

Contents

1	Response	1
	1.1 Review by Joost Beckers	1
	1.2 Review by Kean Foster	3
5	1.3 Comments by Fredrik Wetterhall	6
2	Output of \LaTeXdiff	7

1 Response

1.1 Review by Joost Beckers

The MOS method is presented as an option for streamflow forecasting at the seasonal time scale (Page 1-Line 16, Page 2-Line 18, Page 3-Line 9). However, the results show only forecast skill relative to climatology for the first month ahead. This is usually not what is called seasonal forecasting (rather medium- or extended-range forecasting). So the conclusion must be drawn that no skill was found at the seasonal time scale for any of the models (including ESP and H-TESSSEL). This is indeed concluded for the MOS models (Page 17-Line 15), but the suggestion that the performance may be better for particular calendar months (Page 17-Line 15,16) is unfounded and should be removed. Also, the conclusion that H-TESSSEL is an interesting option for seasonal (i.e. beyond 1 month lead time) streamflow forecasting (Page 1-Line 13, Page 17-Line 26) is not supported by the results shown in Figure 2 (at least not clear to me).

Wherever possible we removed the term 'season'. Its usage is now restricted to the 'seasonal climate predictions' of ECMWF and to paragraphs dealing with seasonal forecasting in general. In addition, we added a short subsection to the results that looks at the variation of forecast skill within the calendar year (this was also proposed by Kean Foster).

Moreover, if the MOS method is to be considered for operational streamflow forecasting, it would need to be tested against the more traditional approach, which uses ESP- or GCM-driven hydrologic models.

We removed the term 'operational' from the article; it remains in the introduction (ESP approach as the de facto standard in operational seasonal streamflow forecasting) and in the 'operational analysis' of numerical weather forecasting.

I am not sure which version of H-TESSSEL you are using, but since you mention it does not include routing (Page 16-Line 26), this must be a relatively limited model. When comparing the MOS method to this H-TESSSEL for zero lead time, the results are not very convincing: according to Table 4 the MAE of the S4* models is worse than for H-TESSSEL at Lobith and only marginally better at Basel. Given these results, would the conclusion not be that the MOS methods are not a viable option for operational streamflow forecasting? This conclusion is missing on Page 17.

We agree, it seems that the simplest approach (linear bias correction of H-TESSSEL runoff) often performs best. This point is now more stressed in the conclusion.

My final concern with this paper is that it is not clear how the skill of the various models at zero lead time is composed. Probably, the forecast skill for the first 5 days is higher than for the last 5 days of the 1-month lead time. By averaging over the entire first month this information is lost. It could very well be that the average positive skill for lead times 1-30 days is entirely due to the positive skill for the first few days. Moreover, this positive skill for the first few days may be a result of the persistence of weather patterns in the GCM, similar to that for a normal short range weather forecast.

Related to this is the ESP-revESP analysis. The paper states that there is no clear difference in skill between the ESP and revESP at zero lead time (Page 16-Line 2). But the zero lead time is actually an average for lead times of 1 to 30 days. A separate analysis for lead times of 5, 10, 15, etc days would probably reveal a cross-over from dominance of initial conditions (higher skill of ESP) for short lead times to dominance of meteorological forcing (higher skill of revESP) for longer lead times. But this cannot be seen in the monthly average. Therefore I encourage the authors to do an skill assessment at higher temporal resolution.

We did an experiment similar to the monthly analysis for five day mean streamflow and lead times of 0, 5, . . . , 175 days.

Page 3-Line 23: What approaches do these earlier studies use? Are these (bias-corrected) hydrologic model forecasting studies or do they use MOS/PP?

We tried to clarify the paragraph – these are all studies using hydrological models forced by subseasonal or seasonal climate predictions.

Page 6-Line 2: I believe the reference for H-TESSSEL should be Balsamo et al., 2008 or 2009.

We added Balsamo et al., 2009 to the references.

Page 7-Line 6: Sample size is 31? 1981-2011 is 31 years, but you leave out the year of forecast and (according to Section 4.1.3) also the two preceding and subsequent years, so n must be 26.

5 We equated n to 26.

*Page 16-Line 14: It is found that the S4PT model outperforms the ESP model for subcatchments with smooth terrain and weak influence of initial conditions. Can you explain why? I would expect the opposite: the S4PT model (which includes forecast temperature) should do well for catchments
10 that are dominated by snowmelt (rough terrain and strong influence of initial conditions).*

We only can speculate: GCM skill for the Rhine basin is on a low level, and thus hard to detect. When the initial conditions are strongly relevant like in the case of a snow dominated catchment, any error in estimating these initial
15 conditions produces larger errors than the GCM skill can reduce. Thus, we suggest that GCM skill is better detectable in catchments where the relevance of the initial conditions is small. For example, we would argue that it is hard to successfully force a hydrological model with seasonal climate predictions in a catchment situated in the Alps – if we get the snow pack wrong, the
20 small skill contained in the precipitation and temperature forecasts vanishes. This point is now included in the discussion.

1.2 Review by Kean Foster

*Is there a reason why the initial hydrological conditions are not included as predictors (page 3-lines 10-12 and table 2)? Predictors related to storages
25 such as soil moisture content, snow, and reservoir/lake levels all impact future streamflow yet only meteorological predictors are used. I agree that many of these initial storages are affected by the antecedent meteorological conditions, but these connections are not necessarily linear or significant depending on the time frame used. For example, if only predictors for the preceding
30 month are used then there is little connection to snow pack size or reservoir levels and therefore little added value. Thus I miss a description of the time period, and to a lesser extent the domain, for the predictors.*

We agree, there exist many other potential predictors. We restricted the set of predictors to precipitation and surface air temperature for practical
35 reasons: These variables are available as gridded products, cover the entire

study region (and thus are present in all subcatchments), and are available for a long time period; the assumption of independence is more or less valid; and the regression strategy stays simple. As long as it is reasonable to include precipitation and temperature of the target season in the model, then it does
5 so for the 'initial conditions' too. In fact, using this restricted set of predictors guarantees a fair comparison of the predictor combinations and spatial levels as they all rely on the same source of data.

In case of the 'preceding' predictors (the predictors that act as a proxy to catch the initial conditions via the antecedent meteorological conditions), the
10 time aggregation is allowed to vary between 10 and 720 days. The predictors are defined as catchment area averages. Two example plots showing the regression coefficients and aggregation periods of the refRun model at Lobith and Basel are now included in the additional materials.

*Similarly, I question whether the use of the terms ESP and revESP are technically correct in this paper as it stands. Without any information regarding the initial conditions at the forecast initialisation one can argue that this is not similar to what Wood and Lettenmeier (2008) meant. If it were possible I would suggest the authors include predictors that represented the initial conditions (soil moisture, snow depth, or even streamflow) otherwise they should
20 add a paragraph explaining why the current approach is still an adaption of the VESPA methodology. I believe that the latter may be difficult to justify especially with respect to revESP.*

We renamed the ESP to preMet and revESP to subMet.

*I echo Joost's point where he suggests that the suggestion that the performance may be better for particular months (page 17-lines15, 16) is unfounded as the
25 article stands now. However, I do expect this to be the case and therefore I disagree with him in that this should be removed. Rather I think it would be of interest to include some results or a section that addresses this variability. This can be done in part in the form of a figure along the lines of the one below (figure 1). Related to this, why are the authors concentrating only
30 on the general performance throughout the year? The usefulness of these forecasts may be much higher, even only, during specific times during the year e.x. during the snow melt period or low flow period.*

We added a subsection to the results that looks at the variation of forecast
35 skill within the calendar year.

With regards to H-TESSSEL, Table 4 shows that it has some skill, at least at the spatial level 1. Have the authors tested using these data as predictors in the MOS approach at levels 2 and 3?

We completed Tab. 4 with the corresponding values.

I am unclear as to whether the S_4^ data is bias corrected. It is now almost common practice for some sort pre-processing or bias correction of the S_4^* forecast data before use in hydrological forecasting studies and work. The authors note that the quality of seasonal climate predictions for the study area are low (page 3-lines 20,21) but it is not clear to me whether any attempt to bias correct the data, and if I did miss it by what method.*

We did not apply any bias correction, since we think it is not useful in case of statistical methods (at least we do not know any study that uses bias corrected predictors for a regression model). The present formulation of the MOS approaches catches any systematic linear error via the regression coefficients. Obviously, this does not hold for nonlinear systematic errors, but we question that e.g. quantile mapping improves the prediction accuracy, since we work with mean values corresponding to at least 5 days.

Lastly, the authors mention how the uncertainties in forecasts can be reduced when the quantity of interest is controlled by teleconnection phenomena (page 1-line 17-19). I don't contest that this is true but rather question how it is relevant to the paper because there does not seem to be any more references to such modulation activity or its importance in the rest of the paper.

We agree, this statement is not strictly necessary for the article. Rather we tried to sketch the basis for environmental seasonal forecasting in order to start somewhere with the article. Please note that the cited 'slowly-varying and predictable phenomena' are not restricted to the thermal coupling of the oceans and the atmosphere (and potential subsequent teleconnections) – a strong cycle of snow accumulation and subsequent melting or persistence in soil moisture are other examples. We tried to clarify the corresponding paragraph.

On page 9-line 12 the authors give a secondary citation where I feel that the original citation, or at least inclusion of the original would be strongly advised. Taylor's original article is: Taylor, K. E. (2001). Summarizing multiple aspects of model performance in a single diagram. Journal of Geophysical Research: Atmospheres, 106(D7), 7183-7192. The authors are encouraged to check their other sources.

We added Taylor, 2001 to the references. Otherwise, there are only minor changes in the list of references. The book published by the National Academies ('Next Generation Earth System Prediction') obviously is grey

literature – however, we retained it in the article since we think it is a good book: It summarises the state-of-the art in research and industry, it looks at seasonal forecasting from a broader perspective (though heavily biased towards climate predictions), and it is written and reviewed by well-known
5 experts in the field.

Lastly, there are some minor grammatical errors in the paper; however these do not detract from the readability or arguments made therein. All the same I do suggest that the authors spend a little time to minimise them if time allows.

10 We tried our best (obviously, we are not native English speakers). However, we also trust the copy-editing skills of the Copernicus team to remove the remaining errors, if the manuscript gets considered for publication.

1.3 Comments by Fredrik Wetterhall

15 *Regarding the analysis of the time aggregation, I would suggest that you add that to the paper since it is worth testing it. If the article gets too lengthy you can remove the skill vs geographical attributes to supplementary material. However, I do not feel that the paper is too long.*

We added the experiment for the five day mean streamflow to the results and retained the skill vs. geographical attributes results.

20 *I am not sure that adding MAE or MSE results from the literature would add anything to this particular study, so I would not recommend that.*

We followed your recommendation.

