1	Seasonal streamflow forecasts for Europe – II. Sources of skill
2	
3	Wouter Greuell, Wietse H. P. Franssen and Ronald W. A. Hutjes
4	
5	Wageningen University and Research
6	
7	all authors:
8 9	Water Systems and Global Change (WSG) group, Wageningen University and Research, Droevendaalsesteeg 3, NL 6708 PB Wageningen, Netherlands
10	
11	correspondence to wouter.greuell@wur.nl
12	
13	
14	

15 Abstract

16

Seasonal forecasts can be exploited to optimize hydropower energy generation, 17 navigability of rivers and irrigation management to decrease crop yield losses. This 18 paper is the second of two papers dealing with a model-based system built to produce 19 seasonal hydrological forecasts (WUSHP: Wageningen University Seamless 20 Hydrological Prediction system), applied here to Europe. In WUSHP, hydrology is 21 simulated by running the Variable Infiltration Capacity (VIC) hydrological model with 22 meteorological forcing from bias-corrected output of ECMWF's Seasonal Forecasting 23 System 4 (S4). WUSHP is probabilistic. For the assessment of skill, hindcasts (1981-24 25 2010) were generated. Whereas the first paper presented the development and the skill evaluation of the system, this paper provides explanations for the skill. 26

27

To that purpose, we first analysed the forcing and found considerable skill in the 28 precipitation forecasts for the first lead month but hardly any significant skill for 29 subsequent lead months. Seasonal forecasts of temperature have more skill. Skill in 30 summer temperature is related to climate change and more or less independent of lead 31 time. Skill in February and March is unrelated to climate change. Sources of skill in 32 runoff were isolated with a suite of specific hindcasts. These revealed that, beyond the 33 second lead month, streamflow hindcasts with meteorological forcing that is identical 34 35 for all years (InitSH) have more skill in runoff than the streamflow hindcasts forced with S4 output (FullSH). This occurs because interannual variability of the S4 forcing 36 has hardly any skill while it adds noise. Other specific hindcasts show that in Europe 37 initial conditions of soil moisture form the dominant source of skill in runoff. From 38 April to July, at the end of the melt season, initial conditions of snow contribute 39 significantly to the skill, provided forecasts do not start earlier than in February. Some 40 remarkable skill features are due to indirect effects, i.e. skill due to forcing or initial 41 42 conditions of snow and soil moisture at an earlier stage is stored in the hydrological state (snow and/or soil moisture) of a later stage, which then contributes to persistence 43 of skill. Finally, predictability of evapotranspiration was analysed in some detail, 44 leading among others to the conclusion that its skill originates mostly in the 45 46 meteorological forcing.

- 48 **1** Introduction
- 49

Society may benefit from seasonal hydrological forecasts (Viel et al., 2016; Soares and Dessai, 2016; Crochemore et al., 2016), i.e. hydrological forecasts for future time periods from more than two weeks up to about a year (Doblas-Reyes et al., 2013). Such predictions can be exploited to optimize e.g. hydropower energy generation (Hamlet et al. 2002), navigability of rivers in low flow conditions (Li, et al., 2008) and irrigation management (Ghile and Schulze 2008; Mushtaq et al. 2012) to decrease crop yield losses.

57

This is the second paper about seasonal hydrological forecasts for Europe produced with WUSHP (Wageningen University Seamless Hydrological Prediction system), a dynamical (i.e. model-based) system. In summary, the forecasts of WUSHP are made with the Variable Infiltration Capacity (VIC) hydrological model, which uses biascorrected output of forecasts from ECMWF's Seasonal Forecast System 4 (S4) as meteorological forcing. The system is probabilistic.

64

In the present and in the companion paper (Greuell et al. 2018), WUSHP is used as a 65 research tool for purposes of academic interest. In the companion paper, the set-up of 66 67 WUSHP has been described and spatial and temporal variations of skill, or lack thereof, in runoff and discharge in Europe have been established by means of hindcasts. 68 69 Significant skill was found for many regions, varying by initialisation and target months. For lead month 2, hot spots of significant skill in runoff are situated in 70 Fennoscandia (for target months from January to October), the southern part of the 71 72 Mediterranean (from June to August), Poland, Northern Germany, Romania and Bulgaria (mainly from November to January) and Western France (from December to 73 74 May). In general, the spatial pattern of significant skill in runoff was found to be fixed 75 in space while the skill decreased in magnitude with increasing lead time. Some significant skill remained even at the end of the hindcasts (7 months). 76

77

The current paper aims to identify the sources of the skill in WUSHP and is structured 78 79 in two main parts. In the first part, an analysis of the skill in the most important meteorological forcing variables (precipitation, two-meter temperature and incoming 80 short-wave radiation from S4) is presented. For S4, this was done earlier by Kim et al. 81 (2012) for the boreal winter months (DJF) with initialisation on the first of November. 82 For that case, they found that in Europe S4 has no skill in the precipitation forecasts 83 and some skill in the temperature forecasts for Southern Sweden, Southern Finland, the 84 region south-east of Saint Petersburg and Northern Germany. Scaife et al. (2014) 85 analysed the skill for the same target months and starting date but with another 86 prediction system, namely the Met Office Global Seasonal forecast System 5 87 (GloSea5). Scaife et al. (2014) found that the GloSea5 temperature forecasts for Europe 88

contain hardly any significant skill but that GloSea5 forecasts of the North Atlantic
Oscillation are correlated significantly with observed temperatures in northern and
southern Europe. This means that there is untapped predictability in the GloSea5
forecasts. We will analyse predictability of the mentioned output variables of S4 for the
whole continent and will consider all combination of lead and target months.

94

The second line of analysis aims to investigate the reasons for presence or absence of 95 skill by means of a series of specific hindcasts that isolate potential sources of skill, 96 namely meteorological forcing, the initial conditions of soil moisture and the initial 97 conditions of snow. The 30 years of standard hindcasts produced by WUSHP, analysed 98 in the companion paper and referred to as Full Streamflow Hindcasts (FullSH; climate-99 model-based hindcasts" according to Yuan et al., 2015) constitute the starting point. 100 Then, a suite of specific hindcasts is carried out and evaluated, an approach explored 101 earlier by Wood et al. (2005), Bierkens and Van Beek (2009) and Koster et al. (2010). 102 Each specific hindcast is largely identical to the FullSH but one or two of the sources of 103 predictability are isolated by eliminating the effect of all of the other sources through 104 removal of their interannual variation. In the ensuing analysis the skills in runoff found 105 in the different specific hindcasts will be compared among themselves and with the 106 skill from the FullSH. 107

108

These specific hindcasts are related to the conventional Ensemble Streamflow 109 110 Prediction (ESP) technique (e.g. Wood and Lettenmaier, 2008, Shukla and Lettenmaier, 2011, Singla et al., 2011), which can, like our specific hindcasts, be used to isolate 111 sources of skill. However, in ESP and its variant reverse-ESP the meteorological 112 forcing is taken from data based on observations, and not from meteorological 113 hindcasts, as in the present study. In fact, we also produced ESPs. In Sect. 4.3 we will 114 compare these with one of the other specific hindcasts and more generally discuss the 115 relation between our specific hindcasts and the ESP suite. 116

117

Since evapotranspiration has a large effect on runoff, the analysis is complemented with an analysis of the skill in this variable. Predictions of evapotranspiration also have independent value because they are useful for planning of water level control in polders and for planning of water use for irrigation and fertiliser application. As for runoff, we will exploit the specific hindcasts to isolate the different sources of predictability in evapotranspiration forecasts.

124

Bierkens and van Beek (2009) investigated the sources of runoff skill for seasonal hydrological forecasting over Europe. They found that in winter initial conditions constitute the dominant source while in summer meteorological forcing and initial conditions are equally important. Singla et al. (2011) assessed the skill of hydrological predictions for France and concluded that over most plains the predictability of hydrological variables primarily depended on forcing, whereas it mainly depended on
snow cover over high mountains. The Seine catchment area was an exception as the
skill mainly came from the initial state of its large and complex aquifers.

133

The version of VIC that we used was only crudely calibrated (by Nijssen et al., 2001). 134 135 Hence, streamflow computed by the present version of the system may be expected to deviate substantially from observations, both in terms of the mean and in terms of the 136 spread of the ensemble of forecasts. Also, within WUSHP no post-processing of 137 discharge is carried out to correct for such deficiencies. This makes the system 138 unsuitable to issue forecasts of absolute amounts of discharge but the system can be 139 140 used to provide information on how likely it is that in a coming month or season discharge will be above or below normal. Consequently, the most important criteria for 141 the selection of skill metrics (see Sect. 2.2) are their ability of discrimination, and their 142 insensitivity to biases and to the spread of the forecasts. 143

144

So, the objective of the present paper is to analyse, at a pan-European and at regional 145 scale, the sources of probabilistic skill of seasonal hydrological forecasts produced by 146 WUSHP. The next section (Sect. 2) will describe the seasonal prediction system itself, 147 the analysis approach as well as details of the various specific hindcast performed. We 148 will present the skill in three variables of the meteorological forcing (Sect. 3.1), 149 followed by skill in runoff found in the various specific hindcasts, which allows 150 151 attribution to either forcing or different types of initial conditions (Sect. 3.2), and 152 finally an analysis of skill in evapotranspiration (Sect. 3.3). We conclude with a discussion (Sect. 4) and conclusions (Sect. 5). 153

- 154
- 155

157

156 2 System and methods

- 158 2.1 The forecast system
- 159

160 The forecasts of WUSHP combine three elements, namely meteorological forcing from ECMWF's Seasonal Forecast System 4 (Molteni et al., 2011), bias correction of the 161 meteorological forcing with the quantile mapping method of Themeßl et al. (2011) and 162 simulations with the Variable Infiltration Capacity (VIC) hydrological model (Liang at 163 al., 1994). The skill of the system was assessed with hindcasts. These cover the period 164 1981-2010, were initialised on the first day of each month and have a length of seven 165 months. The system is probabilistic (15 members), so each set of hindcasts consists a 166 total of 5400 runs (30 years * 12 months * 15 members). In addition a single reference 167 simulation was performed, in which VIC was run with a gridded data set of model-168 assimilated meteorological observations, namely the WATCH Forcing Data Era-169 Interim (WFDEI; Weedon et al., 2014). The reference simulation has a dual aim. The 170

first aim is to create initialisation states for the hindcasts. Secondly, the output of the 171 reference simulation, e.g. discharge, is used for verification of the hindcasts. This 172 output will be named "pseudo-observations" here. 173

