

1 **Seasonal streamflow forecasts for Europe – I. Hindcast verification**  
2 **with pseudo- and real observations**

3

4 Wouter Greuell<sup>1)</sup>, Wietse H. P. Franssen<sup>1)</sup>, Hester Biemans<sup>2)</sup> and Ronald W. A.  
5 Hutjes<sup>1,2)</sup>

6

7 1) Water Systems and Global Change, Wageningen University,  
8 Droevendaalsesteeg 3, NL 6708 PB Wageningen, Netherlands

9 2) Water and Food, Wageningen Environmental Research, Droevendaalsesteeg 3,  
10 NL 6708 PB Wageningen, Netherlands

11

12 correspondence to [ronald.hutjes@wur.nl](mailto:ronald.hutjes@wur.nl)

13

14

15 **Abstract**

16

17 Seasonal predictions can be exploited among others to optimize hydropower energy  
18 generation, navigability of rivers and irrigation management to decrease crop yield  
19 losses. This paper is the first of two papers dealing with a model-based system built to  
20 produce seasonal hydrological forecasts (WUSHP: Wageningen University Seamless  
21 Hydrological Prediction system), applied here to Europe. The paper presents the  
22 development and the skill evaluation of the system. In WUSHP hydrology is simulated  
23 by running the Variable Infiltration Capacity (VIC) hydrological model with forcing  
24 from bias-corrected output of ECMWF's Seasonal Forecasting System 4. The system is  
25 probabilistic. For the assessment of skill, we performed hindcast simulations (1981-  
26 2010) and a reference simulation, in which VIC was forced by gridded meteorological  
27 observations, to generate initial hydrological conditions for the hindcasts and discharge  
28 output for skill assessment (pseudo-observations).

29 Skill in hindcasting runoff and discharge is analysed with monthly temporal resolution,  
30 up to 7 months of lead time, for the entire annual cycle. Using the pseudo-observations  
31 and taking the correlation coefficient as metric, hot spots of significant skill in runoff  
32 were identified in Fennoscandia (from January to October), the southern part of the  
33 Mediterranean (from June to August), Poland, northern Germany, Romania and  
34 Bulgaria (mainly from November to January) and western France (from December to  
35 May). Generally skill decreases with increasing lead time, except in spring in regions  
36 with snow-rich winters. In some areas some skill persists even at the longest lead times  
37 (7 months). On average across the domain, skill in discharge is slightly higher than skill  
38 in runoff. This can be explained by the delay between runoff and discharge and the  
39 general tendency of decreasing skill with lead time.

40 Theoretical skill as determined with the pseudo-observations was compared to actual  
41 skill as determined with real discharge observations from 747 stations. ctual skill is  
42 mostly and often substantially less than theoretical skill. This effect is stronger for small  
43 than for large basins, which is consistent with a conceptual analysis of the structural  
44 differences between the two types of verification. Qualitatively, the use of different skill  
45 metrics (correlation coefficient, ROC area and Ranked Probability Skill Score) lead to  
46 broadly similar spatio-temporal patterns of skill, but the level of skill decreases, and the  
47 area of skill shrinks, in the order correlation coefficient, ROC area below normal tercile,  
48 ROC area above normal tercile, Ranked Probability Skill Score and finally, ROC near  
49 normal tercile.

50

51

52

53 **1 Introduction**

54

55 Society may benefit from seasonal hydrological forecasts, i.e. hydrological forecasts for  
56 future time periods from more than two weeks up to about a year (Doblas-Reyes et al.,  
57 2013). Such predictions can e.g. be exploited to optimize hydropower energy generation  
58 (Hamlet et al. 2002), navigability of rivers in low flow conditions (Li, et al., 2008) and  
59 irrigation management (Mushtaq et al. 2012; Ghile and Schulze 2008) to decrease crop  
60 yield losses. In order to be of any value in decision making processes of such sectors,  
61 forecasts must be credible, i.e. be skilful in predicting anomalous system states, as well  
62 as being relevant and legitimate to the decision making process (e.g. Bruno Soares and  
63 Dessai, 2016). In this paper we will introduce WUSHP (Wageningen University  
64 Seamless Hydrological Prediction system), a dynamical, model-based system (see Yuan  
65 et al., 2015) that was built around the Variable Infiltration Capacity (VIC) hydrological  
66 model and ECMWF's Seasonal Forecast System 4, to produce seasonal hydrological  
67 forecasts. It will be applied to Europe. The usefulness of the system depends partially  
68 on the level of its skill and the paper will therefore focus on an extensive assessment of  
69 the skill of WUSHP. The usual method of assessing skill of predictive systems is by  
70 analysing hindcasts, a strategy that will be adopted here as well.

71

72 During recent years, a number of systems for seasonal hydrological forecasts have been  
73 developed. Examples are the forecasting model suite he University of Washington's  
74 Surface Water Monitor (SWM; Wood and Lettenmaier, 2006) and the African Drought  
75 Monitor (Sheffield et al., 2014). Seasonal hydrological forecast systems for the entire  
76 continent of Europe are scarce (Bierkens and van Beek, 2009; Thober et al., 2015), but  
77 a few more concentrate on smaller domains such as the British Isles (Svensson et al.,  
78 2015), Iberia (Trigo, 2004) or France (Céron et al., 2010; Singla et al., 2012).

79

80 Thober et al. (2015) forced a mesoscale hydrological model (mHM) with meteorological  
81 hindcasts of the North American Multi-Model Ensemble (NMME) to investigate the  
82 predictability of soil moisture in continental Europe, excluding the British Isles and  
83 Fennoscandia. Evaluating a number of forecasting techniques that produced distinct  
84 variations in the magnitude of skill, they found that spatial patterns in skill were  
85 remarkably similar among the different techniques, as well as comparable to the spatial  
86 patterns of the autocorrelation (persistence) of reference soil moisture. High skill was  
87 found in eastern Germany and Poland, Romania, southern Balkans and eastern Ukraine  
88 as well as north-western France. Less skill was found in the mountainous areas of Alps  
89 and Pyrenees, the northern Adriatic and Atlantic Iberia. Most skill was found for winter  
90 months (DJF), least for autumn (SON), this minimum shifting to summer (JJA) at long  
91 lead times (6 months).

92

93 Bierkens and van Beek (2009) developed an analogue events method to select annual  
94 ERA40 meteorological forcing on the basis of annual SST anomalies in the northern  
95 Atlantic and then made hydrological forecasts with a global-scale hydrological model  
96 applied to Europe. Evaluating only winter and summer half year aggregated skill, they

97 found wintertime skill in large parts of Europe with maxima in eastern Spain and a zone  
98 from the southern Balkans and Romania through eastern Poland and western Russia to  
99 the Baltic states and Finland. Summertime skill was lower, generally by about 50% and  
100 even more around the Alps and the Adriatic. NAO based climate forecast added  
101 significant skill only in limited areas, such as Scandinavia, the Iberian Peninsula, the  
102 Balkans, and around the Black Sea.

103  
104 Svensson et al. (2015) found skilful winter river flow forecasts across the whole of the  
105 UK due to a combination of skilful winter rainfall forecasts for the north and west, and  
106 strong persistence of initial hydrological conditions in the south and east. Strong  
107 statistical correlations between NAO and winter precipitation in Iberia lead to skilful  
108 forecasts of JFM river flow and hydropower production (Trigo et al., 2004). Ceron et  
109 al. (2010) and Singla et al. (2012) set up a high resolution river flow forecasting system  
110 (8 km) over France, for which seasonal climate forecast improved MAM skill over  
111 northern France, but worsened it over southern France (compared to a river flow model  
112 with proper initialisation of soil moisture, snow etc., but random atmospheric forcing).  
113 Demirel et al. (2015) found that both two physical models and one neural network over-  
114 predict runoff during low-flow periods using ensemble seasonal meteorological forcing  
115 for the Moselle basin. As a result forecasts of more extreme low flows are less reliable  
116 than forecasts of more moderate ones.

117  
118 It is quite common in seasonal hydrological forecasting (e.g. Shukla and Lettenmaier,  
119 2011, Singla et al., 2012, Mo and Lettenmaier, 2014, and Thober et al., 2015) but also  
120 in medium range forecasting (Alfieri et al., 2014) to determine prediction skill by  
121 comparing the hindcasts with the output from a reference simulation. A reference  
122 simulation is a simulation made with the same hydrological model as the hindcasts,  
123 except that the forcing is taken from meteorological observations or from a gridded  
124 version of meteorological observations. The reference simulation can best be regarded  
125 as a simulation that attempts to make a best estimate of the true conditions (in terms of  
126 e.g. discharge, soil moisture and evapotranspiration), using the modelling system. We  
127 will refer to the output of such a reference simulation as “pseudo-observations”  
128 (misleadingly named “true discharge” in Bierkens and Van Beek, 2009; more  
129 appropriately “synthetic truth” in Shukla and Lettenmaier, 2011; “reanalysis” in Singla  
130 et al., 2012; “a posteriori estimates” in Shukla et al., 2014). We prefer the term “pseudo-  
131 observations” over “re-analysis” since the latter has a meteorological connotation that  
132 often implies the use of some form of (variational) data assimilation. We did not attempt  
133 any form of assimilating observed hydrological variables, such as discharge, in our  
134 reference run.

135  
136 Pseudo-observations have the important advantages of being complete in the spatial and  
137 the temporal domain and to be available for all model variables. Also, they are suitable  
138 for the quantification of small sensitivities, e.g. to bias correction of the meteorological  
139 forcing, which would be hard to detect with real observations. Finally, assessment of  
140 skill based on pseudo observations excludes model errors from the analysis, which is

141 especially useful when addressing various sources of skill (Wood et al., 2016),  
142 something we will do in the companion paper.

143

144 The downside of pseudo-observations is, of course, that they are not equal to real  
145 observations. In this paper we will determine the performance of the prediction system  
146 not only with pseudo-observations, but also with real observations of discharge (like  
147 e.g. Koster et al., 2010, and Yuan et al., 2013) and compare the skill found with the two  
148 different approaches (“theoretical and actual skill”, according to Van Dijk et al., 2013).  
149 Such a comparison was previously made by Bierkens and Van Beek (2009) and Van  
150 Dijk et al. (2013). We will analyse and discuss conceptual differences between using  
151 pseudo- and real observations for verification. We hypothesise that the fact that the  
152 pseudo-observations are obtained with the same model as the hindcasts logically  
153 contributes to an overestimation of the skill when the pseudo-observations are used for  
154 verification.

155

156 This paper aims to analyse to what extent WUSHP is able to predict runoff and discharge  
157 in Europe over the full annual cycle and for lead times up to 7 months. We aim to assess  
158 skill at maximum resolution, i.e. at monthly resolution instead of seasonal or semi-  
159 annual aggregates. Where many studies use correlation coefficient as main skill metric  
160 we will assess skill also for the more probabilistic metrics ROC area and RPSS (see  
161 Sect. 2.3). The second aim of the paper is to get a better understanding of the effects of  
162 using pseudo-observations, as opposed to using actual observations, for the verification  
163 of hindcasts. In the next section we describe the concept and details of our modelling  
164 (Sect. 2.1) and analysis approach (Sect. 2.2 and 2.3). We will start the result section by  
165 assessing theoretical skill of the runoff hindcasts (Sect. 3.1) and then proceed to  
166 theoretical skill of the discharge hindcasts and a comparison between theoretical skill of  
167 discharge and runoff in Sect. 3.2. Differences between theoretical and actual skill of  
168 discharge will be presented (Sect. 3.3) followed by an analysis of differences in skill  
169 determined with various metrics in Sect. 3.4. The discussion starts with a conceptual  
170 analysis of reasons for differences in actual and theoretical skill (Sect. 4.1), followed by  
171 a discussion of uncertainties (Sect. 4.2) and implications (Sect. 4.3).

172

173 In a companion paper (Greuell et al., 2016) we analyse the reasons for the presence or  
174 lack of skill discussed in the present paper, using two different methods. Firstly, skill in  
175 the forcing and other directly related hydrological variables, like evapotranspiration, are  
176 analysed. Secondly, a number of experiments similar to the conventional Ensemble  
177 Streamflow Prediction (ESP) and reverse-ESP, which isolate different causes of  
178 predictability, are discussed. In the present results and discussion sections we will  
179 occasionally look forward to the identified causes of skill.

180

181

182

183

184

185 **2 System, models, data and methods of analysis**

186

187 To assess the forecast quality of our system, two approaches for verification of the  
188 hindcasts are used in this paper. First, we determine the skill of the hindcasts by  
189 comparing predicted discharge with the output of a reference simulation (the “pseudo-  
190 observations” leading to “theoretical skill”), allowing an evaluation that is continuous  
191 in space and time. Secondly, we quantify skill with respect to observations of real  
192 discharge (“real observations” leading to “actual skill”), allowing evaluation at a limited  
193 number of locations (discharge stations) on the river network only. Fig. 1 presents a  
194 flow diagram of the three relevant systems, namely the real world and the two model  
195 systems that generate the hindcasts and the pseudo-observations respectively.

196

197 In each system, confined in the diagram by a box, meteorological and initial conditions  
198 force and initialize hydrology, of which discharge is the relevant component here. There  
199 are three components that cause differences between actual and theoretical skill. First,  
200 the initial conditions are generated by meteorological forcing during the spin up period,  
201 initial conditions at the beginning of the spin up period and hydrology. This is  
202 represented by the upper left branch in each box, omitting initial conditions at the  
203 beginning of the spin up period for simplicity. Second, observations of discharge  
204 generally differ from real discharge (Juston et al., 2014) due to unavoidable  
205 measurement errors as illustrated in the upper right corner of the figure. Third, obviously  
206 a difference exists between real hydrology and model hydrology, central in each box.  
207 Since the hindcasted discharge and pseudo observations share the same model  
208 hydrology and the same initial conditions, and both are free from any observational  
209 errors, theoretical skill will always be larger than actual skill.