25 *Please also do a language check. I would in particular suggest to not use GCM, especially in the terminology of weather forecasting. The term I would suggest here is NWP (Numerical Weather Prediction) or even ESM (Earth System Model). The term GCM is too broad, for example it does not explicitly include the analysis, which is an essential part of weather forecasting.*

We tried to consistently use the term earth system model. 'GCM' remains in the text for a few exceptions, i.e. 'atmospheric GCM' or 'coupled atmosphere-ocean-land GCMs'.
30

2 Output of $\text{\LaTeX}diff$

Below you find the output of $\text{\LaTeX}diff$. Since we had to rewrite the code to gain some speed and we also rerun the complete experiment, there are some minor changes in the results (e.g. MAE values) – these changes are
5 introduced by the way the sequence of random integers is drawn in order to generate the bootstrap replicates. However, these changes neither affect the discussion of the results nor the conclusions.

Please excuse the layout of Table 5 – we have no clue why $\text{\LaTeX}diff$ fails. Eventually the table contains improper tex code, however, we could not solve
10 this issue.

Monthly streamflow forecasting at varying spatial scales in the Rhine basin

Simon Schick^{1,2}, Ole Rössler^{1,2}, and Rolf Weingartner^{1,2}

¹Institute of Geography, University of Bern, Bern, Switzerland

²Oeschger Centre for Climate Change Research, University of Bern, Bern, Switzerland

Correspondence to: Simon Schick (simon.schick@giub.unibe.ch)

Abstract. Model output statistics (MOS) methods can be used to empirically relate an environmental variable of interest to predictions from ~~general circulation models (GCMs)~~ earth system models (ESMs). This variable often belongs to a spatial scale not resolved by the ~~GCM/ESM~~. Here, using the linear model fitted by least squares, we regress monthly mean streamflow of the Rhine River at Lobith and Basel against seasonal predictions of precipitation, surface air temperature, and runoff from the European Centre for Medium-Range Weather Forecasts. To address potential effects of a scale mismatch between the ~~GCM/ESM~~'s horizontal grid resolution and the hydrological application, the MOS method is further tested with an experiment conducted at the subcatchment scale. This experiment applies the MOS method to 133 additional gauging stations located within the Rhine basin and combines the forecasts from the subcatchments to predict streamflow at Lobith and Basel. In so doing, the MOS method is tested for catchments areas covering four orders of magnitude. Using data from the period 1981-10 2011, the results show that skill, with respect to climatology, is restricted on average to the first month ahead. This result holds for both the predictor combination that mimics the initial conditions and the predictor combinations that additionally include the dynamical seasonal predictions. The latter, however, reduces the mean absolute error of the former in the range of 5 to 11 percent, which is consistently reproduced at the subcatchment scale. ~~The results further indicate that bias-corrected runoff from the H-TESSEL land surface model is an interesting option when it comes to seasonal streamflow forecasting in large~~
15 ~~river basins~~ An additional experiment conducted for five day mean streamflow indicates that the dynamical predictions help to reduce uncertainties up to about 20 days ahead, but also reveals some shortcomings of the present MOS method.

1 Introduction

Environmental forecasting at the subseasonal to seasonal time scale promises a basis for planning in e.g. energy production, agriculture, shipping, or water resources management. While the uncertainties of these forecasts are inherently large, they can
20 be reduced when the quantity of interest is controlled by slowly-varying and predictable phenomena, ~~of which~~. For example, the El Niño-Southern Oscillation might be the most prominent one plays an important role in predicting the atmosphere, and snow accumulation and melting often forms the backbone in predicting hydrological variables of the land surface (National Academies, 2016).

In case of streamflow forecasting the ESP-revESP experiment proposed by Wood and Lettenmaier (2008) provides a methodological framework to disentangle forecast uncertainty with respect to the initial conditions and the meteorological forcings. Being a retrospective simulation, the experiment consists of model runs where the initial conditions are assumed to be known and the meteorological forcing series are randomly drawn (ESP, Ensemble Streamflow Prediction) and vice versa (revESP, reverse Ensemble Streamflow Prediction). In this context the initial conditions refer to the spatial distribution, volume, and phase of water in the catchment at the date of prediction.

The framework allows for the estimation of the time range at which the initial conditions control the generation of streamflow: When the prediction error of the ESP simulation exceeds that of the revESP simulation, the meteorological forcings start to dominate the streamflow generation. Similarly, when the prediction error of the ESP simulation approaches the prediction error of the climatology (i.e. average streamflow used as naive prediction strategy), the initial conditions no longer control the streamflow generation.

In both cases this time range depends on the interplay between climatological features (e.g. transitions between wet and dry or cold and warm seasons) and catchment specific hydrological storages (e.g. surface water bodies, soils, aquifers, and snow) and can vary from zero up to several months (van Dijk et al., 2013; Shukla et al., 2013; Yossef et al., 2013). Indeed, this source of predictability is the rationale behind the application of the ESP approach in operational forecast settings, and it can be further exploited by conditioning on climate precursors (e.g. Beckers et al., 2016).

An emerging option for ~~seasonal~~ streamflow forecasting is the integration of seasonal predictions from earth system models (ESMs), i.e. coupled atmosphere-ocean-land general circulation models (Yuan et al., 2015b). Predictions from ~~a general circulation model (GCM)~~ an ESM can be used threefold to the aim of streamflow forecasting by

1. forcing a hydrological model with the predicted evolution of the atmosphere;
2. employing runoff simulated by the land surface model, ~~eventually in combination with a routing model~~;
3. using the predicted states of the atmosphere, ocean, or land surface in a perfect prognosis or model output statistics context with the streamflow as the predictand.

The first approach requires a calibrated hydrological model for the region of interest. In order to correct a potential bias and to match the spatial and temporal resolution of the hydrological model, it further involves a postprocessing of the atmospheric fields. A postprocessing ~~also might~~ might also be applied to the streamflow forecasts to account for deficiencies of the hydrological model. See e.g. Yuan et al. (2015a) or Bennett et al. (2016) for recent implementations of such a model chain.

In the second approach the land surface model takes the hydrological model's place with the difference that the atmosphere and land surface are fully coupled. Since ~~land surface components of coupled GCMs often represent the land surface component~~ of ESMs often represents groundwater dynamics and the river routing in a simplified way (Clark et al., 2015), the simulated runoff might be fed to a routing model as e.g. in Pappenberger et al. (2010). To the best of our knowledge, this approach has not yet been tested with a specific focus on ~~the seasonal time scale~~ subseasonal or seasonal streamflow forecasting.

The third approach deals with developing an empirical prediction rule for streamflow. If the model building procedure is based on observations only, the approach is commonly referred to as perfect prognosis (PP). On the other hand, the model

might be built using the hindcast archive of a particular GCM-ESM (model output statistics, MOS). In both cases the final prediction rule is applied to the actual GCM-ESM outcome to forecast the quantity of interest. Therefore, MOS methods require the presence of a hindcast archive of the involved GCM-ESM, but can take systematic errors of the GCM-ESM into account (Brunet et al., 1988).

5 ~~Only a few studies map GCM-ESM~~ Studies that map ESM output to streamflow with PP or MOS methods ~~, including include~~ multiple linear regression (Marcos et al., 2017), principal components regression and canonical correlation analysis (Foster and Uvo, 2010; Sahu et al., 2016), ~~artificial neural networks (Humphrey et al., 2016), and an ensemble of~~ generalized linear models ~~, locally weighted polynomial regression, and k-nearest-neighbour prediction rules (Chowdhury and Sharma, 2009). By far the most selected predictor is catchment area precipitation, but depending on the study region also surface air temperature,~~
10 ~~sea surface temperature, or wind velocity are used (Slater et al., 2017), or artificial neural networks (Humphrey et al., 2016).~~ Whatever the selected predictors, PP and MOS methods ~~often~~ generally conduct the mapping across spatial scales. For example, if the catchment of interest falls below the grid scale of the GCM-ESM, PP and MOS methods implicitly perform a downscaling step. If the catchment covers several grid points, the method implicitly performs an upscaling.

The present study aims to take up this scale bridging and to test a MOS-based approach for seasonal-monthly mean stream-
15 flow forecasting and a range of catchment areas. To analyse the limits of predictability and to aid interpretation, we first ~~adapt~~ the define predictor combinations motivated by the ESP-revESP ~~framework to the context of regression by defining predictor combinations that conceptually correspond to the ESP and revESP simulations~~ framework. Next, seasonal predictions of precipitation, surface air temperature, and runoff from the European Centre for Medium-Range Weather Forecasts (ECMWF) enter the regression model equation and the resulting forecast skill is estimated with respect to the ~~ESP-like~~ ESP-inspired regression
20 model.

The variation of the catchment area borrows from the concept of the ²‘working scale’ (Blöschl and Sivapalan, 1995): Given a particular target catchment, the regression models are applied at the catchment scale as well as at two levels of subcatchment scales. In case of the ~~subcatchments~~ latter, the resulting forecasts are combined in order to get a forecast at the outlet of the target catchment. By validating the combined forecasts of the subcatchments at the main outlet, any differences in the forecast
25 quality can be attributed to the working scales.

This experiment is conducted for the Rhine River at Lobith and Basel in Western Europe. ~~In general the current quality of seasonal climate predictions is classified to be low for this region (Kim et al., 2012; Doblas-Reyes et al., 2013). Streamflow hindcast experiments with dynamical~~ Studies using subseasonal or seasonal climate predictions ~~, however, indicate indicate~~ for several parts of the Rhine basin moderate skill beyond the lead time of traditional weather forecasts. These studies apply
30 the model chain as outlined above in approach number one: Concerning catchments of the Alpine and High Rhine, Orth and Seneviratne (2013) estimate the skillful lead time for daily mean streamflow to lie between one and two weeks, which increases to about one month when focusing on low flows (Fundel et al., 2013; Jörg-Hess et al., 2015). Also for daily low flow Demirel et al. (2015) report for the Moselle River a sharp decrease in skill after 30 days. For a set of French catchments Crochemore et al. (2016) show that weekly streamflow forecasts are improved for lead times up to about one month when

using postprocessed ~~seasonal~~ precipitation predictions. Singla et al. (2012) advance spring mean streamflow forecasts for the French part of the Rhine basin with seasonal predictions of precipitation and surface air temperature.

~~The above studies show that in case of the Rhine basin current model chains skillfully forecast daily mean streamflows approximately two weeks ahead. When considering low flows only, these two weeks extend to about one month, and by reducing the forecasts temporal resolution even longer forecast ranges seem to be feasible.~~ As a compromise between skillful lead time and temporal resolution, we decide to focus on monthly mean streamflow at lead times of zero, one, and two months. In order to resolve the monthly time scale and to test the MOS method at shorter time intervals, an experiment is further conducted for five day mean streamflow. Here, zero lead time refers to forecasting ~~the next month~~ one time interval ahead, while e.g. ~~the a~~ a one month lead time denotes a temporal gap of one month between the release of a forecast and its time of validity.