174

To spin up discharge, each 7-month hindcast was preceded by a one month simulation 175 176 with WFDEI forcing. All hindcasts and simulations were performed on a $0.5^{\circ} \times 0.5^{\circ}$ grid in natural flow mode, i.e. river regulation, irrigation and other anthropogenic 177 influences were not considered. VIC is run with a time step of 3 hours. More details 178 about the set-up of the system and the hindcasts can be found in the companion paper 179 (Greuell et al., 2018). 180

- 181
- 182

184

183

2.2 Methods of analysis and observations

In this paper we analyse hindcasts of runoff, discharge and evapotranspiration. Runoff 185 is defined as the amount of water leaving the model soil either along the surface or at 186 the bottom, while we define discharge as the flow of water through the largest river in 187 each grid cell. 188

189

Discrimination skill (briefly skill from now on) is measured in terms of the correlation 190 coefficient between the median of the hindcasts and the observations (R). We will 191 192 designate R-values as significant for p-values less than 0.05. We also considered metrics designed for the evaluation of categorical forecasts (terciles), namely the 193 Relative Operating Characteristics (ROC) area and the Ranked Probability Skill Score 194 (RPSS). The thresholds used for assigning individual observations to terciles were 195 determined from the observations themselves. Similarly hindcasts were assigned to 196 terciles by reference to themselves. Due to this strategy metrics are unaffected by 197 biases, a desired property (see Sect. 1). In the companion paper skills in terms of the 198 considered metrics were compared and it was found that for all combinations of target 199 and lead month the skill patterns in the maps were similar to a high degree. For that 200 reason we selected only one of them (R) for this paper. 201

202

Unless mentioned otherwise, prediction skill of the hydrological variables is 203 determined against the pseudo-observations (see Sect. 2.1). These have the advantages 204 of being complete in the spatial and the temporal domain and to be available for all 205 model variables. We will refer to this type of skill as "theoretical skill". In the 206 207 companion paper theoretical skill for discharge was compared to "actual skill", which is the skill assessed with real observations. It was concluded that, in terms of R and on 208 average across all target months and for lead month 2, the ratio of actual to theoretical 209 210 skill was 0.67 for "large basins" and 0.54 for "small basins". These two categories were defined on the basis of the observations of discharge, which were acquired from the 211

Global Runoff Data Centre, 56068 Koblenz, Germany (GRDC) and gridded onto the 0.5° x 0.5° model grid. Large basins are catchments upwards from the monitoring station larger than 9900 km² and small basins are catchments upwards from the station with an area smaller than that of the corresponding grid cell.

216

For the determination of the skill of the meteorological forcing we used the WFDEI data.

219

To investigate the possible contribution of trends to skill, skill in the meteorological 220 forcing and in runoff was determined both before and after removing the trend from 221 222 both the (pseudo-) observations and the hindcasts. Data were detrended by first constructing time series (1981-2010) for each variable, target month, lead month and 223 grid cell (30 values). We then removed the trend from each time series by first fitting a 224 least-squares regression line to the original time series and then subtracting the time 225 series corresponding to the line from the original data. For the hindcasts, time series 226 were constructed for the mean of the ensembles and the resulting best fit was subtracted 227 from each member individually. 228

229

Like in the companion paper, skill was analysed on a monthly and not on a seasonal basis with the aim of achieving a relatively high temporal resolution in the skill analysis. Attention was confined to consistent skill, which we define as skill that persists during at least two consecutive target or lead months. In accordance with Hagedorn et al. (2005), we designated the first month of the hindcasts as lead month zero, so target month number is equal to the number of the month of initialisation plus the lead month number.

237

In most result sections, we will first analyse and explain skill at the level of the entire domain. We will then take out the most noteworthy details of the summary plots and seek an explanation for them.

- 241
- 242

243 2.3 Isolation of sources of skill and surface water initialisation

244

As already pointed out in the introduction, a number of specific hindcasts were carried out with the aim of isolating the contributions of different sources to skill. The Full Streamflow Hindcasts (FullSH), in which skill is due to both meteorological forcing and initial conditions, constitute the starting point. The specific hindcasts can be seen as restricted, in the sense of limiting the types of sources of skill, versions of the FullSH. The following five sets of specific hindcasts, each consisting of 5400 computer runs, were produced: 1) The *InitSH* isolate the skill due to both types of initial conditions considered here
(soil moisture and snow). Like in the FullSH, the annually varying initial
conditions are taken from the reference simulation while for each year the
meteorological forcing is identical and consists of an ensemble of fifteen S4
hindcasts. More specifically, we selected member 1 from the 1981 hindcasts,
member 2 from the 1983 hindcasts, etc. By using identical meteorological forcing
for all of the years of the hindcasts, skill due to skill in the forcing is eliminated.

- 259 2) The *SMInitSH* isolate the skill due to the initial conditions of soil moisture only.
 260 SMInitSH is identical to InitSH but in all SMInitSH snow initial conditions are
 261 taken as the 30 year average of the snow conditions in the reference simulation.
- 3) The *SnInitSH* isolate the skill due to the initial conditions of snow only. SnInitSH
 is identical to InitSH but in all SnInitSH soil moisture initial conditions are taken
 as the 30 year average of the soil moisture conditions in the reference simulation.
- 4) The *MeteoSH* isolate the skill due the meteorological forcing and as such are the full complement of the InitSH. Like in the FullSH, the annually varying forcing is taken from the probabilistic S4 hindcasts while for each year the initial soil moisture and snow conditions are identical and equal to the 30 year average of the soil moisture and snow conditions in the reference simulation. By taking identical initial conditions for all of the years of the hindcasts, skill due to initial conditions is eliminated.
- 5) The *ESP* are identical to the InitSH, both in terms of their construction and in terms of their purpose. However, in the ESP the forcing is not taken from the S4 hindcasts but from the WFDEI data by selecting the 15 uneven years from 1981 to 2009.
- 276

Forcings and initial conditions of all of these hindcasts differ among the calendar months, so that the annual cycle is conserved. Hence,:

- "Identical for all years" means that the forcings (or the initial conditions) for all
 hindcasts starting in January are identical.
- "30 year average" means that the initial conditions for all hindcasts starting in
 January are averaged over all of the January 1st conditions in the reference
 simulation.
- "Annually varying" means that the forcings (or the initial conditions) for all
 hindcasts starting in January vary from year to year.
- These statements also hold for the other calendar months.
- 287

Thus, like the FullSH, all specific hindcasts for a single starting date consist of 15 members, which is important since ensemble size affects skill metrics (Richardson, 2001). Also, in all hindcasts the probabilistic character is exclusively due to the 15 members of the meteorological forcing while initial conditions are deterministic. This consistency is important since the main aim of the various specific hindcasts is to

compare them with each other. A disadvantage of the small ensemble size of the forcing is the sampling uncertainty, see Sect. 4.2 of the companion paper. All of the hindcasts were preceded by a one month run with reference forcing (WFDEI) with the single aim of initialising the amount of discharge in the rivers. So, discharge initialisation, a potential source of skill, is not considered. This has no effect on most of the analyses of the paper, since these are made in terms of runoff. Where discharge is analysed the effect of discharge initialisation is, due to the limited residence time of water in the rivers, restricted to the first lead month of the hindcasts (see Yuan, 2016).

3 Explanations of skill in hydrological variables

3.1 Skill in the meteorological forcing after bias correction

- In this sub-section, the skill of the meteorological forcing will be analysed. Attention will be limited to the three most important input variables of VIC, namely precipitation, two-meter temperature and incoming short-wave radiation. The WFDEI data are used as a reference. Here the data after bias correction are considered. In Appendix A we will discuss the skill of the raw S4 data, which is the meteorological forcing before bias correction. Differences in skill between the bias-corrected and the uncorrected data are negligible for temperature and short-wave radiation and small for precipitation.



Figure 1: Skill of the precipitation hindcasts after bias correction. Fig. 1a shows a map of the correlation coefficient between the observations and the median of the hindcasts (R), for target month January as lead month 0. The threshold of significant skill lies at 0.31, so yellow cells have insignificant skill. Grid cells with other colours have significant skill, with the amount of skill increasing with darkening colours. The legend provides the percentage of cells with significant values of R and the domain-averaged value of R. Fig. 1b depicts the percentage of cells with significant skill in terms of R, as a function of the target and lead month. Each coloured curve represents the

hindcasts starting in a single month of the year and has a length of 7 (lead) months. For better visualisation the parts of the curves that end in the next year are shown twice, namely at the left hand and the right hand side of the graph. Black lines (dashed, dotted and dashed-dotted) connect the results for identical lead times. The horizontal line gives the expected fraction of cells with significant skill due to chance in the case that the hindcasts have no skill at all (5%).

332 333

Fig. 1 shows results of the skill analysis of the precipitation forcing. Fig. 1a provides an 334 335 example of the skill for a single target and lead month (January as lead month 0). A summary of the skill in the precipitation hindcasts is given in Fig. 1b, which plots the 336 fraction of all cells within the domain with statistically significant R values. So, Fig. 1a 337 condenses into a single point in Fig. 1b. During the entire year, there is considerable 338 skill for lead month 0 (on average in 61% of the domain) but skill declines very rapidly 339 to 6% for lead months 1 and 2, just 1% more than the percentage of cells in the case of 340 no true skill at all. Hence, from lead month 1 on, skill is almost negligible. Regarding 341 lead month 0, there is more skill in January, February and March than during the other 342 months. For the same lead month, hot spots of consistent skill, i.e. with a duration of 343 significant skill of at least three target months, are situated on the Iberian Peninsula 344 from November to March, in Western Norway from January to April, in Greece and 345 Western Turkey from December to February and in Scotland from December to March. 346 347 All these occurrences of consistent skill are restricted to the winter half of the year and mostly to coastal regions (see Fig. 1a), suggesting them to be linked to the initial state 348 of the sea surface temperature. 349

350



Figure 2: Skill of the two-meter temperature hindcasts after bias correction. Figures 2a and 2b give the percentage of cells with significant values of R for the un-detrended (a) and the detrended (b) temperature hindcasts (see Fig. 1b for further explanation). Fig. 2c compares annual cycles of skill of undetrended and detrended data for the first three lead months. The three panels in the middle row show maps of R for the un-detrended temperature hindcasts for target months February (Fig. 2d) and March (Fig. 2e) as lead month 1 and July as lead month 5 (Fig. 2f). The bottom three panels depict the correlation coefficient of the trend (not the trend itself) of the observed monthly mean temperature for March (Fig. 2g) and July (Fig. 2i), and of the trend in the median of the hindcasted temperature for July as lead month 5 (Fig. 2h).

