdrawal data and geodetic observations to estimate hydraulic and  
479 geomechanical subsoil properties.

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



- 485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526
- 6) Derive additional new datasets using meta-analysis and/or geospatial analysis such as gaining or losing stream reaches (e.g., from interpolated head measurements close to the streams), springs and groundwater-dependent surface water bodies, or tracers. Each of these new data sources could in principle be developed from available data using methods already applied at regional scales but do not currently have an ‘off the shelf’ global dataset. For example, some large-scale models have been explicitly compared with residence time and tracer data (Maxwell et al., 2016) which have also been recently compiled globally (Gleeson et al., 2016; Jasechko et al., 2017). This could be an important evaluation tool for large-scale models that are capable of simulating flow paths, or can be modified to do so. Future meta-analyses data compilations should report on the quality of the data and include possible uncertainty ranges as well as the mean estimates.
  - 7) Use machine learning to identify spatiotemporal patterns, for example of perennial streams, depth table depths or baseflow fluxes, which might not be obvious in multi-dimensional datasets and could be useful in evaluation. For example, Yang et al. (2019) predicted the state of losing and gaining streams in New Zealand using random forests. A staggering variety of machine learning tools are available and their use is nascent yet rapidly expanding in geoscience and hydrology (Reichstein et al., 2019; Shen, 2018; Shen et al., 2018; Wagener et al., 2020). While large-scale groundwater models are often considered ‘data-poor’, it may seem strange to propose using data-intensive machine learning methods to improve model evaluation. But some of the data sources are large (e.g over 2 million water level measurements in Fan et al. 2013 although biased in distribution) whereas other observations such as evapotranspiration (Jung et al., 2011) and baseflow (Beck et al. 2013) are already interpolated and extrapolated using machine learning.
  - 8) Consider comparing models against hydrologic signatures - indices that provide insight into the functional behavior of the system under study (Wagener et al., 2007; McMillan, 2020). The direct comparison of simulated and observed variables through statistical error metrics has at least two downsides. One, the above mentioned unresolved problem of commensurability, and two, the issue that such error metrics are rather uninformative in a diagnostic sense - simply knowing the size of an error does not tell the modeller how the model needs to be improved, only that it does (Yilmaz et al., 2009). One way to overcome these issues, is to derive hydrologically meaningful signatures from the original data, such as the signatures derived from transient groundwater levels by Heudorfer et al. (2019). For example, recharge ratio (defined as the ratio of groundwater recharge to precipitation) might be hydrologically more informative than recharge alone (Jasechko et al., 2014) or the water table ratio and groundwater response time (Cuthbert et al. 2019) which are spatially-distributed signatures of groundwater systems dynamics. Such signatures might be used to assess model consistency (Wagener & Gupta, 2005; Hrachowitz et al. 2014) by looking at the similarity of patterns or spatial trends rather than the size of the aggregated error, thus reducing the commensurability problem.



527 9) Understand and quantify commensurability error issues better so that a fairer  
528 comparison can be made across scales using existing data. As described above,  
529 commensurability errors will depend on the number and locations of observation  
530 points, the variability structure of the variables being compared such as hydraulic head  
531 and the interpolation or aggregation scheme applied. While to some extent we may  
532 appreciate how each of these factors affect commensurability error in theory, in  
533 practice their combined effects are poorly understood and methods to quantify and  
534 reduce commensurability errors for groundwater model purposes remain largely  
535 undeveloped. As such, quantification of commensurability error in (large-scale)  
536 groundwater studies is regularly overlooked as a source of uncertainty because it cannot  
537 be satisfactorily evaluated (Tregoning et al., 2012). Currently, evaluation of simulated  
538 groundwater heads is plagued by, as yet, poorly quantified uncertainties stemming from  
539 commensurability errors and we therefore recommend future studies focus on  
540 developing solutions to this problem.

541 We recommend evaluating models with a broader range of currently available data sources (with  
542 explicit consideration of data uncertainty and regional differences) while also simultaneously working to  
543 derive new data sets. However, data distribution and commensurability issues will likely still be present,  
544 which underscores the importance of the two following strategies.

### 545 3.2. Model-based model evaluation

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

550  
551 The original MIP concept offers a framework to consistently evaluate and compare models, and  
552 associated model input, structural, and parameter uncertainty under different objectives (e.g., climate  
553 change, model performance, human impacts and developments). Since the Project for the  
554 Intercomparison of Land-Surface Parameterization Schemes (PILPS; Sellers et al., 1993), the first MIP,  
555 the land surface modeling community has used MIPs to deepen understanding of land physical  
556 processes and to improve their numerical implementations at various scales from regional (e.g., Rhône-  
557 aggregation project; Boone et al., 2004) to global (e.g., Global Soil Wetness Project; Dirmeyer, 2011).  
558 Two examples of recent model intercomparison efforts illustrate the general MIP objectives and  
559 practice. First, ISIMIP (Schewe et al., 2014; Warszawski et al., 2014) assessed water scarcity at different  
560 levels of global warming. Second, IH-MIP2 (Kollet et al., 2017) used both synthetic domains and an  
561 actual watershed to assess fully-integrated hydrologic models because these cannot be validated easily  
562 by comparison with analytical solutions and uncertainty remains in the attribution of hydrologic  
563 responses to model structural errors. Model comparisons have revealed differences, but it is often  
564 unclear whether these stem from differences in the model structures, differences in how the  
565 parameters were estimated, or from other modelling choices (Duan et al., 2006). Attempts for modular  
566 modelling frameworks to enable comparisons (Wagener et al., 2001; Leavesley et al., 2002; Clark et al.,



567 2008; Fenicia et al., 2011; Clark et al., 2015) or at least shared explicit modelling protocols and boundary  
568 conditions (Refsgaard et al., 2007; Ceola et al., 2015; Warszawski et al., 2014) have been proposed to  
569 reduce these problems.

570  
571 Inter-scale model comparison - for example, comparing a global model to a regional-scale model - is a  
572 potentially useful approach which is emerging for surface hydrology models (Hattermann et al., 2017;  
573 Huang et al., 2017) and could be applied to large-scale models with groundwater representation. For  
574 example, declining heads and decreasing groundwater discharge have been compared between a  
575 calibrated regional-scale model (RRCA, 2003) and a global model (de Graaf et al., 2019). A challenge to  
576 inter-scale comparisons is that regional-scale models often have more spatially complex subsurface  
577 parameterizations because they have access to local data which can complicate model inter-  
578 comparison. Another approach which may be useful is running large-scale models over smaller  
579 (regional) domains at a higher spatial resolution (same as a regional-scale model) so that model  
580 structure influences the comparison less. In the future, various variables that are hard to directly  
581 observe at large scales but routinely simulated in regional-scale models such as baseflow or recharge  
582 could be used to evaluate large-scale models. In this way, the output fluxes and intermediate spatial  
583 scale of regional models provide a bridge across the “river of incommensurability” between highly  
584 location-specific data such as well observations and the coarse resolution of large-scale models. It is  
585 important to consider that regional-scale models are not necessarily or inherently more accurate than  
586 large-scale models since problems may arise from conceptualization, groundwater-surface water  
587 interactions, scaling issues, parameterization etc.

588  
589 In order for a regional-scale model to provide a useful evaluation of a large-scale model, there are  
590 several important documentation and quality characteristics it should meet. At a bare minimum, the  
591 regional-scale model must be accessible and therefore meet basic replicability requirements including  
592 open and transparent input and output data and model code to allow large-scale modelers to run the  
593 model and interpret its output. Documentation through peer review, either through a scientific journal  
594 or agency such as the US Geological Survey, would be ideal. It is particularly important that the  
595 documentation discusses limitations, assumptions and uncertainties in the regional-scale model so that  
596 a large-scale modeler can be aware of potential weaknesses and guide their comparison accordingly.  
597 Second, the boundary conditions and/or parameters being evaluated need to be reasonably comparable  
598 between the regional- and large-scale models. For example, if the regional-scale model includes human  
599 impacts through groundwater pumping while the large-scale model does not, a comparison of baseflow  
600 between the two models may not be appropriate. Similarly, there needs to be consistency in the time  
601 period simulated between the two models. Finally, as with data-driven model evaluation, the purpose of  
602 the large-scale model needs to be consistent with the model-based evaluation; matching the hydraulic  
603 head of a regional-scale model, for instance, does not indicate that estimates of stream-aquifer  
604 exchange are valid. Ideally, we recommend developing a community database of regional-scale models  
605 that meet this criteria. It is important to note that Rossman & Zlotnik (2014) review 88 regional-scale  
606 models while a good example of such a repository is the California Groundwater Model Archive  
607 ([https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-](https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-modeling.html)  
608 [modeling.html](https://ca.water.usgs.gov/sustainable-groundwater-management/california-groundwater-modeling.html)).



609  
610 In addition to evaluating whether models are similar in terms of their outputs, e.g. whether they  
611 simulate similar groundwater head dynamics, it is also relevant to understand whether the influence of  
612 controlling parameters are similar across models. This type of analysis provides insights into process  
613 controls as well as dominant uncertainties. Sensitivity analysis provides the mathematical tools to  
614 perform this type of model evaluation (Saltelli et al., 2008; Pianosi et al., 2016; Borgonovo et al., 2017).  
615 Recent applications of sensitivity analysis to understand modelled controls on groundwater related  
616 processes include the study by Reinecke et al. (2019b) trying to understand parametric controls on  
617 groundwater heads and flows within a global groundwater model. Maples et al. (2020) demonstrated  
618 that parametric controls on groundwater recharge can be assessed for complex models, though over a  
619 smaller domain. As highlighted by both of these studies, more work is needed to understand how to  
620 best use sensitivity analysis methods to assess computationally expensive, spatially distributed and  
621 complex groundwater models across large domains (Hill et al., 2016). In the future, it would be useful to  
622 go beyond parameter uncertainty analysis (e.g. Reinecke et al. 2019b) to begin to look at all of the  
623 modelling decisions holistically such as the forcing data (Weiland et al., 2015) and digital elevation  
624 models (Hawker et al., 2018). Addressing this problem requires advancements in statistics (more  
625 efficient sensitivity analysis methods), computing (more effective model execution), and access to large-  
626 scale models codes (Hutton et al. 2016), but also better utilization of process understanding, for  
627 example to create process-based groups of parameters which reduces the complexity of the sensitivity  
628 analysis study (e.g. Hartmann et al., 2015; Reinecke et al., 2019b).  
629

