

Response to Reviewer #1

We thank again Reviewer #1 for his/her time spent providing feedback and constructive comments to our manuscript. The comments made by the reviewer are in black, and our responses are in blue.

I believe that the authors did significant improvements to their paper and I appreciated the detailed answers they gave to each reviewer.

I do remain a little puzzled by one of your answers, where you show that providing the 'true' signature value to the parameter selection changes almost nothing. In fact I would have expected that if you give all the 'true' information to the parameter selection procedure, you should have obtained the same performance as in calibration. is it the case? Or did I miss something?

Authors' reply: The reason for obtaining such similar results is the high uncertainty we assume in the 'true' (observed) value. Although we use the observed instead of the regionalized value of the signature, we assume the same uncertainty as identified for the regionalized value. If we derived a lower estimate of the uncertainty in the observed value, then we could expect to obtain results closer to the calibration results. The paragraph at the end of Section 3.3.1 has been clarified:

“It is worth noting that very similar results (not shown here) are obtained when instead of regionalized signatures, 'observed' signatures are used but with the same errors derived from regionalization. This suggests that the uncertainty around the regionalized signatures values, as well as signature information content, are the key factors leading to the results shown in Fig. 4.”

Accounting for dependencies in regionalized signatures for predictions in ungauged catchments

Susana Almeida^{1,2}, Nataliya Le Vine², Neil McIntyre^{2,3}, Thorsten Wagener^{1,4}, and Wouter Buytaert²

¹Department of Civil Engineering, University of Bristol, Bristol, UK

²Department of Civil and Environmental Engineering, Imperial College London, London, UK

³Centre for Water in the Minerals Industry, Sustainable Minerals Institute, The University of Queensland, Brisbane, Australia

⁴Cabot Institute, University of Bristol, Bristol, UK

Correspondence to: Susana Almeida (susana.almeida@bristol.ac.uk)

Abstract. A recurrent problem in hydrology is the absence of streamflow data to calibrate rainfall-runoff models. A commonly used approach in such circumstances conditions model parameters on regionalized response signatures. While several different signatures are often available to be included in this process, an outstanding challenge is the selection of signatures that provide useful and complementary information. Different signatures do not necessarily provide independent information, and this has led to signatures being omitted or included on a subjective basis. This paper presents a method that accounts for the inter-signature error correlation structure so that regional information is neither neglected nor double-counted when multiple signatures are included. Using 84 catchments from the MOPEX database, observed signatures are regressed against physical and climatic catchment attributes. The derived relationships are then utilized to assess the joint probability distribution of the signature regionalization errors that is subsequently used in a Bayesian procedure to condition a rainfall-runoff model. The results show that the consideration of the inter-signature error structure may improve predictions when the error correlations are strong. However, other uncertainties such as model structure and observational error may outweigh the importance of these correlations. Further, these other uncertainties cause some signatures to appear repeatedly to be disinformative.

1 Introduction

In many areas of the world the absence of past observational streamflow time series to calibrate rainfall-runoff models limits the ability to apply such models reliably to predict streamflow and inform effective water resources management. Whilst a large and increasing number of regions across the world are insufficiently gauged (Mishra and Coulibaly, 2009), there are also many highly mon-

itored catchments (Gupta et al., 2014). Transferring the knowledge gained in data-rich areas to ungauged catchments - a process known as regionalization - offers a possible way of overcoming the absence of streamflow observations in data-scarce regions. Several techniques for transferring information are reported in the literature (for an overview of different methods used in continuous streamflow regionalization see He et al. (2011), Peel and Blöschl (2011), and Razavi and Coulibaly (2013), and for a recent comparative assessment of some of the most commonly used methods see Parajka et al. (2013)).

A commonly applied approach is to use response signatures (e.g. the runoff ratio and the base flow index), which can provide insight into the hydrological functional behavior of a catchment (Wagener et al., 2007). Response signatures are calculated from available system output or input-output time series for numerous gauged catchments with known catchment attributes, i.e. physiographic and/or meteorological attributes (e.g. drainage area, latitude and longitude, average annual temperature, average monthly precipitation, etc.). Subsequently, statistical models relating each response signature to a set of catchment attributes can be identified. Given the attributes of an ungauged catchment, the signatures for the ungauged location can then be estimated using the derived statistical models. Numerous regional models of this type can be found in the literature (e.g. Boorman et al., 1995). These regionalized signatures can be used to constrain the prior range of streamflow simulations generated using a pre-selected rainfall-runoff model structure and hence restrict the model parameter space (Yadav et al., 2007; Zhang et al., 2008; Bulygina et al., 2009; Castiglioni et al., 2010). Advantages of this approach include: the flexibility in the selection of the response signatures allowing it to be based on the specific parts of the hydrograph that are of greatest importance for a given application and, if known, on the dominant hydrological processes of the catchment; access to readily available regional models for different signatures in the literature (such as base flow index from the Hydrology of Soil Types system (Boorman et al., 1995) and curve number from the United States Department of Agriculture's Soil Conservation Service soil and land use classification (USDA, 1986)) hence eliminating the need to build new regional regression models; the relationships between response signatures and catchment and climatic characteristics are not specific to any rainfall-runoff model nor to a particular calibration method used in the gauged catchments and are therefore not obscured by model structural error and can be used to condition any model.

Different ways of incorporating the regionalized information into a catchment model have been suggested in the literature. This includes set-theoretic approaches (e.g. Yadav et al., 2007; Winsemius et al., 2009) and more formal Bayesian data assimilation frameworks (e.g. Bulygina et al., 2009, 2011; Castiglioni et al., 2010; Singh et al., 2011). Where probability distributions characterizing regionalization quality have been estimated, a Bayesian conditioning procedure is one of the possibilities (Bulygina et al., 2009, 2011). This provides a framework for combining prior knowledge with the regionalized data and/or other sources of information (e.g. small scale physics-based

knowledge and hydrological measurements as in Bulygina et al., 2012), which has the potential to formally encompass the nature of the errors arising from the regionalization.

Conditioning a rainfall-runoff model on multiple independent signatures would reflect a spectrum of processes and in principle lead to an accurate prediction of flow time series (Parajka et al., 2013). However, regionalized signatures have correlated errors, for example if the signatures have been estimated using a common dataset of catchment attributes or using the same hydro-climatic data; and in general the correlations are expected to be stronger for pairs of signatures that represent similar functional behaviors of the catchment. This raises the questions of, not only how many and which signatures should be used, but also how to avoid double-counting of the information in signatures with correlated error distributions. Previous applications have tended to use a small number of signatures (e.g. Bulygina et al., 2009, 2011) and/or have tended to select signatures that are considered to be independent (e.g. Yadav et al., 2007). When multiple signatures are used, the correlations between the errors in the different sources of information are commonly disregarded (e.g. Bulygina et al., 2012). To make better use of information in available sets of signatures, a formal way of combining them so that information is neither double-counted nor neglected is required. Using formal methods to include autocorrelated data errors in model calibration is well-researched (e.g. Sorooshian and Dracup, 1980); an application of comparable methods in the regionalization context will allow making more formal and rigorous assessments of the value of correlated information sources.

Formally, in a Bayesian context, it is necessary to distinguish between *correlated signatures* and *correlated signature errors*. It is the correlation between the errors that should be accounted for in the likelihood function to avoid double-counting of information. It is possible to have two highly correlated signatures that are derived from independent information sources and therefore they have uncorrelated errors. In that case it would be valid to include both signatures in the likelihood function without accounting for correlation. This principle is well established when considering Bayesian calibration to a time series of flow observations, where flow values are typically strongly autocorrelated - but it is the observation error autocorrelation that is relevant to the likelihood function derivation (e.g. Sorooshian and Dracup, 1980). The same principle applies to adopting signatures as the observations. In the case study below, the signatures are derived from a common dataset using a common approach, so in practice the signature correlations are comparable to the signature error correlations; nevertheless for the sake of formality, we use the term *signature error correlations* (or *covariance*).

In this paper we introduce and test a method that considers multiple regionalized signatures, explicitly accounting for the signature error correlations. By formally accounting for the error covariance, we hypothesize that accuracy of flow predictions will generally improve and a greater number of signatures can usefully be included without introducing avoidable bias related to the duplication of information. This should allow the modeler to use all signatures available without having to select, on a more or less subjective basis, the most relevant (independent) signatures. The objective is thus

to explore how to get fuller value out of a set of regionalized information than has been achieved in
 95 past applications. The method is applied to a set of 84 United States catchments with a broad range
 of hydro-meteorological characteristics, obtained from the Model Parameter Estimation Experiment
 (MOPEX) dataset (Duan et al., 2006; Schaake et al., 2006). The impact of signature error covariance
 is assessed using pairs of signatures to condition a rainfall-runoff model. Along with the real data,
 synthetic streamflow data are used to isolate the effect of model structural error. Further, the model
 100 is conditioned on a variable number of regionalized signatures to evaluate whether an increasing
 number of signatures is justifiable when formally accounting for the error covariance.