210

211 For now we simply accept, and even stress this a-priori ‘superiority’ of theoretical over  
212 actual skill. In Sect. 4.1 we will come back to this and further discuss, at least in  
213 qualitative terms, how each of the differences between the three systems affect skill  
214 assessment.

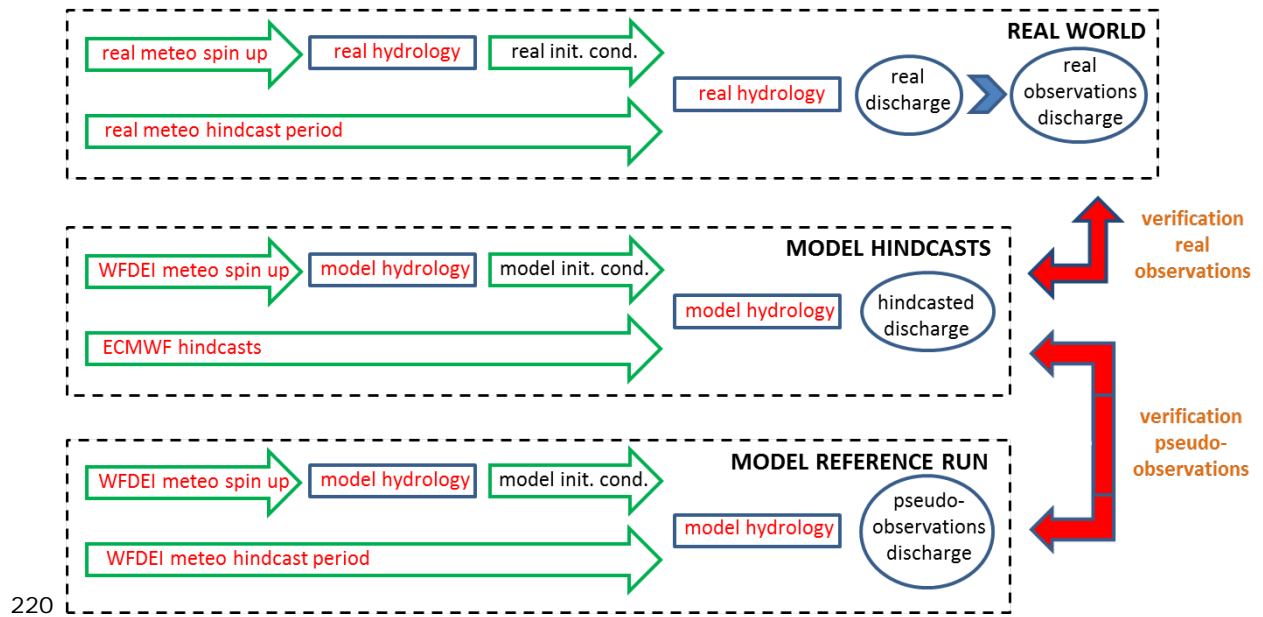
215

216 In the following subsections we will describe each component.

217

218

219



220  
221

222 Figure 1: Conceptual setup of the present study, showing differences between  
223 verification of hindcasts (in the middle) with pseudo observations (bottom)  
224 and with observations of real discharge (top). See the text in this section for  
225 further explanation and Sect. 4.1 for further discussion.

226  
227

## 228 2.1 The hindcasts and the reference simulation

229

230 WUSHP consists of two simulation branches: a single reference simulation and the  
231 hindcasts themselves. In both branches, terrestrial hydrology is simulated with the  
232 Variable Infiltration Capacity model (VIC, see Liang et al., 1994), which runs on a  
233 domain extending from 25 W to 40 E and from 35 to 72 N, including 5200 land based  
234 cells of  $0.5^\circ \times 0.5^\circ$  (see maps in e.g. Fig. 2. VIC is forced by a gridded data set of daily  
235 meteorological data. VIC is run in so-called ‘energy balance mode’ which requires  
236 resolving the diurnal cycle. Therefore, internally the model temporally disaggregates  
237 the daily input to 3-hourly data and runs at 3 hourly time step. Output of all variables is  
238 again at daily resolution. Because snow may contribute significantly to the seasonal  
239 predictability of other hydrological variables, VIC was run with the option of subgrid  
240 elevation bands. This means that for each gridcell calculations were carried out at up to  
241 16 different elevations, with the aim of simulating the elevation gradient of snow. VIC  
242 was run in naturalised flow mode, i.e. river regulation, irrigation and other  
243 anthropogenic influences were not considered.

244

245 In the reference simulation VIC is forced by the WATCH Forcing Data Era-Interim  
246 (WFDEI; Weedon et al., 2014) for the period of 1979-2010, of which the first two years  
247 were used to spin up the states of snow, soil moisture and discharge, and not used in  
248 further analysis. The reference simulation has the dual aim to create the pseudo-

249 observations for verification purposes and to create a best estimate of the temporally  
250 varying model state, which is then used for the initialisation of the hindcasts.

251

252 The second branch, the hindcasts, consists of three steps. Seasonal predictions of daily  
253 meteorological variables are taken from ECMWF's Seasonal Forecast System 4 (S4  
254 hereafter). These are then bias-corrected using WFDEI as the reference data set. Finally,  
255 VIC is run with the bias-corrected S4 hindcasts as forcing, taking initial states from the  
256 reference simulation.

257

258 The S4 hindcasts used in the present study include 15 members, cover the period from  
259 1981 to 2010 and consist of 7 month simulations initialised on the first day of every  
260 month (see Molteni et al., 2011 and ECMWF Seasonal Forecast User Guide, online).  
261 The S4 ensemble is constructed by combining a 5-member ensemble analysis of the  
262 ocean initial state with SST perturbations of that state and with activation of stochastic  
263 physics. The whole system is thus probabilistic.

264

265 The variables taken from the S4 hindcasts are daily values of precipitation, minimum  
266 and maximum temperature, atmospheric humidity, wind speed and incoming short- and  
267 long wave radiation, since these are all needed to force VIC. All of these variables were  
268 regridded with bi-linear interpolation from the  $0.75 \times 0.75^\circ$  lat-lon grid of the S4  
269 hindcasts to a  $0.5^\circ \times 0.5^\circ$  grid. Since bias correction generally improves forecasting skill,  
270 the quantile mapping method of Themeßl et al. (2011) was applied to bias-correct the  
271 forcing variables, taking the WFDEI as reference. For each variable and grid cell, 84  
272 correction functions were established and applied by separating the data according to  
273 target month (12) and lead month (7). Such empirical distribution mapping of daily  
274 values has been successful in improving especially forecast reliability (rather than  
275 sharpness and accuracy; Crochemore et al., 2016).

276

277 VIC was run for the period of the S4 hindcasts (1981 – 2010). Additionally, for the  
278 reference simulation two extra years (1979 – 1980) were simulated to spin up the states  
279 of snow, soil moisture and discharge. The hindcast simulations were initialised with  
280 states of soil moisture and snow from the reference simulation, so for these variables  
281 spin up was not needed. However, due to the set-up of the routing module of VIC, the  
282 state of discharge could not be saved and loaded. Hence to spin up discharge, each 7-  
283 month hindcast simulation was preceded by one month simulation with WFDEI forcing.  
284 Since the hindcasts cover 30 years with 12 initialisation dates each and consist of 15  
285 members, a total of 5400 hindcast simulations was carried out.

286

287 Simulations of historic discharge made with VIC (and four other hydrological models)  
288 were validated with observations from large European rivers by Greuell et al. (2015)  
289 and Roudier et al. (2016). VIC exhibits a fairly small average bias (across 46 stations)  
290 of +23 mm/yr (= 7%) and overall differentiates well between low and high runoff basins  
291 with a spatial correlation coefficient of 0.955. However, specific discharge was  
292 overestimated in the Mediterranean and underestimated in northern Fennoscandia.



293 Annual cycles are fairly well reproduced across Europe, though VIC somewhat  
294 overestimates their amplitude. In northern Fennoscandia the spring peak is too late and  
295 too long. Annual cycles are best reproduced for rain-fed rivers in central Europe while  
296 those for rivers with significant snow dynamics are good (Alps). However, the annual  
297 cycle is more poorly reproduced in basins with strong soil freezing dynamics (northern  
298 Fennoscandia) or strong damping of discharge amplitudes by large lakes (southern  
299 Finland).

300

301 Perhaps more relevant in the present context is the model's capability to reproduce inter-  
302 annual variations in discharge. The standard deviation of simulated annual discharge  
303 was 9% higher than observed and the correlation between the two 0.935. Like most  
304 models, VIC is better in simulating high flows (95 percentile: Q95) than low flows (Q5);  
305 the first is slightly overestimated, the second more seriously underestimated. The inter-  
306 annual variation in Q5 is overestimated in central Europe and the Alps, but  
307 underestimated in Fennoscandia (overall correlation across Europe 0.40). The inter-  
308 annual variation in Q95 shows no clear spatial pattern and the overall correlation is 0.7.

309

310 All validation results discussed in these two paragraphs are for the VIC model forced  
311 by E-OBS (v9, Haylock et al. 2008). Our forcing, WFDEI, shows higher precipitation  
312 (+104 mm/yr) across most of Europe, except the Alps, Scotland and westernmost  
313 Norway. According to Greuell et al. (2015) this leads to higher mean discharge, higher  
314 inter annual variability and higher Q95 (not Q5) of simulated discharge for almost all  
315 stations.

316

317

## 318 **2.2 Discharge observations**

319

320 For the assessment of skill with real discharge observations, two data sets were acquired  
321 from the Global Runoff Data Centre, 56068 Koblenz, Germany (GRDC): the GRDC  
322 data set proper and the European Water Archive (EWA) data set. We mapped these two  
323 station data sets onto the VIC grid with a resolution of  $0.5^\circ \times 0.5^\circ$  and a time step of a  
324 month. To enable the investigation of the effect of size on some of our results, we made  
325 two sub-classes of observations. The first comprised observations for basins larger than  
326  $9900 \text{ km}^2$  ("large basins"), the second basins smaller than the area of the grid cells, i.e.  
327 smaller than about  $2530 \text{ km}^2$  in southern Europe (at  $35^\circ \text{ N}$ ) or  $< 1050 \text{ km}^2$  at  $70^\circ \text{ N}$   
328 ("small basins").

329

330 Initially, in many cases the location of observation stations did not match with the  
331 corresponding river in the digital river network used in the routing calculations  
332 (DDM30, see Döll and Lehner, 2002). We corrected for this issue by matching the  
333 observations with the simulations by means of basin size. The size of the model basins  
334 ("model basin area") was determined by the DDM30 network. The size of the basins  
335 upstream of the observation stations ("station basin area") was taken from the meta data

336 of the observations. First the station basin area was compared to the model basin area of  
337 the cell that is nearest to the station (“nearest model cell basin area”).

338

339 Then, the mapping procedure for each observation varied slightly between the two  
340 classes of basins.

341

342 For large basins we then proceeded as follows:

- 343 - If the station and the nearest model cell basin area differed by less than 15%, the  
344 observations were matched with the model calculations for the nearest model cell.
- 345 - Otherwise, the station basin area was compared with the model basin area of the  
346 eight cells surrounding the nearest model cell.
- 347 - The minimum of the eight differences was determined.
- 348 - If that minimum was less than 15%, the simulations for the corresponding cell were  
349 matched with the observations.
- 350 - Otherwise, the station was discarded.

351

352 For small basins we proceeded as follows:

- 353 - If the nearest model cell did not have an influx from any of the neighbouring cells,  
354 its simulations were matched with the observations.
- 355 - Otherwise, all of the eight neighbouring cells without influx were selected.
- 356 - Their simulations were averaged and matched with the observations.

357

358 We further discarded all observations with less than 21 years of data within the  
359 simulation period (1981-2010) for any of the months of the year. The final data set  
360 within our European domain contained 111 cells with observations for large basins and  
361 636 cells with observations for basins smaller than a model gridcell.

362

363 These data sets do not include any variable or parameter characterising the level of  
364 human impact. To enable analysis of the effect of anthropogenic flow modifications on  
365 predictive skill, we quantified the human impact by performing two model simulations  
366 with the Lund-Potsdam-Jena managed Land (LPJmL) model (Rost et al., 2008;  
367 Schaphoff et al., 2013). This model was operated at the same spatial resolution ( $0.5^\circ \times$   
368  $0.5^\circ$ ) and with the same river network (DDM30) as VIC, but the former does include  
369 dams (GRanD database; Lehner et al., 2011) and associated reservoir management.  
370 From the discharge output of a naturalized LPJmL run and an LPJmL run with reservoir  
371 operation and irrigation, the human impact at cell level was quantified by computing the  
372 so-called Amended Annual Proportional Flow Deviator (AAPFD; see Marchant and  
373 Hehir, 2002). For the analysis in Sect. 3.3, we selected all discharge observations for  
374 large basins with an AAPFD  $< 0.3$ , i.e. basins with a relatively small degree of human  
375 impact (about half of all 111 basins).