Strictly speaking, the present study deals with hindcasts or retrospective forecasts. However, for the sake of readability we use the terms forecast, hindcast, and prediction interchangeably.

Below, Sect. 2 introduces the study region, Sect. 3 describes the data set, Sect. 4 exposes the methodology in more detail, and in Sect. 5 and 6 the results are presented and discussed, respectively.

15 2 Study region

The Rhine River is situated in Western Europe and discharges into the North Sea; in the south its basin is defined by the Alps. About 58 million people use the Rhine water for the purpose of navigation, hydro power, industry, agriculture, drinking water supply, and leisure (ICPR, 2009). The present study focuses on two gauging stations: The first is located in Lobith near the Dutch-German border, the second in Basel in the tri-border region of France, Germany, and Switzerland.

20 Table 1 lists some geographical attributes. The Rhine at Basel covers an area of approximately one fifth of the Rhine at Lobith whereas the mean elevation halves when going from Basel to Lobith. The negative minimum elevation of the Rhine at Lobith is due to a coal mine. Dominant land use classes are farmed areas and forests, but the Rhine at Basel proportionately includes more grass land, wasteland, surface water, and glacier.

Concerning the climatology of the period 1981-2011 (Fig. 1), we observe that streamflow peaks at Lobith in winter and at 25 Basel in early summer. Streamflow at Basel is dominated by snow accumulation in winter, subsequent snow melting in spring, and high precipitation in summer. At Lobith precipitation exhibits less variability and higher surface air temperature intensifies evaporation. Based on recent climate projections, it is expected that streamflow in the Rhine basin ~~is going to increase~~ increases in winter, ~~to decrease~~ decreases in summer, and ~~to slightly decrease~~ slightly decreases in its annual mean in the last third of the 21th century (Bosshard et al., 2014).

Table 1. Geography of the Rhine River at Basel and Lobith according to CORINE (2013), EU-DEM (2013), and GRDC (2016).

	Lobith	Basel
area (km ²)	159700	36000
gauging station (m a. s.)	20	250
elevation min (m a. s.)	-230	250
elevation max (m a. s.)	4060	4060
elevation mean (m a. s.)	490	1050
farmed area (%)	47.7	36.8
forest (%)	35.8	31.6
grass land (%)	3.4	11.4
urban area (%)	9.6	7.0
wasteland (%)	1.8	8.2
surface water (%)	1.4	4.0
glacier (%)	0.3	1.0

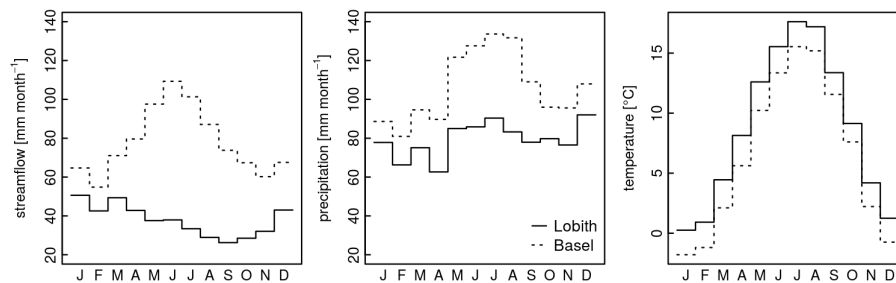


Figure 1. Monthly area averages of streamflow, precipitation, and surface air temperature for the Rhine at Lobith and Basel with respect to the period 1981-2011 (GRDC, 2016; E-OBS, 2016).

3 Data

Observations of river streamflow and gridded ~~runoff, precipitation, and precipitation~~, surface air temperature, ~~and runoff~~ of the period 1981-2011 in daily resolution constitute the data set. Throughout the study gridded quantities get aggregated to (sub)catchment area averages.

3.1 Observations

The streamflow observations consist of a set of 135 time series in $\text{m}^3 \text{s}^{-1}$. These series as well as the corresponding catchment boundaries are provided by several public authorities and the Global Runoff Data Centre (GRDC (2016); see also Sect. 9), and belong to catchments with nearly natural to heavily regulated streamflow.

- 5 The ENSEMBLES gridded observational data set in Europe (E-OBS, version 14.0) provides precipitation and surface air temperature on ~~the~~ a 0.25° regular grid (Haylock et al., 2008; E-OBS, 2016). These fields base upon the interpolation of station data and are subject to inhomogeneities and biases. However, a comparison against meteorological fields derived from denser station networks attests a high correlation (Hofstra et al., 2009). In case of the Rhine basin an E-OBS tile approximately covers an area of 500 km^2 .

10 3.2 Dynamical seasonal predictions

Precipitation, surface air temperature, and runoff from ECMWF's seasonal forecast system 4 (S4) archive are on a 0.75° regular grid. ~~This, amounting in case of the Rhine basin to a tile area of about 4500 km^2 .~~ The hindcast set consists of 15 members of which we take the ensemble mean. Runs of the coupled atmosphere-ocean-land model are initialised on the first day of each month ~~with a lead time of and simulate the subsequent~~ seven months. Up to 2010, initial conditions are out of ERA Interim, and the year 2011 is based on the operational analysis.

The atmospheric model (IFS cycle 36r4) consists of 91 vertical levels with the top level at 0.01 hPa in the mesosphere. The horizontal resolution is truncated at TL255 and the temporal discretisation equals 45 min. The NEMO ocean model has 42 levels with a horizontal resolution of about 1° . Sea ice is considered by using its actual extent from the analysis and relaxing it towards the climatology of the past five years (Molteni et al., 2011).

- 20 The H-TESEL land surface model implements four soil layers with an additional snow layer on the top. Interception, infiltration, surface runoff, and evapotranspiration are dealt with by dynamically separating a grid cell ~~in to into~~ fractions of bare ground, low and high vegetation, intercepted water, and shaded and exposed snow. In contrast, the soil properties of a particular layer are uniformly distributed within one grid cell. Vertical water movement in the soil follows Richards's equation with an additional sink term to allow for water uptake by plants. Runoff per grid cell ~~then finally equals~~ the sum of surface runoff and open drainage at the soil bottom ~~(ECMWF, 2016).~~

~~In case of the Rhine basin an E-OBS tile in the above configuration approximately covers an area of 500, and an S4 tile an area of about 4500 (Balsamo et al., 2009; ECMWF, 2016).~~

4 Method

The following subsections outline the experiment, which is individually conducted for both the Rhine at Lobith and Basel.

- 30 Section 4.1 ~~first~~ details the predictor combinations and the regression strategy. ~~Section, Sect.~~ 4.2 introduces the variation of the catchment area, and Sect. 4.3 illustrates the validation of the resulting hindcasts.

Table 2. Predictor combinations consisting of (with respect to the date of prediction) preceding and subsequent precipitation (p), surface air temperature (t), and runoff (q); the numerical values are either out of [the E-OBS gridded data set](#) or [the ECMWF's S4 hindcast archive](#).

model	preceding		subsequent		
	p^{pre}	t^{pre}	p^{sub}	t^{sub}	q^{sub}
refRun	E-OBS	E-OBS	E-OBS	E-OBS	-
ESP-preMet	E-OBS	E-OBS	-	-	-
revESP-subMet	-	-	E-OBS	E-OBS	-
S4P	E-OBS	E-OBS	S4	-	-
S4T	E-OBS	E-OBS	-	S4	-
S4PT	E-OBS	E-OBS	S4	S4	-
S4Q	E-OBS	E-OBS	-	-	S4

4.1 Model building

~~Let~~ [The predictand \$y_{i,j}\$ denote](#) ~~denotes~~ observations of mean streamflow at a specific gauging site in $\text{m}^3 \text{s}^{-1}$ for $j = 30, 60, 90$ [d, starting the first day of each calendar month \$i = 1, \dots, 12\$ \$i = 1, \dots, 12\$ in the period 1981-2011.](#) ~~Henceforth $y_{i,j}$ is the predictand.~~

5 4.1.1 Predictor combinations

The set of predictors consists of variables that either precede or succeed the date of prediction ~~, i.e. the first day of month i~~ (Tab. 2). The first model refRun (reference run) is aimed to estimate how well the regression works given the best available input data. ~~The second and third combinations imitate the ESP and revESP simulations. The ESP-revESP framework thus is mimicked by constraining the model to observed precipitation and temperature either prior to or following combinations~~ [named preMet \(preceding meteorology\) and subMet \(subsequent meteorology\) are constrained to precipitation and surface air temperature preceding and subsequent to the date of prediction forecast, respectively.](#)

The S4* combinations ~~actually~~ constitute the MOS method and consider the seasonal predictions out of the S4 hindcast archive, where we use the asterisk as wildcard to refer to any of the S4P, S4T, S4PT, and S4Q models. The S4P and S4T models are used to separate the forecast quality with respect to precipitation and temperature. The S4Q model is tested as H-TESSSEL ~~does not implement any groundwater dynamics and preceding precipitation and temperature might tap this source of predictability. Let aside the S4Q model, the preceding and subsequent predictors conceptually approximate the initial conditions and the meteorological forcings, respectively.~~

4.1.2 Regression

For a particular ~~predictor combination and~~ $y_{i,j}$ we first apply a correlation screening to select the optimal aggregation time $a_{i,j}$ for each predictor ~~according to~~

$$a_{i,j} = \underset{k}{\operatorname{argmax}} | \operatorname{cor}(y_{i,j}, x_{i,k}) | \quad (1)$$

5 where $x_{i,k}$ is one of the predictors from Tab. 2 and ~~$k = -10, -20, \dots, -720$~~ $k = -10, -20, \dots, -720$ d in case of p^{pre} and t^{pre} (backward in time relative to the date of prediction) and ~~$k = 5, 10, \dots, j$~~ $k = 5, 10, \dots, j$ d in case of p^{sub} , t^{sub} , and q^{sub} (forward in time relative to the date of prediction). The limit of 720 d is chosen since larger values rarely get selected.

The ordinary least squares hyperplane is then used for prediction without any transformation, basis expansion, or interaction. However, model variance can be an issue: Specifically for the ~~ESP-preMet~~ model from Tab. 2 we expect the signal-to-noise ratio to be low ~~in for~~ most of the ~~seasons~~ ~~predictands~~. In combination with the moderate sample size ~~$n = 31$~~ $n = 26$ for model fitting (with respect to the cross-validation, see Sect. 4.1.3), perturbations in the training set can lead to large changes in the ~~predictors~~ ~~predictor's~~ time lengths $a_{i,j}$ and regression coefficients. In order to ~~reduce~~ ~~stabilise~~ model variance, we draw 100 non-parametric bootstrap replicates of the training set, fit the model to these replicates, and combine the predictions by unweighted averaging (Breiman, 1996; Schick et al., 2016).