Figure 2 shows important aspects of skill in the two-meter temperature hindcasts. One 368 aspect is the possible contribution of a 30-year trend, which could be related to 369 greenhouse warming, to the skill. Figure 2a and 2b provide summaries of the skill of the 370 un-detrended and the detrended data, respectively, whereas Figure 2c compares these 371 two types of data. For lead month 0 the un-detrended hindcasts have significant skill in 372 the largest part of the domain (Fig. 2a). At longer lead times, the percentage of cells 373 with significant skill quickly drops towards the theoretical no skill limit (5%) but there 374 are a few exceptions, namely: 375

- For lead month 1, February and March temperatures are predicted with significant 376 skill in a considerable part of the domain (44% in February; 53% in March). In both 377 378 months the region with skill is more or less contiguous and comprises the Russian part of the domain, the Ukraine and the regions bordering the southern part of the 379 Baltic Sea (Figs. 2d and 2e). In February the region of skill extends towards Central 380 Europe. In March it also comprises northern Fennoscandia. This skill hardly 381 diminishes by detrending the data (Figs. 2b and 2c), suggesting that the skill is not 382 related to climate change. Indeed, in February and March the observed trend (in the 383 WFDEI data set) is insignificant across most of the domain (11% of the domain in 384 February and 18% in March) and, more importantly here, it is insignificant in the 385 regions with significant skill in the temperature hindcasts (Fig. 2g demonstrates this 386 for March). We conclude that the temperature skill in February and March as lead 387 month 1 must be due to initial conditions of the climate model (see also the 388 discussion on Fig. 10). 389
- 390 _ The three summer months (JJA) exhibit significant skill at all lead times in much more than 5% of the domain (a range from 22 to 56% for all combinations of the 391 three summer months and all lead months beyond lead month 0), see Fig. 2a. In this 392 case the fraction of cells with significant skill is not a function of lead time, which 393 is the type of behaviour that Yuan (2016) also found for the Yellow River basin. 394 Since Figs. 2b and 2c demonstrate that the skill more or less vanishes when the 395 temperature hindcasts and observations are detrended, we conclude that this skill is 396 due to trends in the data and hence probably related to greenhouse warming. 397 Another conclusion is that skill might be related to climate change, if the magnitude 398 of the skill does not or hardly varies with lead time. 399
- 400

It should be noted here that trends can only cause correlation between hindcasts and 401 observations, and hence skill in the hindcasts, if they are present in both time series. 402 A random time series of hindcasts is not correlated with a time series of 403 observations with a trend and vice versa. Indeed, time series of both hindcasts and 404 observations have a maximum in significant trends in summer, when trends form 405 the prime source of skill according to our analyses. In the hindcasts and on average 406 over all lead times beyond the first month, the summer months exhibit significant 407 trends in almost the entire domain (95%), versus 79% of the domain in the other 408

months of the year, on average. Similarly, observed trends are significant during the 409 three summer months in 67% of the domain, versus only 24% of the domain in the 410 other months of the year, on average. These percentages also show that significant 411 trends occur in a larger part of the domain in the hindcasts than in the observations. 412 So, the observations, and not the hindcasts, are mostly limiting the occurrence of 413 trend-related skill in the temperature hindcasts. This point is illustrated by Figs. 2g-414 i. Figure 2h shows that the trends of the hindcasts for July as lead month 5 are 415 significant across almost the entire domain (99% of the domain). However, 416 according to Fig. 2i only 69% of the domain has a significant trends in the observed 417 July temperatures. Indeed, the patterns of significance of Fig. 2f (skill in the 418 temperature hindcasts) and Fig. 2i (significance of observed trends) agree to a large 419 420 extent.

- April, May and September combine the behaviour of February and March, which
 have skill due to initial conditions of the climate model, with the skill of the
 summer months, which show skill related to trends (Fig. 2c).
- January has a considerable amount of significant skill but only for lead month 2
 (42% across the domain). This skill occurs in a stroke of land reaching from
 England to Russia, which vaguely coincides with the region in which Kim et al.
 (2012) found skill in the S4 temperature hindcasts for the three winter months.
 However, as this skill is not found in adjacent lead and target months, we speculate
 that this skill is spurious.



431 432

430

Figure 3: Skill of the incoming short-wave radiation hindcasts after bias correction.
Figure 3a gives the percentage of cells with significant values of R (see Fig. 1b for further explanation). Fig. 3b compares annual cycles of skill of undetrended and detrended data for the first three lead months (see Fig. 2c for further explanation).

438

439 Since short-wave incoming radiation is important for evapotranspiration, we finalise440 this sub-section with a short analysis of its predictability (Fig. 3). In terms of R, skill is

considerable during the first lead month with 58% of the cells having significant skill,
on average over the year. There tends to be more skill from March to September than
during the remaining months. Beyond lead month 0 skill settles around the no skill line,
except from April to July, but the fraction of cells with significant skill never exceeds
21% (in May as lead month 0). Trends in the data hardly affect skill (Fig. 3b).

- 446
- 447
- 448 449

3.2 Sources of skill in runoff and discharge

While Sect. 3.1 dealt with predictability of the meteorological forcing, this sub-section 450 451 analyses the effects of skill in the forcing and in the initial conditions on the predictability of runoff and discharge (discharge is only considered in Fig. 4). We first 452 address the question of how much of the skill in the runoff hindcasts could be linked to 453 trends that are possibly related to climate change. To examine this question, the 454 pseudo-observations and the hindcasts of runoff were detrended and the skill was 455 compared to that of the un-detrended data sets. We found that for lead month 2 and 456 averaged over all months of the year, the fraction of cells with a significant R decreased 457 from 58.7 to 57.4% due to detrending, a difference of 1.3%. This difference is much 458 smaller than the decrease for temperature (11.8%). We conclude that trends contribute 459 very little to skill in runoff. All analyses of this sub-section hereafter pertain to un-460 detrended data. Unless indicated otherwise, the pseudo-observations are used for 461 verification. 462

463

464

466

3.2.1 The relative importance of initial hydrological conditions

467 Figure 4 compares the InitSH with the FullSH in terms of the fraction of cells with a significant R. Figs. 4a and 4b show the result for runoff and discharge, respectively, 468 using the pseudo-observations, so calculations for all cells of the domain contribute to 469 the result. While the lumped results hardly differ between runoff and discharge (the 470 companion paper discusses small differences in skill between these two variables), 471 systematic differences in skill between the FullSH and InitSH are revealed. In lead 472 month 0, skill is higher in the FullSH than in the InitSH for all target months of the 473 year. Beyond lead month 1, the reverse occurs for most target months. Lead month 1 is 474 transitional with the order of skill depending on the time of the year. Figs. 4c (for large 475 catchments) and 4d (for small catchments) compare actual discharge skill, i.e. skill 476 determined with real discharge observations, of the FullSH with that of the InitSH. As 477 discussed in the companion paper, domain-average actual skill is less than domain-478 average theoretical skill but, more importantly here, the reversal of skill after lead 479 480 month 1 found with the pseudo-observations is confirmed with real observations, both for large and for small basins. While Fig. 4c is based on the data for all 111 large 481

catchments, a similar graph was produced for a selection of the large catchments with
relatively little human impact (about half of the 111 basins; Fig. S1 in the
supplementary material; see the companion paper for a description of the selection
procedure). Again the reversal occurs after lead month 1, so this phenomenon is also
confirmed by real observations from relatively pristine basins. Figures similar to Fig. 4
(Figs. S2-S5) illustrate that the skill reversal is found for all of the metrics considered in
this study and also for the domain-mean of R.

489



490

491

Figure 4: Comparison of the annual cycles of skill of the InitSH (blue) and the FullSH (red). The top two panels show theoretical skill obtained with the pseudo-observations for runoff (Fig. 4a) and discharge (Fig. 4b) at four different lead times. The bottom two panels compare actual skill of discharge for large (Fig. 4c) and for small (Fig. 4d) basins at three different lead times (months 0, 2 and 4).

498 499

We hypothesize that the reason for the reversal lies in the signal-to-noise ratio of the meteorological forcing. The InitSH forcing is the same for each year, so its interannual variation does not contain a signal nor noise. However, the forcing of the FullSH varies from year to year. During the first lead month this forcing has considerable skill (see Sect. 3.1), so the signal-to-noise ratio of the forcing is relatively high. This enhances

skill in the FullSH with respect to InitSH. At longer lead times the interannual variation 505 506 in the forcing hardly contains a signal, with the exception of some limited skill in the temperature hindcasts (see Sect. 3.1), so the signal-to-noise ratio is low. Noise in the 507 meteorological hindcasts reduces skill in the FullSH with respect to InitSH. Figure 4 508 demonstrates that averaged across the domain, the reversal between skill enhancement 509 due to the signal in the forcing of the FullSH and skill reduction due to noise in the 510 forcing of the FullSH occurs at some time between the first and the third lead month, 511 with the exact timing of the reversal depending on the target month. 512

513



- 514
- 515

516 517

Figure 5: Comparison of the annual cycles of the skill in the runoff hindcasts of four specific hindcasts for lead months 0 and 2. Different colours correspond to different specific hindcasts and different line types to different lead months.

- 518 519
- 520

3.2.2 The relative contributions of soil moisture and snow initial conditions, and of meteorological forcing

523

Figure 5 compares the skill in run-off of the specific hindcasts (except ESP) for two lead months (0 and 2). At both lead times and for all target months, initialisation of soil moisture is the dominant source of skill in Europe. Initialisation of snow and meteorological forcing are less important. This is true for all lead times (not shown here).