376

377

378

379

### 380 2.3 Methods of analysis

381

382 From the model output, consisting of daily means, monthly mean values were computed,  
383 which were then used for the analysis. The analysis is restricted to runoff, defined here  
384 as the amount of water leaving the model soil either along the surface or at the bottom,  
385 and discharge, defined here as the flow of water through the largest river in each grid  
386 cell. Discharge accumulates all runoff from cells that are upstream in the model river  
387 network, with delays due to transport inside cells and through the river network. Hence,  
388 whereas runoff represents only local hydrological processes, discharge aggregates  
389 hydrological processes occurring in the entire basin upstream of a particular cell.

390

391 Instead of analysing skill per target season and/or for a number of consecutive lead  
392 months, we analysed skill for every combination of 12 target and 7 lead months. The  
393 thus achieved higher temporal resolution of the skill metrics enables a more accurate  
394 determination of the beginning and end of periods of skill. Moreover, skill at a monthly  
395 resolution provides the possibility to determine the consistency of the skill where we  
396 define consistent skill as skill that persists during at least two consecutive target or lead  
397 months. In accordance with Hagedorn et al. (2005) we designated the first month of the  
398 hindcasts as lead month zero, so target month number is equal to the number of the  
399 month of initialisation plus the lead month number.

400

401 Three skill metrics (see Mason and Stephensen, 2008, for a good discussion of the why  
402 and how of these) were computed: i) the correlation coefficient between the  
403 observations and the *median* values of the hindcasts (shortly “correlation coefficient” or  
404 R), ii) the area beneath the Relative Operating Characteristics (ROC) curve (shortly  
405 “ROC area”) and iii) the Ranked Probability Skill Score (RPSS). The ROC area is  
406 computed for each month separately and for three categories of the observations and  
407 hindcasts with an equal number of values, with the categories containing the one third  
408 highest, lowest and the remaining values (upper, lower and middle tercile, resp.; above,  
409 below and near-normal, AN, BN and NN categories). The same subdivision of  
410 observations and hindcasts in terciles was made to compute the RPSS. Since none of  
411 these metrics is sensitive to systematic biases in the forecasting system, no attempt was  
412 made to correct simulated runoff or discharge for any such errors prior to computing the  
413 skill metrics. So we focus our evaluation on the models capability to predict river flow  
414 anomalies rather than absolute river flows.

415

416 All three skill metrics quantify, though in different ways, how well the ranking of the  
417 hindcasts matches the ranking of the observations. The correlation coefficient is a  
418 measure of the association between (pseudo-) observation and forecast ensemble  
419 median; we used the Pearson correlation coefficient. The ROC area is a measure of  
420 resolution or discrimination and indicates whether the forecast probability of an event  
421 (i.e. value falling in the considered tercile) is higher when such an event occurs  
422 compared to when not. The RPSS is a measure of accuracy and summarizes in a single  
423 number the skill of a forecast system to make forecasts with the correct percentage of

424 ensemble members falling in any of the defined terciles. Perfect forecasts have values  
425 of 1 for all three skill metrics. Climatological forecasts (probabilistic forecasts that in  
426 our case each year predict a 0.33 chance of a high or low anomaly occurring) lead to  
427 values of 0 for R, 0.5 for the ROC area and 0 for the RPSS. Random forecasts were used  
428 to determine the significance of the metrics. In the case of the RPSS, these were  
429 generated by sampling randomly from the multinomial distribution with  $p = (1/3, 1/3,$   
430  $1/3)$  and  $N = 15$  (the number of ensemble members), which is the distribution of  
431 climatological ensemble forecasts. Each metric will be designated as significant for p-  
432 values less than 0.05. This implies association is significant for  $R > 0.31$ , resolution is  
433 significant for ROC area  $> 0.69$  and accuracy is significant for  $RPSS > 0$ .

434

435 To a large extent, we found that our results and conclusions, in terms of spatio temporal  
436 patterns of skill, are independent of the chosen metric. Hence, and because among the  
437 three metrics the correlation coefficient is the easiest to understand, we will discuss  
438 results mostly in terms of the correlation coefficient, which is in line with Doblas-Reyes  
439 et al. (2013). The sensitivity to the chosen metric and significant differences between  
440 these metrics will be discussed in Sect. 4.2.

441

442 All metrics were computed using the low and high level R packages “SpecsVerification”  
443 (Siegert et al., 2014) and “easyVerification” (Bhend et al., 2016), respectively. Metrics  
444 cannot be computed if observations or hindcasts within the entire 30 year period consist  
445 for more than one third of zeros or one sixth of ties (i.e. equal values). Such skill gaps  
446 (i.e. the white terrestrial cells in Fig. 2 and 3) only occur in the far North due to rivers  
447 that are frozen for at least a month in winter.

448

449

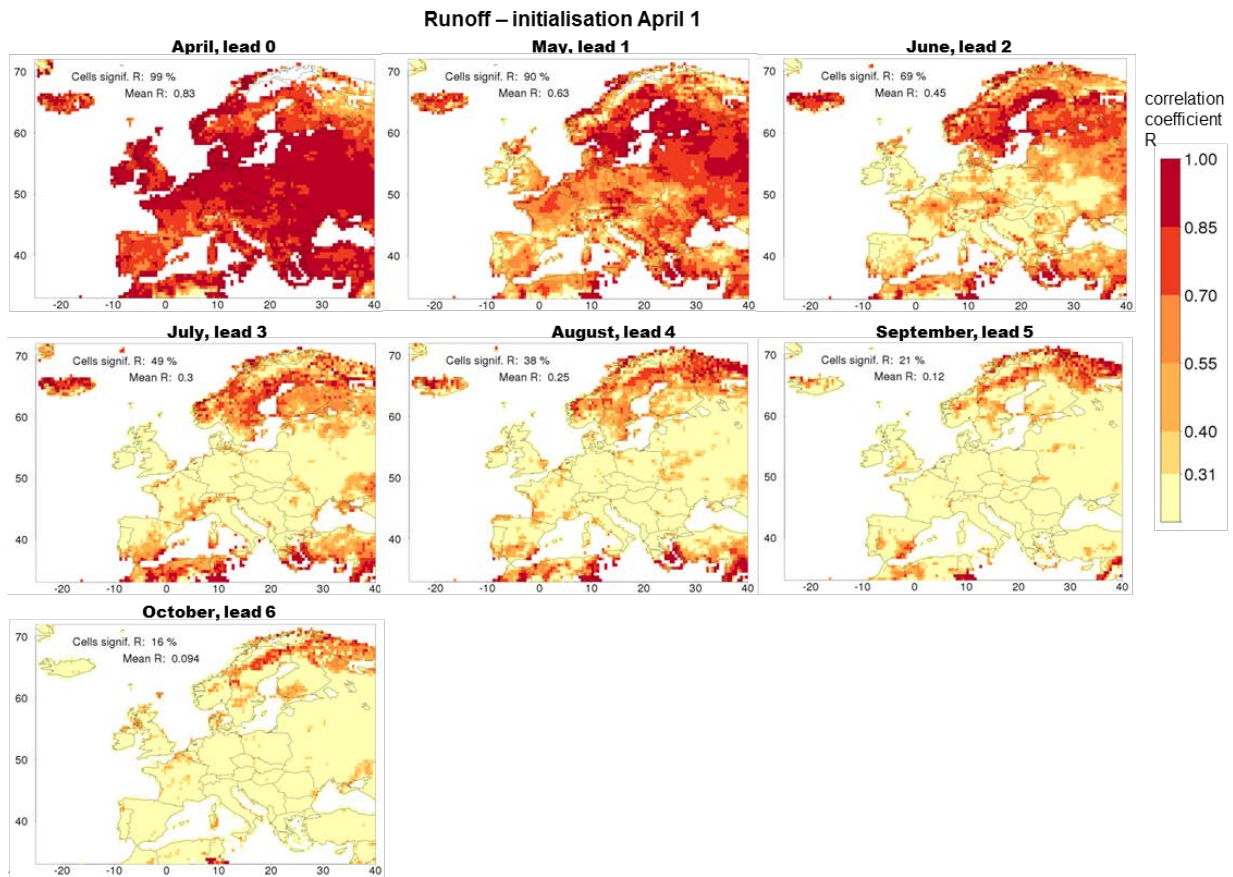
## 450 **3 Results**

451

### 452 **3.1 Spatiotemporal variation of skill in runoff forecasts**

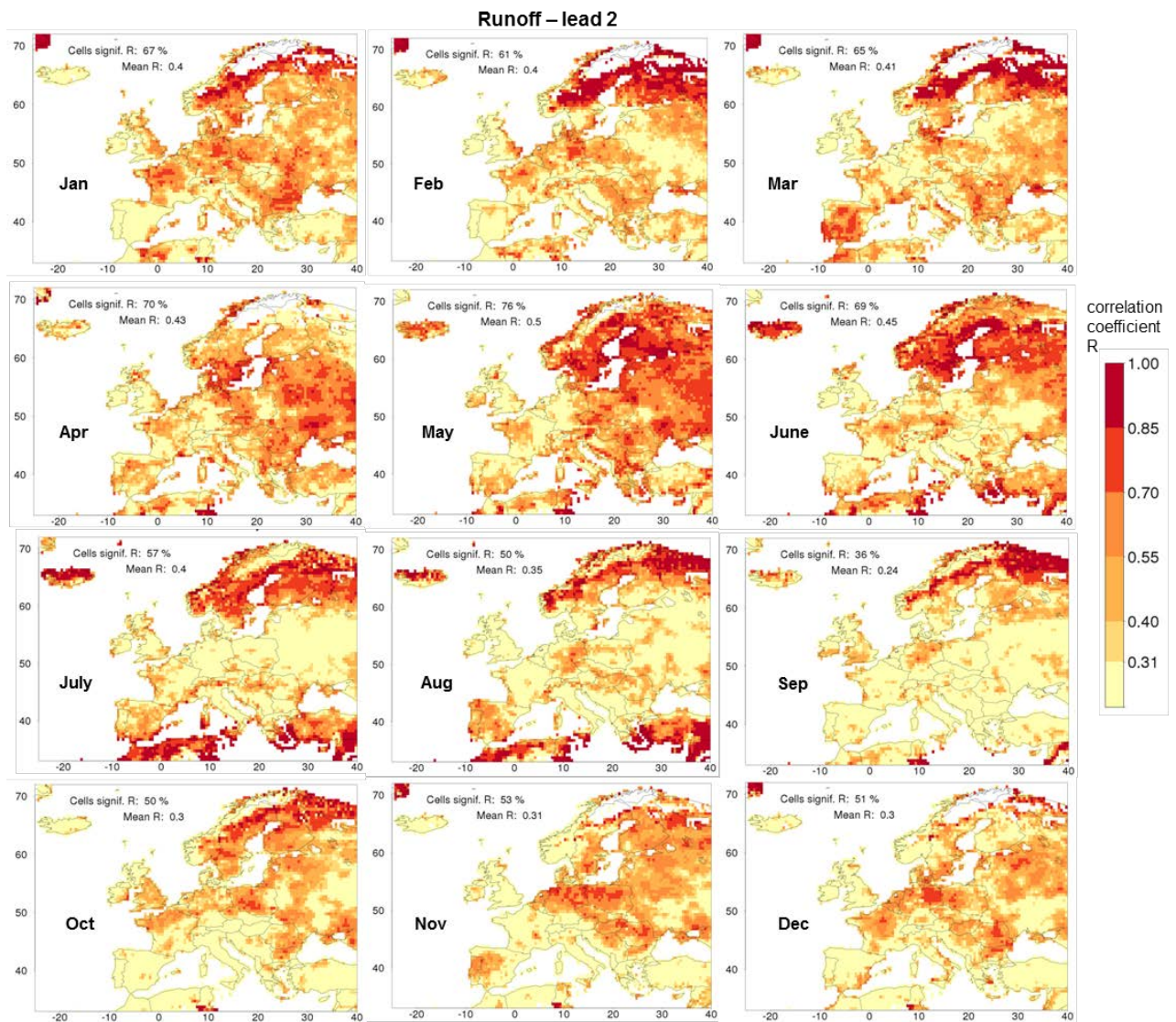
453

454 Eighty-four maps of skill of the runoff hindcasts were drawn for all 12 initialisation  
455 months and all 7 lead months (all are presented in supplementary material S1). Two  
456 cross-cuts through that collection are shown in Fig. 2 (for a single initialisation month)  
457 and 2 (for a single lead month). The seven panels of Fig. 3 show the skill of the hindcasts  
458 initialised on April 1 as a function of lead time. Cells with an insignificant amount of  
459 skill are tinted yellow; cells where no metric could be computed remain white. In lead  
460 month 0, significant skill is found across almost the entire domain (99% of the cells).  
461 After the first lead month, the fraction of cells with significant skill gradually decreases  
462 to reach 16% at the longest lead time (lead month 6). This is more than expected for the  
463 case of completely unskilful simulations (5% of the cells), so at the end of the hindcast  
464 simulations significant skill that does not occur due to chance is still present in some  
465 regions. The general impression is that the pattern of skill does not move in space but  
466 that skill is fading, i.e. for individual grid cells R is mostly decreasing with increasing  
467 lead time.



469

470 Figure 2: Skill of the runoff hindcasts initialised on April 1 for all seven lead months.  
 471 Skill is measured in terms of the Pearson correlation coefficient between the  
 472 median of the hindcasts and the observations (R). The threshold of  
 473 significant skill lies at 0.31, so yellow cells have insignificant skill, (dark)  
 474 red cells have (most) skill. White, terrestrial cells correspond to cells where  
 475 observations or hindcasts consist for more than one third of zeros or one  
 476 sixth of ties. The legend provides the fraction of cells with significant values  
 477 of R (at the 5% level) and the domain-averaged value of R.