2 Method

2.1 Bayesian Method for Signature Assimilation

Using a simple least-squares regression, observed signatures of catchments' functional responses
 105 are related to physical and climatic attributes of the catchments. Assuming that the same catchment
 attributes are available for an ungauged location, it is possible to obtain an estimate of the set of
 signatures for the location. Further, the parametric distribution of regression errors can be directly
 translated to a response signature(s) likelihood function. The likelihood function can then be used to
 update the prior available knowledge about model parameters via Bayes' law, which is expressed as

$$110 \quad p(\Theta|s^*, \mathbf{I}, M) = \frac{L(s(\Theta)|s^*, \mathbf{I}, M) \times p(\Theta|\mathbf{I}, M)}{p(s^*|\mathbf{I}, M)} \quad (1)$$

where, for one catchment, s^* represents the regionalized response signature(s); $p(\Theta|\mathbf{I}, M)$ is the
 prior distribution of parameters Θ for a model structure M and inputs \mathbf{I} ; $L(s(\Theta)|s^*, \mathbf{I}, M)$ is the
 likelihood function of the modeled response signature(s) $s(\Theta)$ given s^* , \mathbf{I} and M ; $p(s^*|\mathbf{I}, M)$ is the
 marginal density of s^* ; and $p(\Theta|s^*, \mathbf{I}, M)$ is the posterior distribution of Θ given s^* , \mathbf{I} and M . For
 115 the purpose of this paper, M is selected in advance and considered to be fixed (as it is the common
 practice in regionalization studies, Wagener and Montanari, 2011), as is \mathbf{I} for any one catchment,
 and so both these terms are dropped from (Eq. (1)) for convenience, resulting in

$$p(\Theta|s^*) = \frac{L(s(\Theta)|s^*) \times p(\Theta)}{p(s^*)} \quad (2)$$

Parameter sets are then sampled from the parameter posterior to allow an ensemble of rainfall-runoff
 120 simulations and a posterior distribution of flow at each time-step to be estimated and evaluated
 against observed flow. This can be repeated using different sets of signatures and different assump-
 tions about their error correlations.

2.2 Prior Distribution and Likelihood Function

2.2.1 Prior Distribution

125 To apply Bayes' law (Eq. (2)) it is necessary to specify the likelihood function ($L(s(\Theta)|s^*)$) in
Eq. (2)) and the prior distribution ($p(\Theta)$ in Eq. (2)). The prior is defined so that it reflects our
initial lack of knowledge. We follow Almeida et al. (2013) and sample sets of signature values
from uniform distributions representing the feasible ranges of signatures. This approach allows the
signatures to be sampled uniformly using a simple amendment to the commonly applied approach
130 of sampling from uniform parameter priors, which avoids highly skewed signature priors that have
undue influence on the posterior likelihood. More specifically, N parameter sets (N is equal to
10000 in our study) are sampled from a uniform distribution using Latin Hypercube sampling, so
that probability of each parameter set is $1/N$ (10^{-4} in our study). Subsequently, to provide parameter
samples that correspond to a uniform in signatures prior distribution, the parameter probabilities are
135 re-weighted (see Almeida et al., 2013), and used in the further posterior distribution approximation.
This allows accounting for correlation among the parameters imposed by the uniform in signatures
prior distribution.

2.2.2 Likelihood Function Approximation

The likelihood functions are defined using joint distributions of respective signature errors obtained
140 from the regionalization model. Errors introduced by the regionalization procedure may come from
at least five sources. First, errors are introduced by the fact that the regression model is estimated
using a specific sample of catchments rather than the entire population; second, differences may exist
between the observed and the true value of the response signature due, for example, to factors such
as the discharge record length and time period of record used in the computation (Kennard et al.,
145 2010); third, errors are present due to errors in the catchment properties data; fourth, errors exist
due to the incomplete set of catchment properties used as explanatory variables in the regression
equations; and, fifth, they exist due to the assumed linear regression structure. It is assumed that the
total error model for the regionalized signature(s) s^* can be estimated using the following procedure:

1. Considering all available gauged catchments stepwise regression is applied to each signature
150 independently to determine the predictors to include. The predictors are then fixed for the
remaining steps.
2. Considering all available gauged catchments, one catchment is left out and the remaining are
used in the fitting of the regression models for each signature.
3. The regression models obtained in Step 2 are used to estimate the signature values for the
155 omitted catchment.

4. The error for each signature is calculated for the omitted catchment by comparing the regionalized and observed signature values.
5. The process is repeated for all catchments.
6. A parametric joint probability distribution is fitted to all the computed errors. Furthermore, the errors are tested for independence that allows (approximately) factorizing a joint distribution into a product of marginal distributions.

The resultant error distribution defines the likelihood function L in Eq. (2). The main assumption here is that the potentially complex nature of errors in the set of signature values can be usefully represented by the fitted error distributions.

2.2.3 Synthetic Case and Likelihood Functions

To avoid masking the potential value of the regionalized signatures with model structure and observational errors, a “perfect model” is first employed. This involves using the pre-selected rainfall-runoff model and the observed forcing data to generate the “observed” catchment signatures. The Nash-Sutcliffe criteria (NSE) (Nash and Sutcliffe, 1970) optimal parameter set is taken to generate a “perfect model” streamflow time series for each catchment. To produce regionalized signature analogues in this case, two types of imposed errors are introduced to these “observed” signatures. The first error type is characterized by a range of standard deviations (1, 5, 10 and 20 % of the signature value range observed over all catchments used in this study) and a range of inter-signature error correlations (Pearson correlation coefficients equal to 0, 0.25, 0.50, 0.75 and 0.90). This allows the sensitivity of the results to the regionalization quality and the regionalization errors’ correlations to be evaluated. The second error type is set to be equal to the observation-based likelihood function (Sect. 2.2.2). These error structures are the likelihoods used in Eq. (2) for the synthetic case when flows are generated by a “perfect model”.

2.3 Case Study and Rainfall-Runoff Model

2.3.1 Study Catchments

A set of 84 medium sized United States catchments (242 to 8657 km²) from the MOPEX database (Schaake et al., 2006; Duan et al., 2006), for which a variety of regional response signature models have been determined in Almeida et al. (2012), namely runoff ratio, base base flow index, streamflow elasticity, slope of slow duration curve and high pulse count, are used to test the method proposed in this paper. It has proven difficult to derive regionalization equations of acceptable prediction quality for all catchments in the MOPEX dataset (Almeida, 2014). This is due to the lack of descriptive power in the set of available catchment attributes, e.g. the attributes do not provide satisfactory information about catchment geology. To isolate the effect of variable geology on the regression

equations, the selected 84 catchments are grouped based on the underlying geology, namely, Middle
190 Paleozoic sedimentary rocks. Use of more catchments from the MOPEX database would require
different regionalization equations due to changing process controls, and would be unnecessary given
that the focus of the study is on signature error correlations given a regionalization model. For more
details on the motivation for choosing these specific 84 catchments the reader is referred to Almeida
et al. (2012) and Almeida (2014).

195 The 84 catchments are hydrologically varied with a selection of properties summarized in Table 1.
Daily time series for the period from 1 October 1949 to 30 September 1959 are employed. As
highlighted in Almeida et al. (2012), these 10 years of data, representing only a subset of all the data
available, are assumed to be of sufficient length to capture climatic variability, but short enough to
avoid effects of long-term climatic trends (Sawicz et al., 2011).

200 **2.3.2 Response Signatures**

Five response signatures are considered: runoff ratio (RR), base flow index (BFI), streamflow elas-
ticity (SE), slope of flow duration curve (SFDC) and high pulse count (HPC) (Table 1). This specific
subset of signatures is selected to cover a wide range of different qualities of regionalized informa-
tion, and also to ensure that some signature errors are largely uncorrelated, whilst others are strongly
205 correlated (see also Sect. 3.1).

RR reflects the amount of precipitation that becomes streamflow over a certain area and time.
It is determined as the ratio of catchment's outlet streamflow and catchment average precipitation
over the 10 years used in this study. BFI gives the proportion of streamflow that is considered to be
base flow. A simple one-parameter single-pass digital filter method is used to derive BFI (Arnold
210 and Allen, 1999). SE provides a measure of the sensitivity of streamflow to changes in precipitation
(Sankarasubramanian et al., 2001). It is calculated as a median of the inter-annual variation in total
annual streamflow to the inter-annual variation in total annual precipitation ratios normalized by
the long-term runoff ratio (Sawicz et al., 2011; Sankarasubramanian et al., 2001). SFDC gives an
indication of the streamflow variability and is calculated as the slope of the flow duration curve
215 between the 33 and 66 % flow exceedance values in a semi-log scale (Sawicz et al., 2011). HPC
reflects aspects of the high flow regime and catchment flashiness, and is calculated as the average
number of events per year that exceed three times the median daily flow (Clausen and Biggs, 2000;
Yadav et al., 2007).