529

530 Meteorological forcing does not only have a relatively small contribution to the 531 domain-averaged skill of Fig. 5 but also to regional skill. We searched for combinations

of a region and target months where the MeteoSH produce consistently equal or more

skill than the SMInitSH but we did not find any combination where this clearly was the 533 case. During the first lead month there is more skill due to the forcing than due to snow 534 initial conditions (SnInitSH). For later lead months this order depends on the target 535 month, mainly because skill due to snow initial conditions varies strongly during the 536 year. Although skill in run-off due to meteorological forcing is relatively small, it does 537 exceed the skill in the forcing variable to which runoff is most sensitive, precipitation 538 (compare Fig. 5 with Fig. 1). Whereas predictability of precipitation is almost limited to 539 the first lead month, significant skill in runoff due to forcing is more widespread for 540 lead months 1 and 2 (on average over the year in 23 and 15 % of the domain, 541 respectively). We explain the enhanced skill in runoff mainly by an indirect effect. Skill 542 543 in the precipitation forcing of the first lead month leads to skill in the states of soil moisture and snow at the end of that month. These model states then serve as the source 544 of skill during the next lead months, when the precipitation forcing has no skill at all. In 545 addition to this indirect effect of precipitation, the skill in the hindcasts of temperature 546 (Fig. 2) contributes to the skill in runoff. 547

548

From April to July, a considerable part of Europe has significant skill derived from 549 snow initialisation provided initialisation does not occur earlier than in February. Skill 550 due to snow initialisation reaches a maximum in May and June, resulting in a 551 maximum in skill in the InitSH-hindcasts and the FullSH for these months and for most 552 553 lead times. When snow contributes considerably to predictability (from April to July), 554 the skill in the InitSH exceeds the skill in the SMInitSH. Because for target months from August to March snow contributes little to predictability, the percentages of cells 555 with significant skill in InitSH and SMInitSH are almost identical for these months. 556 The rapid rise in skill due to snow initialisation at the transition from April to May 557 explains a remarkable feature that we noticed in the companion paper, namely an 558 increase in runoff skill with lead time at this time of year. Another noticeable feature is 559 that the skill due to snow initialisation for lead month 2 exceeds skill due to snow 560 initialisation for lead month 0. This occurs for target months from May to August and 561 562 will be explained in the text corresponding to Fig. 8.

563

Figures similar to Fig. 5 but for all metrics of the present study are included in the supplementary material (Fig. S6). The graphs for the ROC areas for the Above Normal (AN) and Below Normal (BN) terciles are qualitatively similar to the graph for R. This also holds for the RPSS though fractions of the domain with significant RPSS are almost always lower than for the other metrics. An exception is the relatively large amount of significant skill in the SnInitSH when RPSS is used as metric.

- 570
- 571



Figure 6: Example that compares the skill in runoff of three specific hindcasts (SMInitSH (a), SnInitSH (b) and InitSH (c)), for target month May as lead month 2. For more explanation, see Fig. 1a. White, terrestrial cells correspond to cells where observations or hindcasts consist for more than one third of zeros or one sixth of ties.

Figure 6 compares skill maps for the three specific hindcasts that isolate skill due to initial conditions (InitSH, SMInitSH and SnInitSH). It illustrates that at regional scale skill due to snow and soil moisture initialisation are more or less additive. Copies of the patterns of skill due to soil moisture initialisation e.g. in Africa, on the Iberian Peninsula and in Western France (Fig. 5a) are found in the map of skill due to both soil moisture and snow initialisation (Fig. 5c). Small regions with considerable skill due to snow initialisation (Fig. 5b) like those near Stockholm, in South-east Czechia and South-east Austria also stick out as foci of skill on the map of skill due to both soil moisture and snow initialisation (Fig. 5c). Where both soil moisture and snow initialisation cause moderate skill, e.g. in Southern Finland, the combined specific hindcast exhibits more significant skill. The additive behaviour of skill in the two initialisation components is also visible in Fig. 5.



595

596 597

598 599 600

Figure 7: Example showing the variation of skill in runoff as a function of lead time in the SnInitSH, for initialisation on March 1. For more explanation, see Figs. 1a and 6.

Figure 7 zooms in on the specific hindcast that isolates skill due to snow initialisation 601 (SnInitSH), giving the example of a time series of skill as a function of lead time, after 602 initialisation on March 1st. One observation is that skill does not gradually decrease with 603 604 time but has a maximum during the snow melt season. We like to note that locally skill is hardly generated during the part of the melt season when a snow pack covers the 605 surface in each year. The reason is that in VIC the rate of snow melt is almost 606 insensitive to snow pack thickness (Sun et al., 1999). Hence, as long as the surface is 607 covered by snow in each year, inter-annual variation in snow melt is absent or 608 negligible. Skill is only generated towards the end of the melt season, when snow melt 609 differs from year to year because snow stops to be available for melt at different dates 610 due to different initial amounts of snow. So, the initial snow conditions cause skill 611 because of interannual variation in the duration of the period that it takes to melt the 612 613 snow present at the time of initialisation and not because of interannual variation in the melt rate. Of course, the timing of the end of the melt season differs regionally and with 614 elevation, which largely explains the patterns of skill visible in the maps of Fig. 7. A 615 good example is Scandinavia, where the earliest skill (in April; lead month 1) occurs at 616 low elevations near the coasts of Southern Norway and Sweden, at the end of the local 617 snow season. The latest skill (in July; lead month 4) occurs in the Norwegian 618 mountains, again at the end of the local snow season (we ascribe the skill in South-east 619

Sweden in July and August to chance). It is also relevant to note that the skill patterns in the maps of Fig. 7 are influenced by the fact that VIC has higher vertical resolution than its horizontal resolution may suggest, by performing simulations in multiple elevation bands within each grid cell, accounting for sub-grid variations in topography. Therefore, sub-grid topography leads to spreading of the snow skill signal of individual cells over longer periods of time.

626



628

Figure 8: Example illustrating that skill in runoff for a target month may increase with lead time, namely for runoff in May as target month 0 (a) and 1 (b) in the SnInitSH. Skill in the hindcast of soil moisture in the SnInitSH for May as lead month 1 is shown (c) because it provides an explanation for the mechanism causing the increase in skill with lead time. For more explanation, see Figs. 1a and 6.

- 635
- 636

To finish the analysis of the SnInitSH, Fig. 8 analyses a remarkable feature. In 637 SnInitSH, hindcasts for May have less skill when the hindcasts are initialised on May 1 638 (Fig. 8a) compared to initialisation during preceding months (February, March or April, 639 Fig. 8b is for initialisation on April 1). Similar counterintuitive results are found for 640 June and July as target months. This result is counterintuitive because in hindcasts with 641 initialisation on May 1 there is, due to the use of pseudo-observations for verification, 642 perfect knowledge about snow conditions on that date. With initialisation on April 1, 643 snow conditions on May 1 differ from those of the pseudo-observations, which by itself 644 must lead to less skill in May runoff. However, there is compensation for this direct 645 effect by an indirect effect through soil moisture. In SnInitSH, soil moisture has no skill 646 on the date of initialisation, e.g. May 1 in the hindcasts starting on that date. However, 647 in the hindcasts starting on April 1 the perfect knowledge of the snow conditions on that 648 date leads via skill in snow melt in April to some skill in soil moisture on May 1 (Fig. 649 8c), which then leads indirectly to skill in runoff in May. Since we find more skill in 650 May runoff after snow initialisation on April 1 than after snow initialisation on May 1, 651

the gain of skill in the runs starting on April 1 due to the indirect effect
overcompensates for the loss of skill in the same runs due to the direct effect.

To finalise this section, the specific hindcasts were exploited to attribute the hotspots of 655 significant skill in runoff for lead month 2, listed in the companion paper, to the 656 different potential sources of skill. This was done for each of the hotspots by an 657 inspection of the maps of skill (like those of e.g. Fig. 6) for three specific hindcasts that 658 isolate the different sources of skill (SMInitSH, SnInitSH and MeteoSH). If the hotspot 659 was present in e.g. SMInitSH, soil moisture initialisation is one of the sources of skill. 660 Results are summarised in Table 1. Almost all of the significant skill in the hotspot 661 regions is due to the initial conditions of soil moisture. Exceptions are formed by the 662 target months from April to July when skill is caused by a mix of the initial conditions 663 of snow and soil moisture in regions with significant snow melt skill. In these cases the 664 relative contributions of the two sources varies in time and space but soil moisture is 665 more important than snow, except in Fennoscandia where in June snow dominates and 666 in July both sources are of about equal importance. In none of the hotspots of skill, 667 meteorological forcing contributed significantly to this. 668

> Region period source of skill Fennoscandia Jan - Mar SM Apr - Jul SM and snow Aug - Oct SM Poland and Northern Germany Oct - Mar SM Apr - May SM and snow SM Western France Dec - May Romania and Bulgaria Oct - Mar SM Apr - May SM and snow southern Mediterranean Jun - Aug SM

Table 1 Sources of skill for hotspot regions and periods of skill. SM is soil moisture.

671 672

669

3.3 Skill and source of skill in evapotranspiration

674

Because hindcasts of evapotranspiration are useful in themselves, because
evapotranspiration affects runoff (see Sect. 1), and in order to demonstrate the rich
possibilities of the pseudo-observations and the specific hindcasts, this section analyses

skill in the hindcasts of evapotranspiration. In VIC evapotranspiration is computed withthe Penman-Monteith method (see Shuttleworth, 1993).

680





Figure 9: Summary plots of the skill of the hindcasts of evapotranspiration. Figure 9a summarises the FullSH (for more explanation, see Fig. 1b), Fig. 9b depicts the annual cycles of skill for the FullSH and three specific hindcasts (SnInitSH, SMInitSH and MteoSH) for lead months 0 and 2, and Fig. 9c compares the annual cycles of skill of the un-detrended and the detrended FullSH for the first three lead months.

689 690

Figure 9a summarizes skill in evapotranspiration in the FullSH. Levels of predictability 691 are higher than for precipitation (Fig. 1), similar to those for temperature (Fig. 2) and 692 693 lower than those for runoff (Fig. 4a). Figure 9b isolates the diverse contributions to skill 694 for lead months 0 and 2 by showing the skill for the FullSH and three specific hindcasts (SMInitSH, SnInitSH and MeteoSH). Averaged over the year, meteorological forcing 695 contributes more and initial soil moisture less to predictability in evapotranspiration 696 697 than to predictability in runoff. Initial snow is the least important of the three sources of skill. 698

699

Focusing on lead month 2, there is hardly any skill in the evaporation hindcasts from 700 November to March (9% of the domain, on average over these months), with the 701 702 exception of January (18%) when the region of skill (Germany and Benelux) is part of a larger region of skill in the temperature hindcasts for the same target and lead month. 703 We blame the winter minimum of skill in evapotranspiration to the low levels of 704 evapotranspiration and the low levels of skill in the temperature forecasts for the same 705 period. The next month (April) exhibits the highest level of skill of all months (44% of 706 the domain), mainly due to meteorological forcing (MeteoSH) and with smaller 707 contributions by the initial conditions of soil moisture (SMInitSH) and snow 708

(SnInitSH). From May to September there is some significant skill (23% of the domain, 709 710 on average over these months). Whereas in May forcing is still the most important contributor to skill, initial conditions of soil moisture form the main contributor from 711 712 June to October. We speculate that this shift in the order of importance between forcing and soil moisture is due to the amount of variability in soil moisture. In Europe in spring 713 (April, May), soil moisture variations are relatively small and hence hardly contribute to 714 variations in evapotranspiration. Later in the year (June to September), soil moisture is 715 often available in limited amounts, so variations are larger and hence contribute more to 716 variations in evapotranspiration. Snow initial conditions contribute to skill only during 717 the snow melt season from April to July. 718

719

The contribution of trends to predictability of evapotranspiration is summarised in Fig. 720 9c, for lead months 0, 1 and 2. For lead month 2 and averaged over all target months of 721 the year, detrending leads to a decrease in the fraction of cells with a significant R from 722 17.6 to 13.8%, a difference of 3.8%. The contribution of trends to skill in 723 evapotranspiration is less than its contribution to skill in temperature (a difference of 724 11.8%) but larger than its contribution to skill in runoff (a difference of 1.3%). Trends 725 contribute to skill in evapotranspiration during the part of the year when they also 726 contribute to skill in atmospheric temperature (Fig. 2c), namely from April to 727 September and in November (for lead month 0). However, whereas during the three 728 729 summer months the skill in the temperature hindcasts is almost exclusively linked to climate change, a considerable part of the domain still exhibits skill in 730 evapotranspiration after detrending. 731

732 733



735 736

737 738

Figure 10: Explanation of the skill in the hindcasts of evapotranspiration for target month April as lead month 2. The panels map the skill in evapotranspiration of the FullSH (a), of the MeteoSH (b) and of the hindcasts of temperature (c). For more explanation, see Fig. 1a.