479

480 Figure 3: Annual cycle of skill (R) of runoff hindcasts of lead month 2. More  
 481 explanation is given in the caption of Fig. 2.

482

483

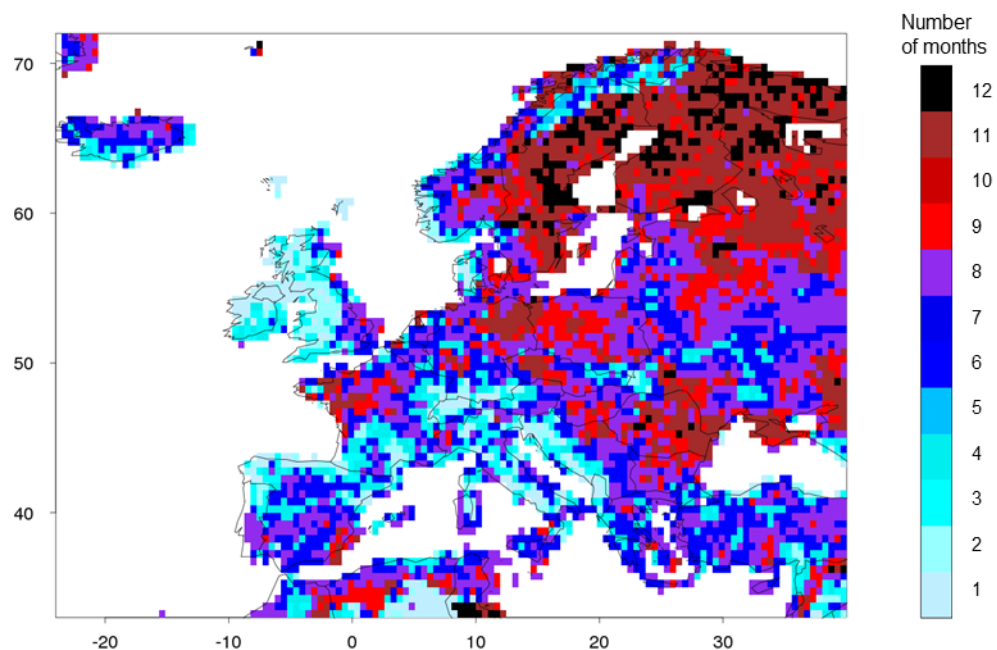
484 The twelve panels of Fig. 3 show the annual cycle of skill of the hindcasts for lead  
 485 month 2. Consistent skill (persistent during at least 2 consecutive target months) is found  
 486 in (causes of skill are reproduced here from the companion paper, Greuell et al., 2016):

- 487 - Fennoscandia. Much skill is present during the entire year, except for November  
 488 and December, and there is a dip in skill in April. On average across the entire  
 489 region, skill reaches a maximum in May and June, i.e. the end of the melting season,  
 490 and –as shown in the companion paper- largely due to initialising snow. Compared  
 491 to the rest of the peninsula, there is generally less skill along the Scandinavian  
 492 Mountain range. The companion paper shows some evidence that this may be due  
 493 to high variability of orographic rain, ill-represented in the S4 re-forecasts.
- 494 - Poland and northern Germany. The core period lasts from November to January,  
 495 but it is extended with periods of less skill into October and the months from

496 February to May. Here both initialisation of soil moisture and snow are important  
 497 for skill .

- 498 - western France, more or less from Paris to Brittany and roughly from December to  
 499 May. Skill derives from initialisation of soil moisture.
- 500 - The eastern side of the British Isles from January to April up to lead month 2. Also  
 501 here skill derives from soil moisture initialisation.
- 502 - Romania and Bulgaria. The core as well as the whole period are the same as that  
 503 for Poland and northern Germany. In addition to causes mentioned there, in this  
 504 part of Europe also summer precipitation and evapotranspiration are forecasted  
 505 fairly well.
- 506 - The southern part of the Mediterranean region from June to August. The high  
 507 amounts of skill are limited to the coastal parts of northern Africa, Sicily, southern  
 508 Greece, Turkey, Syria and Lebanon.
- 509 - The Iberian peninsula from January to March up to lead month 2, and July and  
 510 August like the other parts of the Mediterranean mentioned before. Skill derives  
 511 from soil moisture in initialisation and in winter also from some skill in  
 512 precipitation.

513



514  
 515  
 516 Figure 4: Number of months in a year with significant skill (R) in the runoff  
 517 forecasts of lead month 2.

518  
 519 Figure 4 displays a synthesis of Fig. 3 in the form of a map with the fraction of the 12  
 520 months of the year with significant skill for lead month 2. Many of the regions with very  
 521 little or no skill all over the year are coastal regions (e.g. northern coast of Spain),  
 522 especially coastal regions on the western side of land masses (e.g. western coasts of

523 Denmark, southern Norway, Croatia and the British Isles), and mountain regions (e.g.  
524 the Alps, mountains in northern Norway and Sweden and on the Tatra on the border of  
525 Poland and Slovakia). The British Isles exhibit little skill, except for the eastern coast of  
526 Great Britain in late winter and early spring (JFMA). The companion paper shows that  
527 for regions with skill during a large part of the year, this skill is derived from  
528 complementary periods of skill due to initial conditions of snow and/or soil moisture.

529

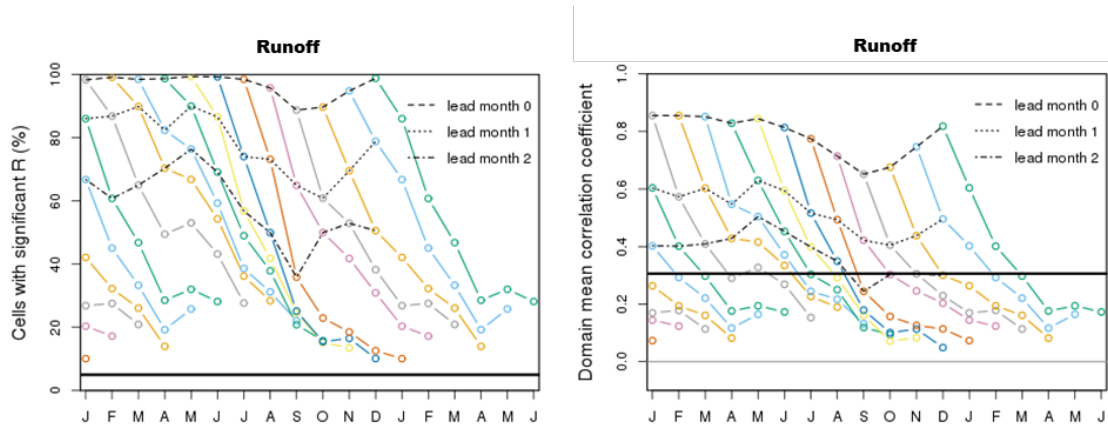
530 These pan-European results can be compared to those of Bierkens and Van Beek (2009).  
531 They found maxima in predictability of winter discharge in Northern Sweden, Finland,  
532 the region between Moscow and the Baltic Sea, Romania and Bulgaria, and Eastern  
533 Spain. For the winter there is crude agreement with the current study about Northern  
534 Sweden, Romania and Bulgaria, but not about the other regions. For the summer,  
535 Bierkens and Van Beek (2009) compute maxima in skill for Southern Spain, Sardinia,  
536 Western Turkey and South-western Finland, a pattern that broadly agrees with the  
537 locations of the summertime maxima in skill (most of Fennoscandia and southern part  
538 of the Mediterranean region) we find.

539

540 Singla et al. (2012) found considerable skill in the Seine basin for low flows from June  
541 – September, a bit more eastern from the region where we found skill. Trigo et al. (2004)  
542 using a statistical model based on December NAO indices found skill for JFM discharge  
543 (and hydropower production) for the Douro, Tejo and Guadiana basins covering most  
544 of central and western Iberia. We confirm this skill which last till about May here, when  
545 initialised in January. In addition (not analysed by Trigo) we find skill beyond lead zero  
546 also in summer but then more concentrated around the south eastern coast of Iberia.  
547 Svensson et al. (2015) using a statistical model, based on NAO indices and river flow  
548 persistence, found good skill for winter river flows on the eastern side of the British  
549 Isles, consistent with our findings, and barely significant skill on its western coast that  
550 we do not reproduce.

551





552

553

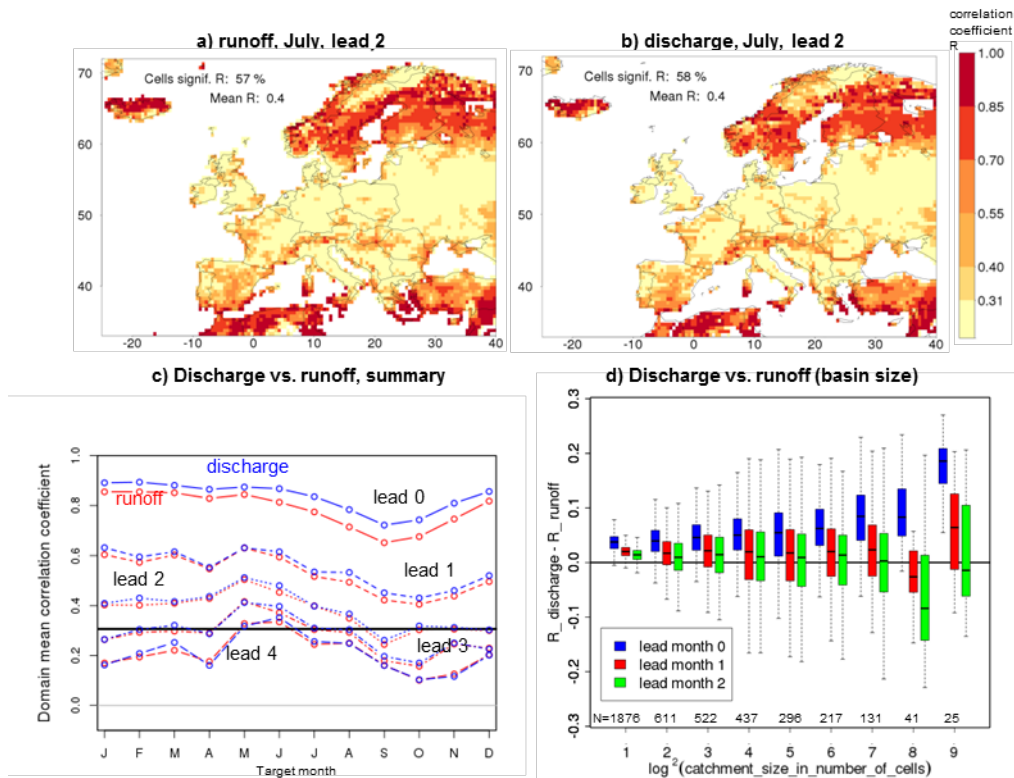
554 Figure 5: At left a) Fraction of cells with significant skill (in terms of R), and at right  
 555 b) domain average correlation in the runoff hindcasts, as a function of  
 556 initialisation month and lead time. Each coloured curve corresponds to the  
 557 hindcasts initialised in a single month. For better visualisation, parts of the  
 558 curves that end in the next year are shown twice, namely at the left hand and  
 559 the right hand side of the graph. Black lines (dashed, dotted and dashed-  
 560 dotted) connect the results for identical lead times. The horizontal line in a)  
 561 shows the expected fraction of cells with significant skill, in the case that  
 562 the hindcasts have no skill at all (5%), in b) the minimal magnitude of the  
 563 correlation of a single cell for it to be statistically significant

564

565

566 Figure 5a summarizes skill across the domain in terms of the fraction of cells with  
 567 significant R for all initialisation and lead months. Overall there is a considerable  
 568 amount of significant skill, with a minimum roughly from August to November and a  
 569 maximum in May. For lead month 2 the fraction of cells with significant skill varies  
 570 between 36% (September) and 76% (May). In all of the 84 combinations of initialisation  
 571 and lead month, the theoretical value of no skill at all (5%) is exceeded, implying there  
 572 are (small) pockets of skill even at lead month seven. Individual curves show the loss of  
 573 skill with increasing lead time. The exception is formed by hindcasts starting in  
 574 November, December and January which gain skill when they progress from April to  
 575 May, a phenomenon caused by initial conditions of snow that takes longer or shorter to  
 576 melt in (late) spring. For details, see the companion paper. Fig. 5b shows decay and  
 577 gain trends of the domain-averaged R. It shows that a forecast initialised in February  
 578 exhibits higher domain average skill into June (5 lead months), than one starting in July  
 579 has for August (2 lead months). Similar summary plots for the other skill metrics are  
 580 presented in the Fig. S2, and discussed in Sect. 3.4.

581



582

583 Figure 6: Comparison of the skill of the hindcasts of discharge and runoff. The two  
 584 maps display R for runoff (a) and discharge (b) for hindcasts initialised on  
 585 May 1 and target month July (see further explanation in Fig. 1). Panel c  
 586 depicts the annual cycle of the domain-averaged R for runoff (red) and  
 587 discharge (blue) for lead months 0 to 4. The horizontal line at 0.31 is the  
 588 threshold of significance for a single cell. Panel d is a box plot of the  
 589 difference between R for discharge and runoff as a function of the basin size.  
 590 Each bin  $i$  contains the results for all basins with a maximum of  $2^i$  cells and  
 591 more than  $2^{(i-1)}$  cells, e.g. bin 4 is for all basins with a size from 10 to 16  
 592 cells. Boxes represent the interquartile range and the median; whiskers  
 593 extend to minimum and maximum values found in the bin. All values are  
 594 average differences over the twelve months of the year and results are shown  
 595 for three different lead times. The value above the abscissa give the number  
 596 of cells in each bin.