2.3.3 Rainfall-Runoff Model Choice

220 The probability distributed moisture (PDM) model (Moore, 2007) together with two parallel lin-
ear routing stores and a simple snow model (Hock, 2003) is selected with two major motivations
(a detailed description of the model is given in Appendix A). First, this type of model has been
shown to have a suitable complexity for modelling daily rainfall-runoff over a large sample of the

MOPEX catchments (Wagener and McIntyre, 2012). Second, the model has been successfully applied in other regionalization studies across a wide range of climate and physiographic conditions, for example Calver et al. (1999), Lamb and Kay (2004), McIntyre et al. (2005), Young (2006), and De Vleeschouwer and Pauwels (2013). Even though other model structures may be better suited for some specific catchments, it is prohibitively difficult to vary model structure between catchments and no single model structure will ever be best for all catchments (Lidén and Harlin, 2000; Clark et al., 2008; van Werkhoven et al., 2008). Consequently, the selected model structure is believed to be a sufficient choice for the purposes of this paper. Most importantly, the general framework is independent of the rainfall-runoff model choice.

2.4 Posterior Distribution and Performance Assessment

Employing Bayes' law (Eq. (2)), the rainfall-runoff model is conditioned on different combinations of signatures: (1) assuming independence between the signature regionalization errors (setting the correlation values to zero in the joint probability function); and (2) accounting for the inter-signature error correlations (using the estimated covariance in the joint probability function).

Two metrics are used to assess the effectiveness of the parameter conditioning procedure: (1) the Bayes factor (Jeffreys, 1961) to assess convergence of the parameter posteriors to known parameter values; (2) the probabilistic Nash-Sutcliffe efficiency (Bulygina et al., 2009) to assess convergence of the flow ensembles to the observed flows.

The Bayes factor BF is defined as the ratio between two marginal distributions of the data \mathbf{y} (e.g. observed streamflow time series) for two competing hypotheses (H_1 and H_2) (Kass and Raftery, 1995) (more detail is given in Appendix B):

$$\text{BF} = \frac{p(\mathbf{y}|H_1)}{p(\mathbf{y}|H_2)} \quad (3)$$

Thus, to test the impact of representing the error correlations, the hypothesis H_1 corresponds to the inter-signature errors being treated as correlated, while the hypothesis H_2 corresponds to the inter-signature errors assumed to be independent. If the resulting Bayes factor is greater than 1, there is more support for hypothesis H_1 , and the inter-signature error correlation is worth considering.

When using synthetic streamflow data ("perfect model" approach), with the streamflow time series generated by a pre-selected parameter set, $p(\mathbf{y}|H)$ in Eq. (3) can be seen as either the posterior probability of the known observed streamflow time series under hypothesis H or the probability of the known parameter set that generated that particular flow time series under hypothesis H . As in a "perfect model" approach there is no observational error, $p(\mathbf{y}|H)$ is the probability estimated for the known value of the parameter set that generated the observed streamflow under each of the hypotheses H_1 and H_2 . Since there is no known parameter value corresponding to the real data, the application of the Bayes factor is less useful in this situation. In this case, defining \mathbf{y} as an NSE-

optimal parameter set allows an indication of the relative degree of convergence around the chosen point.

260 The probabilistic Nash-Sutcliffe efficiency NSE_{prob} (Bulygina et al., 2009) is a probabilistic analogue of the traditional Nash-Sutcliffe efficiency coefficient (Nash and Sutcliffe, 1970), and allows both prediction accuracy and precision to be summarized by a single statistic (Eq. (4)).

$$\text{NSE}_{\text{prob}} = \left\{ 1 - \frac{\sum_{t=1}^T (E[\hat{q}_t] - q_t)^2}{\sum_{t=1}^T (q_t - E[q])^2} \right\} - \frac{\sum_{t=1}^T \text{Var}[\hat{q}_t]}{\sum_{t=1}^T (q_t - E[q])^2} \quad (4)$$

q_t denotes a set of streamflow observations for time $t = 1, \dots, T$, $E[q]$ is the average value for the
 265 q_t time series, \hat{q}_t is the simulated time series of streamflow for time $t = 1, \dots, T$, $\text{Var}[\hat{q}_t]$ is the prediction variance at time t , $E[\hat{q}_t]$ is the mathematical expectation of the predictions at time t , and T is the total number of time steps in the sequence. The first part of Eq. (4) corresponds to the traditional Nash-Sutcliffe efficiency coefficient (Nash and Sutcliffe, 1970) in which expected streamflow values are considered as predictors. The latter part of the equation represents the variance, whereby higher
 270 predictor variance corresponds to less precise predictions (Bulygina et al., 2009). An NSE_{prob} of 1 indicates a perfect fit, i.e. the results are both accurate and precise. The incremental improvement in the NSE_{prob} can be used to measure the value of adding signatures into the conditioning or otherwise changing the likelihood function.

For model validation, we use a jack-knife approach (or leave-one-out strategy), commonly em-
 275 ployed in regionalization studies (e.g. Merz and Blöschl, 2004; Shu and Ouarda, 2012). One catchment at a time is removed as a test “ungauged” catchment and the remaining gauged catchments are used to support the regionalization process, including Steps 2–6 listed in Sect. 2.2.2 The procedure is repeated for each of the available catchments.

3 Results and Discussion

280 3.1 Regionalized Signature Errors and Likelihood Functions

The regionalization error probability distributions (that define the likelihoods) are generated follow-
 ing Steps 2 to 6 in Sect. 2.2.2 and are shown in Fig. 1. The marginal error distributions, shown on the
 Fig. 1 diagonal, are approximated using histograms, and parameters of normal distributions are fitted
 using the method of moments. The univariate Kolmogorov-Smirnov test shows that the marginal dis-
 285 tribution normality cannot be rejected at the 95 % confidence level for each of the five signatures. The
 off-diagonal shows the regionalization errors for different signature pairs (lower off-diagonal), the
 corresponding correlation coefficient values and their statistical significance (upper off-diagonal).
 The joint error distributions are approximated using multivariate normal distributions that are fitted
 using estimates of the marginal normal distribution parameters and the inter-signature error corre-
 290 lations. These marginal and joint distributions define the likelihood functions in Eq. (2). Note that
 Fig. 1 represents the regionalization errors based on all 84 catchments. Meanwhile, the jack-knife

procedure (see Sect. 2.4) utilized in the performance assessment employs only 83 catchments at a time.

3.2 The Impact of Inter-Signature Error Correlations (Pairs of Signatures)

295 This section considers the role of inter-signature error correlation on model parameter estimation when pairs of signatures are used. First, different imposed error variances and correlations together with synthetic streamflow data are employed to test the impact of inter-signature error correlation without the impact of model structural error. Then, the results obtained using the observation-based error structure, for both synthetic and observed data streamflow, are analyzed.

300 3.2.1 Synthetic Streamflow Data (Imposed Likelihoods)

Synthetic streamflow data are generated as described in Sect. 2.2.3, and the imposed likelihood functions are defined as described in Sect. 2.2.3. The imposed likelihoods are considered to have standard deviations equal to 1, 5, 10, 20 % of the signature value range observed over all catchments. A comparison of the imposed error structures under the different levels of variance and the observed error structure is given in Table 2. Furthermore, different inter-signature error correlations are also tested, namely 0 (linear independence), 0.25, 0.50, 0.75 and 0.90.

Ten possible pairs of the five response signatures are used in parameter conditioning, and the median Bayes factor, calculated over the 84 MOPEX catchments, is calculated for each pair. The Bayes factor (Eq. 3) compares the two following hypothesis: H_1) the inter-signature error correlation is to be taken into account, and H_2) the errors between the different sources of information can be assumed independent. The Bayes factor is found to be relatively insensitive to the selection of response signature pairs (Kruskal–Wallis test). Table 3 summarizes the 95 % pooled confidence intervals for the median Bayes factor across all catchments and across all 10 signature pairs, for each choice of the likelihood (i.e. 20 likelihoods). This provides reference values indicative of the error interdependency importance in model regionalization depending on the signature pair correlations and marginal distribution variances. As it would be expected, the median Bayes factor is equal to 1 when signatures errors are not correlated (i.e. $\rho = 0$). However, as correlations between signatures errors increase the median Bayes factor increases noticeably. This suggests that considering error correlations allocates higher likelihoods to parameter sets that capture a considered signature pair. Furthermore, the results shown in Table 3 also imply that the median Bayes factor is relatively insensitive to the precision with which the signatures are regionalized.