740

739

To provide a deeper understanding of the skill in evapotranspiration, the skill in April 742 and July is analysed in some detail. Fig. 10 deals with April as lead month 2, showing 743 the skill in evapotranspiration from the FullSH in Fig. 10a and from the MeteoSH in 744 Fig. 10b. Regions of skill, mainly a stroke of land from southern Fennoscandia to the 745 Black Sea, are the same in the FullSH and in the MeteoSH though skill is somewhat 746 degraded in the MeteoSH. This indicates that meteorological forcing causes most, 747 though not all, of the skill. Indeed, Fig. 2e (skill in temperature for March as lead month 748 1) and Fig. 10c (skill in temperature for April as lead month 2) show that the 749 temperature forecasts of the preceding lead month and the lead month considered 750 contain skill in the same regions. We conclude that much of the skill in 751 752 evapotranspiration is due to skill in the temperature hindcasts. The remaining part of the skill is due to initial conditions (Fig. 9b shows this for the entire domain). We found 753 limited amounts of skill in the SnInitSH and the SMInitSH for April in the same stroke 754 of land from southern Fennoscandia to the Black Sea (not shown here). This means that 755 in that region initial conditions of the hydrological model on February 1 provide some 756 skill to the hindcasts of evapotranspiration for April. We like to note that this could be 757 consistent with the conclusion in Sect. 3.1 that the skill in the temperature hindcasts of 758 759 February and March in this same region are due to the initial conditions of the climate 760 model. These initial conditions could e.g. be sea surface temperatures but also the local state of snow and/or soil conditions. In the latter case, the two types of predictability in 761 762 the mentioned regions would have the same or a similar source. Initial conditions of 763 snow and/or soil conditions in S4 would lead to skill in the temperature hindcasts of S4 while initial conditions of snow and soil moisture in VIC lead to skill in the 764 evapotranspiration hindcasts of VIC. 765

766

During the summer months and for all lead times, skill in evapotranspiration occurs in 767 two regions, namely the southern part of the Mediterranean, and Western and Northern 768 769 Norway. Fig. 11 shows target month July as lead month 5, as an example. Whereas Fig. 11a is for the FullSH, Figs. 11b-d depict the maps for three specific hindcasts 770 771 (SnInitSH, SMInitSH and MeteoSH) and Fig. 11e shows skill for the FullSH after detrending. Since the SnInitSH and the MeteoSH exhibit hardly any skill while 772 773 SMInitSH has considerable skill in the Mediterranean (Figs. 11b-d), it can be concluded 774 that the skill in this region is due to soil moisture initial conditions. So, in this particular case, knowledge of soil moisture conditions on February 1 still yields skill in 775 evapotranspiration in July. This skill in the Mediterranean is not affected by detrending 776 (compare Figs. 11a and 11e), so it does not have a climate change component. 777

778

The skill in Norway has a more complicated origin. The three specific hindcasts show that it is due to a mix of initial snow conditions (Fig. 11c) and meteorological forcing (Fig. 11d). The effect of the initial snow conditions (on February 1) can be understood with the help of the analysis of runoff skill in the SnInitSH (Fig. 7), which led to the

conclusion that runoff skill caused by snow initialisation occurs at the end of the melt 783 season, which is July in much of Norway. Therefore, in this country and in July the 784 timing of the disappearance of snow cover varies from year to year. This then has a 785 considerable effect on evapotranspiration since bare soil has, compared to snow, higher 786 surface temperatures and hence more evapotranspiration in summer. The contribution to 787 skill by forcing (Fig. 11d) fades with but is not removed by detrending (not shown 788 here), so it has a part that is related to climate change and a part that is unrelated to 789 climate change. The climate-change-related skill due to forcing resides in the 790 temperature hindcasts, which have significant skill in this region at all lead times (Fig. 791 2f). The non-climate change related skill in the MeteoSH for July is likely an indirect 792 793 effect of the skill in the forcing (especially precipitation) in the first lead month (February). This leads to skill in snow water equivalent towards the end of February, 794 which fades but has not disappeared completely on July 1 (Fig. 11f) and then causes 795 skill in evapotranspiration at the end of the melt season. 796





Figure 11: Explanation of the skill in the hindcasts of evapotranspiration (ET) for July 801 by taking lead month 5 as an example. The panels map the skill in 802 evapotranspiration of the FullSH (Fig. 11a), SMInitSH (Fig. 11b), SnInitSH 803 (Fig. 11c), MeteoSH (Fig. 11d) and the FullSH after detrending (Fig. 11e). 804 Figure 11f depicts skill of the hindcasts of snow water equivalent (swe) in 805 the MeteoSH. For more explanation, see Fig. 1a. Note that statistics in the 806 807 legends of the panels refer only to that part of the domain for which R was computed, which consists of all coloured cells. 808

809 810

811

4

Discussion

812 813

814

4.1 Comparison of skill with previous studies

A remarkable result of our work is the reduction of the skill in runoff beyond lead 815 month 1, when annually varying S4 forcing is used (FullSH) instead of meteorological 816 forcing that is identical for all years (InitSH), see Fig. 4. This result is counter-intuitive 817 but, as we discussed, a logical consequence of forcing with interannual variation that 818 819 has no or insufficient skill, such as the S4 forcing. Other studies compared FullSH (also called climate-model based hindcasts) with ESP hindcasts, which are slightly different 820 from our InitSH (see Sect. 4.3) but like the InitSH have identical meteorological forcing 821 for each year. Some of these studies (e.g. Singla et al., 2011, and Mackay et al., 2015) 822 found little overall difference in skill between the FullSH and ESP hindcasts. However, 823 in contrast with our results, skill is enhanced when using meteorological hindcasts, also 824 at longer leads, in the studies of Yuan et al. (2013), Thober et al. (2015) and Yuan 825 (2016). This contrast might be explained by more skill in the meteorological hindcasts 826 of the mentioned studies than in the present study, which could be due to the type of 827 meteorological hindcasts (none used S4) or the investigated region (in the mentioned 828 829 studies US, Europe and China, respectively). Indeed, Europe is a region with relatively little skill in meteorological hindcasts (Kim et al., 2012, Scaife et al., 2014, and Baehr et 830 al., 2015). Effects of regional differences in the skill of the forcing on the relative skill 831 of full and ESP hindcasts are mentioned by Wood et al. (2005), who reported that full 832 hindcasts for the Western United States have practically no skill improvement over the 833 ESP, except for some regions and seasons with predictability of the forcing originating 834 in ENSO teleconnections. 835

836

The specific hindcasts of this study show that in Europe initial conditions of soil 837 moisture are the largest source of skill in the seasonal run-off forecasts produced with 838 WUSHP. In terms of domain averages, this is true for all lead and target months. 839 Contributions to skill by the initial conditions of snow and by the meteorological 840 forcing are mostly much smaller. To our knowledge, two other studies analysed sources 841 of skill of hydrological seasonal forecasts for Europe with dynamical systems similar to 842 those of the present study, namely Bierkens and Van Beek (2009) and Singla et al. 843 (2011). Results of these two studies were summarised in the introduction. However, the 844 conclusions of Singla et al. (2011) are not directly comparable with those of the present 845 study as they used ESP (see Sect. 4.3). 846

847

Comparing our results with those of Bierkens and van Beek (2009), both studies agree that initial conditions form the dominant source of skill. However, compared to the

present study, Bierkens and van Beek (2009) find a larger contribution to skill by the 850 meteorological forcing, at least in summer. This difference might be due to the quality 851 of the forcing. Bierkens and van Beek (2009) developed an analogue events method to 852 853 select, on the basis of annual SST anomalies in the North Atlantic, annual ERA40 meteorological forcings, which they used as forcing for their hydrological model. One 854 855 might speculate that in Europe their semi-statistical forcing is more skilful than the S4 forcing used in WUSHP. This suggests that there is room for improvement of climate 856 model seasonal forecasts, so if and when this improvement is realised, the relative 857 contribution of the meteorological forcing to skill in hydrological variables would 858 increase. In any case, that contribution depends and will depend on the climate model 859 860 used (e.g. S4 or GloSea5).