597

598

599

### 3.2 Spatiotemporal variation of skill in discharge forecasts

600

601

602

603

604

605

606

This sub-section compares skill for discharge with skill for runoff. The two maps of Fig. 6, which depict the skill in runoff and discharge hindcasts for July as lead month 2, show a high degree of similarity in terms of the patterns and the magnitude of the skill. The same holds for other target months and lead times (not shown). There are, however, subtle differences because rivers aggregate the skill, or lack of skill, from the whole upstream part of their basin. As a result, cells containing rivers with large basins may contrast against adjacent cells if these contain rivers with a small, local basin. Indeed,

607 some downstream parts of large rivers stick out in the skill map for discharge, but not  
608 in the skill map for runoff. An example in Fig. 6b are the reaches of the Danube along  
609 the Romanian-Bulgarian border, which show more skill than local small rivers in  
610 adjacent cells, because some upstream parts of the Danube have more skill than the  
611 region around the Romanian-Bulgarian border. An example that demonstrates the  
612 opposite is the downstream part of the Loire showing less skill than local small rivers,  
613 because upstream parts of the Loire have less skill than small, local rivers in the  
614 downstream part.

615

616 Domain summary statistics of skill also differ slightly between runoff and discharge.  
617 Figure 6c compares the annual cycle of the skill in discharge with the skill in runoff at  
618 five different lead times. Here we show the difference in the domain-averaged R instead  
619 of the fraction of cells with a significant R because in lead month 0 that fraction is close  
620 to one for both variables. In terms of the domain-averaged R, predictability is higher for  
621 discharge than for runoff for the first lead month. On average over the 12 months of the  
622 year, the difference is 0.049. We ascribe this result to the combined effect of the delay  
623 between runoff and discharge and the general tendency of decreasing skill with lead  
624 time. The curves for the different lead times in Fig. 6c show that the difference in skill  
625 between the two variables gradually disappears with increasing lead time (an annual  
626 average of 0.020 and 0.012 for lead months 1 and 2, respectively). This is compatible  
627 with the given explanation for the difference and the fact that the rate with which skill  
628 is lost gradually decreases with increasing lead time.

629

630 We finally analysed whether the difference in skill between discharge and runoff was a  
631 function of the size of the basin (Fig. 6d). For the first lead month, when on average  
632 there is more skill in discharge than in runoff, the difference increases with the size of  
633 the basin. Again, this can be explained by the combination of the skill decaying with  
634 time and the delay between runoff and discharge, with the delay increasing with the size  
635 of the basin. For longer lead times (from lead month 1 on), when the domain-averaged  
636 difference in skill has become very small, the figure shows no effect of the basin size.  
637 Referring to the comparison between runoff and discharge in panels Fig. 6a and 6b for  
638 lead month 2, cases like the Danube (more skill than local rivers) and the Loire (less  
639 skill than local rivers) tend to cancel when the entire domain is considered. .

640

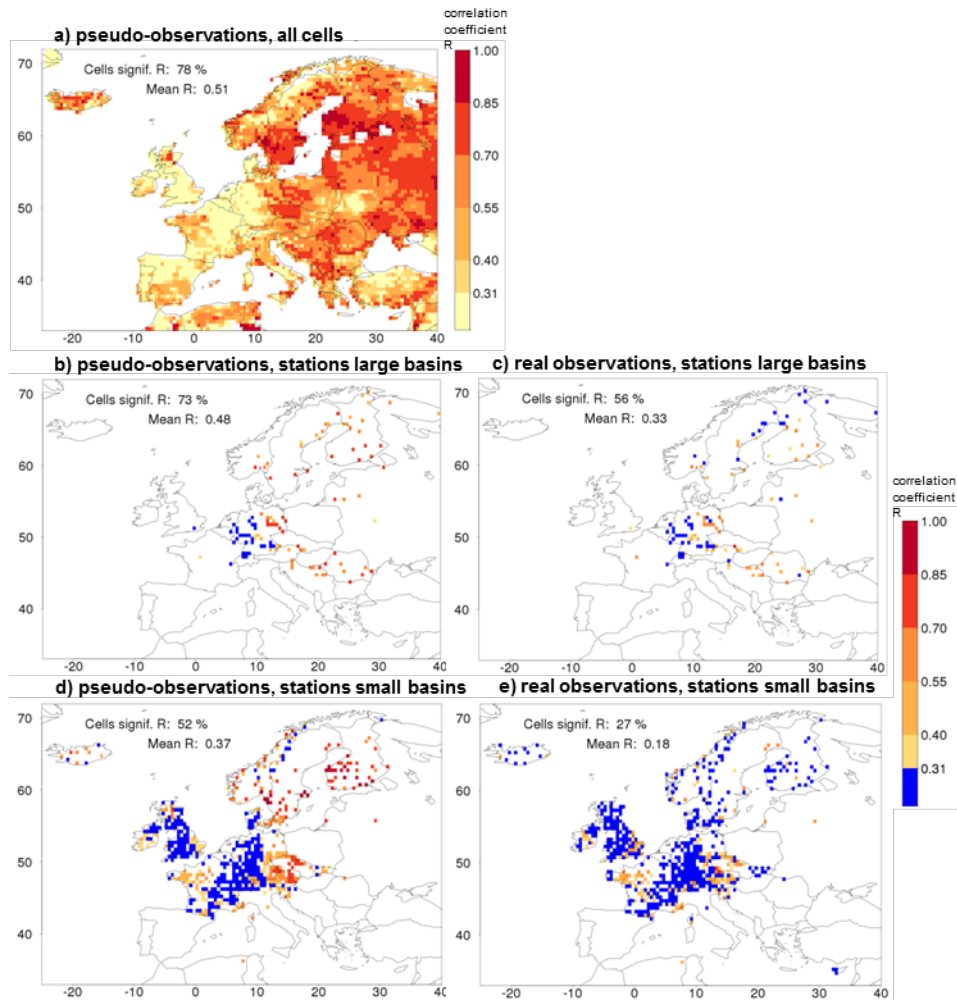
641

### 642 **3.3 Verification of discharge with pseudo- and real observations**

643

644 So far, all skill was determined by using the discharge generated with the reference  
645 simulation. i.e. with pseudo-observations. In this section, this “theoretical skill” will be  
646 compared with the skill determined with real discharge as observed at gauging stations  
647 (“actual skill”) from the GRDC and EWA data bases. Figure 7 compares the theoretical  
648 skill (Fig. 7b and 7d for large and small basins, respectively) with actual skill (Fig. 7c  
649 and 7e for large and small basins, respectively) for a single combination of a target  
650 month (May) with a lead month (2).

## Discharge May as lead 2



651

652

653 Figure 7: Skill ( $R$ ) of the discharge hindcasts for May as lead month 2 (initialisation  
 654 on March 1). In sequence: a) discharge verified with pseudo-observations,  
 655 b) as a but for cells representing large basins only, c) discharge verified with  
 656 real observations for large basins. Panels d) and e) are identical to b) and c),  
 657 respectively, but for cells representing small basins. More explanation is  
 658 given in the caption of Fig. 1 but in panels d) and e) cells with insignificant  
 659 skill are coloured blue instead of yellow for better contrast.

660

661

662 For this combination of May forecasts initialised in March, a substantial degradation in  
 663 skill is found when the pseudo-observations are replaced by real observations. In terms  
 664 of the fraction of cells with significant skill, the reduction is from 73 to 56 % for large  
 665 basins and from 52 to 27 % for small basins and the domain-averaged  $R$  decreases from  
 666 0.48 to 0.33 for large basins and from 0.37 to 0.18 for small basins. The larger basins,  
 667 especially those in northern Fennoscandia lose all skill when using actual observations,  
 668 a region where VIC also performed poorly in reproducing historic flows. There the  
 669 specific discharge was underestimated and the annual cycle was poorly reproduced;  
 670 especially the spring peak occurred too late and too long (Greuell et al., 2015). In central

671 Europe useful skill remains when using real observations, a region where VIC well  
 672 reproduced annual cycles, though interannual variation in low flows where  
 673 overestimated in that area.

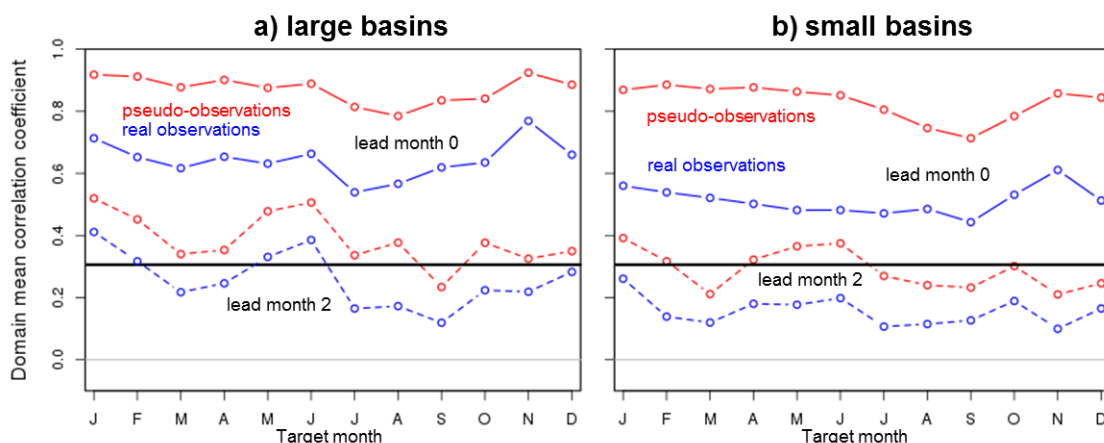
674

675 Figure 8 compares actual with theoretical skill for all target months and two lead times  
 676 by considering the domain-mean R. Similar figures for the other skill metrics are  
 677 presented in Fig. S4 and discussed in the next Sect. 3.4. The reduction in skill occurs for  
 678 all combinations of target and lead months and does not exhibit a clear annual cycle. On  
 679 average across all target months and for lead month 2, the ratio of actual to theoretical  
 680 skill is 0.667 (0.258 divided by 0.387) for large basins and 0.538 (0.156 divided by  
 681 0.290) for small basins. This is comparable to Van Dijk et al. (2013), who found a ratio  
 682 of actual to theoretical skill of 0.54 for 6192 basins worldwide in terms of the ranked  
 683 correlation coefficient.

684

685 Comparing skill for small basins with skill for large basins in Fig. 8, we notice two  
 686 differences. Firstly, in terms of the domain mean R, theoretical skill is higher for large  
 687 basins than for small basins (0.39 and 0.29, respectively, for the annual mean and lead  
 688 month 2). However, this result holds for the cells with observations. If all cells of the  
 689 domain are considered, this difference becomes insignificantly small. So, the apparent  
 690 difference in theoretical skill between large and small basins can be attributed almost  
 691 entirely to the geographical distribution of the discharge monitoring stations, with  
 692 stations on small basins being relatively more often located in regions with relatively  
 693 little skill like Germany, France and the British Isles than large basin stations.

694



695

696

697 Figure 8: Comparison between verification of discharge with pseudo- (red) and real  
 698 (blue) observations in terms of the annual cycle of the domain mean R. The  
 699 horizontal line at 0.31 is the threshold of significance for a single cell.  
 700 Results are shown for cells representing large basins (left) and cells  
 701 representing small basins (right). Both panels depict cycles for lead months  
 702 0 and 2 only.

703

704 The second effect of the size of basins is that reduction between theoretical and actual  
705 skill is larger for small basins than for large basins. This is perhaps even more clear from  
706 Fig. S4 in the supplementary material. We speculate that this is due to a combination of  
707 two effects. Firstly, there is more skill in simulations of historic streamflow in large  
708 basins than in small basins (Van Dijk and Warren, 2010, confirmed for VIC in Europe  
709 by Greuell et al. 2015). Secondly, as Van Dijk et al. (2013) demonstrated, the ratio of  
710 actual to theoretical skill is almost linear in the skill of simulating historic streamflow.  
711 Combining these two relationships confirms the relationship that we found, namely an  
712 increase in the ratio of actual to theoretical skill with basin size.

713

714 Finally, we investigated to what extent these results are affected by human interference,  
715 keeping in mind that the simulations are naturalized, while the observations include  
716 human impacts to a variable but unknown degree. Human interference is expected to  
717 have a negative effect on actual skill and hence on the ratio of actual to theoretical skill.  
718 For relatively natural basins ( $AAPFD < 0.3$ ; see end of Sect. 2.2), the ratio of actual to  
719 theoretical skill was computed in terms of the domain mean  $R$ , averaged across all target  
720 months and for lead month 2. We found a ratio of 0.686, which should be compared to  
721 a ratio of 0.667 for the entire set of large basins (see above). So, as expected the ratio is  
722 larger for basins with less impact. However, since the difference between the two ratios  
723 is small we conclude that the effect of evaluating naturalised runs against observations  
724 that are obviously affected by human interference, contributes only little to the  
725 difference between actual and theoretical skill. A similar analysis was not applied to the  
726 collection of small basins with observations, since these are smaller than the spatial  
727 resolution of the simulations.