### 630 3.3 Expert-based model evaluation

631 A path much less traveled is expert-based model evaluation which would develop hypotheses of  
632 phenomena (and related behaviors, patterns or signatures) we expect to emerge from large-scale  
633 groundwater systems based on expert knowledge, intuition, or experience. In essence, this model  
634 evaluation approach flips the traditional scientific method around by using hypotheses to test the  
635 simulation of emergent processes from large-scale models, rather than using large-scale models to test  
636 our hypotheses about environmental phenomena. This might be an important path forward for regions  
637 where available data is very sparse or unreliable. The recent discussion by Fan et al. (2019) shows how  
638 hypotheses about large-scale behavior might be derived from expert knowledge gained through the  
639 study of smaller scale systems such as critical zone observatories. While there has been much effort to  
640 improve our ability to make hydrologic predictions in ungauged locations through the regionalization of  
641 hydrologic variables or of model parameters (Bloeschl et al., 2013), there has been much less effort to  
642 directly derive expectations of hydrologic behavior based on our perception of the systems under study.  
643  
644 Large-scale models could then be evaluated against such hypotheses, thus providing a general  
645 opportunity to advance how we connect hydrologic understanding with large-scale modeling - a strategy  
646 that could also potentially reduce epistemic uncertainty (Beven et al., 2019), and which may be  
647 especially useful for groundwater systems given the data limitations described above. Developing  
648 appropriate and effective hypotheses is crucial and should likely focus on large-scale controlling factors



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

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

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

672  
673 Formal approaches are available to gather the opinions of experts and to integrate them into a joint  
674 result, often called expert elicitation (Aspinall, 2010; Cooke, 1991; O’Hagan, 2019). Expert elicitation  
675 strategies have been used widely to describe the expected behavior of environmental or man-made  
676 systems for which we have insufficient data or knowledge to build models directly. Examples include  
677 aspects of future sea-level rise (Bamber and Aspinall, 2013), tipping points in the Earth system (Lenton  
678 et al., 2018), or the vulnerability of bridges to scour due to flooding (Lamb et al., 2017). In the  
679 groundwater community, expert opinion is already widely used to develop system conceptualizations  
680 and related model structures (Krueger et al., 2012; Rajabi et al., 2018; Refsgaard et al., 2006), or to  
681 define parameter priors (Ross et al., 2009; Doherty and Christensen, 2011; Brunner et al., 2012;  
682 Knowing and Werner, 2016; Rajabi and Ataie-Ashtiani, 2016). The term expert opinion may be  
683 preferable to the term expert knowledge because it emphasizes a preliminary state of knowledge  
684 (Krueger et al., 2012).

685  
686 A critical benefit of expert elicitation is the opportunity to bring together researchers who have  
687 experienced very different groundwater systems around the world. It is infeasible to expect that a single  
688 person could have gained in-depth experience in modelling groundwater in semi-arid regions, in cold  
689 regions, in tropical regions etc. Being able to bring together different experts who have studied one or a  
690 few of these systems to form a group would certainly create a whole that is bigger than the sum of its



691 parts. If captured, it would be a tremendous source of knowledge for the evaluation of large-scale  
692 groundwater models. A challenge though is to formalize this knowledge in such a way that it is still  
693 usable by third parties that did not attend the expert workshop itself.

694  
695 So, while expert opinion and judgment play a role in any scientific investigation (O'Hagan, 2019),  
696 including that of groundwater systems, we rarely use formal strategies to elicit this opinion. It is also less  
697 common to use expert opinion to develop hypotheses about the dynamic behavior of groundwater  
698 systems, rather than just priors on its physical characteristics. Yet, it is intuitive that information about  
699 system behavior can help in evaluating the plausibility of model outputs (and thus of the model itself).  
700 This is what we call expert-based evaluation herein. Expert elicitation is typically done in workshops with  
701 groups of a dozen or so experts (e.g. Lamb et al., 2018). Upscaling such expert elicitation in support of  
702 global modeling would require some web-based strategy and a formalized protocol to engage a  
703 sufficiently large number of people. Contributors could potentially be incentivized to contribute to the  
704 web platform by publishing a data paper with all contributors as co-authors and a secondary analysis  
705 paper with just the core team as coauthors. We recommend the community develop expert elicitation  
706 strategies to identify effective hypotheses that directly link to the relevant large-scale hydrologic  
707 processes of interest.

#### 708 **4. TOWARDS A HOLISTIC EVALUATION OF GROUNDWATER REPRESENTATION IN LARGE-SCALE** 709 **MODELS**

710  
711 Ideally, all three strategies (observation-based, model-based, expert-based) should be pursued  
712 simultaneously because the strengths of one strategy might further improve others. For example,  
713 expert- or model-based evaluation may highlight and motivate the need for new observations in certain  
714 regions or at new resolutions. Or observation-based model evaluation could highlight and motivate  
715 further model development or lead to refined or additional hypotheses. We thus recommend the  
716 community significantly strengthens efforts to evaluate large-scale models using all three strategies.  
717 Implementing these three model evaluation strategies may require a significant effort from the scientific  
718 community, so we therefore conclude with two tangible community-level initiatives that would be  
719 excellent first steps that can be pursued simultaneously with efforts by individual research groups or  
720 collaborations of multiple research groups.

721  
722 First, we need to develop a 'Groundwater Modeling Data Portal' that would both facilitate and  
723 accelerate the evaluation of groundwater representation in continental to global scale models (Bierkens,  
724 2015). Existing initiatives such as IGRAC's Global Groundwater Monitoring Network ([https://www-  
725 igrac.org/special-project/ggm-global-groundwater-monitoring-network](https://www-igrac.org/special-project/ggm-global-groundwater-monitoring-network)) and HydroFrame  
726 ([www.hydroframe.org](http://www.hydroframe.org)), are an important first step but were not designed to improve the evaluation of  
727 large-scale models and the synthesized data remains very heterogeneous - unfortunately, even  
728 groundwater level time series data often remains either hidden or inaccessible for various reasons. This  
729 open and well documented data portal should include:

730 a) observations for evaluation (Table 1) as well as derived signatures (Section 3.1);



- 731           b) regional-scale models that meet the standards described above and could facilitate inter-scale  
732           comparison (Section 3.2);  
733           c) Schematizations, conceptual or perceptual models of large-scale models since these are the  
734           basis of computational models; and  
735           d) Hypothesis and other results derived from expert elicitation (Section 3.3).

736 Meta-data documentation, data tagging, aggregation and services as well as consistent data structures  
737 using well-known formats (netCDF, .csv, .txt) will be critical to developing a useful, dynamic and evolving  
738 community resource. The data portal should be directly linked to harmonized input data such as forcings  
739 (climate, land and water use etc.) and parameters (topography, subsurface parameters etc.), model  
740 codes, and harmonized output data. Where possible, the portal should follow established protocols,  
741 such as the Dublin Core Standards for metadata (<https://dublincore.org>) and ISIMIP protocols for  
742 harmonizing data and modeling approach, and would ideally be linked to or contained within an existing  
743 disciplinary repository such as HydroShare (<https://www.hydroshare.org/>) to facilitate discovery,  
744 maintenance, and long-term support. Additionally, an emphasis on model objective, uncertainty and  
745 regional differences as highlighted (Section 3) will be important in developing the data portal. Like  
746 expert-elicitation, contribution to the data portal could be incentivized through co-authorship in data  
747 papers and by providing digital object identifiers (DOIs) to submitted data and models so that they are  
748 citable. By synthesizing and sharing groundwater observations, models, and hypotheses, this portal  
749 would be broadly useful to the hydrogeological community beyond just improving global model  
750 evaluation.

751  
752 Second, we suggest ISIMIP, or a similar model intercomparison project, could be harnessed as a  
753 platform to improve the evaluation of groundwater representation in continental to global scale models.  
754 For example, in ISIMIP (Warszawski et al., 2014), modelling protocols have been developed with an  
755 international network of climate-impact modellers across different sectors (e.g. water, agriculture,  
756 energy, forestry, marine ecosystems) and spatial scales. Originally, ISIMIP started with multi-model  
757 comparison (model-based model evaluation), with a focus on understanding how model projections  
758 vary across different sectors and different climate change scenarios (ISIMIP Fast Track). However, more  
759 rigorous model evaluation came to attention more recently with ISIMIP2a, and various observation data,  
760 such as river discharge (Global Runoff Data Center), terrestrial water storage (GRACE), and water use  
761 (national statistics), have been used to evaluate historical model simulation (observation-based model  
762 evaluation). To better understand model differences and to quantify the associated uncertainty sources,  
763 ISIMIP2b includes evaluating scenarios (land use, groundwater use, human impacts, etc) and key  
764 assumptions (no explicit groundwater representation, groundwater availability for the future, water  
765 allocation between surface water and groundwater), highlighting that different types of hypothesis  
766 derived as part of the expert-based model evaluation could possibly be simulated as part of the ISIMIP  
767 process in the future. While there has been a significant amount of research and publications on MIPs  
768 including surface water availability, limited multi-model assessments for large-scale groundwater  
769 studies exist. Important aspects of MIPs in general could facilitate all three model evaluation strategies:  
770 community-building and cooperation with various scientific communities and research groups, and  
771 making the model input and output publicly available in a standardized format.  
772



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

788

789 **Acknowledgements:**

790 The commentary is based on a workshop at the University of Bristol and significant debate and  
791 discussion before and after. This community project was directly supported by a Benjamin Meaker  
792 Visiting Professorship at the Bristol University to TG and a Royal Society Wolfson Award to TW  
793 (WM170042). We thank many members of the community who contributed to the discussions,  
794 especially at the IGEM (Impact of Groundwater in Earth System Models) workshop in Taiwan.

795

796 **Author Contributions:** (using the [CRediT taxonomy](#) which offers standardized descriptions of author  
797 contributions) conceptualization and writing original draft: TG, TW and PD; writing - review and  
798 editing: all co-authors. Authors are ordered by contribution for the first three coauthors (TG, TW and PD)  
799 and then ordered in reverse alphabetical order for all remaining coauthors.