The dominance of soil moisture initial conditions in terms of domain-lumped skill also

861

862

864

865

4.2 Understanding the skill due to initial soil moisture

extends to the hotspot regions and periods of skill (Table 1). The understanding of the 866 skill linked to soil moisture can be deepened by another level as in Shukla and 867 Lettenmaier (2011). The underlying idea is that this type of skill increases with the 868 interannual variability of soil moisture at the date of initialisation and that this skill is 869 870 gradually eliminated during the course of the hindcasts by interannual variability in 871 processes like rain fall and snow melt. The question is to what extent hotspots of skill (see Table 1) linked to soil moisture initialisation are due to the cause of the skill and to 872 873 874

what extent they are due to a lack of interannual variability in the processes that eliminate the skill? Figure 12 helps answering this question for the skill found in the runoff hindcasts of August as lead month 2 with a simple method of analysis. Figure 12a shows the standard deviation of total modelled soil moisture (σ_{SM}) on the day of initialisation (June 1), taken from the reference simulation. Figure 12b depicts the standard deviation of total rain fall (σ_{RF}) during the course of the hindcast (June – August), taking from the WFDEI data set, which is the investigated skill-eliminating factor. These two quantities were combined into an estimate of the skill (S_{est}):

$$S_{est} = \exp\left(-\frac{\sigma_{RF}^2}{\sigma_{SM}^2}\right)$$
 (1)



883 884

Figure 12: Illustration of a simple method that partly explains skill in runoff due to 885 initial soil moisture, exemplified for target month August as lead month 2. 886 Figure 12a is a map of the standard deviation in soil moisture at the date of 887 initialisation (June 1). Similarly, Fig. 12b maps the standard deviation of 888 observed rain fall during the course of the hindcasts (June-August). These 889 two standard deviations are combined into an estimate of the skill (Eq. 1) in 890 Fig. 12c, which is compared with the skill of the FullSH (Fig. 12d). Note 891 that the colour scales of Figs. 12c and 12d differ from each other and differ 892 from scales of other figures (e.g. Fig. 1a). 893

894 895

This estimate (Fig. 12 c) needs to be compared with the skill of the hindcasts, mapped in 896 Fig. 12d in terms of R. The two maps are not expected to be exactly equal, not only 897 because of the simplicity of the estimation method but also because Sest is not a 898 correlation coefficient. However, in the limits Sest has the desired properties. It is equal 899 to zero for the cases of constant initial amounts of soil moisture or infinite variability in 900 901 rain fall. It is equal to one for the cases of infinite variability in soil moisture or constant rain fall. The correlation coefficient between the patterns in Figs. 12c and d is highly 902 significant (0.67) and the hotspot regions of skill are the same in both panels, namely 903 the northern part of Fennoscandia and the southern part of the Mediterranean. So, in the 904 case of August as lead month 2 the estimation method is reasonably successful in 905 computing the pattern of skill in the hindcasts with the simple means of the WFDEI data 906 set and model calculations from the reference simulation. The additional merit of the 907

estimation method is the deeper understanding of the cause of the skill in the two
hotspot regions. Northern Fennoscandia is a hotspot because the amount of interannual
variability in initial soil moisture is larger than elsewhere (Fig. 12a). The southern part
of the Mediterranean is a hotspot because the amount of interannual variability in
rainfall is lower than elsewhere (Fig. 12b).

913

This simple method of analysis helped to bring the understanding of the skill in northern Fennoscandia and the southern Mediterranean to a deeper level but it was less successful for the other hotspots. A more thorough analysis along these lines and a deeper understanding of skill in the hindcasts is left for future work.

- 918
- 919 920

4.3 Relation of the present specific hindcasts with conventional ESP

921

The specific hindcasts of this study are related to the well-known Ensemble Streamflow Predictions (ESP) (e.g. Wood and Lettenmaier, 2008, Shukla and Lettenmaier, 2011, Singla et al., 2011, and Van Dijk et al., 2013). ESP are not only used as an experimental tool in science but are also widely used to produce forecasts in operational mode (Day, 1985). ESP used for scientific purposes can be subdivided into ESP proper (called ESP from now on) and reverse-ESP.

928

929 ESP (hindcasts) are similar to the InitSH of this study. In both types of hindcasts the 930 initial conditions vary from year to year and are quasi-perfect, i.e. they are taken from a simulation like our reference simulation, while the meteorological forcing does not vary 931 from year to year. This eliminates skill due to the meteorological forcing, so skill can 932 only be due to the initial conditions. However, while in ESP the forcing is selected 933 from historic observations, it is selected from the S4 hindcasts in InitSH in order to 934 retain an inter-member variability and other statistical characteristics of the time series 935 similar to that in the FullSH. An advantage of ESP is that its production is relatively 936 cheap because no climate model forecasts are needed. 937

938

Similarly, reverse-ESP (see Wood and Lettenmaier, 2008) resemble the MeteoSH of this study. In both types of hindcasts the meteorological forcing varies from year to year while the initial conditions are identical for each year. This eliminates skill due to the initial conditions, so skill can only be due to the forcing. However, while in reverse-ESP the forcing of each year is made up of the observations of that year, it is made up of the S4 hindcasts in the MeteoSH.

- 945
- 946



947 948

949 950

951 952

Figure 13 Comparison of the annual cycles of skill of the FullSH (red), the InitSH (blue) and the ESP (green) for three different lead times. Where blue or green symbols seem to be missing, they coincide.

- So, both in ESP and in the InitSH the meteorological forcing is identical for all years of 953 954 the hindcasts with the aim of eliminating the skill due to the forcing. If indeed all skill due to the forcing is removed, the remaining skill, which is due to the annually varying 955 initial conditions, should logically be the same in ESP and InitSH since the initial 956 conditions in both types of hindcasts are the same. To test this expectation we produced 957 ESP and compared their skill with the skills of the FullSH and the InitSH. Figure 13 958 959 shows results for three different lead times. In the graph most of the points for the ESP are indistinguishable from their counterparts for the InitSH. So, all conclusions that 960 were drawn from Fig. 4 and especially the reversal with lead time of the ranking of 961 predictability for the FullSH and the InitSH, are equally true for the ranking of the 962 FullSH and the ESP. We also conclude that, though forcings in the InitSH and the ESP 963 differ, skills from both types of specific hindcasts are, as expected, virtually identical, 964 which can be ascribed to the fact that the forcings do not vary from year to year. We 965 speculate that this result would also hold for other plausible forcings that do not vary 966 967 from year to year.
- 968

This behaviour is in sharp contrast with the skill resulting from reverse-ESP and MeteoSH, which are expected to be totally different. Keeping in mind that in both types of hindcasts skill is caused only by skill of the meteorological forcing, this is the skill of the S4 hindcasts in the MeteoSH. The present study showed that in Europe there is a small contribution to skill in the streamflow hindcasts by the forcing and that this contribution tends to decrease with time. This differs from reverse-ESP, in which skill
is small at the beginning and then increases with lead time (see Wood and Lettenmaier,
2008) because the meteorological forcing is quasi-perfect (i.e. identical to the forcing in
the reference simulation) while the influence of the initial conditions, which in reverseESP eliminate skill, decreases with time.

979

In summary, amounts of skill are almost the same for ESP and InitSH, and totally 980 different for reverse-ESP and MeteoSH. Also, interpretations of reverse-ESP and 981 MeteoSH differ. MeteoSH can be used to assess skill in the streamflow hindcasts due 982 exclusively to skill in the meteorological hindcasts. Reverse-ESP can be used to 983 quantify skill due to prescribing meteorological observations, i.e. the skill if we had 984 perfect knowledge about the meteorological forcing during the forecast period and no 985 knowledge about the initial conditions. Such a specific hindcast would not fit into the 986 present study. In fact, ESP, reverse-ESP and a reference simulation were produced by 987 Wood and Lettenmaier (2008) and Shukla and Lettenmaier (2011) for purposes that 988 differ from those of the present study. Their aim was to quantify what can be gained if 989 the meteorological forcing (in ESP) or the initial conditions (in reverse-ESP) are 990 improved from containing climatological information to being quasi-perfect, i.e. when 991 992 they are equal to the meteorological observations and the initial state of the reference simulation. Wood et al. (2016) extended this type of analysis by determining 993 sensitivities of the streamflow to changes in the information of the meteorological 994 forcing and the initial conditions. 995

996

999

997998 4.4 Towards an operational system

We plan to launch an operational version of WUSHP. That version might include a post-processing procedure with the aims of removing biases in discharge and making the system more reliable. This could perhaps be done with statistical calibration (e.g. Gneiting et al., 2005, and Schepen et al., 2014), a technique that, contrary to quantile mapping, considers information that is available from correlations between hindcasts and observations (see Wood and Schaake, 2008, and Madadgar et al., 2014).

1006

The superiority of the InitSH (and the ESP) with respect to the FullSH for hindcasts beyond the first two lead months raises the question whether one should, in an operational version of WUSHP and for these lead months, issue forecasts like the InitSH (or ESP) and not forecasts like the FullSH. The logical answer is "yes" but such a strategy should then be reconsidered when the meteorological forcing is taken from a new, possibly improved version of the climate model, or from another, possibly better type of climate model.

The applied methods of analysis are not suitable for giving quantitative advice on what 1015 would be the best investment for increasing the amount of skill of WUSHP. However, 1016 since initial soil moisture is the dominant source of predictability, a large gain of skill 1017 could possibly be made by assimilation of soil moisture observations into the modelled 1018 state of soil moisture (see e.g. Draper and Reichle, 2015). In addition, observations of 1019 1020 snow water equivalent could be assimilated into the modelled state of snow (see e.g. Griessinger et al., 2016). Improving the calibration of VIC would be another obvious 1021 road towards improvement of the seasonal predictions discussed in this paper. This 1022 should lead to higher actual skill but not necessarily to more theoretical skill, see the 1023 discussion section of the companion paper. 1024

1025

In this study we analysed the effect of trends on predictability by comparing skill for 1026 un-detrended and detrended observations and hindcasts. The effect of trends was almost 1027 negligible for runoff, considerable and in summer even dominant for atmospheric 1028 temperature and of intermediate magnitude for evapotranspiration. While for academic 1029 purposes the distinction between climate-change related and non-climate change related 1030 skill contributes to a better understanding of the skill, this is not the case for practical 1031 applications. To our knowledge, operational forecast systems issue their predictions 1032 without considering trends, i.e. predictions are issued by taking a historic period 1033 covered by observations or hindcasts as reference. 1034

1035

1036 This study demonstrates the power of using pseudo-observations for verification. These 1037 data cover all cells of the entire domain and they are available for all hydrological model variables, e.g. runoff and evapotranspiration, which have no equivalent (runoff) 1038 or are sparse (evapotranspiration) in the realm of real observations. Many features of 1039 skill would not have been detectable with real observations. At the same time, we like to 1040 stress that actual skill obtained with real observations is generally less than theoretical 1041 1042 skill obtained with pseudo-observations, as discussed extensively in the companion paper. 1043

1044 1045

1046 5 Conclusions

1047

The present paper explains skill in the hindcasts of WUSHP, a seasonal hydrological forecast system, applied to Europe. We first analysed the meteorological forcing, which consists of bias-corrected output from a climate model (S4), and found considerable skill in the precipitation forecasts of the first lead month but negligible skill for later lead times. Seasonal forecasts for temperature have more skill. Skill in summer temperature was found to be related to climate change occurring in both the observations and the hindcasts, and to be more or less independent of lead time. Skill in North-East Europe in February and March is unrelated to climate change and musthence be due to initial conditions of the climate model.