728

729

### 730 **3.4 Results for other skill metrics**

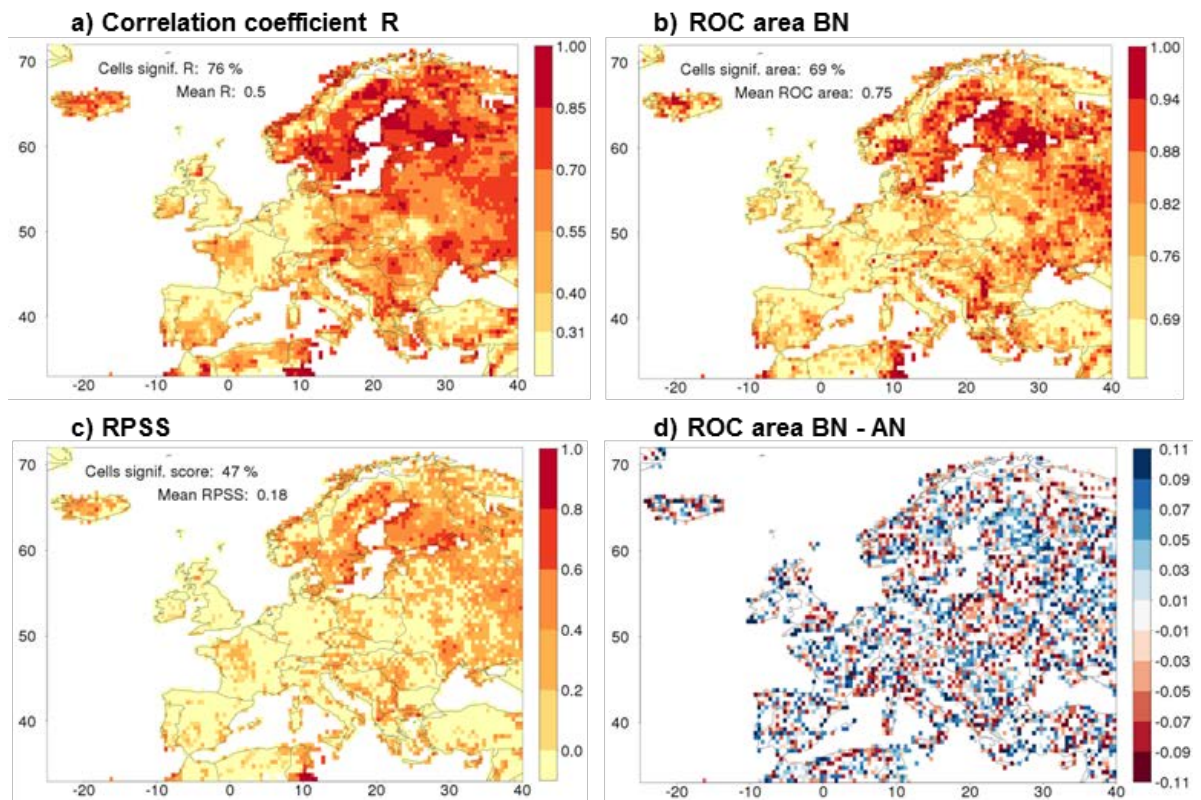
731

732 So far, skill was measured in terms of the correlation coefficient between the median of  
733 the hindcasts and the observations ( $R$ ) only. This section compares those results, for  
734 runoff, with results in terms of other skill metrics. Figure 9 gives an example for one  
735 particular target month and lead month, i.e. target May initialised in March (lead 2). Fig.  
736 9a, 9b and 9c show the skill patterns for  $R$ , for the ROC area for Below Normal (BN)  
737 years and for the RPSS. The three patterns are spatially similar to a large degree, though  
738 the magnitudes and number of significant cells do differ. The pattern of the map of the  
739 ROC area for Above Normal (AN) years (see Fig. S1) is also similar to the patterns of  
740 the three maps shown. On average, across all lead and target months, 89% of the cells  
741 that have significant  $R$  also have significant ROC scores for the BN tercile, 84% also  
742 for the ROC scores for the AN tercile. Finally, 65% of the cells that have significant  $R$   
743 also have significant RPSS scores. The fraction of cells with no significant  $R$ , but with  
744 significant ROC or RPSS remains below the 5% level across all target and lead months,  
745 and thus such cases are likely due to chance.

746

747 The agreement that we find between the patterns of the different metrics is in accordance  
 748 with a result mentioned in a global analysis of seasonal streamflow predictions by Van  
 749 Dijk et al. (2013) who found high spatial correlation between the different skill metrics  
 750 they used (among which R, the RPSS and the ranked correlation coefficient).  
 751  
 752

### Runoff May as lead 2



753  
 754 Figure 9: Maps of different skill metrics for one combination of a target month (May)  
 755 and a lead month (2) of the runoff hindcasts. Panels show a) R, b) the ROC  
 756 area for the below normal tercile, c) the Ranked Probability Skill Score  
 757 (RPSS) and d) the difference in ROC area between the BN and AN terciles.  
 758 In panels a, b and c skill is not significant in cells with a yellow colour.  
 759 Legends provide the fraction of cells with significant values of the metric  
 760 and the domain-averaged value of the metric.

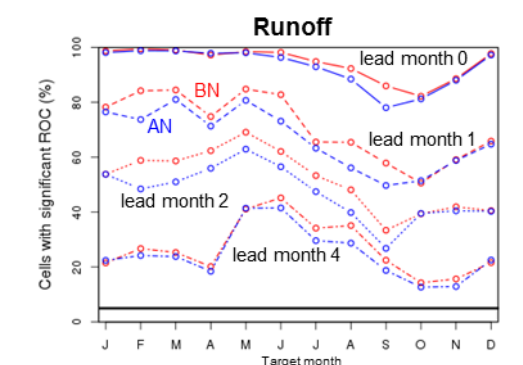
761  
 762 Although the different nature of the different metrics does not enable a quantitative  
 763 comparison of the metrics, ROC areas for the different terciles can be compared among  
 764 each other. For the particular combination of May target month and lead month two  
 765 shown in Fig. 9, the domain-mean ROC area is largest for the BN tercile (0.75), slightly  
 766 smaller for the AN tercile (0.73) and much lower for the near-normal (NN) tercile (0.58,  
 767 see Fig. S2c and d; 0.5 corresponds to climatological forecasts). A similar tendency is  
 768 found in the fraction of cells with a significant ROC area (69%, 63% and 21%,  
 769 respectively). The fraction of cells with a significant value of the RPSS is 47%, which  
 770 is somewhere between the fractions for ROC areas of the three terciles because the RPSS

771 represents the skill across all terciles. All metrics show a minimum value in the annual  
 772 cycles in either September or in October, irrespective of lead time; maxima are attained  
 773 in February for lead month 0 shifting to May at longer lead times (Fig. S2). Finally, Fig.  
 774 9d presents a map of the difference between the BN and the AN ROC area. BN ROC  
 775 values are larger than AN (blue colours) in southern Finland and central Sweden,  
 776 western France, Hungary and Serbia and large parts of Russia. The reverse (ROC AN >  
 777 ROC BN, red colours) is true in eastern Poland and the Baltic states, southern eastern  
 778 France (Rhône basin) and eastern UK.

779 For other combinations of target and lead months the results of this analysis are similar,  
 780 though numbers may vary. See supplementary figures.

781 Figure 10 compares the BN with the AN tercile in terms of the fraction of cells with a  
 782 significant ROC area across all target and initialisation months. The main finding is that  
 783 in all combinations of lead and target month the fraction significant cells is larger for  
 784 the BN than for the AN tercile. This is perhaps not as expected from the VIC  
 785 performance in reproducing historic flows, which is better for high flows than for low  
 786 flows (Greuell et al., 2015; recall that their high/low flows are defined as Q95 and Q5,  
 787 respectively, while here they are Q67 and Q33; see also Sect. 2.1). However, the AN  
 788 and BN fractions tend to become equal (i) when these ROC areas approach 1.0, (ii) when  
 789 they approach the limit of no skill (5%) and (iii) during target months from October to  
 790 January.

791



792

793 Figure 10: Skill of the runoff hindcasts in the Below Normal (BN) compared to the skill  
 794 of the runoff hindcasts in the Above Normal (AN) tercile. The plot depicts  
 795 annual cycles of the fraction of cells with a significant ROC area for the two  
 796 terciles and for four lead months.

797

798

799



## 800 4 Discussion

801

### 802 4.1 Theoretical versus actual skill

803

804 The two essential questions are: 1) What are the conceptual differences between the  
805 physical systems that generate the pseudo- and the real discharge observations, i.e.  
806 between the model reference run and the real world. To answer this question, the  
807 components in the upper and the lower box of the diagram need to be compared. 2) What  
808 are the expected effects of these differences on skill, i.e. on the comparison with the  
809 hindcasts. To answer this question, the components that differ between the real world  
810 and the model reference run need to be compared with the model hindcasts. The rule  
811 then is that skill decreases with increasing disagreement between a component of the  
812 hindcast system and the corresponding component of one of the other systems. The  
813 following components (red text in Fig. 1) differ between the real world and the model  
814 reference simulation, and their expected effect on skill are:

- 815 1. Real meteorology differs from the meteorology assumed in the reference  
816 simulation (WFDEI), both during the spin up period and during the hindcast  
817 period. During spin up, model reference run and hindcasts have identical  
818 meteorological forcing (WFDEI), which differs from real meteorology.  
819 Therefore, this difference is expected to lead to more theoretical than actual skill.  
820 During the hindcast period, all three systems have different meteorological  
821 forcing. For cases with skill in the meteorological hindcasts, one would need to  
822 have an expectation about the agreement between the skilful hindcasts and  
823 reality, on one side, and the skilful hindcasts and the WFDEI data set, on the  
824 other side. Unfortunately, we do not have a well-founded expectation about such  
825 a disagreement and, hence, we have no expectation about its effect on the  
826 difference between theoretical and actual skill. However, in Europe and beyond  
827 the first lead month almost all skill in the seasonal forecasts is due to the initial  
828 conditions (see the companion paper). Therefore, beyond the first lead month  
829 and in Europe differences in forcing during the hindcast period have a negligible  
830 effect on skill.
- 831 2. Models are imperfect, in terms of physics and in terms of spatial and temporal  
832 discretisation, so model hydrology differs from real world hydrology. Hindcasts  
833 and the pseudo-observations are produced with the same model, so  
834 imperfections in model hydrology are expected to lead to more theoretical than  
835 actual skill. One assumption implicitly made in the diagram is that the basin of  
836 the observation station and the model basin are identical. This is not the case  
837 (see Sect. 2.2), so differences between observation and model basin form an  
838 additional cause of disagreements between theoretical and actual skill. Again,  
839 this will favour theoretical skill with respect to actual skill since basins are  
840 identical in the hindcasts and the reference simulation. In particular, differences  
841 in meteorological forcing between the basin of the observation station and the  
842 model basin reduce actual skill. Van Dijk et al. (2013) investigated this aspect  
843 by making simulations for Australia at different spatial resolutions and verifying

844 with networks of observations with different spatial densities. They found that  
845 the resolution and perhaps the quality of the forcing data contributed at least half  
846 to the difference between theoretical and actual skill.

- 847 3. In the real world discharge observations are subject to measurement errors.  
848 Measurement errors of discharge are not constant over time (due to varying cross  
849 sectional areas, following erosion and sedimentation) and therefore add noise to  
850 the data; noise always reduces skill. There is no equivalent of this error in the  
851 model environment. Hence, as for differences 1) and 2) this difference is  
852 expected to lead to more theoretical than to actual skill.
- 853 4. Initial conditions are absent in this list of differences since in WUSHP they are  
854 not independent components but entirely determined by two components of the  
855 system listed above, namely meteorology and hydrology. Alternatively, initial  
856 hydrological conditions could be taken from observations or by assimilation of  
857 observations into model calculations. In that case, initial conditions would  
858 become an independent or semi-dependent component of the system. However,  
859 while model initial conditions would, of course, differ from real initial  
860 conditions, the two model system had identical initial conditions. Hence, this  
861 difference would again be expected to lead to more theoretical than to actual  
862 skill.

863

864 In summary, all of the conceptual differences between the generation of pseudo- and  
865 real observations are expected to lead to more theoretical skill than actual skill, except  
866 for the difference in meteorology during the hindcast period, which has, in the case of  
867 Europe beyond the first lead month, a neutral effect, and otherwise an unknown effect.

868

869 Our data analysis, Sect. 3.3, broadly confirms that theoretical skill exceeds actual skill.

870

871 It is interesting to discuss what would happen in the utopian case that the system of the  
872 model reference run would converge with the real world, i.e. if model meteorological  
873 forcing and hydrology would approach perfection and if measurement errors would  
874 approach zero. Equality of the two systems would, according to the analysis above, lead  
875 to equality of theoretical and actual skill. However, we like to note that at the same time  
876 optimisation of the model system can lead to a degradation of the theoretical skill if the  
877 hydrological models have unrealistic memory time scales in their storage compartments.  
878 If this memory, from stored water in either snow, soil or aquifer (or man-made reservoirs  
879 behind dams) , is too strong then skill will reduce with calibrating the model towards  
880 more realistic storage accumulation. However, if this memory is too small before  
881 improving the model, then, of course, the reverse may happen and skill increases with  
882 optimization.