3.2.2 Synthetic and Observed Streamflow Data (Observation-Based Likelihoods)

Figure 2 shows the distribution of the Bayes factor values obtained across the 84 catchments for each of the 10 possible different pairs of signatures, when the observation-based error structure is used for each catchment. Figure 2a shows the results for the observed streamflow data with regionalized

signatures calculated from the derived regressions; Fig. 2b shows the results for the synthetic streamflow data with regionalized signatures calculated by adding noise to the exact signature values. The Tukey boxplots in red correspond to pairs of signatures whose errors are statistically significantly correlated (see Fig. 1). The upper whisker represents the upper quartile plus one and a half times the interquartile range, and the lower whisker represents the lower quartile minus one and a half times the interquartile range. The matrix below Fig. 2b shows the pairs of signatures used.

The signature pair [SFDC, HPC] shows the strongest correlation between errors ($\rho = 0.65$, Fig. 1). A likelihood function with a standard deviation equal to 10 % of the observed signature ranges and $\rho = 0.75$ in Table 3 is comparable to the observation-based likelihood of the pair [SFDC, HPC] (Table 2), with Table 3 indicating [1.45, 1.53] as a 95 % confidence interval for the median Bayes factor. However, a median Bayes factor of 2.17 is obtained for the observed streamflow data (Fig. 2a). Similar differences are found for the other pairs of signatures, although the comparison with the reference table (Table 3) becomes challenging, as the individual signatures have not been regionalized necessarily with similar quality. On the other hand, Fig. 2b shows that the Bayes factors for the synthetic study (when there is no model structural error) are consistent with the values provided in the look-up Table 3. The difference between the median Bayes factor for the two cases is likely to be caused by the model structure error, or may be related to the location of the NSE-optimal in the parameter space.

Nevertheless, it is clear from Fig. 2 that those pairs of signatures whose errors are significantly correlated (i.e. [SFDC, HPC], [BFI, HPC], [BFI, SFDC] and [BFI, SE]) have wider interquartile ranges. Furthermore, the pair of signatures with the strongest correlation between errors [SFDC, HPC] presents the greatest interquartile range. Therefore the inclusion of significant correlations in the likelihood function matters, but whether or not it is beneficial to conditioning the parameters seems to depend on the interplay between model structure error, parameter space and likelihood function. Only strong correlations (as in the [SFDC, HPC] case) can be expected to result in a median Bayes factor clearly above 1.

3.3 The Impact of Inter-Signature Error Correlations (Multiple Signatures)

Multiple signatures are used for parameter constraining and flow prediction. The information value of multiple signatures and its dependence on inter-signature error correlations is explored in this section.

3.3.1 Synthetic Streamflow Data (Observation-Based Likelihood)

Figure 3 shows Bayes factors derived for the synthetic streamflow data (generated using the NSE-optimal parameter set) when the observation-based likelihood is used. The Bayes factor considers $p(\cdot|H_2)$ to be the prior parameter distribution, and $p(\cdot|H_1)$ to be one of the parameter posteriors that includes or ignores the inter-signature error correlations. Figure 3 summarizes the variability

in the Bayes factor for the different combinations of signatures for all 84 catchments. Boxplots are color coded by the total number of signatures combined, when the inter-signatures error correlation is considered in the likelihood function definition. The grey dashed boxplots correspond to the results obtained assuming that the inter-signature errors are independent when defining the likelihood function. Although the colored boxplots visually seem to have higher values than the grey dashed boxplots, these differences are not statistically significant at a 95 % confidence level (Kolmogorov–Smirnov two-sided tests).

To better evaluate whether the incorporation of additional sources of information improves parameter identification, one-sided Kolmogorov–Smirnov tests are applied between any combination of certain signatures (e.g. [SE, SFDC]) and any other combination that contains the same signatures and a new one (e.g. [SE, SFDC, HPC]). It is found that adding more signatures improves parameter identification in 82.5 % of the cases (66 out of 80 cases) at a 95 % confidence level).

Figure 4 summarizes the variability in the analog Nash-Sutcliffe efficiency measure NSEprob for different combinations of signatures for all 84 catchments. The colored boxplots correspond to the results obtained when the inter-signature error correlations are considered in the likelihood definition, and the grey dashed boxplots correspond to the results when the inter-signature errors are assumed to be independent. There is no visual or statistical (two-sided Kolmogorov–Smirnov tests) difference between the colored boxplots and the grey dashed boxplots in Fig. 4. Moreover, visually, adding more response signatures seems to improve streamflow predictions in terms of accuracy and precision when no model structure error exists. However, only in 59 % of the cases (47 out of 80 cases) more signatures contribute to improved streamflow predictions at a 95 % confidence level (one-sided Kolmogorov–Smirnov test). The other 33 cases always involve the inclusion of the most poorly regionalized signatures (with the highest variance from the five regionalized signatures) - SE, SFDC or HPC - as additional sources of information (see Table 2).

It is worth noting that very similar results (not shown here) are obtained when instead of regionalized signatures (~~calculated by adding noise to the exact signature value~~), “observed” signatures (~~the exact signature value~~) are used are used but with the same error derived from regionalization. This suggests that the uncertainty around the regionalized signatures values, as well as signature information content, are the key factors leading to the results shown in Fig. 4.

3.3.2 Observed Streamflow Data (Observation-Based Likelihood)

Figure 5 shows the results when the same methodology as in the Sect. 3.3.1 is applied using the observed streamflow data. As in the synthetic streamflow case, the differences between the Bayes factor distributions when inter-signature error correlations are considered and when inter-signature errors are assumed to be independent are not statistically significant at a 95 % confidence level (Kolmogorov–Smirnov two-sided tests).

Further, by comparing Fig. 5 with Fig. 3, it becomes clear that the signatures contribute less information, and there is a smaller increase in performance as more signatures are added. It is found that adding more signatures tends to improve parameter identification only in half of the cases when compared to the synthetic streamflow case at a 95 % confidence level (42.5 % versus 82.5 % in the synthetic streamflow case). Furthermore, and contrastingly to the case where no structural error exists, in five situations adding more signatures contributes to a decrease in performance. These five cases always involve adding either SFDC or HPC as an additional source of information. This performance deterioration can be attributed to model structure and observational error. Overall, a statistically significant drop in performance with regard to the Bayes factor is observed most of the time when model structural error is present.

Figure 6 presents the results in terms of NSEprob using the observed streamflow data. As in the synthetic study in Sect. 3.3.1, there is no statistically significant difference at a 95 % confidence difference between the NSEprob distributions when the inter-signature error correlations are considered and when the errors are treated independently (Kolmogorov–Smirnov two-sided tests).

Figure 6 shows that better results in terms of NSEprob are not necessarily achieved when all five signatures are used simultaneously. It is found that adding more signatures tends to improve parameter identification only in 36 % of the cases at a 95 % confidence level (compared to 59 % when there is no model structure error). Furthermore, and contrasting the case where no model structure error exists, in two situations, adding more signatures may contribute to a decrease in performance (when we start with [RR, BFI] and add HPC, and when we start with [RR, BFI] and add SFDC). This might be due to regionalization biases in SFDC and HPC and/or due to the inability of the PDM model to maintain a satisfactory overall performance when conditioned on high peak flow and medium flow information. This negative impact is not observed when synthetic streamflow data are used (Fig. 4), indicating that the decrease in performance may be due to model structural deficiencies. Moreover, a statistically significant drop in performance with regard to NSEprob is observed most of the time when there is model structural error.

In summary, unless there is no model structural error, an all-round performance improvement is not guaranteed by adding more signatures. Furthermore, model structure uncertainty seems to have a much bigger effect on the performance than the explicit inclusion of the inter-signature error correlations.

3.4 Limitations and Applicability

The main feature of the method suggested in this paper lies in the possibility of allowing a large number of signatures to be added to the conditioning process, without worrying about double-counting of information or degree of uncertainty in signature estimates, and avoiding subjective decisions about removal of possibly nonindependent information. Although the proposed framework can be applied to any number of signatures, the limited sample size (i.e. number of gauged catchments available)

can have an impact on the definition of the likelihood distribution. For this specific study 83 samples were available to define that distribution. When a single response signature is used to condition the hydrological model this sample size is likely to be sufficient to confidently judge whether the normal distribution assumption is sufficient. However, when moving to multidimensional problems, in which various signatures may be used simultaneously to condition the hydrological model, it is increasingly difficult to judge the adequacy of any multivariate parametric distribution and to judge which catchments are outliers. This implies that as more signatures are used simultaneously in the conditioning of the hydrological model, the more gauged catchments should be used to define the likelihood function. As stressed by Gupta et al. (2014), large samples are of great importance to support statistical regionalization of uncertainty estimates, and this is particularly the case if dependencies between information sources are to be specified.