800

801 **Code/Data availability:**

802 No code or data were used in the writing of this manuscript

803

804 **Competing interests:**

805 The authors declare no competing interests.

806

807

808



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

810

Data type	Strengths	Limitations	Data availability and spatial resolution
<b>Available observations already used to evaluate large-scale models</b>			
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	<a href="#">IGRAC Global Groundwater Monitoring Network</a> ; 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 <a href="#">European Groundwater Drought Initiative</a> 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 $\approx 100,000 \text{ km}^2$ (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 $\text{km}^2$ )
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	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 <a href="#">data sources</a> ; large to small basin.
Evapotranspiration	Widely available; related to groundwater recharge	Not a direct groundwater observations	Various datasets (Miralles et al., 2016); gridded

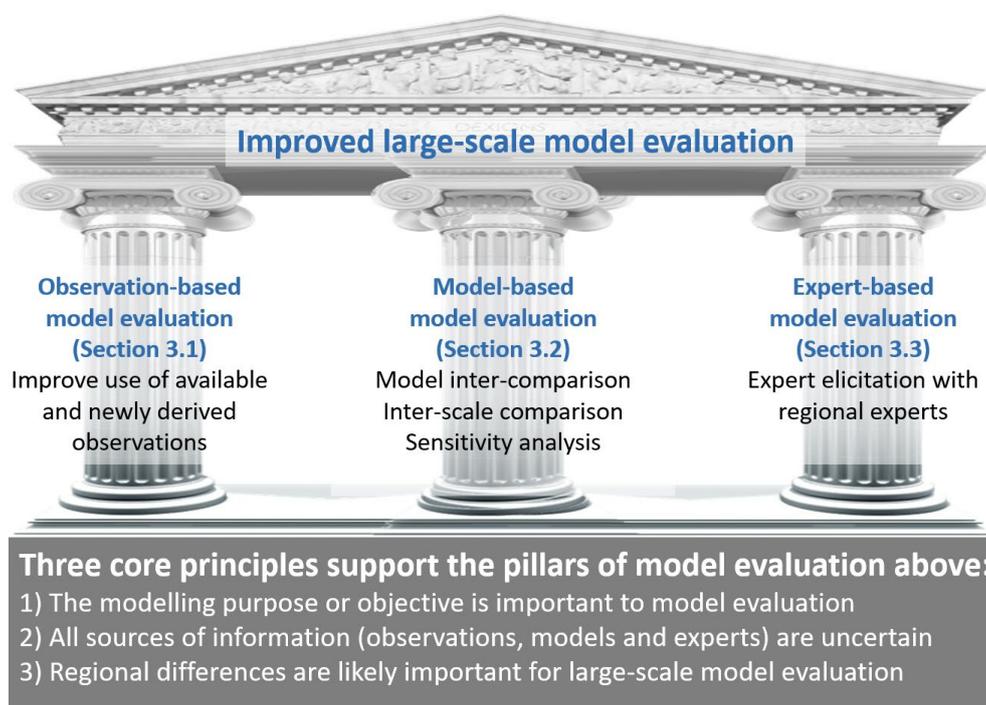


	or discharge (for shallow water tables)		
<b>Available observations not being used to evaluate large-scale models</b>			
Baseflow index (BFI) or baseflow recession (k)	Possible integrator of groundwater contribution to streamflow over a basin	BFI and k values vary with method; baseflow may be dominated 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,	Provides information about changes in surface	Provides indirect information and needs a geomechanical	Leveling data, GPS data and lidar observations mostly limited to



radar/lidar) an in particular land subsidence observations	elevation that are related to groundwater head variations or groundwater head decline	model to translate to head. Introduces additional uncertainty of geomechanical properties.	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.
--	---	--	--

811  
 812  
 813  
 814



815  
 816 **Figure 1: Improved large-scale model evaluation rests on three pillars: observation-, model-, and**  
 817 **expert-based model evaluation. We argue that each pillar is an essential strategy so that all three**  
 818 **should be simultaneously pursued by the scientific community. The three pillars of model evaluation**  
 819 **all rest on three core principles related to 1) model objectives, 2) uncertainty and 3) regional**  
 820 **differences.**



821

822

823 **References**

824 Addor, N., & Melsen, L. A. (2018). Legacy, Rather Than Adequacy, Drives the Selection of Hydrological Models.  
825 *Water Resources Research*, 0(0). <https://doi.org/10.1029/2018WR022958>

826

827 Al-Yaari, A., Ducharne, A., Cheruy, F., Crow, W.T. & Wigneron, J.P. (2019). Satellite-based soil moisture provides  
828 missing link between summertime precipitation and surface temperature biases in CMIP5 simulations over  
829 conterminous United States. *Scientific Reports*, 9, article number 1657, doi:10.1038/s41598-018-38309-5

830

831 Anderson, M. P., Woessner, W. W. & Hunt, R. (2015a). *Applied groundwater modeling- 2nd Edition*. San Diego:  
832 Academic Press.

833 Anderson, R. G., Min-Hui Lo, Swenson, S., Famiglietti, J. S., Tang, Q., Skaggs, T. H., Lin, Y.-H., and Wu, R.-J. (2015b),  
834 Using satellite-based estimates of evapotranspiration and groundwater changes to determine anthropogenic  
835 water fluxes in land surface models, *Geosci. Model Dev.*, 8, 3021-3031, doi:10.5194/gmd-8-3021-2015. Alley, W.M.  
836 and LF Konikow (2015) Bringing GRACE down to earth. *Groundwater* 53 (6): 826–829

837 Anyah, R. O., Weaver, C. P., Miguez-Macho, G., Fan, Y., & Robock, A. (2008). Incorporating water table dynamics in  
838 climate modeling: 3. Simulated groundwater influence on coupled land-atmosphere variability. *J. Geophys. Res.*,  
839 113. Retrieved from <http://dx.doi.org/10.1029/2007JD009087>

840 Archfield, S. A., Clark, M., Arheimer, B., Hay, L. E., McMillan, H., Kiang, J. E., et al. (2015). Accelerating advances in  
841 continental domain hydrologic modeling. *Water Resources Research*, 51(12), 10078–10091.  
842 <https://doi.org/10.1002/2015WR017498>

843 Aspinall, W. (2010). A route to more tractable expert advice. *Nature*, 463, 294–295.

844 <https://doi.org/10.1038/463294a>

845 Bamber, J.L. and Aspinall, W.P. (2013). An expert judgement assessment of future sea level rise from the ice  
846 sheets. *Nature Climate Change*. 3(4), 424-427.

847 Barthel, R. (2014). HESS Opinions “Integration of groundwater and surface water research: an interdisciplinary  
848 problem?” *Hydrology and Earth System Sciences*, 18(7), 2615–2628.

849 Beck, H. et al (2013). Global patterns in base flow index and recession based on streamflow observations from  
850 3394 catchments. *Water Resources Research*.

851 Befus, K., Jasechko, S., Luijendijk, E., Gleeson, T., Cardenas, M.B. (2017) The rapid yet uneven turnover of Earth's  
852 groundwater. (2017) *Geophysical Research Letters* 11: 5511-5520 doi: 10.1002/2017GL073322

853 Best, M. J., Pryor, M., Clark, D. B., Rooney, G. G., Essery, R. L. H., Ménard, C. B., Edwards, J. M., Hendry, M. A.,  
854 Porson, A., Gedney, N., Mercado, L. M., Sitch, S., Blyth, E., Boucher, O., Cox, P. M., Grimmond, C. S. B., & Harding,  
855 R. J. (2011). The Joint UK Land Environment Simulator (JULES), model description – Part 1: Energy and water fluxes,  
856 *Geosci. Model Dev.*, 4, 677-699. <https://doi.org/10.5194/gmd-4-677-2011>



- 857 Beven, K. (2000). Uniqueness of place and process representations in hydrological modelling. *Hydrology and Earth*  
858 *System Sciences*, 4(2), 203–213.
- 859 Beven, K. (2005). On the concept of model structural error. *Water Science & Technology*, 52(6), 167–175.
- 860 Beven, K. (2016). Facets of uncertainty: epistemic uncertainty, nonstationarity, likelihood, hypothesis testing, and  
861 communication. *Hydrological Sciences Journal*, 61(9), 1652–1665, DOI: 10.1080/02626667.2015.1031761
- 862 Beven, K. (2019) How to make advances in hydrological modelling. In: *Hydrology Research*. 50, 6, p. 1481–1494. 14  
863 p.
- 864 Beven, K. J., and H. L. Cloke (2012), Comment on “Hyperresolution global land surface modeling: Meeting a grand  
865 challenge for monitoring Earth’s terrestrial water” by Eric F. Wood et al., *Water Resour.Res.*, 48, W01801,  
866 doi:10.1029/2011WR010982.
- 867 Beven, K.J., Aspinall, W.P., Bates, P.D., Borgomeo, E., Goda, K., Hall, J.W., Page, T., Phillips, J.C., Simpson, M., Smith,  
868 P.J., Wagener, T. and Watson, M. 2018. Epistemic uncertainties and natural hazard risk assessment – Part 2: What  
869 should constitute good practice? *Natural Hazards and Earth System Sciences*, 18, 10.5194/nhess-18-1-2018
- 870 Bierkens, M. F. P. (2015). Global hydrology 2015: State, trends, and directions. *Water Resources Research*, 51(7),  
871 4923–4947. <https://doi.org/10.1002/2015WR017173>
- 872 Bierkens, M. F.P. & Wada, Y. (2019). Non-renewable groundwater use and groundwater depletion: A review.  
873 *Environmental Research Letters*, 14(6), 063002
- 874 Boone, A. A., Habets, F., Noilhan, J., Clark, D., Dirmeyer, P., Fox, S., Gusev, Y., Haddeland, I., Koster, R., Lohmann,  
875 D., Mahanama, S., Mitchell, K., Nasonova, O., Niu, G. Y., Pitman, A., Polcher, J., Shmakin, A. B., Tanaka, K., Van Den  
876 Hurk, B., Vérant, S., Verseghy, D., Viterbo, P. and Yang, Z. L.: The Rhône-aggregation land surface scheme  
877 intercomparison project: An overview, *J. Clim.*, 17(1), 187–208, doi:10.1175/1520-  
878 0442(2004)017<0187:TRLSSI>2.0.CO;2, 2004.
- 879 Borgonovo, E. Lu, X. Plischke, E. Rakovec, O. and Hill, M. C. (2017). Making the most out of a hydrological model  
880 data set: Sensitivity analyses to open the model black-box. *Water Resources Research*.  
881 DOI:10.1002/2017WR020767
- 882 Bresciani, E., P. Goderniaux, and O. Batelaan (2016), Hydrogeological controls of water table-land surface  
883 interactions, *Geophysical Research Letters*, 43, 9653–9661.
- 884 Bresciani, E., Cranswick, R. H., Banks, E. W., Batlle-Aguilar, J., et al. (2018). Using hydraulic head, chloride and  
885 electrical conductivity data to distinguish between mountain-front and mountain-block recharge to basin aquifers.  
886 *Hydrology and Earth System Sciences*, 22(2), 1629–1648.
- 887 Brunner, P., J. Doherty, and C. T. Simmons (2012), Uncertainty assessment and implications for data acquisition in  
888 support of integrated hydrologic models, *Water Resources Research*, 48.  
889
- 890 Burgess, W. G., Shamsudduha, M., Taylor, R. G., Zahid, A., Ahmed, K. M., Mukherjee, A., et al. (2017). Terrestrial  
891 water load and groundwater fluctuation in the Bengal Basin. *Scientific Reports*, 7(1), 3872.