883 An example proving this statement is a model that accumulates too much snow. The  
884 model will do so both in the initial state of the reference simulation and the initial state  
885 of the hindcasts and since more snow leads, at some stage of the melting season, to more  
886 predictive skill, theoretical skill will be overestimated. A perfect model, accumulating

887 less but more realistic amounts of snow, would exhibit less skill. Another example is  
888 predictive skill caused by interannual variations in the initial amount of soil moisture  
889 and/or groundwater. A model that is imperfect because it overestimates the transport  
890 speed of water through the soil and the groundwater reservoirs will do so both in the  
891 reference simulation and the hindcasts. Predictive skill due to soil moisture initial  
892 conditions will then occur too early. Compared to the model that overestimates transport  
893 speed, a perfect model with smaller, realistic transport speed would yield less theoretical  
894 skill at the early lead times.

895 Hence, theoretical skill is not equal to the maximum that could be accomplished if  
896 hydrological model and meteorological forcing during the reference simulation were  
897 perfect.

898 The version of VIC used in this study was calibrated by Nijssen et al. (2001) in a crude  
899 way, in the sense that they assumed no spatial variation of the parameters set by  
900 calibration within almost the entire European continent. Improving the calibration of  
901 VIC would be an obvious candidate for trying to improve the seasonal predictions  
902 discussed in this paper. This should lead to higher actual skill. However, the two  
903 examples discussed in the previous paragraph show that theoretical skill may actually,  
904 for certain locations, months of initialisation and lead months, decline due to the  
905 recalibration.

906

## 907 **4.2 Results and uncertainties**

908 There seems to be a broad correspondence between the probabilistic forecast  
909 verification presented here and the model validation presented in Greuell et al. 2015;  
910 and Roudier et al. 2016. These studies found that average discharge and inter-annual  
911 variations therein are well reproduced against observations, consistent with our result  
912 that all skill scores -also against real observations (see Fig. S4 for the lead 0 results)- are  
913 good for large parts of Europe in the first lead month. Their finding that high flows are  
914 generally better reproduced than low flows seems to contradict with our fact that BN  
915 forecasts are more skilful than AN forecasts (although by a small margin, and for lead  
916 0 mostly so in southern Europe, Fig. S1 e.g. Initialisation April). This discrepancy may  
917 be due to different definitions of high or low flows between these studies and the present  
918 one. They define high and low flows by Q95 and Q5 based on daily discharge,  
919 respectively, while here we use Q66 and Q33 based on monthly discharge, much less  
920 extreme values. Also, their study showed that the variability in Q5 was more  
921 overestimated than the variability in Q95, which may be a reason for the higher skill we  
922 find in the lower tercile (skill requires variability, see discussion of companion paper),  
923 though this inference is hard to prove. This prior work also invokes some warnings.  
924 Greuell et al. (2105) found that seasonal flow cycles show a too late and too broad spring  
925 peak in (northern) Fennoscandia. This suggests that our theoretical forecast skills may  
926 also be too high at too long lead times in that region and season, (as was also already  
927 revealed by comparing Figure 7b vs 7c).

928 In a future extension of our work, an objective method like cluster analysis could reveal  
929 regions where skill has a similar signature. This could lead to an improved assessment  
930 of the physical and climatological factors that are responsible for the spatial variations  
931 in skill found in this and its companion paper.

932 There also seems to be a broad correspondence between the regions and seasons with  
933 skill identified in the present work, with that from more spatially or temporally confined  
934 studies based on entirely different physical or even statistical models. Without repeating  
935 the more detailed description in the Introduction and closer comparison in Sect. 3.1, we  
936 restate here that the results of Bierkens and van Beek, (2009) and Thober et al. (2015)  
937 were similar at the European domain. These pan-European studies, like ours, confirm  
938 more regional studies such as for the British Isles (Svensson et al., 2015), Iberia (Trigo,  
939 2004) or France (Céron et al., 2010; Singla et al., 2012). Though a high resolution study  
940 like the latter may add much spatial detail, this does not change the region and season  
941 of skill.

942 Our results are based on a forcing with the 15 member, monthly initialized, 7 month  
943 forecast version of ECMWF System 4, basically because at the start of this work that  
944 hindcast was the only one accessible to us but also because it allows verification at the  
945 highest temporal resolution. Alternatively, we could have used the 51 member  
946 seasonally initialised (4 times per year), 7 month forecast version of the same model.  
947 That would have provided us with better constrained, more precise statistics (larger  
948 sample size), or would have allowed assessment of more percentiles (e.g. quintiles  
949 instead of terciles) at similar precision. However, the variation of skill over a year would  
950 not have been resolved with such detail as in the present work. Finally, a 15 member,  
951 seasonally initialized, 12 month forecast version is available. Our results show that for  
952 some regions at lead month 6 still a few, small pockets of persistent skill remain,  
953 suggesting that extending the forecast for our domain might be worth exploring.

954 Other seasonal forecasting systems, based on different coupled ocean-climate models,  
955 exist that could have been used, such as CFSv2 (Saha et al., 2014), GloSea5  
956 (MacLachlan et al., 2014). Given that, at least at large scales, multi model ensembles  
957 exhibit better climate forecast skill, it is interesting to investigate if that additional skill  
958 also propagates into river flow forecasts. While this seems to be true for the Eastern  
959 United States (Luo & Wood, 2008) it is not known if similar conclusions could be drawn  
960 for Europe. A similar reasoning can also be extended to the hydrological models: using  
961 a multi climate model ensemble to force a multi hydrological model ensemble might  
962 also provide improved skill, as the latter models may be complementary in the regions  
963 and seasons of best model performance. Bohn et al. (2010) showed some advantage of  
964 using an ensemble of three hydrological models (but with a single forcing), over using  
965 only the best of the three, but only after bias correcting the hydrological output and  
966 making a linear combination of them with monthly varying weights.

967

968

### 969 **4.3 Implications and recommendations**

970

971 Many conclusions drawn from this work are valid at the scale of our domain and not  
972 necessarily at the scale of river basins. Only in some parts of our analysis, especially  
973 where we focused on the annual cycle of the skill (Fig. 2), regional patterns at a scale  
974 smaller than that of the domain were discussed. This was done in a qualitative way.

975

976 For applications of these seasonal forecasts in decision making processes at (sub) basin  
977 level, a more detailed skill analysis is recommended for that specific (sub)basin,  
978 preferably after a better model calibration for that same basin. The facts presented in  
979 this study that anomaly correlations and ROC scores for the AN and BN terciles are  
980 significant for large parts of the domain several lead months in advance, supported by  
981 (fairly) positive validation results for interannual variability of high and low flows  
982 (Greuell et al., 2015; Roudier et al. 2016), suggest these anomaly forecasts are good  
983 enough to be used as such. However, areas of significant RPSS are much smaller and  
984 remain significant for shorter lead times. Spatially distributed calibration of VIC model  
985 parameters, or distribution based calibration of modelled discharge to observed, or both,  
986 might also increase the RPSS. This might then allow forecasting of absolute discharge  
987 magnitudes and thus inform decision making processes that involve certain absolute  
988 discharge thresholds.

989 In the respective Result sections we already discussed the probable reasons for skill,  
990 which are much elaborated on in the companion paper. In general that paper shows that  
991 for most areas skill in runoff is caused by initialising snow and /or soil moisture  
992 properly, only in few areas and seasons skill in precipitation or skill in temperature and  
993 evapotranspiration adds to that beyond the first lead month. This has two implications:  
994 one is that, if ever the skill of seasonal climate forecasts improves for Europe, this may  
995 well translate to improved seasonal river flow forecast too. The second is that better  
996 initial conditions of snow water equivalent and soil moisture from observations may do  
997 the same, but the latter only if the spatial distribution of the soil moisture storage  
998 capacity is more realistic too (see Sect. 4.1).

999

1000 Overall the present analysis shows that especially in winter, spring and early summer,  
1001 there is potentially good skill to forecast runoff and discharge in large parts of Europe,  
1002 with considerable lead time. While this broadly confirms previously published work,  
1003 the present study (while being specific to our model setup) gives much more spatial and  
1004 temporal (season and lead time) details. As such it provides a good basis to support  
1005 operational forecasts and to add information about skill to seasonal forecasts, which is  
1006 very important for proper value assessment and decision making.

1007

1008

1009

1010

1011 **5 Conclusions**

1012

1013 This paper is the first of two papers dealing with a model-based system built to produce  
1014 seasonal hydrological forecasts (WUSHP: Wageningen University Seamless  
1015 Hydrological Predictions). The present paper presents the development and the skill  
1016 evaluation of the system for Europe, the companion paper provides an explanation of  
1017 the skill or the lack of skill.

1018

1019 First, “theoretical skill” of the runoff hindcasts was determined taking the output of the  
1020 reference simulation as “pseudo-observations”. Using the correlation coefficient (R) as  
1021 metric, hot spots of significant skill were found in Fennoscandia (from January to  
1022 October), the southern part of the Mediterranean (from June to August), Poland,  
1023 northern Germany, Romania and Bulgaria (mainly from November to January) and  
1024 western France (from December to May). There is very little or no significant skill all  
1025 over the year in some coastal and mountain regions. The entire British Isles exhibit very  
1026 little skill, except for the eastern coast of Great Britain. If the entire domain is  
1027 considered, the annual cycle of skill has a minimum roughly from August to November  
1028 and a maximum in May.

1029

1030 Runoff and discharge show a high degree of similarity in terms of the spatial patterns  
1031 and the magnitude of the skill. However, when averaged over the domain and the year,  
1032 predictability is slightly higher for discharge than for runoff for the first lead month (by  
1033 0.049 in terms of R). The difference then decreases with increasing lead time. These  
1034 tendencies can be ascribed to the combined effect of the delay between runoff and  
1035 discharge and the fact that skill decreases with lead time. We also found that the  
1036 difference between discharge and runoff skill increases with the size of the basin.

1037

1038 Theoretical skill as determined with the pseudo-observations was compared to actual  
1039 skill as determined with real discharge observations. On average across all target months  
1040 and for lead month 2, the ratio of actual to theoretical skill in terms of the domain-mean  
1041 R is 0.67 (0.26 divided by 0.39) for large basins and 0.54 (0.16 divided by 0.29) for  
1042 small basins. So, skill reduction due to replacing pseudo- by real observations is larger  
1043 for small basins than for large basins. For 10 day flow forecasts Alfieri et al. (2014) also  
1044 found that, especially in mountain areas, performance drops significantly in river basins  
1045 with upstream area smaller than 300 km<sup>2</sup>.

1046

1047 Spatio-temporal patterns for the different skill metrics considered in this study  
1048 (correlation coefficient, ROC area and Ranked Probability Skill Score) are similar to a  
1049 large degree. ROC areas tend to be slightly larger for the below normal than for the  
1050 above normal tercile but not during target months from October to January.

1051

1052

1053 **6 Author Contributions**

1054

1055 Greuell and Hutjes designed the experiments, with suggestions from the other co-  
1056 authors. Franssen and Greuell developed the workflow scripts and performed all the  
1057 simulations. Greuell and Franssen developed the analyses and plotting scripts in R.  
1058 Biemans did the LPJmL work on AAPFD. All co-authors participated in repeated  
1059 discussions on interpretations of results and suggested ways forward in the analysis.  
1060 Greuell prepared the first version of the manuscript with contributions from all co-  
1061 authors. Hutjes prepared the revision of the manuscript with contributions from all co-  
1062 authors.

1063

1064 **7 Conflicting Interests**

1065 The authors declare that they have no conflict of interest.

1066

1067 **8 Acknowledgments**

1068 This study was financially supported by the EUPORIAS project (EUropean Provision  
1069 of Regional Impact Assessment on Seasonal-to-decadal timescale); grant agreement No.  
1070 308291, funded by the European Commission (EU) project in the Seventh Framework  
1071 Programme. We thank the valued suggestions and insightful comments from two  
1072 (anonymous) reviewers that contributed to an improved version of the manuscript.

1073

1074 **9 References**

1075

1076 Alfieri, L., Pappenberger, F., Wetterhall, F., Haiden, T., Richardson, D., & Salamon, P.  
1077 (2014). Evaluation of ensemble streamflow predictions in Europe. *Journal of*  
1078 *Hydrology*, 517, 913-922.

1079 Bhend, J., Ripoldi, J., Mignani, C., Mahlstein, I., Hiller, R., Spirig, C., Liniger, M.,  
1080 Weigel, A., Bedia Jimenez, J., De Felice, M., Siegert, S., (2016) easyVerification:  
1081 Ensemble Forecast Verification for Large Data Sets. [https://CRAN.R-](https://CRAN.R-project.org/package=easyVerification)  
1082 [project.org/package=easyVerification](https://CRAN.R-project.org/package=easyVerification)

1083 Bierkens, M. F. P., & Van Beek, L. P. H. (2009). Seasonal predictability of European  
1084 discharge: NAO and hydrological response time. *Journal of Hydrometeorology*, 10(4),  
1085 953-968.

1086 Bohn, T. J., Sonessa, M. Y., Lettenmaier, D. P. (2010). "Seasonal Hydrologic  
1087 Forecasting: Do Multimodel Ensemble Averages Always Yield Improvements in  
1088 Forecast Skill?" *Journal of Hydrometeorology* 11(6): 1358-1372.