15 4.1.3 Cross-validation

Each year with a buffer of two years (i.e. the two preceding and subsequent years) is left out and the regression outlined in Sect. 4.1.2 is applied to the remaining years. The fitted models then predict the central left-out years. Buffering is used to avoid artificial forecast quality due to hydrometeorological persistence (Michaelsen, 1987).

4.1.4 Lead time

20 Lead time is introduced by integrating the predicted $\hat{y}_{i,j}$ in time and taking differences with respect to j . For example monthly mean streamflow z_i in July ($i = 7$) is predicted with a lead time of one month according to

$$\hat{z}_7 = (\hat{y}_{6,60} \cdot (30 + 31) \cdot \underline{sb} - \hat{y}_{6,30} \cdot 30 \cdot \underline{sb}) / (31 \cdot \underline{sb}) \quad (2)$$

where ~~$s = 24 \cdot 60 \cdot 60$~~ $b = 24 \cdot 60 \cdot 60$ s equals the number of seconds of one day and both \hat{y} and \hat{z} have unit $\text{m}^3 \text{s}^{-1}$. For zero lead time, we set $\hat{z}_i = \hat{y}_{i,30}$. Please note that the year 1981 needs to be dropped from the validation (Sect. 4.3) since the length
25 of the streamflow series prevents to forecast e.g. January 1981 with a lead time of one month.

4.2 Spatial levels

Contrasting the forecast quality of a given model for ~~individual~~ catchments separated in space inevitably implies a large number of factors, e.g. the geographic location (and thus the involved ~~GCM-ESM~~ grid points), the orography, or the degree to which streamflow is regulated. In order to hold these factors ~~whilst~~ ~~while~~ screening through a range of catchment areas, we propose
30 to vary the working scale within a particular target catchment.

Table 3. Subcatchment division of the Rhine at Lobith and Basel. The median area covers four orders of magnitude.

	number of subcatchments	area km ²		
		min	median	max
Lobith level 1	1	-	159700	-
Lobith level 2	5	19690	33220	43550
Lobith level 3	12	8284	13040	17610
Basel level 1	1	-	36000	-
Basel level 2	10	1871	2946	6346
Basel level 3	124	6	187	2654

Following this line of argumentation we apply the model building procedure from Sect. 4.1 to three distinct sets of subcatchments, which we term ‘spatial levels’ (Tab. 3). Spatial level 1 simply consists of the target catchment itself, i.e. the Rhine at Lobith and Basel. At spatial levels 2 and 3 we take additional gauging stations from within the Rhine basin, which naturally divide the basin into subcatchments.

5 For these subcatchments we have streamflow observations belonging to the entire upstream area, but not the actual subcatchment area itself. To arrive at an estimate of the water volume generated by the subcatchment, we equate the predictand $y_{i,j}$ to the difference of outflow and inflow of that subcatchment. For a particular date of prediction and spatial level, the sum of the resulting subcatchment forecasts \hat{z}_i then constitutes the final forecast for the Rhine at Lobith and Basel, respectively.

~~A drawback of this procedure is~~ This procedure implies that we ignore the water travel time: First when taking the differences of outflows and inflows and second when summing up the subcatchment forecasts. While the former increases the observational noise, the latter does not affect the regression itself, but adds a noise term to the final forecast at Lobith and Basel. As the statistical properties of the noise introduced by the water travel time is unknown, we only can argue that the results ~~below~~ provide a lower bound of the forecast quality due to this methodological constraint.

4.3 Validation

15 The forecast quality of the regression models is analysed with the pairs of cross-validated monthly mean streamflow forecasts and observations (\hat{z}, z) . These series cover the period 1982-2011 and have a sample size of $n = 360$. In general the validation is based on the mean absolute error (MAE) and Pearson’s correlation coefficient (ρ).

The first validation steps focus on the forecasts at Lobith and Basel and thus consider the sum of the ~~subcatchments~~ subcatchment forecasts \hat{z} per spatial level. The forecasts in the subcatchments itself are addressed in Sect. 4.3.5. Finally,
20 the validation of the five day mean streamflow forecasts (Sect. 4.3.6) complements the monthly analysis.

4.3.1 Benchmarks

Climatology and runoff simulated by H-TESSSEL serve as benchmarks. The ~~monthly~~ climatology is estimated with the arithmetic mean from the daily streamflow observations. ~~The monthly basin averages of~~ After averaging in time, runoff from H-TESSSEL ~~get gets~~ post-calibrated via linear regression against the ~~monthly mean streamflow observations at Lobith and~~ Basel, respectively streamflow observations per spatial level. For both benchmarks the cross-validation scheme from Sect. 4.1.3 is applied.

4.3.2 Taylor diagram

Taylor diagrams (~~Jolliffe and Stephenson, 2012~~) are employed to ~~get a global overview. For a particular model, let ρ be the~~ Pearson correlation coefficient of the forecasts \hat{z} and the corresponding observations z (Taylor, 2001) provide an instrument to contrast model performances. The plotting position of ~~the a particular~~ model has a distance from the origin equal to the standard deviation of its forecasts \hat{z} and is located on the line having an angle of incline $\phi = \arccos(\rho)$. The plotting position of the observations z has a distance from the origin equal to the standard deviation of z and is located on the abscissa. The distance between ~~the these~~ two plotting positions equals the root mean squared error with the unconditional bias $E(\hat{Z} - Z)$ removed.

15 4.3.3 ~~Mean absolute error~~ Statistical significance

~~The statistical significance of the difference in forecast accuracy between the ESP and a S4* model is tested in terms of the mean absolute error (MAE). As~~ In case of the monthly analysis it turns out ~~that~~ the paired differences of absolute errors for a given lead time and spatial level, spatial level, and reference model r

$$d = | \hat{z}^{\text{ESPr}} - z | - | \hat{z}^{\text{S4}^*} - z | \quad (3)$$

20 no longer exhibit serial correlation and approximately follow a Gaussian distribution. Using the mean difference \bar{d} , we then report the p-values of the two-sided t-test with null hypothesis $\bar{d} = 0$ and alternative hypothesis $\bar{d} \neq 0$. The sample autocorrelation functions and quantile plots against the Gaussian distribution of d for zero lead time and r being the preMet model are included in the additional materials (Sect. 10).

4.3.4 Skill

25 To evaluate whether a particular model m has skill with respect to a reference model r the MAE ratio

$$s = 1 - \frac{\text{MAE}_m}{\text{MAE}_r} \quad (4)$$

is employed. For example, m could be a S4* model and r the ESP-preMet model. $s = 0.1$ means that the model m lowers the MAE of model r by 10 %.

4.3.5 Subcatchments

To help in the interpretation of the forecast quality of the MOS method regarding the spatial levels at Lobith and Basel, we ~~finally have a look at the subcatchments itself, which are up to now only implicitly addressed. In plot in~~ a qualitative manner ~~we plot~~ the MAE skill score (Eq. 4) of the S4* and ~~ESP_preMet~~ models in space as well as against the subcatchment area, the median of the terrain roughness, the MAE skill score of the ~~revESP with the ESP_subMet with the preMet~~ model as reference, and the MAE skill score of the refRun model with the climatology as reference.

The terrain roughness is included since the atmospheric flow in complex terrain is challenging to simulate and atmospheric ~~GCMs-general circulation models~~ need to filter the topography according to their spatial resolution (Maraun and Widmann, 2015; Torma et al., 2015). The terrain roughness is defined as the difference of the maximum and minimum elevation value within a 3 times 3 pixel window (Wilson et al., 2007). It is derived here from the digital elevation model EU-DEM (2013), which has a horizontal resolution of 25 m.

4.3.6 Five day mean streamflow

In order to predict five day mean streamflow, Eq. 2 is used with a step size of five days. However, the monthly date of predictions impose some restrictions to the validation: First, it is not possible to derive regular time series at different lead times as in the monthly analysis. Furthermore, the distributional assumptions required for the statistical test from Sect. 4.3.3 are not valid. The results of the five day mean streamflow experiment thus are restricted to a qualitative interpretation.

5 Results

The experiment spans several dimensions (i.e. Lobith versus Basel, date of prediction, lead times, predictor combinations, spatial levels), so we frequently need to collapse one or several dimensions. The additional materials as listed in Sect. 10 try to complete the results as presented bellow.

5.1 Taylor diagram

Figure 2 shows the Taylor diagrams for Lobith and Basel to get a global overview regarding the lead times, predictor combinations, and spatial levels. Accurate forecasts reproduce the standard deviation of the observations (thus lie on the circle with radius equal to the ~~the~~ standard deviation of the observations), and also exhibit high correlation (so travel on this circle towards the observations on the abscissa). At a first glimpse the spatial levels do not introduce clear differences and most of the models mass at the same spots.

The benchmark climatology is outperformed at zero lead time by all models. At longer lead times the ~~revESP_subMet~~ model pops up besides the refRun model and the remaining models approach climatology. ~~H-TESEL-stays close to the regression models and tends to score a higher correlation in case of Lobith, but not Basel.~~ For the refRun model we note a correlation of about 0.9 independently of the lead time while the ~~observations-observation's~~ variability generally is underestimated.

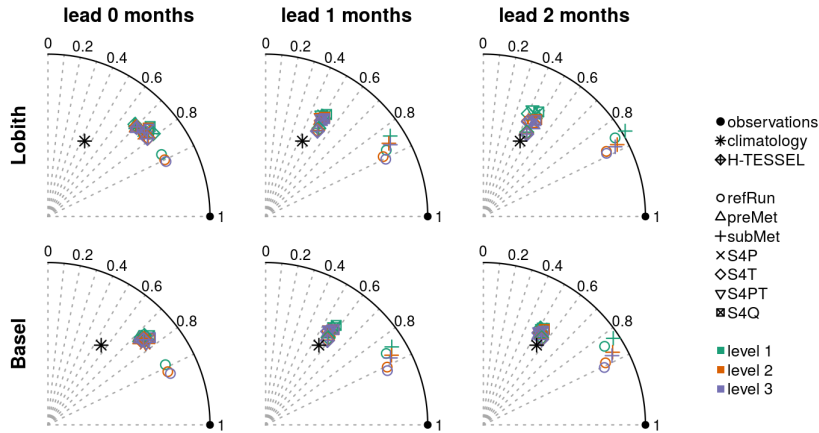


Figure 2. Taylor diagrams for the benchmarks climatology and H-TESESEL and the predictor combinations from Tab. 2 at Lobith (top row) and Basel (bottom row); $n = 360$.

For Lobith and zero lead time we observe an elongated cluster, which comprises all models but the climatology and the refRun model. Some models score a higher correlation – zooming in would reveal that these are the S4P, S4PT, and S4Q models with H-TESESEL standing at the forefront. ~~In the following we focus on the forecasts with~~

5.2 Date of prediction versus lead time

- 5 Figure 3 takes a closer look at the clusters in Fig. 2 at hand of the S4PT model and in addition breaks down the prediction skill into the different calendar months. Please note that the ordinate lists the calendar month and not the date of prediction – e.g. the top rows show the skill in predicting January’s mean streamflow for lead times of zero up to two months. Crosses indicate p-values smaller than 0.05 when Eq. 3 is applied to the individual calendar months.