- 892 Caceres, D., Marzeion, B., Malles, J.H., Gutknecht, B., Müller Schmied, H., Döll, P. (2020): Assessing global water  
893 mass transfers from continents to oceans over the period 1948–2016. *Hydrol. Earth Syst. Sci. Discuss.*  
894 doi:10.5194/hess-2019-664
- 895 Ceola, S., Arheimer, B., Baratti, E., Blöschl, G., Capell, R., Castellarin, A., et al. (2015). Virtual laboratories: new  
896 opportunities for collaborative water science. *Hydrology and Earth System Sciences*, 19(4), 2101–2117.
- 897 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)  
898 Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between  
899 hydrological models, *Water Resour. Res.*, 44, W00B02, doi:10.1029/2007WR006735.
- 900 Clark, M. P., et al. (2015), A unified approach for process-based hydrologic modeling: 1. Modeling concept, *Water*  
901 *Resources Research*, 51, 2498–2514, doi:10.1002/2015WR017198
- 902 Condon, L. E., & Maxwell, R. M. (2019). Simulating the sensitivity of evapotranspiration and streamflow to large-  
903 scale groundwater depletion. *Science Advances*, 5(6), eaav4574. <https://doi.org/10.1126/sciadv.aav4574>
- 904 Condon, LE et al Evapotranspiration depletes groundwater under warming over the contiguous United States  
905 *Nature Comm*, 2020, <https://doi.org/10.1038/s41467-020-14688-0>
- 906 Condon, L. E., Markovich, K. H., Kelleher, C. A., McDonnell, J. J., Ferguson, G., & McIntosh, J. C. (2020). Where Is the  
907 Bottom of a Watershed? *Water Resources Research*, 56(3). <https://doi.org/10.1029/2019wr026010>
- 908 Cooke, R. (1991). *Experts in uncertainty: opinion and subjective probability in science*. Oxford University Press on  
909 Demand.
- 910 Cuthbert, M. O., Gleeson, T., Moosdorf, N., Befus, K. M., Schneider, A., Hartmann, J., & Lehner, B. (2019). Global  
911 patterns and dynamics of climate–groundwater interactions. *Nature Climate Change*, 9, 137–141  
912 <https://doi.org/10.1038/s41558-018-0386-4>
- 913 Cuthbert, M. O., et al. (2019) Observed controls on resilience of groundwater to climate variability in sub-Saharan  
914 Africa. *Nature* 572: 230–234
- 915 Dalin, C., Wada, Y., Kastner, T., & Puma, M. J. (2017). Groundwater depletion embedded in international food  
916 trade. *Nature*, 543(7647), 700–704. <https://doi.org/10.1038/nature21403>
- 917 Dirmeyer, P. A.: A History and Review of the Global Soil Wetness Project (GSWP), *J. Hydrometeorol.*, 12(5),  
918 110404091221083, doi:10.1175/jhm-d-10-05010, 2011
- 919 Doherty, J., and S. Christensen (2011), Use of paired simple and complex models to reduce predictive bias and  
920 quantify uncertainty, *Water Resources Research*, 47(12),
- 921 Döll, P., Fiedler, K. (2008): Global-scale modeling of groundwater recharge. *Hydrol. Earth Syst. Sci.*, 12, 863–885,  
922 doi: 10.5194/hess-12-863-2008
- 923 Döll, P., Douville, H., Güntner, A., Müller Schmied, H., Wada, Y. (2016): Modelling freshwater resources at the  
924 global scale: Challenges and prospects. *Surveys in Geophysics*, 37(2), 195–221. doi: 10.1007/s10712-015-9343-1
- 925 Döll, P., Müller Schmied, H., Schuh, C., Portmann, F. T., & Eicker, A. (2014a). Global-scale assessment of  
926 groundwater depletion and related groundwater abstractions: Combining hydrological modeling with information



- 927 from well observations and GRACE satellites. *Water Resources Research*, 50(7), 5698–5720.  
928 <https://doi.org/10.1002/2014WR015595>
- 929 Döll, P., Fritsche, M., Eicker, A., Müller Schmied, H. (2014b): Seasonal water storage variations as impacted by  
930 water abstractions: Comparing the output of a global hydrological model with GRACE and GPS observations.  
931 *Surveys in Geophysics*, 35(6), 1311-1331, doi: 10.1007/s10712-014-9282-2.
- 932 Döll, P., Hoffmann-Dobrev, H., Portmann, F.T., Siebert, S., Eicker, A., Rodell, M., Strassberg, G., Scanlon, B. (2012):  
933 Impact of water withdrawals from groundwater and surface water on continental water storage variations. J.  
934 *Geodyn.* 59-60, 143-156, doi:10.1016/j.jog.2011.05.001.
- 935 Duan Q., Schaake, J., Andreassian, V., Franks, S., Gupta, H.V., Gusev, Y.M., Habets, F., Hall, A., Hay, L., Hogue, T.S.,  
936 Huang, M., Leavesley, G., Liang, X., Nasonova, O.N., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T. and Wood,  
937 E.F. (2006). Model Parameter Estimation Experiment (MOPEX): Overview and Summary of the Second and Third  
938 Workshop Results. *Journal of Hydrology*, 320(1-2), 3-17.
- 939 Enemark, T., Peeters, L. J. M., Mallants, D., & Batelaan, O. (2019). Hydrogeological conceptual model building and  
940 testing: A review. *Journal of Hydrology*, 569, 310–329. <https://doi.org/10.1016/j.jhydrol.2018.12.007>
- 941 Erban L E, Gorelick S M and Zebker H A 2014 Groundwater extraction, land subsidence, and sea-level rise in the  
942 Mekong Delta, Vietnam *Environ. Res. Lett.* 9 084010
- 943 Famiglietti, J. S., & E. F. Wood (1994). Multiscale modeling of spatially variable water and energy balance  
944 processes, *Water Resour. Res.*, 30(11), 3061–3078, <https://doi.org/10.1029/94WR01498>
- 945 Fan, Y. et al., (2019) Hillslope hydrology in global change research and Earth System modeling. *Water Resources*  
946 *Research*, doi.org/10.1029/2018WR023903
- 947 Fan, Y. (2015). Groundwater in the Earth’s critical zone: Relevance to large-scale patterns and processes. *Water*  
948 *Resources Research*, 51(5), 3052–3069. <https://doi.org/10.1002/2015WR017037>
- 949 Fan, Y., & Miguez-Macho, G. (2011). A simple hydrologic framework for simulating wetlands in climate and earth  
950 system models. *Climate Dynamics*, 37(1–2), 253–278.
- 951 Fan, Y., Li, H., & Miguez-Macho, G. (2013). Global patterns of groundwater table depth. *Science*, 339(6122), 940–  
952 943.
- 953 Fenicia, F., D. Kavetski, and H. H. G. Savenije (2011), Elements of a flexible approach for conceptual hydrological  
954 modeling: 1. Motivation and theoretical development, *Water Resources Research*, 47(11), W11510,  
955 10.1029/2010wr010174.
- 956 Forrester, M.M. and Maxwell, R.M. Impact of lateral groundwater flow and subsurface lower boundary conditions  
957 on atmospheric boundary layer development over complex terrain. *Journal of Hydrometeorology*,  
958 doi:10.1175/JHM-D-19-0029.1, 2020.
- 959 Forrester, M.M., Maxwell, R.M., Bearup, L.A., and Gochis, D.J. Forest Disturbance Feedbacks from Bedrock to  
960 Atmosphere Using Coupled Hydro-Meteorological Simulations Over the Rocky Mountain Headwaters. *Journal of*  
961 *Geophysical Research-Atmospheres*, 123:9026-9046, doi:10.1029/2018JD028380 2018.



- 962 Freeze, R. A., & Witherspoon, P. A. (1966). Theoretical analysis of regional groundwater flow, 1. Analytical and  
963 numerical solutions to a mathematical model. *Water Resources Research*, 2, 641–656.
- 964 Foster, S., Chilton, J., Nijsten, G.-J., & Richts, A. (2013). Groundwater — a global focus on the ‘local resource.’  
965 *Current Opinion in Environmental Sustainability*, 5(6), 685–695. doi.org/10.1016/j.cosust.2013.10.010
- 966 Garven, G. (1995). Continental-scale groundwater flow and geologic processes. *Annual Review of Earth and*  
967 *Planetary Sciences*, 23, 89–117.
- 968 Gascoïn, S., Ducharne, A., Ribstein, P., Carli, M., Habets, F. (2009). Adaptation of a catchment-based land surface  
969 model to the hydrogeological setting of the Somme River basin (France). *Journal of Hydrology*, 368(1-4), 105-116.  
970 <https://doi.org/10.1016/j.jhydrol.2009.01.039>
- 971 Genereux, D. (1998). Quantifying uncertainty in tracer-based hydrograph separations. *Water Resources Research*,  
972 34(4), 915–919.
- 973 Gilbert, J.M., Maxwell, R.M. and Gochis, D.J. Effects of water table configuration on the planetary boundary layer  
974 over the San Joaquin River watershed, California. *Journal of Hydrometeorology*, 18:1471-1488, doi:10.1175/JHM-  
975 D-16-0134.1, 2017.
- 976 Gleeson, T., & Manning, A. H. (2008). Regional groundwater flow in mountainous terrain: Three-dimensional  
977 simulations of topographic and hydrogeologic controls. *Water Resources Research*, 44. Retrieved from  
978 <http://dx.doi.org/10.1029/2008WR006848>
- 979 Gleeson, T., Befus, K. M., Jasechko, S., Luijendijk, E., & Cardenas, M. B. (2016). The global volume and distribution  
980 of modern groundwater. *Nature Geosci*, 9(2), 161–167.
- 981 de Graaf, I. E. M., van Beek, L. P. H., Wada, Y., & Bierkens, M. F. P. (2014). Dynamic attribution of global water  
982 demand to surface water and groundwater resources: Effects of abstractions and return flows on river discharges.  
983 *Advances in Water Resources*, 64(0), 21–33. <https://doi.org/10.1016/j.advwatres.2013.12.002>
- 984 de Graaf, I. E. M., Sutanudjaja, E. H., Van Beek, L. P. H., & Bierkens, M. F. P. (2015). A high-resolution global-scale  
985 groundwater model. *Hydrology and Earth System Sciences*, 19(2), 823–837.
- 986 de Graaf, I. E. M., van Beek, L. P. H., Gleeson, T., Moosdorf, N., Schmitz, O., Sutanudjaja, E. H., & Bierkens, M. F. P.  
987 (2017). A global-scale two-layer transient groundwater model: Development and application to groundwater  
988 depletion. *Advances in Water Resources*, 102, 53–67. <https://doi.org/10.1016/j.advwatres.2017.01.011>
- 989 de Graaf, I. E. M., Gleeson, T., Beek, L. P. H. (Rens) van, Sutanudjaja, E. H., & Bierkens, M. F. P. (2019).  
990 Environmental flow limits to global groundwater pumping. *Nature*, 574(7776), 90–94.  
991 <https://doi.org/10.1038/s41586-019-1594-4>
- 992 Gnann, S. J., Woods, R. A., & Howden, N. J. (2019). Is there a baseflow Budyko curve? *Water Resources Research*,  
993 55(4), 2838–2855.
- 994 Goderniaux, P., P. Davy, E. Bresciani, J.-R. de Dreuzy, and T. Le Borgne (2013), Partitioning a regional groundwater  
995 flow system into shallow local and deep regional flow compartments, *Water Resources Research*, 49(4), 2274-  
996 2286.