1057

1058 Sources of skill in runoff were isolated with specific hindcasts, namely SMInitSH (soil moisture initialisation), SnInitSH (snow initialisation), InitSH (a combination of soil 1059 1060 moisture and snow initialisation) and MeteoSH (meteorological forcing). These hindcasts revealed that, beyond the second lead month, hindcasts with forcing that is 1061 identical for all years but with "perfect" initial conditions (InitSH) produce, averaged 1062 across the model domain, more skill in runoff than the hindcasts forced with S4 output 1063 (FullSH). This occurs because interannual variability of the S4 forcing adds noise 1064 1065 while it has hardly any skill. The other specific hindcasts showed that in Europe initial conditions of soil moisture form the dominant source of skill in runoff. For target 1066 months from April to July, initial conditions of snow contribute significantly, with a 1067 domain-mean maximum in May and June. The timing of that maximum varies 1068 spatially and coincides with the end of the melt season, when snow melt differs from 1069 year to year because snow stops to be available for melt at different dates. All regional 1070 and temporal hotspots of skill in runoff found in the companion paper are due to initial 1071 conditions of soil moisture, with smaller or larger contributions by the initial 1072 conditions of snow for target months from April to July in hotspot regions with snow 1073 1074 fall in earlier months. We further showed that skill due to snow and soil moisture 1075 initialisation is more or less additive.

1076

1077 Some remarkable skill features are due to indirect effects, i.e. skill due to forcing or 1078 initial conditions of snow and/or soil moisture is, during the course of the model 1079 simulation, stored in the hydrological state (snow and/or soil moisture), which then by 1080 itself acts as a source of skill. Examples occur in the skill of run-off for target months 1081 from May to July in the SnInitSH and in the skill of run-off in the MeteoSH.

1082

Predictability of evapotranspiration was analysed in some detail. Levels of predictability and the annual cycle of skill are similar to those for temperature. For most combinations of target and lead months, forcing forms the most important contributor to skill but for lead month 2 initial conditions of soil moisture dominate from June to October. The sources of some regional and temporal hotspots of skill in evapotranspiration were analysed with the specific hindcasts.

- 1089
- 1090
- 1091
- 1092

1093 Appendix A Skill in the meteorological forcing before bias correction

1094

1095



1098Figure A1Skill, in terms of the percentage of cells with significant values of R, for1099three components of the raw S4 forcing. Figure A1a shows precipitation1100skill, as a function of target and lead month. The other three panels compare1101the skill of the raw S4 output (noBC) with its bias-corrected version (BC) as1102a function of the target month and for the first three lead months.1103Precipitation is plotted in Fig. A1b, temperature in Fig. A1c and incoming1104short-wave radiation in Fig. A1d.

- 1105
- 1106 1107 Sect. 3.1 is an analysis of the skill of the meteorological forcing after bias correction. Because predictability of the meteorological forcing is an interesting topic by itself, we 1108 here present an analysis of the skill of the meteorological forcing before bias correction, 1109 i.e. of the raw S4 output, limiting attention again to the three variables considered in 1110 Sect. 3.1. Fig. A1a summarizes the skill of the raw precipitation hindcasts, which 1111 should be compared with the summary for the bias-corrected hindcasts precipitation in 1112 Fig. 1b. Such a comparison is made for lead months 0, 1 and 2 in Fig. A1b. Similar 1113 1114 comparisons are made for the two-meter temperature and incoming short-wave radiation in Figs. A1c and A1d, respectively. At this level of summarizing the 1115

differences in skill between the two types of data, differences are small for precipitation 1116 and negligible for temperature and short-wave radiation. Also, patterns of skill for all 1117 three variables, such as those shown in the maps of Figs. 1 and 2, are almost identical 1118 for the bias-corrected and the raw data. The fact that differences are small is not 1119 surprising because the bias corrections hardly change the ranking of the values while 1120 the value of the correlation coefficient largely depends on the ranking of the hindcasts 1121 relative to the ranking of the observations. Results, in terms of differences in skill 1122 between raw and bias-corrected meteorological forcing, are essentially the same for the 1123 other metrics used (ROCarea and RPSS). 1124

- 1125
- 1126
- 1127 Appendix B Reliability of the hindcasts
- 1128

To complement the analysis of discrimination skill of WUSHP published in the companion paper, this appendix presents a short evaluation of the reliability of the system. Per definition forecasts are considered "reliable" when the forecast probability is an accurate estimation of the relative frequency of the predicted outcome (Mason and Stephenson, 2008). We assessed the reliability of the discharge hindcasts of the FullSH by means of so-called reliability diagrams (see Mason and Stephenson, 2008), which we produced and evaluated as follows:

1136

1139

1140 1141

- For each grid cell and combination of a category (or tercile; AN, NN and BN),
 lead month and target month we proceeded as follows:
 - Divide the 30 (number of years) observations into terciles and give them a binary number (1 if the event falls in the considered category, 0 otherwise).
- Divide the 450 (number of years x number of ensemble members)
 forecasts into terciles.
- Determine for each of the 30 years the forecast probability of the event occurring (forecast falling in the considered tercile).
 - Pair the binary observations with the forecast probabilities.
- Sort the paired data into eight bins stratified by the forecast probabilities of the event.
- Compute bin averages of the forecast probability and of the binary observations.
- Pool the results for two consecutive lead months and the three target months of
 the same season.
- 1153 The results were further processed as follows:
- They were aggregated for the entire domain and then plotted. Examples for the BN tercile and the spring months (MAM) as target are shown in

Figs. B1a-c and B2a-c with lead month number increasing from left to 1156 right. In each diagram a linear regression is applied to the data points, 1157 weighing individual points by the number of data pairs in the bins. 1158 Because tercile thresholds are set independently for observations and 1159 forecasts, the resulting line always goes through the climatological 1160 intersection (one-third in our case; see Weisheimer and Palmer, 2013) 1161 and it is insensitive to biases. As in Weisheimer and Palmer (2013) we 1162 use the slope of the line as a measure of reliability. A slope equal to 1 1163 corresponds to perfect reliability and a slope equal to 0 indicates no 1164 reliability at all. 1165

> Reliability diagrams similar to those in Figs. B1a-c and B2a-c were produced for each terrestrial grid cell, and best-fit lines and their slopes were computed. The slopes were plotted in maps, of which examples for the BN tercile and the spring months (MAM) as target are shown in Figs. B1d-f and B2d-f.



1173

1166

1167

1168 1169

1170

Figure B1 Reliability of the FullSH discharge hindcasts for the BN tercile in spring 1174 (target months MAM). Pseudo-observations were used for verification. 1175 Lead time increases from left to right. Figures B1a-c depict aggregated 1176 reliability diagrams for the full domain. The forecasted probabilities of BN 1177 discharge (horizontal axis) are collected in eight bins. The vertical co-1178 ordinate is the relative frequency of BN discharge observations for all of the 1179 forecasts in a specific bin. The solid line is the 1:1 line. The dashed line 1180 shows the best fit to the eight data points, each weighted by the number of 1181

1182observations contributing to the bin (N_{bin}). The area of the symbols is1183proportional to N_{bin}. The dotted lines are the averages of the variables along1184the two axes (one-third). Similar reliability diagrams were made for all grid1185cells individually and the slopes of the best-fit lines are plotted in Figs.1186B1d-f.

- 1187
- 1188

For the analysis it is helpful to first consider the value of the slope in two extreme 1189 cases. If pseudo-observations are used for verification and lead time approaches zero, 1190 all members of the hindcasts for a specific year approach the pseudo-observation of that 1191 year. Hence, all hindcasts fall in the same category as the observation, so the reliability 1192 diagram condenses to two points at the coordinates [0,0] and [1,1], which represent, 1193 respectively, two-third and one-third of all contributing data. In this case the hindcasts 1194 1195 are utterly sharp. The second case is when the hindcasts have no discrimination skill at all, i.e. forecast probabilities of an event are randomly paired with the outcome 1196 (whether the event occurs or not). In this case, the slope of the fitted line is equal to 1197 zero and sharpness is minimal, i.e. forecast probabilities tend to approach one-third for 1198 each of the terciles. 1199

1200

In Fig. B1 reliability is evaluated for the case of verification with pseudo-observations. 1201 For the first two lead months, the slope of the line in the diagram of the aggregated data 1202 1203 (Fig. B1a) is 0.916. Hence, during these two lead months the system is not far from 1204 being perfectly reliable and it is rather sharp with relative maxima in forecast probability in the lowest and the highest bin. Then, with progressing lead time, 1205 reliability is reduced, i.e. the slope of the aggregated data decreases to 0.767 (for lead 1206 months 2 and 3; Fig. B1b) and 0.469 (for lead months 4 and 5; Fig. B1c). Moreover, 1207 with increasing lead time sharpness is reduced, with gradually more ensemble forecasts 1208 approaching the climatological forecast, i.e. a probability of one-third for each of the 1209 terciles. 1210

1211

1212 The maps of Figs. B1d-f show the geographical distribution of the slope from the 1213 reliability diagrams. For the first two lead months most values of the slope for individual grid cells lie between 0.7 and 1.1 (Fig. B1d) and the domain-averaged slope 1214 is 0.910. At longer leads, the highest values are found in some regions with 1215 considerable amounts of discrimination skill, such as Poland and Northern Germany, 1216 Western France, and Romania and Bulgaria (see Table 1). Reliability also tends to 1217 increase towards the northeast of the continent. Domain mean values of the grid level 1218 slope are generally somewhat lower than the slope of the aggregated data. This can, at 1219 least partly, be ascribed to more scatter of individual points around the best-fit line 1220 1221 because of the much smaller sample size for individual grid cells.



1223 1224

Figure B2 Reliability of the FullSH discharge hindcasts for the BN tercile in spring (target month MAM). Real observations for large basins were used for verification. See Figure B1 for more explanation.

1228

1229

Figure B2 is analogous to Fig. B1 but instead of using the pseudo-observations, the real 1230 observations for large basins (Sect. 2.2) are taken for verification. Compared to 1231 verification with pseudo-observations, slopes are closer to zero and therefore the 1232 forecasts seem to be less reliable and more overconfident. Independent of the type of 1233 verification data, reliability for the AN tercile is almost equal to that for the BN tercile 1234 while slopes are much closer to one for the NN tercile (not shown here). Finally, levels 1235 of reliability show little variation during the year, except for the autumn (SON), when 1236 1237 slopes are smaller (not shown here).

1238

1239 Strikingly, discrimination skill and reliability have similar characteristics. Both 1240 decrease with increasing lead time, differences between the AN and BN terciles are 1241 relatively small while scores for the NN tercile are clearly inferior to those for the two 1242 outer terciles. Also, regional maxima in discrimination skill and reliability tend to 1243 coincide, scores of discrimination skill and reliability are smallest in autumn and higher 1244 for verification with pseudo-observations than for verification with real observations.