In general, the patterns repeat more or less along the spatial levels and the S4PT model beats the reference models in the denominator of Eq. 4 only at zero lead time~~since at longer lead times we virtually do not have any improvements relative to the climatology.~~ An exception can be observed for June, for which the S4PT model most likely outperforms the climatology at one month lead time.

For May, the S4PT model outscores both the preMet and the subMet model. While significant differences between the S4PT and the preMet models are rare, the subMet model starts to outperform the S4PT model already at a lead time of one month. The comparison against the bias corrected H-TESESEL runoff shows that the S4PT model might provide more accurate predictions for late spring and early summer, but not otherwise.

5.3 Mean absolute error

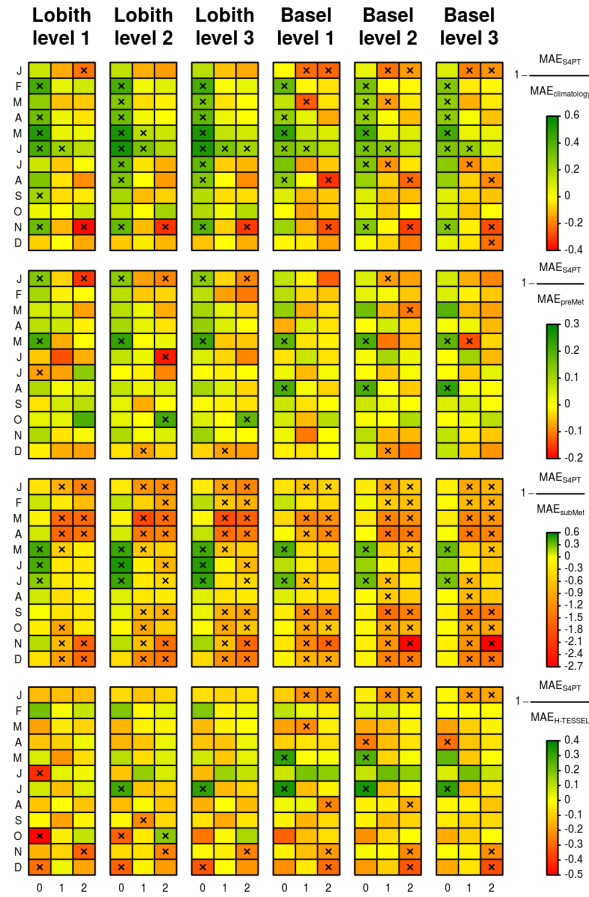


Figure 3. Taylor diagrams for MAE skill score of the S4PT model with respect to the benchmarks climatology, the preMet and subMet models, and bias corrected H-TESEL runoff. The ordinate depicts the calendar month and the predictor-combinations from Tababscissa the monthly lead time. \times at Lobith (top-row) Crosses indicate p-values smaller than 0.05 for the null hypothesis 'the reference model in the denominator and Basel (bottom-row) the S4PT model score an equal mean absolute error'; $n = 30$.

Table In order to conclude the analysis of the monthly predictions at Lobith and Basel, Tab. 4 reports the mean absolute error (MAE) MAE at zero lead time. Reading Tab. 4 along the rows reveals a more or less consistent pattern: The refRun model approximately halves the MAE of the climatology; differences between the ESP, revESP, preMet, subMet, and S4T models are small; compared to the ESP-preMet model, the S4P, S4PT, and S4Q models lower the MAE by about $40 \text{ m}^3 \text{ s}^{-1}$ for Lobith and by about $15 \text{ m}^3 \text{ s}^{-1}$ for Basel; and H-TESEL outperforms the S4* models in case of Lobith, but not Basel. When reading Tab. 4 along the columns, we generally note at Lobith a decreasing MAE when going from spatial level 1 to spatial level 3. In case of Basel, the MAE remains more or less constant except for the refRun model.

Focusing on the MOS method, Tab. 4.3.3 contains the 5 lists the corresponding MAE skill score (Eq. 4) of the S4* models using the preMet model as the reference. The p-values for the null hypothesis 'the ESP-preMet and S4* models score an equal

Table 4. Mean absolute error at zero lead time of the benchmarks climatology and H-TESEL and the predictor combinations from Tab. 2, rounded to integers. All values have unit $\text{m}^3 \text{s}^{-1}$; $n = 360$.

	climatology	H-TESEL	refRun	ESP-preMet	revESP-subMet	S4P	S4T	S4PT	S4Q
Lobith level 1	633	419	334	499-500	499-498	460-459	506-503	464-467	446-445
Lobith level 2	-633	-417	299-295	484	497-494	440	484	445	442
Lobith level 3	-633	-417	288-287	480-482	495	437-436	479-481	442-441	439
Basel level 1	239	191	130-131	201-199	194-195	188-189	195-196	187-188	190-189
Basel level 2	-239	-186	118-117	199	192	185-184	194	184	187-186
Basel level 3	-239	-184	113-112	199	193	184	195	183	187

Table 5. MAE skill score of the S4* models relative to the preMet model (Eq. 4, expressed in percent) at zero month lead time. p-values for the null hypothesis 'the ESP-preMet and S4* models score an equal mean absolute error' at zero lead time are enclosed in brackets; $n = 360$.

	S4P	S4P
Lobith level 1		<u>8</u>
Lobith level 2		<u>9</u>
Lobith level 3		<u>10</u>
Basel level 1	<0.01-0.03-5	
Basel level 2	<0.01-0.08-2	<0.01-7 (<0.01 MAE skill score of the S4* models relative to the ESP model (Eq. 4), expressed in percent; $n = 360$)
Basel level 3	+6-3	<u>7</u>

mean absolute error'. Apart from the S4T model the results among the spatial levels agree. While at Lobith the null hypothesis for the S4T model should not be rejected, at Basel one might do so.

Table 5 shows the corresponding MAE skill score (Eq. 4) using the ESP model as reference. The are listed in brackets. We see that the S4P, S4PT, and S4Q models score an error reduction ranging from 5 to 11 %. In case of the S4T model an error reduction is either not existent (Lobith) or small (Basel), supporting the which comes along with high p-values from Tab. 4.3.3. The MAE reduction generally tends to increase along the spatial levels, however, on a rather low level.

In order to reduce the number of models, we drop the S4P, S4T, and S4Q models and focus in the next section on the S4PT model. Temperature is retained as predictor because the S4T model might not be rejected at Basel (Tab. 4.3.3). Among the similar performing S4P, S4PT, and S4Q models, the S4PT model is selected for ease of interpretation as the refRun, ESP, revESP, and S4PT models share the same predictors. For the sake of completeness Fig. 4 and 5 below are included in the additional materials for the dropped S4* models (Sect. 10).

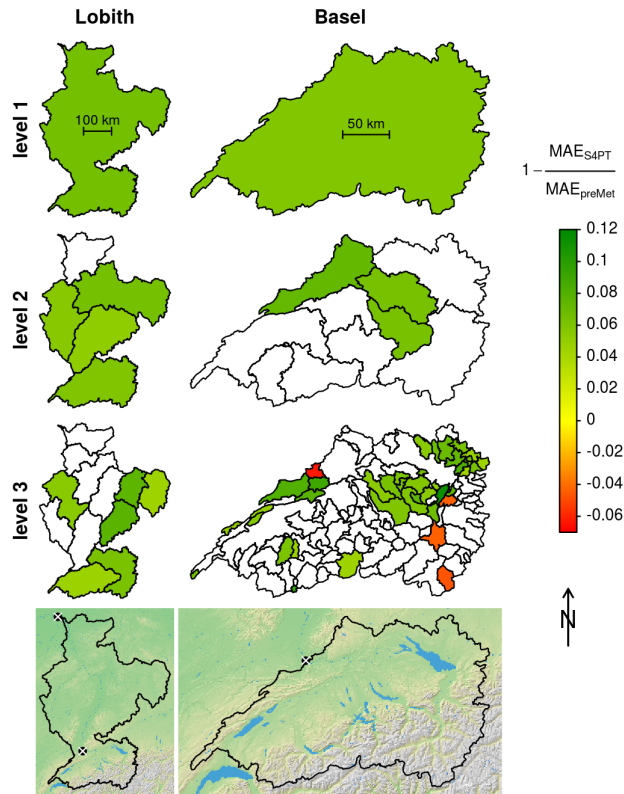


Figure 4. MAE skill score of the S4PT model with respect to the [ESP-preMet](#) model for each subcatchment and zero lead time. Subcatchments are coloured only when the p-value for the null hypothesis ‘the [ESP-preMet](#) and S4PT models score an equal mean absolute error’ is smaller than 0.05. In the bottom maps the main outlets at Lobith and Basel are marked with a [black-circle-white cross](#) and open water surfaces are coloured in blue (CORINE, 2013; EU-DEM, 2013); $n = 360$.

5.4 Subcatchments

Figure 4 depicts the MAE skill score (Eq. 4) for the S4PT model relative to the [ESP-preMet](#) model for each subcatchment at zero lead time. If the MAE difference does not exhibit a p-value smaller than 0.05 (Eq. 3), the subcatchment is coloured in white. We observe that the MAE skill score takes values in the range of [about](#) -0.06 to [0.12-0.11](#) and both the lowest and highest scores occur at Basel and spatial level 3. Negative scores can only be found at Basel and spatial level 3, and positive skill tends to cluster in space.

The same skill scores from Fig. 4 are contrasted in Fig. 5 with the subcatchment area, the median of the terrain roughness, the MAE skill score of the [revESP-subMet](#) model relative to the [ESP-preMet](#) model, and the MAE skill score of the refRun model relative to the climatology. [If the MAE difference of the S4PT and the preMet models does not exhibit a p-value smaller than 0.05, the symbol is drawn with a reduced size.](#) The horizontal lines depict the MAE skill scores from Tab. 5.

While the first two attributes concern the geography of the subcatchment, the third attribute indicates the relevance of the initial conditions for the subsequent generation of streamflow. The fourth attribute shows how well the S4PT model performs relative to the climatology as benchmark, when it has access to the best available input data.

~~In addition to the MAE skill scores of the subcatchments, the horizontal lines in Fig. 5 depict the MAE skill scores for each spatial level at Lobith and Basel (i.e. the values from Tab. 5). If the MAE difference does not exhibit a p-value smaller than 0.05 (Eq. 3), the symbol is drawn with a reduced size.~~

The resulting patterns suggest that positive skill does not depend on the subcatchment area. On the other hand, a low terrain roughness and a weak relevance of the initial conditions seem to favour positive skill. The last row finally indicates that positive skill is restricted to subcatchments where the refRun model outperforms climatology. Roughly, a hypothetical ~~linear~~ relationship seems relationship appears to strengthen from the top to the bottom plots.