- 997 Gosling, S. N., Zaherpour, J., Mount, N. J., Hattermann, F. F., Dankers, R., Arheimer, B., et al. (2017). A comparison  
998 of changes in river runoff from multiple global and catchment-scale hydrological models under global warming  
999 scenarios of 1 °C, 2 °C and 3 °C. *Climatic Change*, 141(3), 577–595. <https://doi.org/10.1007/s10584-016-1773-3>
- 1000 Guimberteau, M., Ducharne, A., Ciais, P., Boisier, J. P., Peng, S., De Weirtdt, M., & Verbeeck, H. (2014). Testing  
1001 conceptual and physically based soil hydrology schemes against observations for the Amazon Basin, *Geosci. Model*  
1002 *Dev.*, 7, 1115–1136. <https://doi.org/10.5194/gmd-7-1115-2014>
- 1003 Habets, F., Boé, J., Déqué, M., Ducharne, A., Gascoin, S., Hachour, A., Martin, E., Pagé, C., Sauquet, E., Terray, L.,  
1004 Thiéry, D., Oudin, L. & Viennot, P. (2013). Impact of climate change on surface water and ground water of two  
1005 basins in Northern France: analysis of the uncertainties associated with climate and hydrological models, emission  
1006 scenarios and downscaling methods. *Climatic Change*, 121, 771–785. <https://doi.org/10.1007/s10584-013-0934-x>
- 1007 Hartmann, A., Gleeson, T., Rosolem, R., Pianosi, F., Wada, Y., & Wagener, T. (2015). A large-scale simulation model  
1008 to assess karstic groundwater recharge over Europe and the Mediterranean. *Geosci. Model Dev.*, 8(6), 1729–1746.  
1009 <https://doi.org/10.5194/gmd-8-1729-2015>
- 1010 Hartmann, Andreas, Gleeson, T., Wada, Y., & Wagener, T. (2017). Enhanced groundwater recharge rates and  
1011 altered recharge sensitivity to climate variability through subsurface heterogeneity. *Proceedings of the National*  
1012 *Academy of Sciences*, 114(11), 2842–2847. <https://doi.org/10.1073/pnas.1614941114>
- 1013 Hattermann, F. F., Krysanova, V., Gosling, S. N., Dankers, R., Daggupati, P., Donnelly, C., et al. (2017). Cross-scale  
1014 intercomparison of climate change impacts simulated by regional and global hydrological models in eleven large  
1015 river basins. *Climatic Change*, 141(3), 561–576. <https://doi.org/10.1007/s10584-016-1829-4>
- 1016 Hay, L., Norton, P., Viger, R., Markstrom, S., Regan, R. S., & Vanderhoof, M. (2018). Modelling surface-water  
1017 depression storage in a Prairie Pothole Region. *Hydrological Processes*, 32(4), 462–479.  
1018 <https://doi.org/10.1002/hyp.11416>
- 1019 Henderson-Sellers, A., Z. L. Yang, and R. E. Dickinson: The Project for Intercomparison of Land-Surface Schemes  
1020 (PILPS). *Bull. Amer. Meteor. Soc.*, 74, 1335–1349, 1993
- 1021 Herbert, C., & Döll, P. (2019). Global assessment of current and future groundwater stress with a focus on  
1022 transboundary aquifers. *Water Resources Research*, 55, 4760–4784. <https://doi.org/10.1029/2018WR023321>
- 1023 Heudorfer, B., Haaf, E., Stahl, K., & Barthel, R. (2019). Index-based characterization and quantification of  
1024 groundwater dynamics. *Water Resources Research*, 55, 5575–5592. <https://doi.org/10.1029/2018WR024418>
- 1025 Hill, M. C. (2006). The practical use of simplicity in developing ground water models. *Ground Water*, 44(6), 775–  
1026 781. <https://doi.org/10.1111/j.1745-6584.2006.00227.x>
- 1027 Hill, M. C., & Tiedeman, C. R. (2007). *Effective groundwater model calibration*. Wiley.
- 1028 Hill, M. C., Kavetski, D. Clark, M. Ye, M. Arabi, M. Lu, D. Foglia, L. & Mehl, S. (2016). Practical use of computationally  
1029 frugal model analysis methods. *Groundwater*. DOI:10.1111/gwat.12330  
1030
- 1031 Hiscock, K. M., & Bense, V. F. (2014). *Hydrogeology—principles and practice* (2nd edition). Blackwell.



- 1032 Huang, S., Kumar, R., Flörke, M., Yang, T., Hundecha, Y., Kraft, P., et al. (2017). Evaluation of an ensemble of  
1033 regional hydrological models in 12 large-scale river basins worldwide. *Climatic Change*, 141(3), 381–397.  
1034 <https://doi.org/10.1007/s10584-016-1841-8>
- 1035 Hrachowitz, M., Fovet, O., Ruiz, L., Euser, T., Gharari, S., Nijzink, R., Freer, J., Savenije, H.H.G. and Gascuel-Oudou, C.  
1036 (2014). Process Consistency in Models: the Importance of System Signatures, Expert Knowledge and Process  
1037 Complexity. *Water Resources Research* 50:7445–7469.
- 1038 Hunt, R. J., Walker, J. F., Selbig, W. R., Westenbroek, S. M., & Regan, R. S. (2013). Simulation of climate-change  
1039 effects on streamflow, lake water budgets, and stream temperature using GSFLOW and SNTMP, Trout Lake  
1040 Watershed, Wisconsin. USGS Scientific Investigations Report No. 2013–5159. Reston, VA: U.S. Geological Survey.
- 1041 Hutton, C., Wagener, T., Freer, J., Han, D., Duffy, C., & Arheimer, B. (2016). Most computational hydrology is not  
1042 reproducible, so is it really science? *Water Resources Research*, 52(10), 7548–7555.  
1043 <https://doi.org/10.1002/2016WR019285>
- 1044 Jasechko, S., Birks, S.J., Gleeson, T., Wada, Y., Sharp, Z.D., Fawcett, P.J., McDonnell, J.J., Welker, J.M. (2014)  
1045 Pronounced seasonality in the global groundwater recharge. *Water Resources Research*. 50, 8845–8867 doi:  
1046 10.1002/2014WR015809
- 1047 Jasechko, S., Perrone, D., Befus, K. M., Bayani Cardenas, M., Ferguson, G., Gleeson, T., et al. (2017). Global aquifers  
1048 dominated by fossil groundwaters but wells vulnerable to modern contamination. *Nature Geoscience*, 10(6), 425–  
1049 429. <https://doi.org/10.1038/ngeo2943>
- 1050 Jung, M., et al. (2011). Global patterns of land-atmosphere fluxes of carbon dioxide, latent heat, and sensible  
1051 heat derived from eddy covariance, satellite, and meteorological observations. *J. Geophys. Res.*, 116,  
1052 G00J07, doi:10.1029/2010JG001566.
- 1053 Keune, J., Sulis, M., Kollet, S., Siebert, S., & Wada, Y. (n.d.). Human Water Use Impacts on the Strength of the  
1054 Continental Sink for Atmospheric Water. *Geophysical Research Letters*, 45(9), 4068–4076.  
1055 <https://doi.org/10.1029/2018GL077621>
- 1056 Keune, J., F. Gasper, K. Goergen, A. Hense, P. Shrestha, M. Sulis, and S. Kollet, 2016, Studying the influence of  
1057 groundwater representations on land surface-atmosphere feedbacks during the European heat wave in 2003, *J.*  
1058 *Geophys. Res. Atmos.*, 121, 13, 301–13,325, doi:10.1002/2016JD025426. doi:10.1002/2016JD025426.
- 1059 Knowling, M. J., and A. D. Werner (2016), Estimability of recharge through groundwater model calibration: Insights  
1060 from a field-scale steady-state example, *Journal of Hydrology*, 540, 973–987.
- 1061 Koirala et al. (2013) Global-scale land surface hydrologic modeling with the representation of water table  
1062 dynamics, *JGR Atmospheres* <https://doi.org/10.1002/2013JD020398>
- 1063 Koirala, S., Kim, H., Hirabayashi, Y., Kanae, S. and Oki, T. (2019) Sensitivity of Global Hydrological Simulations to  
1064 Groundwater Capillary Flux Parameterizations, *Water Resour. Res.*, 55(1), 402–425, doi:10.1029/2018WR023434,
- 1065 Kollet, S. J., & Maxwell, R. M. (2008). Capturing the influence of groundwater dynamics on land surface processes  
1066 using an integrated, distributed watershed model. *Water Resources Research*, 44(2).



- 1067 Kollet, S., Sulis, M., Maxwell, R. M., Paniconi, C., Putti, M., Bertoldi, G., et al. (2017). The integrated hydrologic  
1068 model intercomparison project, IH-MIP2: A second set of benchmark results to diagnose integrated hydrology and  
1069 feedbacks. *Water Resources Research*, 53(1), 867–890.
- 1070 Konikow, L. F. (2011), Contribution of global groundwater depletion since 1900 to sea-level rise, *Geophys. Res.*  
1071 *Lett.*, 38, L17401, doi: 10.1029/2011GL048604.
- 1072 Koster, R.D., Suarez, M.J., Ducharne, A., Praveen, K., & Stieglitz, M. (2000). A catchment-based approach to  
1073 modeling land surface processes in a GCM - Part 1: Model structure. *Journal of Geophysical Research*, 105 (D20),  
1074 24809-24822.
- 1075 Konikow, L.F. (2011) Contribution of global groundwater depletion since 1900 to sea-level rise. *Geophysical*  
1076 *Research Letters* <https://doi.org/10.1029/2011GL048604>
- 1077 Krakauer, N. Y., Li, H., & Fan, Y. (2014). Groundwater flow across spatial scales: importance for climate modeling.  
1078 *Environmental Research Letters*, 9(3), 034003.
- 1079 Kresic, N. (2009). *Groundwater resources: sustainability, management and restoration*. McGraw-Hill.
- 1080 Krueger, T., T. Page, K. Hubacek, L. Smith, and K. Hiscock (2012), The role of expert opinion in environmental  
1081 modelling, *Environmental Modelling & Software*, 36, 4-18.
- 1082
- 1083 Lamb, R., Aspinall, W., Odbert, H. and Wagener, T. (2017). Vulnerability of bridges to scour: Insights from an  
1084 international expert elicitation workshop. *Natural Hazards and Earth System Sciences*. 17(8), 1393-1409.
- 1085 Leaf, A. T., Fienen, M. N., Hunt, R. J., & Buchwald, C. A. (2015). Groundwater/surface-water interactions in the Bad  
1086 River Watershed, Wisconsin. USGS Numbered Series No. 2015–5162. Reston, VA: U.S. Geological Survey.
- 1087 Leavesley, G. H., S. L. Markstrom, P. J. Restrepo, and R. J. Viger (2002), A modular approach for addressing model  
1088 design, scale, and parameter estimation issues in distributed hydrological modeling, *Hydrol. Processes*, 16, 173–  
1089 187, doi:10.1002/hyp.344.
- 1090 Lemieux, J. M., Sudicky, E. A., Peltier, W. R., & Tarasov, L. (2008). Dynamics of groundwater recharge and seepage  
1091 over the Canadian landscape during the Wisconsinian glaciation. *J. Geophys. Res.*, 113. Retrieved from  
1092 <http://dx.doi.org/10.1029/2007JF000838>
- 1093 Lenton, T.M. et al. (2008). Tipping elements in the Earth's climate system. *Proceedings of the National Academy of*  
1094 *Sciences* 105 (6), 1786-1793.
- 1095 Liang, X., Z. Xie, and M. Huang (2003). A new parameterization for surface and groundwater interactions and its  
1096 impact on water budgets with the variable infiltration capacity (VIC) land surface model, *J. Geophys. Res.*, 108,  
1097 8613, D16. <https://doi.org/10.1029/2002JD003090>
- 1098 Lo, M.-H., Famiglietti, J. S., Reager, J. T., Rodell, M., Swenson, S., & Wu, W.-Y. (2016). GRACE-Based Estimates of  
1099 Global Groundwater Depletion. In Q. Tang & T. Oki (Eds.), *Terrestrial Water Cycle and Climate Change* (pp. 135–  
1100 146). John Wiley & Sons, Inc. <https://doi.org/10.1002/9781118971772.ch7>
- 1101 Lo, M.-H., Yeh, P. J.-F., & Famiglietti, J. S. (2008). Constraining water table depth simulations in a land surface  
1102 model using estimated baseflow. *Advances in Water Resources*, 31(12), 1552–1564.