- 1245
- 1246
- 1247

1248 **References**

1249

Baehr, J., Fröhlich, K., Botzet, M., Domeisen, D. I., Kornblueh, L., Notz, D., ... &
Müller, W. A. (2015). The prediction of surface temperature in the new seasonal
prediction system based on the MPI-ESM coupled climate model. Climate Dynamics,
44(9-10), 2723-2735.

- Bierkens, M. F. P., & Van Beek, L. P. H. (2009). Seasonal predictability of European
 discharge: NAO and hydrological response time. Journal of Hydrometeorology, 10(4),
 953-968.
- Crochemore, L., Ramos, M. H., Pappenberger, F., Andel, S. J. V., & Wood, A. W.
 (2016). An experiment on risk-based decision-making in water management using
 monthly probabilistic forecasts. Bulletin of the American Meteorological Society, 97(4),
 541-551.
- Day, G. N. (1985). Extended streamflow forecasting using NWSRFS. Journal of WaterResources Planning and Management, 111(2), 157-170.
- Doblas-Reyes, F. J., García-Serrano, J., Lienert, F., Biescas, A. P., & Rodrigues, L. R.
 (2013). Seasonal climate predictability and forecasting: status and prospects. Wiley
 Interdisciplinary Reviews: Climate Change, 4(4), 245-268.
- Draper, C., & Reichle, R. (2015). The impact of near-surface soil moisture assimilation
 at subseasonal, seasonal, and inter-annual timescales. Hydrology and Earth System
 Sciences, 19(12), 4831.
- Ghile, Y. B., & Schulze, R. E. (2008). Development of a framework for an integrated time-varying agrohydrological forecast system for Southern Africa: Initial results for seasonal forecasts. Water SA, 34(3), 315-322.
- Gneiting, T., Raftery, A. E., Westveld III, A. H., & Goldman, T. (2005). Calibrated
 probabilistic forecasting using ensemble model output statistics and minimum CRPS
 estimation. Monthly Weather Review, 133(5), 1098-1118.
- Greuell, W., Franssen, W. H., Biemans, H., & Hutjes, R. W. (2018). Seasonal streamflow forecasts for Europe–Part I: Hindcast verification with pseudo-and real observations. Hydrology and Earth System Sciences, 22(6), 3453-3472.
- Griessinger, N., Seibert, J., Magnusson, J., & Jonas, T. (2016). Assessing the benefit of
 snow data assimilation for runoff modelling in Alpine catchments. Hydrol. Earth Syst.
 Sci., 20, 3895-3905.

- Hagedorn, R., Doblas-Reyes, F. J., & Palmer, T. N. (2005). The rationale behind the
 success of multi-model ensembles in seasonal forecasting–I. Basic concept. Tellus A,
 57(3), 219-233.
- Hamlet, A. F., Huppert, D., & Lettenmaier, D. P. (2002). Economic value of long-lead
 streamflow forecasts for Columbia River hydropower. Journal of Water Resources
 Planning and Management, 128(2), 91-101.
- 1287 Kim, H. M., Webster, P. J., & Curry, J. A. (2012). Seasonal prediction skill of ECMWF
- 1288 System 4 and NCEP CFSv2 retrospective forecast for the Northern Hemisphere Winter.
- 1289 Climate Dynamics, 39(12), 2957-2973.
- Koster, R. D., Mahanama, S. P., Livneh, B., Lettenmaier, D. P., & Reichle, R. H.
 (2010). Skill in streamflow forecasts derived from large-scale estimates of soil moisture
 and snow. Nature Geoscience, 3(9), 613-616.
- Li, H.,, Luo, L. and Wood, E.F. (2008). Seasonal hydrologic predictions of low-flow
 conditions over eastern USA during the 2007 drought. <u>Atmospheric Science Letters</u>
 9(2): 61-66.
- Liang, X., Lettenmaier, D. P., Wood, E. F., & Burges, S. J. (1994). A simple hydrologically based model of land surface water and energy fluxes for general circulation models. Journal of Geophysical Research: Atmospheres (1984–2012), 99(D7), 14415-14428.
- Mackay, J. D., Jackson, C. R., Brookshaw, A., Scaife, A. A., Cook, J., & Ward, R. S.
 (2015). Seasonal forecasting of groundwater levels in principal aquifers of the United
 Kingdom. Journal of Hydrology, 530, 815-828.
- Madadgar, S., Moradkhani, H., & Garen, D. (2014). Towards improved post-processing
 of hydrologic forecast ensembles. Hydrological Processes, 28(1), 104-122.
- Mason, S. J., & Stephenson, D. B. (2008). How do we know whether seasonal climate
 forecasts are any good?. In Seasonal Climate: Forecasting and Managing Risk (pp. 259289). Springer Netherlands.
- Molteni, F, Stockdale, T., Balmaseda, M., Balsamo, G., Buizza, R., Ferranti, L.,
 Magnusson, L., Mogensen, K., Palmer, T., Vitart, F. (2011). The new ECMWF seasonal
 forecast system (System 4). ECMWF Technical Memorandum 656.
- 1311 Mushtaq, S., Chen, C., Hafeez, M., Maroulis, J. and Gabriel, H., 2012: The economic
- value of improved agrometeorological information to irrigators amid climate variability.
- 1313 Int. J. Climatol., 32, 567–581.

- 1314 Nijssen, B., O'Donnell, G. M., Lettenmaier, D. P., Lohmann, D., & Wood, E. F. (2001).
- 1315 Predicting the discharge of global rivers. Journal of Climate, 14(15), 3307-3323.
- 1316 Richardson, D. S. (2001). Measures of skill and value of ensemble prediction systems,
- 1317 their interrelationship and the effect of ensemble size. Quarterly Journal of the Royal
- 1318 Meteorological Society, 127(577), 2473-2489.
- 1319 Scaife, A. A., Arribas, A., Blockley, E., Brookshaw, A., Clark, R. T., Dunstone, N., ...
- 4 Hermanson, L. (2014). Skillful long-range prediction of European and North
 American winters. Geophysical Research Letters, 41(7), 2514-2519.
- Schepen, A., Wang, Q.J. and Robertson, D.E., 2014. Seasonal forecasts of Australian
 rainfall through calibration and bridging of coupled GCM outputs. Monthly Weather
 Review, 142(5), pp.1758-1770.
- Shukla, S., & Lettenmaier, D. P. (2011). Seasonal hydrologic prediction in the United
 States: understanding the role of initial hydrologic conditions and seasonal climate
 forecast skill. Hydrology and Earth System Sciences, 15(11), 3529-3538.
- 1328 Shuttleworth, J. S. (1993), Evaporation, in Handbook of Hydrology, 1992 (D. R.
- 1329 Maidment, Ed.), McGraw-Hill, New York.
- 1330 Singla, S., Céron, J. P., Martin, E., Regimbeau, F., Déqué, M., Habets, F., & Vidal, J. P.
- 1331 (2011). Predictability of soil moisture and river flows over France for the spring season.
- 1332 Hydrology & Earth System Sciences Discussions, 8(4).
- Soares, M. B., & Dessai, S. (2016). Barriers and enablers to the use of seasonal climate
 forecasts amongst organisations in Europe. Climatic Change, 137(1-2), 89-103.
- Sun, S., Jin, J., & Xue, Y. (1999). A simple snow-atmosphere-soil transfer model.
 Journal of Geophysical Research: Atmospheres, 104(D16), 19587-19597.
- Themeßl, M. J., Gobiet, A., & Leuprecht, A. (2011). Empirical-statistical downscaling
 and error correction of daily precipitation from regional climate models. International
 Journal of Climatology, 31(10), 1530-1544.
- Thober, S., Kumar, R., Sheffield, J., Mai, J., Schäfer, D., & Samaniego, L. (2015).
 Seasonal soil moisture drought prediction over Europe using the North American MultiModel Ensemble (NMME). Journal of Hydrometeorology, 16(6), 2329-2344.
- Van Dijk, A. I., Peña-Arancibia, J. L., Wood, E. F., Sheffield, J., & Beck, H. E. (2013).
 Global analysis of seasonal streamflow predictability using an ensemble prediction
 system and observations from 6192 small catchments worldwide. Water Resources
 Research, 49(5), 2729-2746.

- Viel, C., Beaulant, A. L., Soubeyroux, J. M., & Céron, J. P. (2016). How seasonal
 forecast could help a decision maker: an example of climate service for water resource
 management. Advances in Science and Research, 13, 51-55.
- 1350 Weedon, G. P., Balsamo, G., Bellouin, N., Gomes, S., Best, M. J., & Viterbo, P. (2014).
- 1351 The WFDEI meteorological forcing data set: WATCH Forcing Data methodology
- applied to ERA-Interim reanalysis data. Water Resources Research, 50(9), 7505-7514.
- Weisheimer, A., & Palmer, T. N. (2014). On the reliability of seasonal climate
 forecasts. Journal of the Royal Society Interface, 11(96), 20131162.
- Wood, A. W., Kumar, A., & Lettenmaier, D. P. (2005). A retrospective assessment of
 National Centers for Environmental Prediction climate model–based ensemble
 hydrologic forecasting in the western United States. Journal of Geophysical Research:
 Atmospheres, 110(D4).
- Wood, A. W., & Lettenmaier, D. P. (2008). An ensemble approach for attribution ofhydrologic prediction uncertainty. Geophysical Research Letters, 35(14).
- Wood, A. W., & Schaake, J. C. (2008). Correcting errors in streamflow forecastensemble mean and spread. Journal of Hydrometeorology, 9(1), 132-148.
- Wood, A. W., Hopson, T., Newman, A., Brekke, L., Arnold, J., & Clark, M. (2016).
 Quantifying Streamflow Forecast Skill Elasticity to Initial Condition and Climate
- 1365 Prediction Skill. Journal of Hydrometeorology, 17(2), 651-668.
- Yuan, X., Wood, E. F., Luo, L., & Pan, M. (2013). CFSv2-based seasonal hydroclimatic
 forecasts over the conterminous United States. Journal of Climate, 26, 4828-4847.
- Yuan, X., Wood, E. F., & Ma, Z. (2015). A review on climate-model-based seasonal
 hydrologic forecasting: physical understanding and system development. Wiley
 Interdisciplinary Reviews: Water, 2(5), 523-536.
- Yuan, X. (2016). An experimental seasonal hydrological forecasting system over the Yellow River basin – Part 2: The added value from climate forecast models, Hydrol.
- 1373 Earth Syst. Sci., 20, 2453-2466, doi:10.5194/hess-20-2453-2016.