- 1089 Bruno Soares, M. and S. Dessai (2016). Barriers and enablers to the use of seasonal  
1090 climate forecasts amongst organisations in Europe. *Climatic Change* 137(1): 89-103.
- 1091 Céron, J. P., Tanguy, G., Franchistéguy, L., Martin, E., Regimbeau, F., Vidal, J. -P.  
1092 (2010). "Hydrological seasonal forecast over France: feasibility and prospects."  
1093 *Atmospheric Science Letters* 11(2): 78-82.
- 1094 Crochemore, L., Ramos, M. H. and Pappenberger, F.. (2016). Bias correcting  
1095 precipitation forecasts to improve the skill of seasonal streamflow forecasts. *Hydrol.*  
1096 *Earth Syst. Sci.* 20(9): 3601-3618.
- 1097 Demirel, M.C., Booij, M.J. and Hoekstra, A.Y. (2015). The skill of seasonal ensemble  
1098 low-flow forecasts in the Moselle River for three different hydrological models. *Hydrol.*  
1099 *Earth Syst. Sci.*, 19, 275–291, 2015
- 1100 Doblas-Reyes, F. J., García-Serrano, J., Lienert, F., Biescas, A. P., & Rodrigues, L. R.  
1101 (2013). Seasonal climate predictability and forecasting: status and prospects. *Wiley*  
1102 *Interdisciplinary Reviews: Climate Change*, 4(4), 245-268.
- 1103 Döll, P., & Lehner, B. (2002). Validation of a new global 30-min drainage direction  
1104 map. *Journal of Hydrology*, 258(1), 214-231.
- 1105 ECMWF Seasonal Forecast User Guide, retrieved from:  
1106 [http://www.ecmwf.int/en/forecasts/documentation-and-support/long-range/seasonal-](http://www.ecmwf.int/en/forecasts/documentation-and-support/long-range/seasonal-forecast-documentation/user-guide/introduction)  
1107 [forecast-documentation/user-guide/introduction](http://www.ecmwf.int/en/forecasts/documentation-and-support/long-range/seasonal-forecast-documentation/user-guide/introduction)
- 1108 Ghile, Y. B., and Schulze, R. E., 2008: Development of a framework for an integrated  
1109 time-varying agrohydrological forecast system for southern Africa: Initial results
- 1110 Greuell, W., Andersson, J. C. M., Donnelly, C., Feyen, L., Gerten, D., Ludwig, F., ... &  
1111 Schaphoff, S. (2015). Evaluation of five hydrological models across Europe and their  
1112 suitability for making projections of climate change. *Hydrol. Earth Syst. Sci. Discuss.*,  
1113 12, 10289-10330, doi: 10.5194/hessd-12-10289-2015
- 1114 Greuell, W., W. H. P. Franssen, H. Biemans and R. W. A. Hutjes. Seasonal streamflow  
1115 forecasts for Europe – II. Explanation of the skill. (2016, in revision) to *Hydrol. Earth*  
1116 *Syst. Sci.*, doi: 10.5194/hess-2016-604
- 1117 Hagedorn, R., Doblas-Reyes, F. J., & Palmer, T. N. (2005). The rationale behind the  
1118 success of multi-model ensembles in seasonal forecasting–I. Basic concept. *Tellus A*,  
1119 57(3), 219-233.
- 1120 Hamlet, A. F., Huppert, D., & Lettenmaier, D. P. (2002). Economic value of long-lead  
1121 streamflow forecasts for Columbia River hydropower. *Journal of Water Resources*  
1122 *Planning and Management*, 128(2), 91-101.
- 1123 Haylock, M. R., Hofstra, N., Klein Tank, A. M. G., Klok, E. J., Jones, P. D., & New, M.  
1124 (2008). A European daily high-resolution gridded data set of surface temperature and



- 1125 precipitation for 1950–2006. *Journal of Geophysical Research: Atmospheres* (1984–  
1126 2012), 113(D20).
- 1127 Juston, J., Jansson, P. E., & Gustafsson, D. (2014). Rating curve uncertainty and change  
1128 detection in discharge time series: case study with 44-year historic data from the  
1129 Nyangores River, Kenya. *Hydrological Processes*, 28(4), 2509-2523.
- 1130 Koster, R. D., Mahanama, S. P., Livneh, B., Lettenmaier, D. P., & Reichle, R. H. (2010).  
1131 Skill in streamflow forecasts derived from large-scale estimates of soil moisture and  
1132 snow. *Nature Geoscience*, 3(9), 613-616.
- 1133 Lehner, B., Reidy Liermann, C., Revenga, C., Vörösmarty, C., Fekete, B., Crouzet, P.,  
1134 Döll, P., Endejan, M., Frenken, K., Magome, J., Nilsson, C., Robertson, J., Rödel, R.,  
1135 Sindorf, N., Wisser, D. (2011): High resolution mapping of the world’s reservoirs  
1136 and dams for sustainable river flow management. *Frontiers in Ecology and the*  
1137 *Environment* 9(9): 494–502.
- 1138
- 1139 Li, H., Luo, L. and Wood, E.F. (2008). Seasonal hydrologic predictions of low-flow  
1140 conditions over eastern USA during the 2007 drought. *Atmospheric Science Letters*  
1141 **9**(2): 61-66.
- 1142 Liang, X., Lettenmaier, D. P., Wood, E. F., & Burges, S. J. (1994). A simple  
1143 hydrologically based model of land surface water and energy fluxes for general  
1144 circulation models. *Journal of Geophysical Research: Atmospheres* (1984–2012),  
1145 99(D7), 14415-14428.
- 1146 Luo, L. and E.F. Wood, 2008: Use of Bayesian Merging Techniques in a Multimodel  
1147 Seasonal Hydrologic Ensemble Prediction System for the Eastern United States. *J.*  
1148 *Hydrometeor.*, 9, 866–884, doi: 10.1175/2008JHM980.1
- 1149 MacLachlan, C., A. Arribas, K. A. Peterson, A. Maidens, D. Fereday, A. A. Scaife, M.  
1150 Gordon, M. Vellinga, A. Williams, R. E. Comer, J. Camp, P. Xavier and G. Madec,  
1151 2014. Global Seasonal forecast system version 5 (GloSea5): a high-resolution seasonal  
1152 forecast system. *QJR Meteorol Soc*, doi:10.1002/qj.2396.
- 1153 Marchant, R., & Hehir, G. (2002). The use of AUSRIVAS predictive models to assess  
1154 the response of lotic macroinvertebrates to dams in south-east Australia. *Freshwater*  
1155 *Biology*, 47(5), 1033-1050.
- 1156 Mason, S. J., & Stephenson, D. B. (2008). How do we know whether seasonal climate  
1157 forecasts are any good?. In *Seasonal Climate: Forecasting and Managing Risk* (pp. 259-  
1158 289). Springer Netherlands.
- 1159 Mo, K. C., & Lettenmaier, D. P. (2014). Hydrologic prediction over the conterminous  
1160 United States using the national multi-model ensemble. *Journal of Hydrometeorology*,  
1161 15(4), 1457-1472.

- 1162 Molteni, F., Stockdale, T., Balmaseda, M., Balsamo, G., Buizza, R., Ferranti, L.,  
1163 Magnusson, L., Mogensen, K., Palmer, T., Vitart, F. (2011). The new ECMWF seasonal  
1164 forecast system (System 4). ECMWF Technical Memorandum 656.
- 1165 Mushtaq, S., Chen, C., Hafeez, M., Maroulis, J. and Gabriel, H., 2012: The economic  
1166 value of improved agrometeorological information to irrigators amid climate variability.  
1167 *Int. J. Climatol.*, 32, 567–581.
- 1168 Nijssen, B., O'Donnell, G. M., Lettenmaier, D. P., Lohmann, D., & Wood, E. F. (2001).  
1169 Predicting the discharge of global rivers. *Journal of Climate*, 14(15), 3307-3323.
- 1170 Rost, S., Gerten, D., Bondeau, A., Lucht, W., Rohwer, J., & Schaphoff, S. (2008).  
1171 Agricultural green and blue water consumption and its influence on the global water  
1172 system. *Water Resources Research*, 44(9), doi 10.1029/2007WR006331.
- 1173 Roudier, P., Andersson, J.C.M., Donnelly, C., Feyen, L., Greuell, W. and Ludwig, F.,  
1174 (2016). "Projections of future floods and hydrological droughts in Europe under a +2°C  
1175 global warming." *Climatic Change* 135(2): 341-355.
- 1176 Saha, Suranjana and Coauthors, 2014: The NCEP Climate Forecast System Version 2  
1177 *Journal of Climate* J. Climate, 27, 2185–2208. doi: 10.1175/JCLI-D-12-00823.1
- 1178 Schaphoff, S., Heyder, U., Ostberg, S., Gerten, D., Heinke, J., & Lucht, W. (2013).  
1179 Contribution of permafrost soils to the global carbon budget. *Environmental Research*  
1180 *Letters*, 8(1), 014026, doi:10.1088/1748-9326/8/1/014026.
- 1181 Sheffield, J., Wood, E. F., Chaney, N., Guan, K., Sadri, S., Yuan, X., ... & Ogallo, L.  
1182 (2014). A drought monitoring and forecasting system for sub-Saharan African water  
1183 resources and food security. *Bulletin of the American Meteorological Society*, 95(6),  
1184 861-882.
- 1185 Shukla, S. and Lettenmaier, D. P. (2011). "Seasonal hydrologic prediction in the United  
1186 States: understanding the role of initial hydrologic conditions and seasonal climate  
1187 forecast skill." *Hydrol. Earth Syst. Sci.* 15(11): 3529-3538.
- 1188 Shukla, S., McNally, A., Husak, G., & Funk, C. (2014). A seasonal agricultural drought  
1189 forecast system for food-insecure regions of East Africa. *Hydrology and Earth System*  
1190 *Sciences*, 18(10), 3907-3921.
- 1191 Siegart, S., Bhend, J., Kroener, I., De Felice, M. (2014). SpecsVerification: Forecast  
1192 Verification Routines for Ensemble Forecasts of Weather and Climate. [https://CRAN.R-](https://CRAN.R-project.org/package=SpecsVerification)  
1193 [project.org/package=SpecsVerification](https://CRAN.R-project.org/package=SpecsVerification)
- 1194 Singla, S., Céron, J. P., Martin, E., Regimbeau, F., Déqué, M., Habets, F., & Vidal, J.  
1195 P. (2012). Predictability of soil moisture and river flows over France for the spring  
1196 season. *Hydrology & Earth System Sciences*, 16: 201-216.

- 1197 Svensson, C., Brookshaw, A., Scaife, A. A., Bell, V. A., Mackay, J. D., Jackson, C. R.,  
1198 Hannaford, J., Davies, H. N., Arribas A., Stanley, S. (2015). "Long-range forecasts of  
1199 UK winter hydrology." *Environmental Research Letters* 10(6): 064006.
- 1200 Themeßl, M. J., Gobiet, A., & Leuprecht, A. (2011). Empirical-statistical downscaling  
1201 and error correction of daily precipitation from regional climate models. *International*  
1202 *Journal of Climatology*, 31(10), 1530-1544.
- 1203 Thober, S., Kumar, R., Sheffield, J., Mai, J., Schäfer, D., and Samaniegoet, L. (2015).  
1204 "Seasonal Soil Moisture Drought Prediction over Europe Using the North American  
1205 Multi-Model Ensemble (NMME)." *Journal of Hydrometeorology* 16(6): 2329-2344.
- 1206 Trigo, R. M., Pozo-Vázquez, D., Osborn, T.J., Castro-Díez, Y., Gámiz-Fortis, S.,  
1207 Esteban-Parra, M.J. (2004). "North Atlantic oscillation influence on precipitation, river  
1208 flow and water resources in the Iberian Peninsula." *International Journal of Climatology*  
1209 24(8): 925-944.
- 1210 Van Dijk, A. I., Peña-Arancibia, J. L., Wood, E. F., Sheffield, J., & Beck, H. E. (2013).  
1211 Global analysis of seasonal streamflow predictability using an ensemble prediction  
1212 system and observations from 6192 small basins worldwide. *Water Resources Research*,  
1213 49(5), 2729-2746.
- 1214 Van Dijk, A. I. J. M., and G. A. Warren (2010), AWRA Technical Report 4, Evaluation  
1215 Against Observations, WIRADA/CSIROWater for a Healthy Country Flagship,  
1216 Canberra. [http://www.clw.csiro.au/publications/waterforahealthycountry/2010/wfhc-  
1217 awras-evaluationagainstobservations.pdf](http://www.clw.csiro.au/publications/waterforahealthycountry/2010/wfhc-awras-evaluationagainstobservations.pdf)
- 1218 Weedon, G. P., Balsamo, G., Bellouin, N., Gomes, S., Best, M. J., & Viterbo, P. (2014).  
1219 The WFDEI meteorological forcing data set: WATCH Forcing Data methodology  
1220 applied to ERA-Interim reanalysis data. *Water Resources Research*, 50(9), 7505-7514.
- 1221 Wood, A. W., & Lettenmaier, D. P. (2006). A test bed for new seasonal hydrologic  
1222 forecasting approaches in the western United States. *Bulletin of the American*  
1223 *Meteorological Society*, 87(12), 1699.
- 1224 Wood, A. W., Hopson, T., Newman, A., Brekke, L., Arnold, J. and Clark, M. (2016).  
1225 Quantifying Streamflow Forecast Skill Elasticity to Initial Condition and Climate  
1226 Prediction Skill. *Journal of Hydrometeorology* 17(2): 651-668.
- 1227 Yuan, X., Wood, E. F., Luo, L., & Pan, M. (2013). CFSv2-based seasonal hydroclimatic  
1228 forecasts over the conterminous United States. *Journal of Climate*, 26, 4828-4847.
- 1229 Yuan, X., Wood, E. F., & Ma, Z. (2015). A review on climate-model-based seasonal  
1230 hydrologic forecasting: physical understanding and system development. *Wiley*  
1231 *Interdisciplinary Reviews: Water*, 2(5), 523-536.
- 1232