While the work presented in this paper addresses a number of issues associated with model regionalization, it is important to highlight some additional areas for future research. An important source of uncertainty comes from model structure error (Gupta et al., 1998; Kuczera et al., 2006). The conditioning framework suggested here is independent of the selected model, and, in principle, Figs. 5 and 6 could be created by using the model structure that is considered suitable for each catchment rather than using a model structure that we consider good for generalizing. Further research is needed to diagnose the relative importance of different model structures in various climate regimes and for different catchment characteristics (Clark et al., 2008; Hrachowitz et al., 2013). This is crucial to both identifying the most appropriate model structure for an ungauged location and quantifying the uncertainty in the model structure that should be integrated into the likelihood, thus allowing virtually any model choice. Similarly, other sources of uncertainty, namely observational error (e.g. rainfall error), should ideally be evaluated and integrated into the likelihood function. By accounting for all the important sources of uncertainty, further insight should be achieved into the information value of sets of signatures and the value of including their dependencies in the likelihood function.

Some of the results presented may be sensitive to the response signatures used. The relationship between value of signatures and catchment type remains ambiguous and an interesting aspect for posterior evaluation would be how the value of signatures depends on catchment type. Other aspects that are worth further research include whether a similar framework could be applied to different types of information source, e.g. can some discharge measurements be added into the model conditioning process? While Bulygina et al. (2012) suggests a framework capable of combining multiple sources of knowledge, namely physically based information, regionalized signatures and spot observations to identify parameters for models of ungauged catchments, the errors between them were assumed to be independent in their case study. A combination of the framework suggested by Bulygina et al. (2012) and the method proposed in this paper may be the way forward to maximizing the value of the available information within a framework of uncertainty reduction.

4 Conclusions

470 Uncertainty in streamflow estimation in ungauged catchments originates not only from the traditional sources of error generally identified in rainfall-runoff modelling (i.e. model structural, parameter and data errors), but also by errors introduced by the transposition of information from data-rich areas and use of this information to condition model simulations. To identify which and how many types of signatures can usefully be included in model conditioning, it is critical to understand the effects of all
475 these uncertainties. Moreover, when multiple signatures are used simultaneously to condition model simulations, inter-signature error dependencies may also introduce uncertainty and affect decisions about the value of information. While error and uncertainty analyses are quite common in regionalization studies, the question of how much information can be taken from a set of uncertain signatures and determining how many and which signatures should be used given their error dependencies has
480 not been extensively studied.

The method suggested in this paper allows the specification of a signature error structure. A common reason for not including large numbers of signatures in regionalization studies is the potential for under-estimation of uncertainty due to duplication of information. This study helps to justify the inclusion of larger sets of signatures in the regionalization procedure if their error correlations are
485 formally accounted for and thus enables a more complete use of all available information. The results show that adding response signatures to constrain the hydrological model, while accounting for inter-signature error correlations, can contribute to a stronger identification of the optimum parameter set when the error correlations between different sources of information are strong. Furthermore, the results show that assuming independency of errors does not result in significant deterioration in
490 model performance, unless the error correlation is very strong. The results also show that the effect of error correlations is likely to be overwhelmed by model structure and observation errors. The method suggested here can therefore become more relevant if observational and structural errors are reduced. In addition, it is illustrated that using more signatures, with and without considering their error correlations, may lead to deterioration in performance. In our case, there were particular
495 problems when adding the slope of the flow duration curve and/or the high pulse count. As this is likely to be specific to the rainfall-runoff model used, the selected performance criteria and the set of catchments, it is recommended that the disinformative information sources are identified as part of any regionalization study, in a similar manner as has been done here.

Appendix A: Model Structure

500 A schematic representation of the model structure used in this study is shown in Fig. A1. The snowmelt routine is based on the degree-day method. Precipitation accumulates as snow or rain depending whether the air temperature is above or below a threshold temperature (T_{th}). When the air temperature is above the temperature threshold for snowmelt (T_m), snowmelt occurs at a rate

that is proportional to the degree-day factor (DDF). The soil moisture storage component describes
505 the water balance at the soil level. The PDM model uses a probability density function to represent
changes in the catchment storage capacity, defined by the maximum soil moisture storage (c_{max}) -
the maximum soil water storage capacity within the modelled element - and a shape parameter (b)
that controls the degree of spatial variability of storage capacity over the catchment. Interception
is not explicitly modelled. Transpiration and evaporation are lumped into a single term. The actual
510 evapotranspiration (AE) is determined based on a relationship between evapotranspiration and soil
moisture deficit (Moore, 2007). After evapotranspiration, the remaining available water is used to
fill the soil moisture store. When effective rainfall is produced through overflow of the storage ele-
ments, excess water is passed to the routing stores. The routing module channels this water into two
reservoirs, according to a fraction split coefficient (α). A proportion α of the water excess goes to the
515 quick flow reservoir, controlled by the quick flow residence time (k_q), and $(1 - \alpha)$ of the water excess
goes to the slow flow reservoir, controlled by the slow flow residence time (k_s). The streamflow at
the catchment outlet is the sum of the outputs from each of these quick and slow flow reservoirs.

The parameter ranges (Table A1) are selected after Kollat et al. (2012) based largely on the maxi-
mum range sampled from several recent studies, such that only sufficiently extreme values are ruled
520 out.

Appendix B: The Bayes Factor

When evaluating the impact of inter-signature error correlations on model parameter identification,
results are assessed in terms of Bayes factor (Jeffreys, 1961). This form of assessment is preferred
to the most commonly used QQ plots (Laio and Tamea, 2007), due to the particular nature of the
525 problem under analysis. When signature(s) (either regionalized for the case of an ungauged catch-
ment, or derived from actual observations for the case of gauged catchments) is employed to reduce
uncertainty beyond what is possible by defining the priors on model parameters, QQ plots may not
be the most effective form of assessment. Although response signatures are measures of theoretically
relevant system process behaviors (Gupta et al., 2008; Wagener et al., 2007), they reflect fragmented
530 knowledge as different signatures capture different catchment processes. Consequently, the quan-
tiles of observed flows are not conditioned to follow a uniform distribution, as QQ plots assess.
Rather, quantiles of response signatures should follow this condition (for all catchments considered
– Almeida et al., 2013). Therefore, an alternative performance measure that more adequately reflects
the aim of this particular application (i.e. the reproduction of certain aspects of the hydrograph) is
535 used. The Bayes factor BF is particularly relevant in the current context as it allows comparison of
predictions based on two competing theories (Jeffreys, 1961). It is defined as the ratio between the
marginal distributions of the data y for the two hypotheses (H_1 and H_2) being compared (Kass and

Raftery, 1995):

$$\text{BF} = \frac{p(\mathbf{y}|H_1)}{p(\mathbf{y}|H_2)} \quad (\text{B1})$$

540 When the two hypotheses are equally likely a priori, the Bayes factor is the posterior odds in favor of H_1 (Kass and Raftery, 1995). In other words, a value of BF greater than 1 means that H_1 is more strongly supported by the data than H_2 . For example, a Bayes factor equal to 2 implies that H_1 is favored over H_2 with 2 : 1 odds given the evidence provided by the data.

For a given hypothesis H , parameterized by model parameter set Θ , the marginal density $p(\mathbf{y}|H)$
545 represents the likelihood of the data and it is given by

$$p(\mathbf{y}|H) = \int p(\mathbf{y}|\Theta, H)p(\Theta|H)d\Theta \quad (\text{B2})$$

where $p(\mathbf{y}|\Theta, H)$ is the conditional density function given parameters Θ under hypothesis H and $p(\Theta|H)$ is the distribution of parameters under H . Hypothesis H may represent different model and parameter distributions. In this paper, the same model structure is considered. However, different
550 parameter distributions are used in Eq. (B2) to enable prediction comparison based on two theories about parameter distributions.

The above integral can be numerically approximated as,

$$\int p(\mathbf{y}|\Theta, H)p(\Theta|H)d\Theta \approx \frac{1}{N} \sum_{i=1}^N p(\mathbf{y}|\Theta^{(i)}, H)p(\Theta^{(i)}|H) \quad (\text{B3})$$

where $\Theta^{(i)}$ is the i th of N draws from $p(\cdot|\Theta)$, and N is the size of the Monte Carlo sample (in this
555 paper N is equal to 10 000).