6 Discussion

~~The following discussion is valid only for predicting monthly mean streamflow throughout the complete calendar year. An evaluation of~~

5.1 Five day mean streamflow

Figure 6 shows the correlation coefficient of the five day mean streamflow observations and corresponding predictions for all models and benchmarks up to a lead time of 45 days. We observe that the refRun model scores a correlation of about 0.8 with a slowly decreasing tendency towards longer lead times. Furthermore, the subMet model crosses the preMet model approximately in the second week; the preMet model approaches climatology within about three weeks; and the subMet model comes close to the ~~forecast quality with respect to particular calendar months goes beyond the scope of the study~~ refRun model in about three weeks.

In addition, we see that the bias corrected H-TESSSEL runoff starts rather cautious, but seems to slightly outperform the S4* models at longer lead times. While the S4T model is hardly distinguishable from the preMet model, the S4P, S4PT, and S4Q models appear to outperform the preMet model within the first 20 days (Lobith) and 15 days (Basel).

For the full range of lead times, the spatial levels introduce some clear differences (Fig. 7): The refRun and subMet models get improved at longer lead times along the spatial levels. For lead times longer than about 50 days, the bias corrected H-TESSSEL runoff stays in close harmony to the climatology, while the S4* and preMet models instead start to score a smaller correlation. This effect seems to be mitigated at spatial levels 2 and 3.

6 Discussion

6.1 Model building

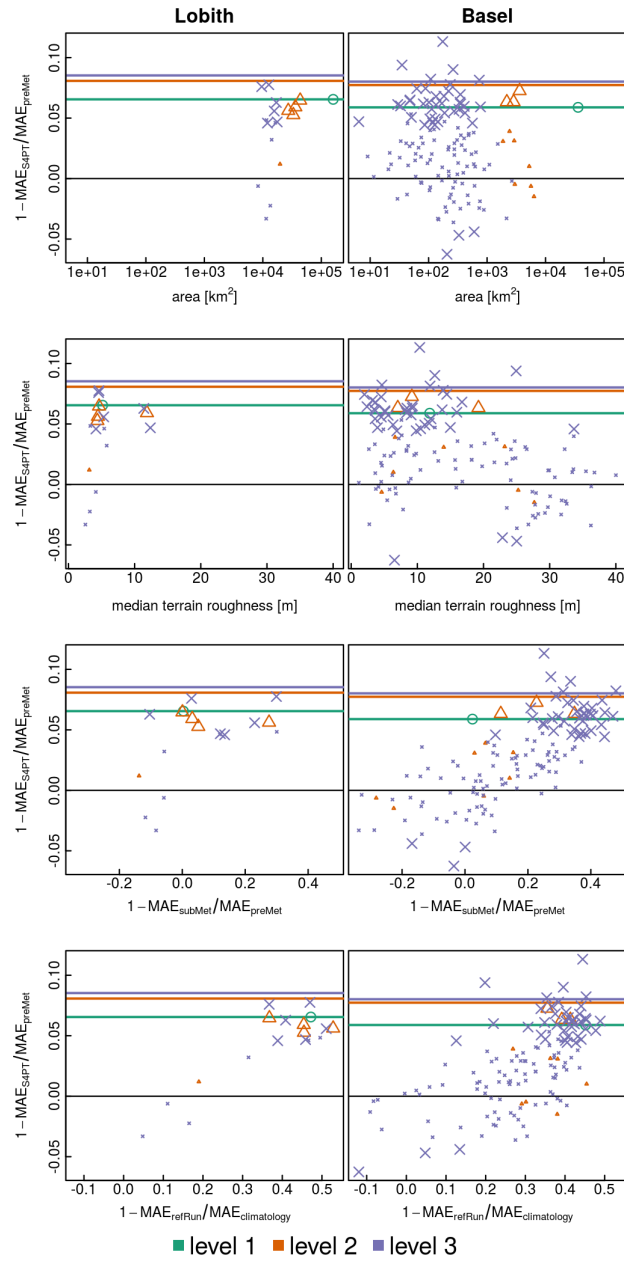


Figure 5. MAE skill score of the S4PT model with respect to the [ESP-preMet](#) model for each subcatchment and zero lead time, plotted against subcatchment attributes (see Sect. 4.3.5 for details). [Lines indicate the corresponding skill per spatial level at Lobith and Basel.](#) Large symbols note a p-value smaller than 0.05 for the null hypothesis ‘the [ESP-preMet](#) and S4PT models score an equal mean absolute error’. [The horizontal lines indicate the corresponding skill per spatial level at Lobith and Basel; \$n = 360\$.](#)

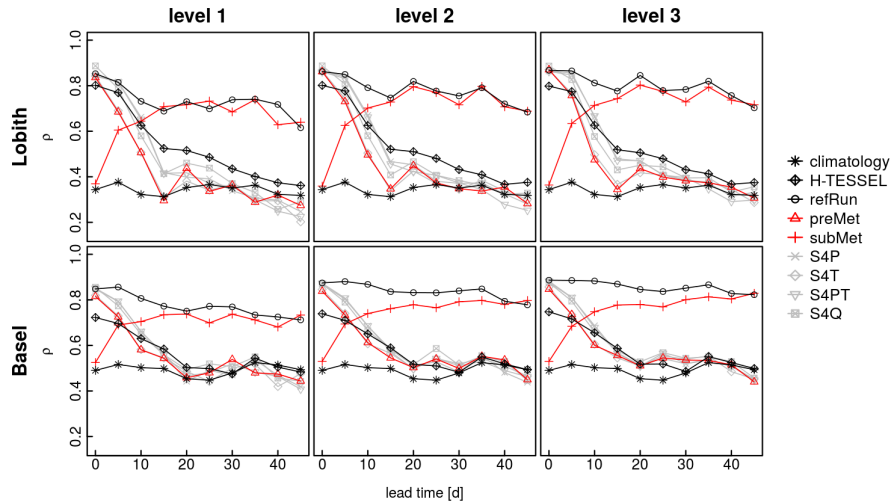


Figure 6. Correlation coefficient of five day mean streamflow observations and predictions for lead times up to 45 days; $n = 360$.

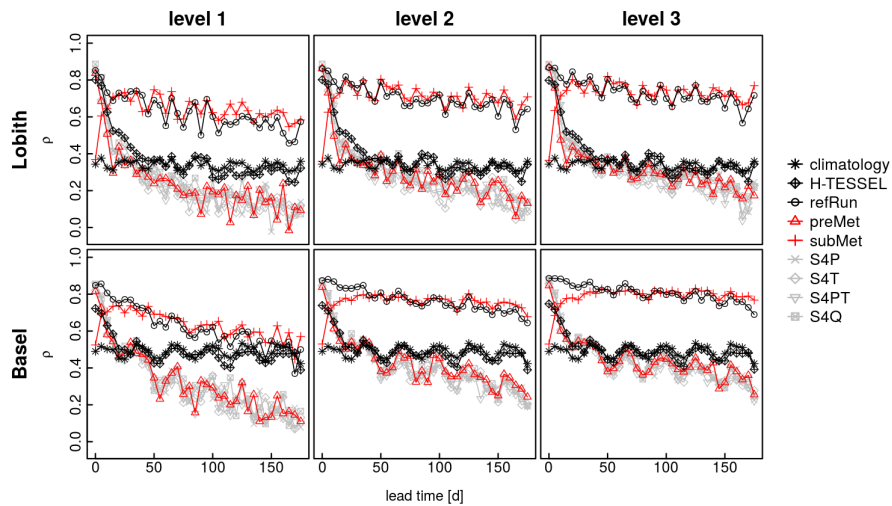


Figure 7. Correlation coefficient of five day mean streamflow observations and predictions for lead times up to 175 days; $n = 360$.

The refRun model, which has access to the best available input data, In case of the monthly streamflow, the refRun model ends up with a correlation of about 0.9 for all lead times, spatial levels, and both Lobith and Basel (Fig. 2). Part of this correlation is also the annual cycle (Fig. 1), which already leads to a correlation of about 0.5 when using the climatology as prediction rule. The forecasts from the refRun model do not fully reproduce the observations' variance, what might be improved with a transformation of the predictand (Wang et al., 2012). This option – along with predictors that more explicitly represent the

initial conditions, e.g. lake levels, soil moisture content, or snow courses – preferably should be tested in a future study with a small number of catchments and longer time series.

For the five day mean streamflow the refRun model gets degraded. At short lead times the correlation amounts to about 0.8, while for longer lead times the correlation exhibits a decreasing trend. Either the present model formulation is less valid (especially for small values of j , say 5 or 10 days, the assumption of linearity might fail) or the scheme to introduce the lead time (Eq. 2) is not appropriate for mean values of small time windows (e.g. the subtraction of streamflow volumes of 155 and 150 days only allows for small prediction errors). Since the final forecast values are not part of the regression equation, it is even possible to perform worse than climatology (Fig. 7).

6.2 Spatial levels

Besides the ignorance of the water travel time (Sect. 4.2), the spatial levels basically can degrade. The spatial levels can affect the forecast quality in three ways. For a particular subcatchment, either

- via the ignorance of the assumption of a linear relationship between the predictors and the predictand might not be valid;
- ~~the present variables precipitation, surface air temperature, and runoff simply do not contain any relevant information (for example due to heavily regulated streamflow);~~ water travel time (Sect. 4.2)
- or the aggregation of the E-OBS and S4 fields at the catchment scale is not the appropriate spatial resolution (e.g. large scale grid averages cancel any spatial variability, and for catchment areas below the grid scale a grid point does not necessarily contain information valid at the local scale).

~~Despite these three sources of uncertainty and the ignorance of the water travel time, we only observe a small gradual improvement of the forecast accuracy along~~

However, clear differences between the spatial levels (Tab. 4). ~~While this result does not allow to relate the forecast accuracy to these uncertainties, it supports at least the robustness of the estimated MAE skill score for the forecasts at Lobith and Basel (Tab. 4): Applying the regression models can only be observed for the five day streamflow predictions, where~~ at spatial levels 2 and 3 ~~virtually does not include streamflow information at Lobith and Basel (with the exception of the subcatchments that include these gauges), so artificial skill can hardly be an issue~~ the forecast quality gets improved. Using local information of precipitation, surface air temperature, or runoff appears to compensate for the ignorance of the water travel time.

6.3 ~~ESP-revESP~~ preMet-subMet

In Yossef et al. (2013) the ESP-revESP framework is applied to the ~~worlds~~ world's largest river basins using the global hydrological model PCRaster Global Water Balance (PCR-GLOBWB). Considering all calendar months and the Rhine at Lobith, the ESP simulation outperforms the climatology only at zero lead time; the revESP simulation is outperformed at zero lead time by both the ESP simulation and climatology; and at longer lead times the revESP simulation clearly outperforms both the

ESP simulation and climatology. Therefore, the results of Yossef et al. (2013) and those of the present study are ~~partly in line—initial conditions are skillful at zero lead time, but for unknown reasons a clear difference between the ESP and revESP model at zero lead time does not exist in our results~~mostly in line.