- 1103 Lo, M. and J. S. Famiglietti, (2010) Effect of water table dynamics on land surface hydrologic memory, *J. Geophys.*  
1104 *Res.*, 115, D22118, doi:10.1029/2010JD014191
- 1105 Lo, M.-H., J. S. Famiglietti, P. J.-F. Yeh, and T. H. Syed (2010), Improving Parameter Estimation and Water Table  
1106 Depth Simulation in a Land Surface Model Using GRACE Water Storage and Estimated Baseflow Data, *Water*  
1107 *Resour. Res.*, 46, W05517, doi:10.1029/2009WR007855.
- 1108 Loheide, S. P., Butler Jr, J. J., & Gorelick, S. M. (2005). Estimation of groundwater consumption by phreatophytes  
1109 using diurnal water table fluctuations: A saturated-unsaturated flow assessment. *Water Resources Research*, 41(7).
- 1110 Luijendijk, E., Gleeson, T. and Moosdorf, N. (2020) Fresh groundwater discharge insignificant for the world's oceans  
1111 but important for coastal ecosystems *Nature Communications*, 11, 1260 (2020). doi: 10.1038/s41467-020-15064-8  
1112
- 1113 Maples, S., Foglia, L., Fogg, G.E. and Maxwell, R.M. (2020). Sensitivity of Hydrologic and Geologic Parameters on  
1114 Recharge Processes in a Highly-Heterogeneous, Semi-Confined Aquifer System. *Hydrology and Earth Systems*  
1115 *Sciences*, in press.  
1116
- 1117 Margat, J., & Van der Gun, J. (2013). *Groundwater around the world: a geographic synopsis*. London: CRC Press
- 1118 Maxwell, R. M., Condon, L. E., and Kollet, S. J. (2015) A high-resolution simulation of groundwater and surface  
1119 water over most of the continental US with the integrated hydrologic model ParFlow v3, *Geosci. Model Dev.*, 8,  
1120 923–937, <https://doi.org/10.5194/gmd-8-923-2015>.
- 1121 Maxwell, R.M., Chow, F.K. and Kollet, S.J., The groundwater-land-surface-atmosphere connection: soil moisture  
1122 effects on the atmospheric boundary layer in fully-coupled simulations. *Advances in Water Resources* 30(12),  
1123 doi:10.1016/j.advwatres.2007.05.018, 2007.
- 1124 Maxwell, R. M., & Condon, L. E. (2016). Connections between groundwater flow and transpiration partitioning.  
1125 *Science*, 353(6297), 377–380.
- 1126 Maxwell, R. M., Condon, L. E., Kollet, S. J., Maher, K., Haggerty, R., & Forrester, M. M. (2016). The imprint of  
1127 climate and geology on the residence times of groundwater. *Geophysical Research Letters*, 43(2), 701–708.  
1128 <https://doi.org/10.1002/2015GL066916>
- 1129 McMilan, H. (2020). Linking hydrologic signatures to hydrologic processes: A review. *Hydrological Processes*. 34,  
1130 1393– 1409.
- 1131 Meixner, T., Manning, A. H., Stonestrom, D. A., Allen, D. M., Ajami, H., Blasch, K. W., et al. (2016). Implications of  
1132 projected climate change for groundwater recharge in the western United States. *Journal of Hydrology*, 534, 124–  
1133 138.
- 1134 Melsen, L. A., A. J. Teuling, P. J. J. F. Torfs, R. Uijlenhoet, N. Mizukami, and M. P. Clark, 2016a: HESS Opinions: The  
1135 need for process-based evaluation of large-domain hyper-resolution models. *Hydrology and Earth System*  
1136 *Sciences*, doi:10.5194/hess-20-1069-2016.
- 1137 Meriano, M., & Eyles, N. (2003). Groundwater flow through Pleistocene glacial deposits in the rapidly urbanizing  
1138 Rouge River-Highland Creek watershed, City of Scarborough, southern Ontario, Canada. *Hydrogeology Journal*,  
1139 11(2), 288–303. <https://doi.org/10.1007/s10040-002-0226-4>



- 1140 Milly, P.C., S.L. Malyshev, E. Shevliakova, K.A. Dunne, K.L. Findell, T. Gleeson, Z. Liang, P. Philipps, R.J. Stouffer, & S.  
1141 Swenson (2014). An Enhanced Model of Land Water and Energy for Global Hydrologic and Earth-System Studies. *J.*  
1142 *Hydrometeor.*, 15, 1739–1761. <https://doi.org/10.1175/JHM-D-13-0162.1>
- 1143 Minderhoud P S J, Erkens G, Pham Van H, Bui Tran V, Erban L E, Kooi, H and Stouthamer E (2017) Impacts of 25  
1144 years of groundwater extraction on subsidence in the Mekong delta, Vietnam *Environ. Res. Lett.* 12 064006
- 1145 Minderhoud, P.S.J., Coumou, L., Erkens, G., Middelkoop, H. & Stouthamer, E. (2019). Mekong delta much lower  
1146 than previously assumed in sea-level rise impact assessments. *Nature Communications* 10, 3847.
- 1147 Minderhoud, P.S.J., Middelkoop, H., Erkens, G. and Stouthamer, E. Groundwater (2020). extraction may drown  
1148 mega-delta: projections of extraction-induced subsidence and elevation of the Mekong delta for the 21st century.  
1149 *Environ. Res. Commun.* 2, 011005.
- 1150 Miralles, D. G., Jimenez, C., Jung, M., Michel, D., Ershadi, A., McCabe, M. F., et al. (2016). The WACMOS-ET project -  
1151 Part 2: Evaluation of global terrestrial evaporation data sets. *Hydrology and Earth System Sciences*, 20(2), 823-842.  
1152 doi:10.5194/hess-20-823-2016.
- 1153 Moeck, C. Nicolas Grech-Cumbo, Joel Podgorski, Anja Bretzler, Jason J. Gurdak ,Michael Berg, Mario Schirmer  
1154 (2020) A global-scale dataset of direct natural groundwater recharge rates: A review of variables, processes and  
1155 relationships. *Science of The Total Environment* <https://doi.org/10.1016/j.scitotenv.2020.137042>
- 1156 Mohan, C., Wei, Y., & Saft, M. (2018). Predicting groundwater recharge for varying land cover and climate  
1157 conditions—a global meta-study. *Hydrology and Earth System Sciences*, 22(5), 2689–2703.
- 1158 Montanari, A., Young, G., Savenije, H.H.G., Hughes, D., Wagener, T., Ren, L.L., Koutsoyiannis, D., Cudennec, C.,  
1159 Toth, E., Grimaldi, S., et al. (2013). “Panta Rhei—Everything Flows”: Change in hydrology and society—The IAHS  
1160 Scientific Decade 2013–2022. *Hydrological Sciences Journal* 58, 1256–1275.
- 1161 Moore, W. S. (2010). The effect of submarine groundwater discharge on the ocean. *Annual Review of Marine*  
1162 *Science*, 2, 59–88.
- 1163 Morris, M. D. (1991). Factorial sampling plans for preliminary computational experiments. *Technometrics*, 33(2),  
1164 161–174.
- 1165 Müller Schmied, H., Eisner, S., Franz, D., Wattenbach, M., Portmann, F.T., Flörke, M., Döll, P. (2014):  
1166 Sensitivity of simulated global-scale freshwater fluxes and storages to input data, hydrological model  
1167 structure, human water use and calibration. *Hydrol. Earth Syst. Sci.*, 18, 3511-3538, doi: 10.5194/hess-  
1168 18-3511-2014.
- 1169 Niu, G.-Y., Z.-L. Yang, R. E. Dickinson, and L. E. Gulden (2005), A simple TOPMODEL-based runoff parameterization  
1170 (SIMTOP) for use in global climate models. *J. Geophys. Res.*, 110, D21106, doi:10.1029/2005JD006111
- 1171 Niu GY, Yang ZL, Dickinson RE, Gulden LE, Su H (2007) Development of a simple groundwater model for use in  
1172 climate models and evaluation with Gravity Recovery and Climate Experiment data. *J Geophys Res* 112:D07103.  
1173 doi:10.1029/2006JD007522