In a ‘‘perfect model’’ study, data \mathbf{y} are generated by a model with parameter set Θ^* , so that there is no model structural or observational error. This means that $p(\mathbf{y}|\Theta^{(i)}, H)$ is always equal to zero, except when $\Theta^{(i)} = \Theta^*$. Mathematically this is expressed as $p(\mathbf{y}|\Theta^{(i)}, H) = \delta_{\Theta^{(i)} = \Theta^*}$, where δ is the Dirac delta function. Therefore Eq. (B3) is equal to $1/N$ times $p(\Theta^{(i)} = \Theta^*|H)$ and the Bayes
560 factor is given by

$$\text{BF} = \frac{\frac{1}{N} \sum_{i=1}^N \delta_{\Theta^{(i)} = \Theta^*} p(\Theta^{(i)}|H_1)}{\frac{1}{N} \sum_{i=1}^N \delta_{\Theta^{(i)} = \Theta^*} p(\Theta^{(i)}|H_2)} = \frac{p(\Theta^{(i)} = \Theta^*|H_1)}{p(\Theta^{(i)} = \Theta^*|H_2)} \quad (\text{B4})$$

While other choices can be made, two cases are considered in this paper. First, the two distributions in Eq. (B4) are posterior distributions, but with different assumptions about the likelihood functions. Given that we are particularly interested in evaluating the impact of considering the inter-
565 signature error correlations versus ignoring them, H_1 will correspond to the joint likelihood defined such that inter-signature error correlations are considered, while H_2 corresponds to the likelihood when inter-signature error correlations are ignored. For the Bayes factor defined in this way, a value greater than 1 supports the idea that considering inter-signature error correlations contributes to an improved specification of the optimum parameter set. In this paper we are also interested in the value

570 of adding/not adding more signatures in model conditioning, and so the Bayes factor will be also
calculated for $p(\cdot|H_2)$ set to be the prior parameter distribution, and $p(\cdot|H_1)$ set to one of the derived
parameter posteriors. For the Bayes factor defined in this way, a value greater than 1 supports the
idea that additional sources of information contribute to a stronger identification of the optimum
parameter set.

575 *Acknowledgements.* The authors would like to acknowledge the support of Fundação para a Ciência e a Tec-
nologia (FCT), Portugal, sponsor of the PhD program of S. Almeida at Imperial College London, under the
grant SFRH/BD/65522/2009. This work was also partially supported by the Natural Environment Research
Council [Consortium on Risk in the Environment: Diagnostics, Integration, Benchmarking, Learning and Elic-
itation (CREDIBLE); grant number NE/J017450/1]. The authors would like to thank Keith Sawicz for advice
580 and support relating to the data used in this study. The authors also thank the editor Vazken Andréassian, who
handled the manuscript, and the three anonymous reviewers for their useful comments.

References

- Almeida, S. M. C. L.: The Value of Regionalised Information for Hydrological Modelling, PhD thesis, Imperial College London, London, UK, 2014.
- 585 Almeida, S., Bulygina, N., McIntyre, N., Wagener, T., and Buytaert, W.: Predicting flows in ungauged catchments using correlated information sources, in: British Hydrological Society's Eleventh National Hydrology Symposium, Hydrology for a Changing World, Dundee, UK, doi:10.7558/bhs.2012.ns02, 2012.
- Almeida, S., Bulygina, N., McIntyre, N., Wagener, T., and Buytaert, W.: Improving parameter priors for data-scarce estimation problems, *Water Resour. Res.*, 49, 6090–6095, doi:10.1002/wrcr.20437, 2013.
- 590 Arnold, J. G. and Allen, P. M.: Automated methods for estimating baseflow and ground water recharge from streamflow records, *J. Am. Water Resour. As.*, 35, 411–424, doi:10.1111/j.1752-1688.1999.tb03599.x, 1999.
- Boorman, D. B., Hollis, J. M., and Lilly, A.: Hydrology of soil types: a hydrologically-based classification of the soils of the United Kingdom, Tech. rep., Institute of Hydrology, Wallingford, UK, 1995.
- Bulygina, N., McIntyre, N., and Wheater, H.: Conditioning rainfall-runoff model parameters for ungauged catchments and land management impacts analysis, *Hydrol. Earth Syst. Sci.*, 13, 893–904, doi:10.5194/hess-13-893-2009, 2009.
- 595 Bulygina, N., McIntyre, N., and Wheater, H.: Bayesian conditioning of a rainfall-runoff model for predicting flows in ungauged catchments and under land use changes, *Water Resour. Res.*, 47, W02503, doi:10.1029/2010wr009240, 2011.
- 600 Bulygina, N., Ballard, C., McIntyre, N., O'Donnell, G., and Wheater, H.: Integrating different types of information into hydrological model parameter estimation: application to ungauged catchments and land use scenario analysis, *Water Resour. Res.*, 48, W06519, doi:10.1029/2011wr011207, 2012.
- Calver, A., Lamb, R., and Morris, S. E.: River flood frequency estimation using continuous runoff modelling, *P. I. Civil Eng.-Water*, 136, 225–234, doi:10.1680/iwtme.1999.31986, 1999.
- 605 Castiglioni, S., Lombardi, L., Toth, E., Castellarin, A., and Montanari, A.: Calibration of rainfall-runoff models in ungauged basins: a regional maximum likelihood approach, *Adv. Water Resour.*, 33, 1235–1242, doi:10.1016/j.advwatres.2010.04.009, 2010.
- Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T., and Hay, L. E.: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose differences between hydrological models, *Water Resour. Res.*, 44, W00B02, doi:10.1029/2007wr006735, 2008.
- 610 Clausen, B. and Biggs, B. J. F.: Flow variables for ecological studies in temperate streams: groupings based on covariance, *J. Hydrol.*, 237, 184–197, doi:10.1016/S0022-1694(00)00306-1, 2000.
- De Vleeschouwer, N. and Pauwels, V. R. N.: Assessment of the indirect calibration of a rainfall-runoff model for ungauged catchments in Flanders, *Hydrol. Earth Syst. Sci.*, 17, 2001–2016, doi:10.5194/hess-17-2001-2013, 2013.
- 615 Duan, Q., Schaake, J., Andréassian, V., Franks, S., Goteti, G., Gupta, H., Gusev, Y., Habets, F., Hall, A., Hay, L., Hogue, T., Huang, M., Leavesley, G., Liang, X., Nasonova, O., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T., and Wood, E.: Model Parameter Estimation Experiment (MOPEX): An overview of science strategy and major results from the second and third workshops, *J. Hydrol.*, 320, 3–17, doi:10.1016/j.jhydrol.2005.07.031, 2006.
- 620

- Gupta, H. V., Sorooshian, S., and Yapo, P. O.: Toward improved calibration of hydrologic models: Multiple and noncommensurable measures of information, *Water Resour. Res.*, 34, 751–763, doi:10.1029/97WR03495, 1998.
- 625 Gupta, H. V., Wagener, T., and Liu, Y.: Reconciling theory with observations: elements of a diagnostic approach to model evaluation, *Hydrol. Process.*, 22, 3802–3813, doi:10.1002/hyp.6989, 2008.
- Gupta, H. V., Perrin, C., Blöschl, G., Montanari, A., Kumar, R., Clark, M., and Andréassian, V.: Large-sample hydrology: a need to balance depth with breadth, *Hydrol. Earth Syst. Sci.*, 18, 463–477, doi:10.5194/hess-18-463-2014, 2014.
- 630 He, Y., Bárdossy, A., and Zehe, E.: A review of regionalisation for continuous streamflow simulation, *Hydrol. Earth Syst. Sci.*, 15, 3539–3553, doi:10.5194/hess-15-3539-2011, 2011.
- Hock, R.: Temperature index melt modelling in mountain areas, *J. Hydrol.*, 282, 104–115, doi:10.1016/S0022-1694(03)00257-9, 2003.
- Hrachowitz, M., Savenije, H., Blöschl, G., McDonnell, J., Sivapalan, M., Pomeroy, J., Arheimer, B., Blume, T., Clark, M., Ehret, U., Fenicia, F., Freer, J., Gelfan, A., Gupta, H., Hughes, D., Hut, R., Montanari, A., 635 Pande, S., Tetzlaff, D., Troch, P., Uhlenbrook, S., Wagener, T., Winsemius, H., Woods, R., Zehe, E., and Cudennec, C.: A decade of Predictions in Ungauged Basins (PUB) – a review, *Hydrol. Sci. J.*, 58, 1198–1255, doi:10.1080/02626667.2013.803183, 2013.
- Jeffreys, H.: *Theory of Probability*, Oxford University Press, Oxford, 1961.
- Kass, R. E. and Raftery, A. E.: Bayes Factors, *J. Am. Stat. Assoc.*, 90, 773–795, doi:10.2307/2291091, 1995.
- 640 Kennard, M. J., Mackay, S. J., Pusey, B. J., Olden, J. D., and Marsh, N.: Quantifying uncertainty in estimation of hydrologic metrics for ecohydrological studies, *River Res. Appl.*, 26, 137–156, doi:10.1002/rra.1249, 2010.
- Kollat, J. B., Reed, P. M., and Wagener, T.: When are multiobjective calibration trade-offs in hydrologic models meaningful?, *Water Resour. Res.*, 48, W03520, doi:10.1029/2011wr011534, 2012.
- 645 Kuczera, G., Kavetski, D., Franks, S., and Thyer, M.: Towards a Bayesian total error analysis of conceptual rainfall-runoff models: Characterising model error using storm-dependent parameters, *J. Hydrol.*, 331, 161–177, doi:10.1016/j.jhydrol.2006.05.010, 2006.
- Laio, F. and Tamea, S.: Verification tools for probabilistic forecasts of continuous hydrological variables, *Hydrol. Earth Syst. Sci.*, 11, 1267–1277, doi:10.5194/hess-11-1267-2007, 2007.
- 650 Lamb, R. and Kay, A. L.: Confidence intervals for a spatially generalized, continuous simulation flood frequency model for Great Britain, *Water Resour. Res.*, 40, W07501, doi:10.1029/2003WR002428, 2004.
- Lidén, R. and Harlin, J.: Analysis of conceptual rainfall-runoff modelling performance in different climates, *J. Hydrol.*, 238, 231–247, doi:10.1016/S0022-1694(00)00330-9, 2000.
- McIntyre, N., Lee, H., Wheater, H., Young, A., and Wagener, T.: Ensemble predictions of runoff in ungauged catchments, *Water Resour. Res.*, 41, doi:10.1029/2005WR004289, 2005.
- 655 Merz, R. and Blöschl, G.: Regionalisation of catchment model parameters, *J. Hydrol.*, 287, 95–123, doi:10.1016/j.jhydrol.2003.09.028, 2004.
- Mishra, A. K. and Coulibaly, P.: Developments in hydrometric network design: A review, *Rev. Geophys.*, 47, RG2001, doi:10.1029/2007RG000243, 2009.
- 660 Moore, R. J.: The PDM rainfall-runoff model, *Hydrol. Earth Syst. Sci.*, 11, 483–499, doi:10.5194/hess-11-483-2007, 2007.