5 The analysis of the five day mean streamflow forecasts (Sect. 5.1) further reveals that the crossover of the preMet and subMet models occurs approximately in the second week. However, this estimate ignores variations within the calendar year and should be considered as a rough guess, since the regression method is far from being perfect in case of the five day mean streamflow.

6.4 MOS method

In case of the monthly mean streamflow forecasts at zero lead time, the MOS method based on precipitation or runoff provides a smaller mean absolute error than the ~~ESP-preMet~~ model (Tab. 5). Figure 6 suggests that this error reduction at the monthly
10 time scale arises from the predictions of the first 15 to 20 days. Here, it must be stressed that for the present regression strategy ~~subsequent temperature—temperature subsequent to the date of prediction~~ often is a weak predictor (~~not shown~~regression coefficients of the refRun model at spatial level 1 are included in the additional materials, see Sect. 10). Thus, a ~~possible~~ rejection of the S4T model does not allow any inference about the forecast quality of surface air temperature itself.

15 ~~While the variation of the MAE skill score along the spatial levels is small (Tab. 5), the skill in the subcatchments itself varies considerably (Fig. 4 and 5). The integration of the seasonal predictions from S4 frequently leads to negative MAE skill scores. Negative scores arise when the model catches spurious relationships, which subsequently get penalised during cross-validation. These negative scores need to be compensated in order to outperform the ESP model at Lobith and Basel—~~

Figure 5 indicates that the subcatchment area most likely is not relevant to score positive skill. Rather the S4PT model outperforms the ~~ESP-preMet~~ model in subcatchments where the terrain roughness and the relevance of the initial conditions
20 ~~is—are~~ low. However, the terrain roughness and the relevance of the initial conditions are not independent attributes: Fig. 4 shows that for small subcatchments in the alpine region positive skill is sparsely present (spatial levels 2 and 3 at Basel). These subcatchments generally exhibit a high terrain roughness as well as a high relevance of the initial conditions due to snow accumulation in winter and subsequent melting in spring and summer.

25 ~~Somewhat trivial, Fig. 5 also shows that skill of the S4PT model is restricted to subcatchments where the refRun model outperforms climatology. If the refRun model downgrades to the climatology, precipitation and temperature do not contain any relevant information to predict streamflow. Consequently also the dynamical seasonal predictions are, however accurate, useless~~A possible explanation could be that errors in the initial condition estimates outweigh the moderate skill contained in the seasonal climate predictions.

6.5 H-TESEL

30 ~~An interesting result finally is the performance of H-TESEL.~~ Within ECMWF's seasonal forecasting system S4, H-TESEL is aimed to provide a lower boundary condition for the simulation of the atmosphere and consequently does neither implement streamflow routing nor ~~ground water storage (ECMWF, 2016). According to Tab. 4~~groundwater storage (Balsamo et al., 2009; ECMWF, 2

However, H-TESEL in combination with a linear bias correction ~~best translates the seasonal predictions in case of Lobith among the models that could be used in an operational forecast setting often performs best (Tab. 4).~~

The S4Q model, which has access to the same input data and in addition conditions on preceding precipitation and temperature, scores a lower forecast accuracy than H-TESEL in case of Lobith (Tab. 4). This most likely is related to overfitting, which is not sufficiently smoothed by the model averaging (Sect. 4.1.2). ~~The question remains whether a more advanced postprocessing instead of the simple linear bias correction leads to further improvements, e.g. by conditioning on other variables or by using a river routing model.~~

7 Conclusions

The present study tests a model output statistics (MOS) method for monthly and five day mean streamflow forecasts in the Rhine basin. The method relies on the linear regression model fitted by least squares and uses predictions of precipitation and surface air temperature from the seasonal forecast system S4 of the European Centre for Medium-Range Weather Forecasts. Observations of precipitation and surface air temperature prior to the date of prediction are employed ~~to estimate as a surrogate for~~ the initial conditions. In addition, runoff simulated by the S4 land surface component, the H-TESEL land surface model, is evaluated for its predictive power.

MOS methods often bridge the grid resolution of the ~~general circulation model (GCM) dynamical model~~ and the spatial scale of the actual predictand. In order to estimate how the forecast quality depends on the catchment area, a hindcast experiment for the period 1981-2011 is conducted ~~where that varies~~ the working scale ~~is varied~~ within the Rhine basin at Lobith and Basel. This variation is implemented by applying the MOS method to subcatchments and combining the resulting forecasts to predict streamflow at the main outlets at Lobith and Basel.

~~The~~ On average, the monthly mean streamflow forecasts based on the initial conditions are skillful with respect to the climatology at zero lead time for both the Rhine at Lobith and Basel. The MOS method, which ~~additionally in addition~~ has access to the dynamical seasonal predictions, further reduces the mean absolute error by about 5 to 11 % compared to the model that is constrained to the initial conditions. ~~When the lead time is increased~~ For lead times of one and two months the forecasts virtually reduce to climatology. ~~However, for a particular calendar month these findings can substantially deviate.~~

~~The above~~ These results hold for the entire range of tested subcatchment scales. ~~Neither do, meaning that~~ effects of a scale mismatch between the ~~GCM's~~ horizontal grid resolution and the catchment area ~~emerge, nor can a subcatchment scale be detected at which do not emerge.~~ Applying the MOS method ~~clearly works best.~~ Moreover, the results indicate that a skillful ~~integration of the dynamical seasonal predictions requires catchments where the initial conditions are less relevant than the meteorological forcings~~ finally for five day mean streamflow results in a rather moderate forecast quality.

The adaptation of the ESP-revESP framework proposed by Wood and Lettenmaier (2008) to the context of regression pays off in that it provides a reference model against which the MOS method can be tested. Clearly, when using regression the ESP-revESP framework does not provide the same insights as when using a hydrological simulation model, but nevertheless it can help in the interpretation of We conclude that the results.

Given the present forecast quality of H-TESEL in combination with present model formulation – in particular the assumption of linearity – is valid for the monthly time scale, catchments with areas up to 160000 km², and water travel times similar to the Rhine river. However, the results also show that a simple linear bias correction, we also conclude that runoff simulated by the land surface component of coupled GCMs is an interesting option when it comes to operational forecasting in large river basins. In addition of the runoff predicted by the H-TESEL land surface model is hard to beat. Given the simplicity of a linear bias correction, we think that it could be interesting to establish such runoff simulations as a common benchmark in studies that use seasonal predictions from GCMs to forecast streamflow. Doing so could reveal where and why model chains, routing algorithms, MOS, and postprocessing techniques reduce uncertainties, and which hydrological processes can be implemented in a simplified manner to forecast at the seasonal time scale worth to further investigate runoff simulations from land surface components of earth system models for subseasonal to seasonal streamflow forecasting.

8 Code availability

The regression approach from Sect. 4.1.2 is compiled implemented in an R package, which is included in the additional materials maintained on github.com/schiggo.

9 Data availability

E-OBS (2016), CORINE (2013), and EU-DEM (2013) are public data sets. Access to the ECMWF and GRDC archive must be requested. Data from the various public authorities as listed in the Acknowledgements is are partly public.

10 Additional materials

Besides the R package and its vignette, the The additional materials include Fig. 4-3, Fig. 4, and Fig. 5 for the S4P, S4Q, and S4T S4* models. Figure 5-7 shows per spatial level at Lobith and Basel and for each S4* model at zero months lead time: The sample autocorrelation function and quantile plots a quantile plot against the Gaussian distribution of the (paired differences of absolute residuals with respect to the ESP model (preMet model, Eq. 3), and scatterplots as well as a scatterplot of predictions and observations. Figure 8 shows for the $y_{i,30}$ predictand the regression coefficients (for predictors standardised to mean zero and standard deviation one) and the aggregation periods $a_{i,j}$ (Eq. 1) of the refRun model at spatial level one ($n = 100$ due to the bootstrap resampling).

Acknowledgements. Streamflow series and catchment boundaries are provided by the following public authorities: State Institute for the Environment, Measurements and Conservation Baden Wuerttemberg; Bavarian Environmental Agency; State of Vorarlberg; Austrian Federal Ministry of Agriculture, Forestry, Environment and Water; and Swiss Federal Office for the Environment. Further we acknowledge the E-OBS data set from the EU-FP6 project ENSEMBLES (ensembles-eu.metoffice.com) and the data providers in the ECA&D project (www.ecad.eu) as well as the Copernicus data and information funded by the European Union (EU-DEM and CORINE). We also thank the Global Runoff

Data Centre and the European Centre for Medium-Range Weather Forecasts for the access to the data archives. The ~~study is~~ [reviews and comments by Joost Beckers, Kean Foster, and Fredrik Wetterhall substantially improved the manuscript. The study was](#) funded by the Group of Hydrology, which is part of the Institute of Geography at the University of Bern, Bern, Switzerland.