- 1174 Ngo-Duc, T., Laval, K. Ramillien, G., Polcher, J. & Cazenave, A. (2007). Validation of the land water storage  
1175 simulated by Organising Carbon and Hydrology in Dynamic Ecosystems (ORCHIDEE) with Gravity Recovery and  
1176 Climate Experiment (GRACE) data. *Water Resour. Res.*, 43, W04427. <https://doi.org/10.1029/2006WR004941>
- 1177 O’Hagan, A. (2019). Expert Knowledge Elicitation: Subjective but Scientific. *The American Statistician*, 73,  
1178 [doi.org/10.1080/00031305.2018.1518265](https://doi.org/10.1080/00031305.2018.1518265)
- 1179 Ortega-Guerrero A, Rudolph D L and Cherry J A 1999 Analysis of long-term land subsidence near Mexico City: field  
1180 investigations and predictive modeling *Water Resour. Res.* 353327–41
- 1181 Pan, M., Sahoo, A. K., Troy, T. J., Vinukollu, R. K., Sheffield, J., & Wood, F. E. (2012). Multisource estimation of long-  
1182 term terrestrial water budget for major global river basins. *J. Climate*, 25, 3191–3206.  
1183 <https://doi.org/10.1175/JCLI-D-11-00300.1>
- 1184 Pappenberger, F., Ghelli, A., Buizza, R. and Bodis, K. (2009). The Skill of Probabilistic Precipitation Forecasts under  
1185 Observational Uncertainties within the Generalized Likelihood Uncertainty Estimation Framework for Hydrological  
1186 Applications. *Journal of Hydrometeorology*, DOI: 10.1175/2008JHM956.1
- 1188 Pellet, V., Aires, F., Munier, S., Fernández Prieto, D., Jordá, G., Dorigo, W. A., Polcher, J., & Brocca, L. (2019).  
1189 Integrating multiple satellite observations into a coherent dataset to monitor the full water cycle – application to  
1190 the Mediterranean region. *Hydrol. Earth Syst. Sci.*, 23, 465-491. <https://doi.org/10.5194/hess-23-465-2019>
- 1191 Perrone, D. and Jasechko (2019). Deeper well drilling an unsustainable stopgap to groundwater depletion. *Nature*  
1192 *Sustain.* 2, 773-782.
- 1193 Person, M. A., Raffensperger, J. P., Ge, S., & Garven, G. (1996). Basin-scale hydrogeologic modeling. *Reviews of*  
1194 *Geophysics*, 34(1), 61–87.
- 1195 Pianosi, F., Beven, K., Freer, J., Hall, J. W., Rougier, J., Stephenson, D. B., & Wagener, T. (2016). Sensitivity analysis  
1196 of environmental models: A systematic review with practical workflow. *Environmental Modelling & Software*, 79,  
1197 214–232.
- 1198 Post, V. E., & von Asmuth, J. R. (2013). Hydraulic head measurements—new technologies, classic pitfalls.  
1199 *Hydrogeology Journal*, 21(4), 737–750.
- 1200 Qiu J. Q., Zipper, S.C., Motew M., Booth, E.G., Kucharik, C.J., & Loheide, S.P. (2019). Nonlinear groundwater  
1201 influence on biophysical indicators of ecosystem services. *Nature Sustainability*, in press, doi: 10.1038/s41893-019-  
1202 0278-2
- 1203 Rajabi, M. M., and B. Ataie-Ashtiani (2016), Efficient fuzzy Bayesian inference algorithms for incorporating expert  
1204 knowledge in parameter estimation, *Journal of Hydrology*, 536, 255-272.
- 1205 Rajabi, M. M., B. Ataie-Ashtiani, and C. T. Simmons (2018), Model-data interaction in groundwater studies: Review  
1206 of methods, applications and future directions, *Journal of Hydrology*, 567, 457-477.
- 1207 Rashid, M., Chien, R.Y., Ducharne, A., Kim, H., Yeh, P.J.F., Peugeot, C., Boone, A., He, X., Séguis, L., Yabu, Y., Boukari,  
1208 M. & Lo, M.H. (2019). Evaluation of groundwater simulations in Benin from the ALMIP2 project. *J. Hydromet.*,  
1209 accepted.



- 1213 Refsgaard, J.C., van der Sluijs, J.P., Højberg, A.L., and Vanrolleghem, P.A. (2007). Uncertainty in the environmental  
1214 modelling process—a framework and guidance. *Environmental Modelling & Software*, 22(11), 1543-1556
- 1215 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais, N., & Prabhat. (2019). Deep learning  
1216 and process understanding for data-driven Earth system science. *Nature*, 566(7743), 195–204.
- 1217 Reinecke, R., Foglia, L., Mehl, S., Trautmann, T., Cáceres, D., & Döll, P. (2019a). Challenges in developing a global  
1218 gradient-based groundwater model. (G<sup>3</sup>M v1.0) for the integration into a global hydrological model. *Geosci. Model*  
1219 *Dev.*, 12, 2401-2418. doi: 10.5194/gmd-12-2401-2019
- 1220 Reinecke, R., Foglia, L., Mehl, S., Herman, J., Wachholz, A., Trautmann, T., and Döll, P. (2019b) Spatially distributed  
1221 sensitivity of simulated global groundwater heads and flows to hydraulic conductivity, groundwater recharge and  
1222 surface water body parameterization, *Hydrology and Earth System Sciences*, (23) 4561–4582. 2019.
- 1223 Reinecke, R., Wachholz, A., Mehl, S., Foglia, L., Niemann, C., Döll, P. (2020). Importance of spatial resolution in  
1224 global groundwater modeling. *Groundwater*. doi: 10.1111/gwat.12996
- 1225 Rodell, M., Velicogna, I., & Famiglietti, J. S. (2009). Satellite-based estimates of groundwater depletion in India.  
1226 *Nature*, 460(7258), 999–1002.
- 1227 Rodell, M., Famiglietti, J. S., Wiese, D. N., Reager, J. T., Beaudoin, H. K., Landerer, F. W., & Lo, M.-H. (2018).  
1228 Emerging trends in global freshwater availability. *Nature*, 557(7707), 651.
- 1229 Rosolem, R., Hoar, T., Arellano, A., Anderson, J. L., Shuttleworth, W. J., Zeng, X., and Franz, T. E.: Translating  
1230 aboveground cosmic-ray neutron intensity to high-frequency soil moisture profiles at sub-kilometer scale, *Hydrol.*  
1231 *Earth Syst. Sci.*, 18, 4363-4379
- 1232 Ross, J. L., M. M. Ozbek, and G. F. Pinder (2009), Aleatoric and epistemic uncertainty in groundwater flow and  
1233 transport simulation, *Water Resources Research*, 45(12).  
1234
- 1235 Rossman, N., & Zlotnik, V. (2013). Review: Regional groundwater flow modeling in heavily irrigated basins of  
1236 selected states in the western United States. *Hydrogeology Journal*, 21(6), 1173–1192.  
1237 <https://doi.org/10.1007/s10040-013-1010-3>
- 1238 RRCA. (2003). Republican River Compact Administration Ground Water Model. Retrieved from  
1239 <http://www.republicanrivercompact.org/>
- 1240 Saltelli, A., Chan, K., & Scott, E. M. (Eds.). (2000). *Sensitivity analysis*. Wiley.
- 1241 Salvucci, G. D., & Entekhabi, D. (1995). Hillslope and climatic controls on hydrologic fluxes. *Water Resources*  
1242 *Research*, 31(7), 1725–1739.
- 1243 Sawyer, A. H., David, C. H., & Famiglietti, J. S. (2016). Continental patterns of submarine groundwater discharge  
1244 reveal coastal vulnerabilities. *Science*, 353(6300), 705–707.
- 1245 Scanlon, B., Healy, R., & Cook, P. (2002). Choosing appropriate techniques for quantifying groundwater recharge.  
1246 *Hydrogeology Journal*, 10(1), 18–39.
- 1247 Scanlon, B. R., Keese, K. E., Flint, A. L., Flint, L. E., Gaye, C. B., Edmunds, W. M., & Simmers, I. (2006). Global  
1248 synthesis of groundwater recharge in semiarid and arid regions. *Hydrological Processes*, 20, 3335–3370.



- 1249 Scanlon, B. R., Faunt, C. C., Longuevergne, L., Reedy, R. C., Alley, W. M., McGuire, V. L., & McMahon, P. B. (2012).  
1250 Groundwater depletion and sustainability of irrigation in the US High Plains and Central Valley. *Proceedings of the*  
1251 *National Academy of Sciences*, 109(24), 9320–9325. <https://doi.org/10.1073/pnas.1200311109>
- 1252 Scanlon, B. R., Zhang, Z., Save, H., Wiese, D. N., Landerer, F. W., Long, D., et al. (2016). Global evaluation of new  
1253 GRACE mascon products for hydrologic applications. *Water Resources Research*, 52(12), 9412–9429.
- 1254 Scanlon, B. R., Zhang, Z., Save, H., Sun, A. Y., Müller Schmied, H., van Beek, L. P., et al. (2018). Global models  
1255 underestimate large decadal declining and rising water storage trends relative to GRACE satellite data. *Proceedings*  
1256 *of the National Academy of Sciences*, 201704665.
- 1257 Schaller, M., and Y. Fan (2009) River basins as groundwater exporters and importers: Implications for water cycle  
1258 and climate modeling. *Journal of Geophysical Research-Atm*, 114, D04103, doi: 10.1029/2008 JD010636
- 1259 Schewe, J., Heinke, J., Gerten, D., Haddeland, I., Arnell, N. W., Clark, D. B., et al. (2014). Multimodel assessment of  
1260 water scarcity under climate change. *Proceedings of the National Academy of Sciences*, 111(9), 3245–3250.  
1261 <https://doi.org/10.1073/pnas.1222460110>
- 1262 Schilling, O. S., Doherty, J., Kinzelbach, W., Wang, H., Yang, P. N., & Brunner, P. (2014). Using tree ring data as a  
1263 proxy for transpiration to reduce predictive uncertainty of a model simulating groundwater–surface water–  
1264 vegetation interactions. *Journal of Hydrology*, 519, Part B, 2258–2271.  
1265 <https://doi.org/10.1016/j.jhydrol.2014.08.063>
- 1266 Schilling, O.S., Cook, P.G., Brunner, P., 2019. Beyond classical observations in hydrogeology: The advantages of  
1267 including exchange flux, temperature, tracer concentration, residence time, and soil moisture observations in  
1268 groundwater model calibration. *Reviews of Geophysics*, 57(1): 146-182.
- 1269 Schneider, A.S., Jost, A., Coulon, C., Silvestre, M., Théry, S., & Ducharne, A. (2017). Global scale river network  
1270 extraction based on high-resolution topography, constrained by lithology, climate, slope, and observed drainage  
1271 density. *Geophysical Research Letters*, 44, 2773–2781. <https://doi.org/10.1002/2016GL071844>
- 1272 Shen, C. (2018). A transdisciplinary review of deep learning research and its relevance for water resources  
1273 scientists. *Water Resources Research*, 54(11), 8558–8593.
- 1274 Shen, C., Laloy, E., Elshorbagy, A., Albert, A., Bales, J., Chang, F.-J., et al. (2018). HESS Opinions: Incubating deep-  
1275 learning-powered hydrologic science advances as a community. *Hydrology and Earth System Sciences*, 22(11).
- 1276 Springer, A., & Stevens, L. (2009). Spheres of discharge of springs. *Hydrogeology Journal*, 17(1), 83–93.  
1277 <https://doi.org/10.1007/s10040-008-0341-y>
- 1278 Steffen, W., Broadgate, W., Deutsch, L., Gaffney, O., & Ludwig, C. (2015). The trajectory of the Anthropocene: the  
1279 great acceleration. *The Anthropocene Review*, 2(1), 81–98.
- 1280 Sutanudjaja, E. H., Beek, R. van, Wanders, N., Wada, Y., Bosmans, J. H., Drost, N., et al. (2018). PCR-GLOBWB 2: a 5  
1281 arcmin global hydrological and water resources model. *Geoscientific Model Development*, 11(6), 2429–2453.
- 1282 Takata, K., Emori, S. and Watanabe, T.: Development of the minimal advanced treatments of surface interaction  
1283 and runoff, *Glob. Planet. Change*, 38(1–2), 209–222, doi:10.1016/S0921-8181(03)00030-4, 2003.