- Nash, J. E. and Sutcliffe, J. V.: River flow forecasting through conceptual models. Part I. A discussion of principles, *J. Hydrol.*, 10, 282–290, doi:10.1016/j.bbr.2011.03.031, 1970.
- Parajka, J., Andréassian, V., Archfield, S. A., Bárdossy, A., Blöschl, G., Chiew, F., Duan, Q., Gelfan, A., Hlavcova, K., Merz, R., McIntyre, N., Oudin, L., Perrin, C., Rogger, M., Salinas, J. L., Savenije, H. G., Skøien, J. O., Wagener, T., Zehe, E., and Zhang, Y.: Prediction of runoff hydrographs in ungauged basins, in: *Runoff Prediction in Ungauged Basins: Synthesis across Processes, Places and Scales*, edited by: Blöschl, G., Sivapalan, M., Wagener, T., Viglione, A., and Savenije, H., 53–69, Cambridge University Press, Cambridge, 2013.
- 665 Peel, M. C. and Blöschl, G.: Hydrological modelling in a changing world, *Prog. Phys. Geog.*, 35, 249–261, doi:10.1177/0309133311402550, 2011.
- 670 Razavi, T. and Coulibaly, P.: Streamflow Prediction in Ungauged Basins: Review of Regionalization Methods, *J. Hydrol. Eng.*, 18, 958–975, doi:10.1061/(ASCE)HE.1943-5584.0000690, 2013.
- Sankarasubramanian, A., Vogel, R. M., and Limbrunner, J. F.: Climate elasticity of streamflow in the United States, *Water Resour. Res.*, 37, 1771–1781, doi:10.1029/2000WR900330, 2001.
- 675 Sawicz, K., Wagener, T., Sivapalan, M., Troch, P. A., and Carrillo, G.: Catchment classification: empirical analysis of hydrologic similarity based on catchment function in the eastern USA, *Hydrol. Earth Syst. Sci.*, 15, 2895–2911, doi:10.5194/hess-15-2895-2011, 2011.
- Schaake, J., Cong, S., and Duan, Q.: U.S. MOPEX data set, Tech. Rep. UCRL-JRNL-221228, Lawrence Livermore National Laboratory, 2006.
- 680 Shu, C. and Ouarda, T. B. M. J.: Improved methods for daily streamflow estimates at ungauged sites, *Water Resour. Res.*, 48, W02523, doi:10.1029/2011WR011501, 2012.
- Singh, R., Wagener, T., van Werkhoven, K., Mann, M. E., and Crane, R.: A trading-space-for-time approach to probabilistic continuous streamflow predictions in a changing climate – accounting for changing watershed behavior, *Hydrol. Earth Syst. Sci.*, 15, 3591–3603, doi:10.5194/hess-15-3591-2011, 2011.
- 685 Sorooshian, S. and Dracup, J. A.: Stochastic parameter estimation procedures for hydrologic rainfall-runoff models: Correlated and heteroscedastic error cases, *Water Resour. Res.*, 16, 430–442, doi:10.1029/WR016i002p00430, 1980.
- United States Department of Agriculture (USDA): Urban hydrology for small watersheds, Technical Release 55, United States Department of Agriculture, Washington, D.C., 1986.
- 690 van Werkhoven, K., Wagener, T., Reed, P., and Tang, Y.: Characterization of watershed model behavior across a hydroclimatic gradient, *Water Resour. Res.*, 44, W01429, doi:10.1029/2007WR006271, 2008.
- Wagener, T. and McIntyre, N.: Hydrological catchment classification using a data-based mechanistic strategy, in: *System Identification, Environmental Modelling, and Control System Design*, edited by: Wang, L. and Garnier, H., 483–500, Springer, London, 2012.
- 695 Wagener, T. and Montanari, A.: Convergence of approaches toward reducing uncertainty in predictions in ungauged basins, *Water Resour. Res.*, 47, W06301, doi:10.1029/2010WR009469, 2011.
- Wagener, T., Sivapalan, M., Troch, P., and Woods, R.: Catchment Classification and Hydrologic Similarity, *Geogr. Compass*, 1, 901–931, doi:10.1111/j.1749-8198.2007.00039.x, 2007.

- Winsemius, H. C., Schaefli, B., Montanari, A., and Savenije, H. H. G.: On the calibration of hydrological models
700 in ungauged basins: A framework for integrating hard and soft hydrological information, *Water Resour. Res.*,
45, W12422, doi:10.1029/2009WR007706, 2009.
- Yadav, M., Wagener, T., and Gupta, H.: Regionalization of constraints on expected watershed response behavior
for improved predictions in ungauged basins, *Adv. Water Resour.*, 30, 1756–1774, doi:10.5194/hess-15-
3539-2011, 2007.
- 705 Young, A. R.: Stream flow simulation within UK ungauged catchments using a daily rainfall-runoff model, *J.*
Hydrol., 320, 155–172, doi:http://dx.doi.org/10.1016/j.jhydrol.2005.07.017, 2006.
- Zhang, Z., Wagener, T., Reed, P., and Bhushan, R.: Reducing uncertainty in predictions in ungauged basins
by combining hydrologic indices regionalization and multiobjective optimization, *Water Resour. Res.*, 44,
doi:10.1029/2008WR006833, 2008.

Table 1. Summary of general catchment properties and response signatures of the 84 MOPEX catchments.

Catchment property	Units	Range
Average annual streamflow	(mm yr ⁻¹)	208–896
Average annual precipitation	(mm yr ⁻¹)	758–1495
Average annual maximum temperature	(°C)	12–23
Average annual minimum temperature	(°C)	0–10
Average annual potential evaporation	(mm yr ⁻¹)	679–1112
Aridity index*	(-)	0.5–1.2
Average elevation	(m)	176–1056
Runoff ratio	(-)	0.16–0.76
Base flow index	(-)	0.36–0.90
Streamflow elasticity	(-)	0.02–4.34
Slope of flow duration curve	(-)	0.01–0.08
High pulse count	(yr ⁻¹)	2.10–120.80

* Long-term ratio of potential evaporation over precipitation.

Table 2. Tested variance values for the data-based and imposed error structures.

	Observed error structure	1 % observed	5 % observed	10 % observed	20 % observed
		signature ranges	signature ranges	signature ranges	signature ranges
RR residuals	0.054 ²	0.005 ²	0.027 ²	0.055 ²	0.109 ²
BFI residuals	0.044 ²	0.006 ²	0.030 ²	0.060 ²	0.121 ²
SE residuals	0.635 ²	0.023 ²	0.116 ²	0.232 ²	0.464 ²
SFDC residuals	0.006 ²	0.0005 ²	0.002 ²	0.005 ²	0.010 ²
HPC residuals	10.687 ²	0.977 ²	4.883 ²	9.767 ²	19.533 ²

Table 3. Reference table showing the 95% confidence interval for the median Bayes factor. The correlation coefficient ρ and the standard deviation of the marginal distributions σ are shown.