References

- Balsamo, G., Beljaars, A., Scipal, K., Viterbo, P., van den Hurk, B., Hirschi, M., and Betts, A. K.: A Revised Hydrology for the ECMWF Model: Verification from Field Site to Terrestrial Water Storage and Impact in the Integrated Forecast System, *J. Hydrometeorol.*, 10, 623–643, doi:10.1175/2008JHM1068.1, 2009.
- 5 Beckers, J. V. L., Weerts, A. H., Tijdeman, E., and Welles, E.: ENSO-conditioned weather resampling method for seasonal ensemble streamflow prediction, *Hydrol. Earth. Syst. Sc.*, 20, 3277–3287, doi:10.5194/hess-20-3277-2016, 2016.
- Bennett, J. C., Wang, Q. J., Li, M., Robertson, D. E., and Schepen, A.: Reliable long-range ensemble streamflow forecasts: Combining calibrated climate forecasts with a conceptual runoff model and a staged error model, *Water Resour. Res.*, 52, 8238–8259, doi:10.1002/2016WR019193, 2016.
- 10 Blöschl, G. and Sivapalan, M.: Scale issues in hydrological modelling: A review, *Hydrol. Process.*, 9, 251–290, doi:10.1002/hyp.3360090305, 1995.
- Bosshard, T., Kotlarski, S., Zappa, M., and Schär, C.: Hydrological Climate-Impact Projections for the Rhine River: GCM–RCM Uncertainty and Separate Temperature and Precipitation Effects, *J. Hydrometeorol.*, 15, 697–713, doi:10.1175/JHM-D-12-098.1, 2014.
- Breiman, L.: Bagging Predictors, *Mach. Learn.*, 24, 123–140, doi:10.1023/A:1018054314350, 1996.
- 15 Brunet, N., Verret, R., and Yacowar, N.: An Objective Comparison of Model Output Statistics and “Perfect Prog” Systems in Producing Numerical Weather Element Forecasts, *Weather Forecast.*, 3, 273–283, doi:10.1175/1520-0434(1988)003<0273:AOCOMO>2.0.CO;2, 1988.
- Chowdhury, S. and Sharma, A.: Multisite seasonal forecast of arid river flows using a dynamic model combination approach, *Water Resour. Res.*, 45, doi:10.1029/2008WR007510, 2009.
- 20 Clark, M. P., Fan, Y., Lawrence, D. M., Adam, J. C., Bolster, D., Gochis, D. J., Hooper, R. P., Kumar, M., Leung, L. R., Mackay, D. S., Maxwell, R. M., Shen, C., Swenson, S. C., and Zeng, X.: Improving the representation of hydrologic processes in Earth System Models, *Water Resour. Res.*, 51, 5929–5956, doi:10.1002/2015WR017096, 2015.
- CORINE: Corine Land Cover 2006 raster data, <http://www.eea.europa.eu/data-and-maps/data/corine-land-cover-2006-raster-3>, last access: 26. October 2017, 2013.
- 25 Crochemore, L., Ramos, M.-H., and Pappenberger, F.: Bias correcting precipitation forecasts to improve the skill of seasonal streamflow forecasts, *Hydrol. Earth. Syst. Sc.*, 20, 3601–3618, doi:10.5194/hess-20-3601-2016, 2016.
- Demirel, M. C., Booij, M. J., and Hoekstra, A. Y.: The skill of seasonal ensemble low-flow forecasts in the Moselle River for three different hydrological models, *Hydrol. Earth. Syst. Sc.*, 19, 275–291, doi:10.5194/hess-19-275-2015, 2015.
- Doblas-Reyes, F. J., García-Serrano, J., Lienert, F., Biescas, A. P., and Rodrigues, L. R. L.: Seasonal climate predictability and forecasting: status and prospects, *Wires Clim. Change*, 4, 245–268, doi:10.1002/wcc.217, 2013.
- 30 E-OBS: Daily temperature and precipitation fields in Europe, <http://www.ecad.eu/download/ensembles/ensembles.php>, last access: 26. October 2017, 2016.
- ECMWF: IFS Documentation, <http://www.ecmwf.int/en/forecasts/documentation-and-support/changes-ecmwf-model/ifs-documentation>, last access: 26. October 2017, 2016.
- 35 EU-DEM: Digital Elevation Model over Europe, <http://www.eea.europa.eu/data-and-maps/data/eu-dem>, last access: 26. October 2017, 2013.
- Foster, K. L. and Uvo, C. B.: Seasonal streamflow forecast: a GCM multi-model downscaling approach, *Hydrol. Res.*, 41, 503–507, doi:10.2166/nh.2010.143, 2010.

- Fundel, F., Jörg-Hess, S., and Zappa, M.: Monthly hydrometeorological ensemble prediction of streamflow droughts and corresponding drought indices, *Hydrol. Earth. Syst. Sc.*, 17, 395–407, doi:10.5194/hess-17-395-2013, 2013.
- GRDC: The Global Runoff Data Centre, http://www.bafg.de/GRDC/EN/Home/homepage_node.html, last access: 26. October 2017, 2016.
- Haylock, M. R., Hofstra, N., Klein Tank, A. M. G., Klok, E. J., Jones, P. D., and New, M.: A European daily high-resolution gridded data set of surface temperature and precipitation for 1950–2006, *J. Geophys. Res.-Atmos.*, 113, doi:10.1029/2008JD010201, 2008.
- Hofstra, N., Haylock, M., New, M., and Jones, P. D.: Testing E-OBS European high-resolution gridded data set of daily precipitation and surface temperature, *J. Geophys. Res.-Atmos.*, 114, doi:10.1029/2009JD011799, 2009.
- Humphrey, G. B., Gibbs, M. S., Dandy, G. C., and Maier, H. R.: A hybrid approach to monthly streamflow forecasting: Integrating hydrological model outputs into a Bayesian artificial neural network, *J. Hydrol.*, 540, 623– 640, doi:10.1016/j.jhydrol.2016.06.026, 2016.
- 10 ICPR: Internationally Coordinated Management Plan for the International River Basin District of the Rhine, International Commission for the Protection of the Rhine, 2009.
- Jolliffe, I. T. and Stephenson, D. B.: *Forecast Verification*, John Wiley & Sons, Ltd, 2 edn., doi:10.1002/9781119960003.ch1, 2012.
- Jörg-Hess, S., Griessinger, N., and Zappa, M.: Probabilistic Forecasts of Snow Water Equivalent and Runoff in Mountainous Areas, *J. Hydrometeorol.*, 16, 2169–2186, doi:10.1175/JHM-D-14-0193.1, 2015.
- 15 Kim, H.-M., Webster, P. J., Curry, J. A., and Toma, V. E.: Asian summer monsoon prediction in ECMWF System 4 and NCEP CFSv2 retrospective seasonal forecasts, *Clim. Dynam.*, 39, 2975–2991, doi:10.1007/s00382-012-1470-5, 2012.
- Maraun, D. and Widmann, M.: The representation of location by a regional climate model in complex terrain, *Hydrol. Earth. Syst. Sc.*, 19, 3449–3456, doi:10.5194/hess-19-3449-2015, 2015.
- Marcos, R., Llasat, M. C., Quintana-Seguí, P., and Turco, M.: Seasonal predictability of water resources in a Mediterranean freshwater reservoir and assessment of its utility for end-users, *Sci. Total Environ.*, 575, 681–691, doi:10.1016/j.scitotenv.2016.09.080, 2017.
- 20 Michaelsen, J.: Cross-Validation in Statistical Climate Forecast Models, *J. Clim. Appl. Meteorol.*, 26, 1589–1600, doi:10.1175/1520-0450(1987)026<1589:CVISCF>2.0.CO;2, 1987.
- Molteni, F., Stockdale, T., Balmaseda, M., Balsamo, G., Buizza, R., Ferranti, L., Magnusson, L., Mogensen, K., Palmer, T., and Vitart, F.: The new ECMWF seasonal forecast system (System 4), 2011.
- 25 National Academies: *Next Generation Earth System Prediction*, The National Academies Press, 1 edn., doi:10.17226/21873, 2016.
- Orth, R. and Seneviratne, S. I.: Predictability of soil moisture and streamflow on subseasonal timescales: A case study, *J. Geophys. Res.-Atmos.*, 118, 10,963–10,979, doi:10.1002/jgrd.50846, 2013.
- Pappenberger, F., Cloke, H. L., Balsamo, G., Ngo-Duc, T., and Oki, T.: Global runoff routing with the hydrological component of the ECMWF NWP system, *Int. J. Climatol.*, 30, 2155–2174, doi:10.1002/joc.2028, 2010.
- 30 Sahu, N., Robertson, A. W., Boer, R., Behera, S., DeWitt, D. G., Takara, K., Kumar, M., and Singh, R. B.: Probabilistic seasonal streamflow forecasts of the Citarum River, Indonesia, based on general circulation models, *Stoch. Env. Res. Risk A.*, pp. 1–12, doi:10.1007/s00477-016-1297-4, 2016.
- Schick, S., Rössler, O., and Weingartner, R.: Comparison of cross-validation and bootstrap aggregating for building a seasonal streamflow forecast model, *Proceedings of the International Association of Hydrological Sciences*, 374, 159–163, doi:10.5194/piahs-374-159-2016, 2016.
- 35 Shukla, S., Sheffield, J., Wood, E. F., and Lettenmaier, D. P.: On the sources of global land surface hydrologic predictability, *Hydrol. Earth. Syst. Sc.*, 17, 2781–2796, doi:10.5194/hess-17-2781-2013, 2013.

- Singla, S., Céron, J.-P., Martin, E., Regimbeau, F., Déqué, M., Habets, F., and Vidal, J.-P.: Predictability of soil moisture and river flows over France for the spring season, *Hydrol. Earth. Syst. Sc.*, 16, 201–216, doi:10.5194/hess-16-201-2012, 2012.
- Slater, L. J., Villarini, G., Bradley, A. A., and Vecchi, G. A.: A dynamical statistical framework for seasonal streamflow forecasting in an agricultural watershed, *Clim. Dynam.*, doi:10.1007/s00382-017-3794-7, 2017.
- 5 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys. Res.-Atmos.*, 106, 7183–7192, doi:10.1029/2000JD900719, 2001.
- Torma, C., Giorgi, F., and Coppola, E.: Added value of regional climate modeling over areas characterized by complex terrain–Precipitation over the Alps, *J. Geophys. Res.-Atmos.*, 120, 3957–3972, doi:10.1002/2014JD022781, 2015.
- van Dijk, A. I. J. M., Peña Arancibia, J. L., Wood, E. F., Sheffield, J., and Beck, H. E.: Global analysis of seasonal streamflow predictability using an ensemble prediction system and observations from 6192 small catchments worldwide, *Water Resour. Res.*, 49, 2729–2746, 10 doi:10.1002/wrcr.20251, 2013.
- Wang, Q. J., Shrestha, D. L., Robertson, D. E., and Pokhrel, P.: A log-sinh transformation for data normalization and variance stabilization, *Water Resour. Res.*, 48, doi:10.1029/2011WR010973, 2012.
- Wilson, M. F. J., O’Connell, B., Brown, C., Guinan, J. C., and Grehan, A. J.: Multiscale Terrain Analysis of Multibeam Bathymetry Data for 15 Habitat Mapping on the Continental Slope, *Mar. Geod.*, 30, 3–35, doi:10.1080/01490410701295962, 2007.
- Wood, A. W. and Lettenmaier, D. P.: An ensemble approach for attribution of hydrologic prediction uncertainty, *Geophys. Res. Lett.*, 35, doi:10.1029/2008GL034648, 2008.
- Yossef, N. C., Winsemius, H., Weerts, A., van Beek, R., and Bierkens, M. F. P.: Skill of a global seasonal streamflow forecasting system, relative roles of initial conditions and meteorological forcing, *Water Resour. Res.*, 49, 4687–4699, doi:10.1002/wrcr.20350, 2013.
- 20 Yuan, X., Roundy, J. K., Wood, E. F., and Sheffield, J.: Seasonal Forecasting of Global Hydrologic Extremes: System Development and Evaluation over GEWEX Basins, *B. Am. Meteorol. Soc.*, 96, 1895–1912, doi:10.1175/BAMS-D-14-00003.1, 2015a.
- Yuan, X., Wood, E. F., and Ma, Z.: A review on climate-model-based seasonal hydrologic forecasting: physical understanding and system development, *Wiley Interdisciplinary Reviews: Water*, 2, 523–536, doi:10.1002/wat2.1088, 2015b.