- 1284 Tallaksen, L. M. (1995). A review of baseflow recession analysis. *Journal of Hydrology*, 165(1–4), 349–370.  
1285 [https://doi.org/10.1016/0022-1694\(94\)02540-R](https://doi.org/10.1016/0022-1694(94)02540-R)
- 1286 Taylor, R. G., Todd, M. C., Kongola, L., Maurice, L., Nahozya, E., Sanga, H., & MacDonald, A. M. (2013). Evidence of  
1287 the dependence of groundwater resources on extreme rainfall in East Africa. *Nature Clim. Change*, 3(4), 374–378.  
1288 <https://doi.org/10.1038/nclimate1731>
- 1289 Taylor, R. G., Scanlon, B., Doll, P., Rodell, M., van Beek, R., Wada, Y., et al. (2013). Groundwater and climate  
1290 change. *Nature Clim. Change*, 3(4), 322–329. <https://doi.org/10.1038/nclimate1744>
- 1291 Thatch, L. M., Gilbert, J. M., & Maxwell, R. M. (2020). Integrated hydrologic modeling to untangle the impacts of  
1292 water management during drought. *Groundwater*, 58(3), 377–391.
- 1293 Thomas, Z., Rousseau-Gueutin, P., Kolbe, T., Abbott, B.W., Marçais, J., Peiffer, S., Frei, S., Bishop, K., Pichelin, P.,  
1294 Pinay, G., de Dreuzy, J.R. (2016). Constitution of a catchment virtual observatory for sharing flow and transport  
1295 models outputs. *Journal of Hydrology*, 543, Pages 59-66. <https://doi.org/10.1016/j.jhydrol.2016.04.067>
- 1296 Tolley, D., Foglia, L., & Harter, T. (2019). Sensitivity Analysis and Calibration of an Integrated Hydrologic Model in  
1297 an Irrigated Agricultural Basin with a Groundwater-Dependent Ecosystem. *Water Resources Research*.  
1298 <https://doi.org/10.1029/2018WR024209>
- 1299 Tóth, J. (1963). A theoretical analysis of groundwater flow in small drainage basins. *Journal of Geophysical*  
1300 *Research*, 68(16), 4795–4812.
- 1301 Tran, H., Jun Zhang, Jean-Martial Cohard, Laura E. Condon, Reed M. Maxwell (2020) Simulating groundwater-  
1302 Streamflow Connections in the Upper Colorado River Basin Groundwater, 2020  
1303 <https://doi.org/10.1111/gwat.13000>
- 1304 Tregoning, P., McClusky, S., van Dijk, A.I.J.M. and Crosbie, R.S. (2012). Assessment of GRACE satellites for  
1305 groundwater estimation in Australia. *Waterlines Report Series No 71*, National Water Commission, Canberra
- 1306 Tustison, B., Harris, D. and Foufoula-Georgiou, E. (2001). Scale issues in verification of precipitation  
1307 forecasts. *Journal of geophysical Research*, 106(D11), 11775–11784.
- 1308 UNESCO. (1978). *World water balance and water resources of the earth* (Vol. USSR committee for the international  
1309 hydrologic decade). Paris: UNESCO.
- 1310 Van Werkhoven, K., Wagener, T., Tang, Y., and Reed, P. 2008. Understanding watershed model behavior across  
1311 hydro-climatic gradients using global sensitivity analysis. *Water Resources Research*, 44, W01429,  
1312 doi:10.1029/2007WR006271.
- 1313 Van Loon, A.F. et al. (2016) [Drought in the Anthropocene](#). *Nature Geoscience* 9: 89-91 doi: 10.1038/ngeo2646.
- 1314 van Loon, Anne F.; Kumar, Rohini; Mishra, Vimal (2017): Testing the use of standardised indices and GRACE  
1315 satellite data to estimate the European 2015 groundwater drought in near-real time. In *Hydrol. Earth Syst. Sci.* 21  
1316 (4), pp. 1947–1971. DOI: 10.5194/hess-21-1947-2017.
- 1317 Vergnes, J.-P., & Decharme, B. (2012). A simple groundwater scheme in the TRIP river routing model: global off-line  
1318 evaluation against GRACE terrestrial water storage estimates and observed river discharges. *Hydrol. Earth Syst.*  
1319 *Sci.*, 16, 3889–3908. <https://doi.org/10.5194/hess-16-3889-2012>



- 1320 Vergnes, J.-P., B. Decharme, & F. Habets (2014). Introduction of groundwater capillary rises using subgrid spatial  
1321 variability of topography into the ISBA land surface model, *J. Geophys. Res. Atmos.*, 119, 11,065–11,086.  
1322 <https://doi.org/10.1002/2014JD021573>
- 1323 Visser, W. C. (1959). Crop growth and availability of moisture. *Journal of the Science of Food and Agriculture*, 10(1),  
1324 1–11.
- 1325 Wada, Y., L. P. H. van Beek, C. M. van Kempen, J. W. T. M. Reckman, S. Vasak, M. F. P. Bierkens, (2010) Global  
1326 depletion of groundwater resources. *Geophys. Res. Lett.* 37, L20402.
- 1327 Wada, Y.; Wisser, D.; Bierkens, M. F. P. (2014). Global modeling of withdrawal, allocation and consumptive use of  
1328 surface water and groundwater resources. *Earth System Dynamics Discussions*, volume 5, issue 1, pp. 15 - 40
- 1329 Wada, Y. (2016). Modeling Groundwater Depletion at Regional and Global Scales: Present State and Future  
1330 Prospects. *Surveys in Geophysics*, 37(2), 419–451. <https://doi.org/10.1007/s10712-015-9347-x>
- 1331 Wada, Y., & Heinrich, L. (2013). Assessment of transboundary aquifers of the world—vulnerability arising from  
1332 human water use. *Environmental Research Letters*, 8(2), 024003.
- 1333 Wagener, T. 2003. Evaluation of catchment models. *Hydrological Processes*, 17, 3375–3378.
- 1334 Wagener, T., & Gupta, H. V. (2005). Model identification for hydrological forecasting under uncertainty. *Stochastic*  
1335 *Environmental Research and Risk Assessment*, 19(6), 378–387.
- 1336 Wagener, T., Sivapalan, M., Troch, P. and Woods, R. (2007). Catchment classification and hydrologic similarity.  
1337 *Geography Compass*, 1(4), 901, doi:10.1111/j.1749-8198.2007.00039.x
- 1338 Wagener, T. and Pianosi, F. (2019) What has Global Sensitivity Analysis ever done for us? A systematic review to  
1339 support scientific advancement and to inform policy-making in earth system modelling. *Earth-Science Reviews*,  
1340 194, 1-18. [doi.org/10.1016/j.earscirev.2019.04.006](https://doi.org/10.1016/j.earscirev.2019.04.006)
- 1341 Wagener, T., Boyle, D.P., Lees, M.J., Wheeler, H.S., Gupta, H.V. and Sorooshian, S. (2001). A framework for  
1342 development and application of hydrological models. *Hydrology and Earth System Sciences*, 5(1), 13-26.
- 1343 Wagener, T., Sivapalan, M., Troch, P. A., McGlynn, B. L., Harman, C. J., Gupta, H. V., et al. (2010). The future of  
1344 hydrology: An evolving science for a changing world. *Water Resources Research*, 46(5).
- 1345 Wagener, T., Gleeson, T., et al. On doing large-scale hydrology with lions: perceptual models and knowledge  
1346 accumulation. submitted to *Water Wires and preprint*: <https://eartharxiv.org/zdy5n/>
- 1347 Wang, F., Ducharme, A., Cheruy, F., Lo, M.H., & Grandpeix, J.L. (2018). Impact of a shallow groundwater table on  
1348 the global water cycle in the IPSL land-atmosphere coupled model, *Climate Dynamics*, 50, 3505-3522,  
1349 <https://doi.org/10.1007/s00382-017-3820-9>
- 1350 Warszawski, L., Frieler, K., Huber, V., Piontek, F., Serdeczny, O., & Schewe, J. (2014). The Inter-Sectoral Impact  
1351 Model Intercomparison Project (ISI-MIP): Project framework. *Proceedings of the National Academy of Sciences*,  
1352 111(9), 3228–3232. <https://doi.org/10.1073/pnas.1312330110>
- 1353 Winter, T. C., Harvey, J. W., Franke, O. L., & Alley, W. M. (1998). *Ground water and surface water: a single resource*  
1354 (p. 79). U.S. Geological Survey circular 1139



- 1355 Woolfenden, L. R., & Nishikawa, T. (2014). Simulation of groundwater and surface-water resources of the Santa  
1356 Rosa Plain watershed, Sonoma County, California. USGS Scientific Investigations Report 2014–5052). Reston, VA:  
1357 U.S. Geological Survey.
- 1358 Yang, J., Griffiths, J., & Zammit, C. (2019). National classification of surface–groundwater interaction using random  
1359 forest machine learning technique. *River Research and Applications*, 35(7), 932–943.  
1360 <https://doi.org/10.1002/rra.3449>
- 1361 Yeh, P. J.-F. and J. Famiglietti, Regional groundwater evapotranspiration in Illinois, *J. Hydrometeorology*, 10(2),  
1362 464–478, 2010
- 1363 Yilmaz, K., Gupta, H.V. and Wagener, T. 2009. Towards improved distributed modeling of watersheds: A process  
1364 based diagnostic approach to model evaluation. *Water Resources Research*, 44, W09417,  
1365 doi:10.1029/2007WR006716.
- 1366 Young, P., Parkinson, S. and Lees, M. (1996). Simplicity out of complexity in environmental modelling: Occam's  
1367 razor revisited. *Journal of Applied Statistics*, 23(2-3), 165-210. <https://doi.org/10.1080/02664769624206>
- 1368 Zipper, S. C., Soylu, M. E., Booth, E. G., & Loheide, S. P. (2015). Untangling the effects of shallow groundwater and  
1369 soil texture as drivers of subfield-scale yield variability. *Water Resources Research*, 51(8), 6338–6358.
- 1370 Zipper, S. C., Soylu, M. E., Kucharik, C. J., & Loheide, S. P. (2017). Quantifying indirect groundwater-mediated  
1371 effects of urbanization on agroecosystem productivity using MODFLOW-AgroIBIS (MAGI), a complete critical zone  
1372 model. *Ecological Modelling*, 359, 201-219
- 1373 Zhang, M and Burbey T J 2016 Inverse modelling using PS-InSAR data for improved land subsidence simulation in  
1374 Las Vegas Valley, Nevada *Hydrol. Process.* 30 4494–516
- 1375 Zhou, Y., Li, W., 2011. A review of regional groundwater flow modeling. *Geoscience Frontiers*, 2(2): 205-214.