		σ			
		1 %	5 %	10 %	20 %
	0	1	1	1	1
	0.25	1.01–1.03	1.03–1.04	1.02–1.04	1.04–1.05
ρ	0.50	1.09–1.15	1.16–1.19	1.14–1.17	1.14–1.18
	0.75	1.41–1.51	1.50–1.57	1.45–1.53	1.40–1.49
	0.90	1.94–2.11	2.11–2.32	2.12–2.26	2.20–2.34

Table A1. Conceptual model prior parameter ranges.

Parameter	Description	Units	Range
DDF	Degree day factor	(mm day ⁻¹ °C ⁻¹)	0–20
T_m	Base temperature for melting	(°C)	0–5
T_{th}	Threshold temperature for snow formation	(°C)	-5–5
c_{max}	Maximum storage capacity within the catchment	(mm)	0–2000
b	Shape Pareto distribution	(-)	0–4
b_e	Evaporation reduction parameter	(-)	0–4
k_q	Time constant for fast routing store	(days)	0–7
k_s	Time constant for slow routing store	(days)	7–20000
α	Fraction of slow through fast routing store	(-)	0–1

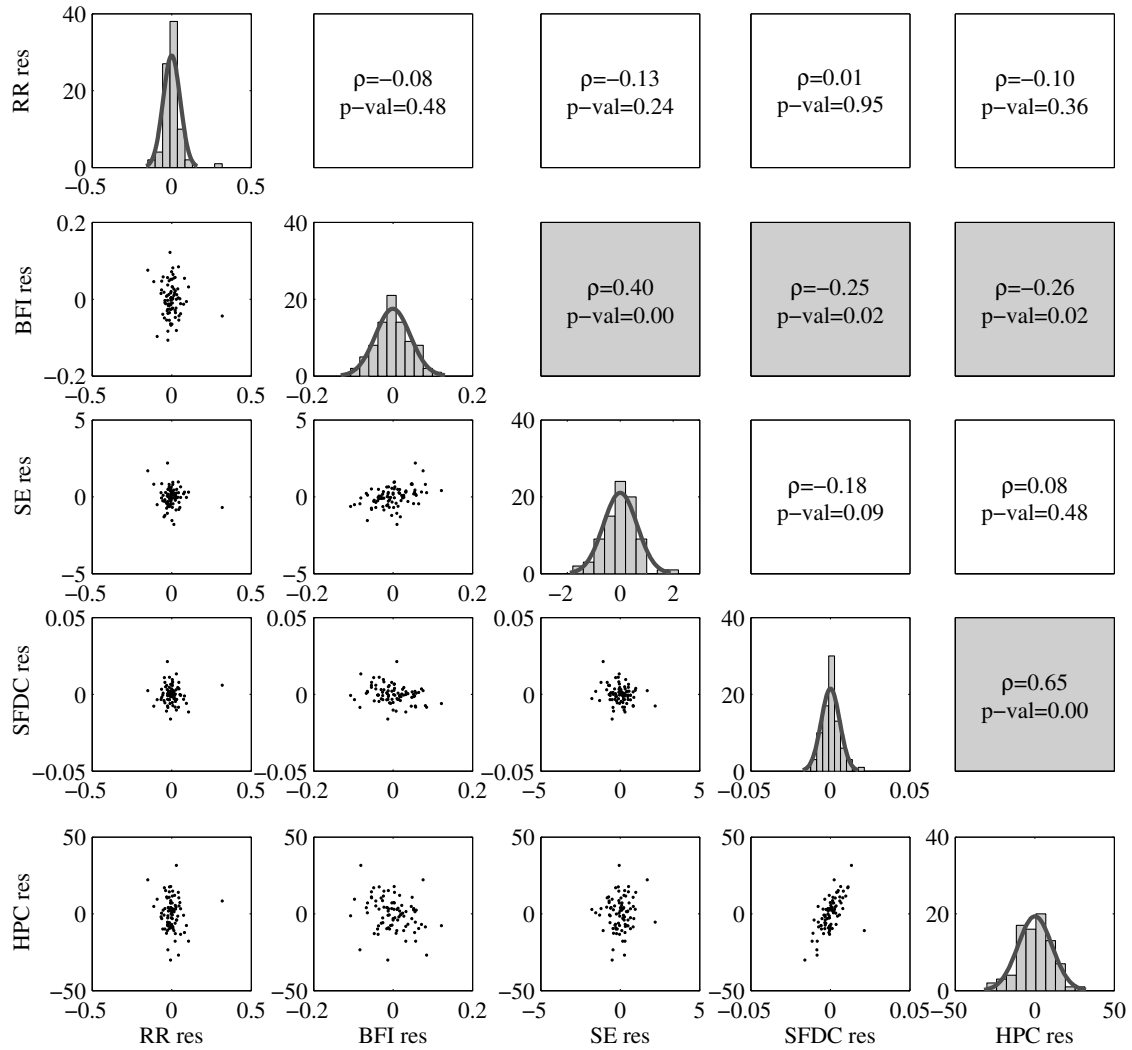


Figure 1. Distribution of individual signature residuals (res) are approximated as histograms and normal distributions. The scatterplots and correlation coefficients (ρ) show correlation between the signature residuals.

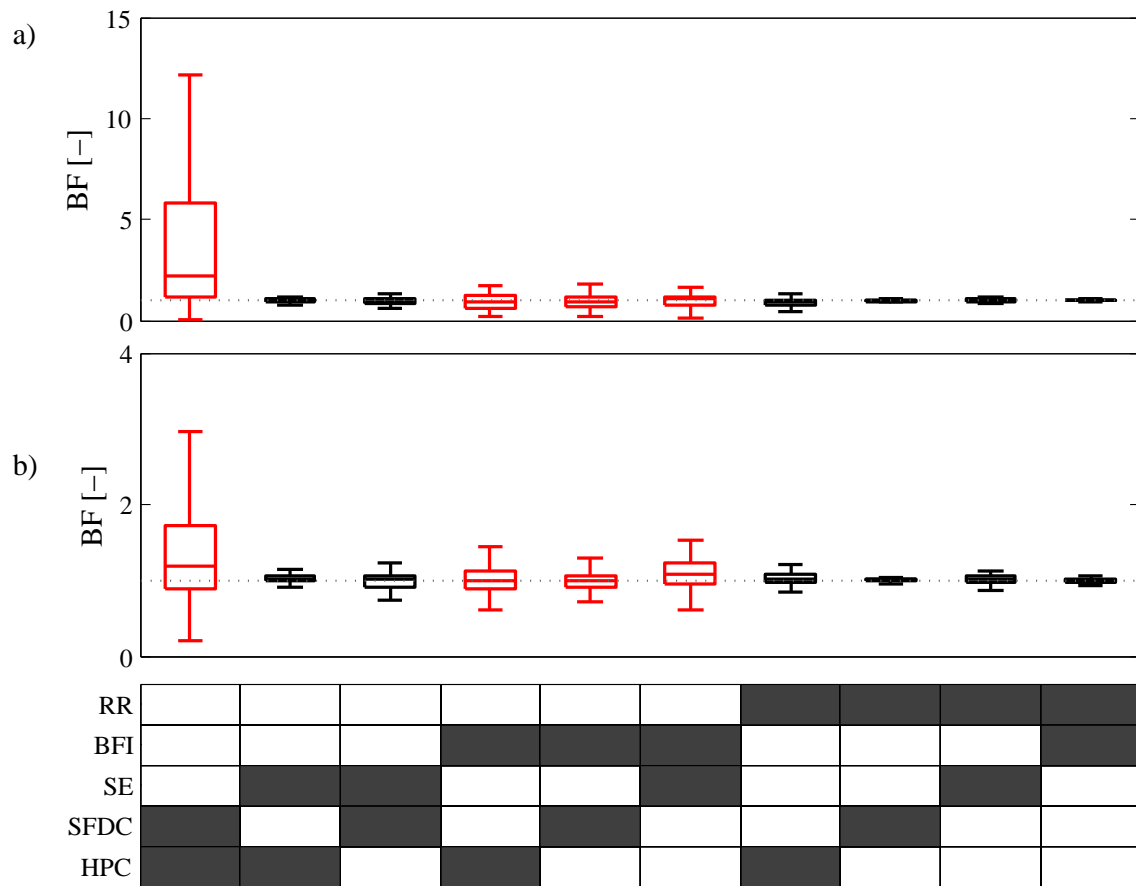


Figure 2. The Bayes factor for the 10 pairs of signatures over the 84 catchments when the observation-based error structure is used with (a) observed streamflow data, (b) synthetic streamflow data. The upper whisker represents the upper quartile plus one and a half times the interquartile range, and the lower whisker represents the lower quartile minus one and a half times the interquartile range. The dashed line represents $BF = 1$.

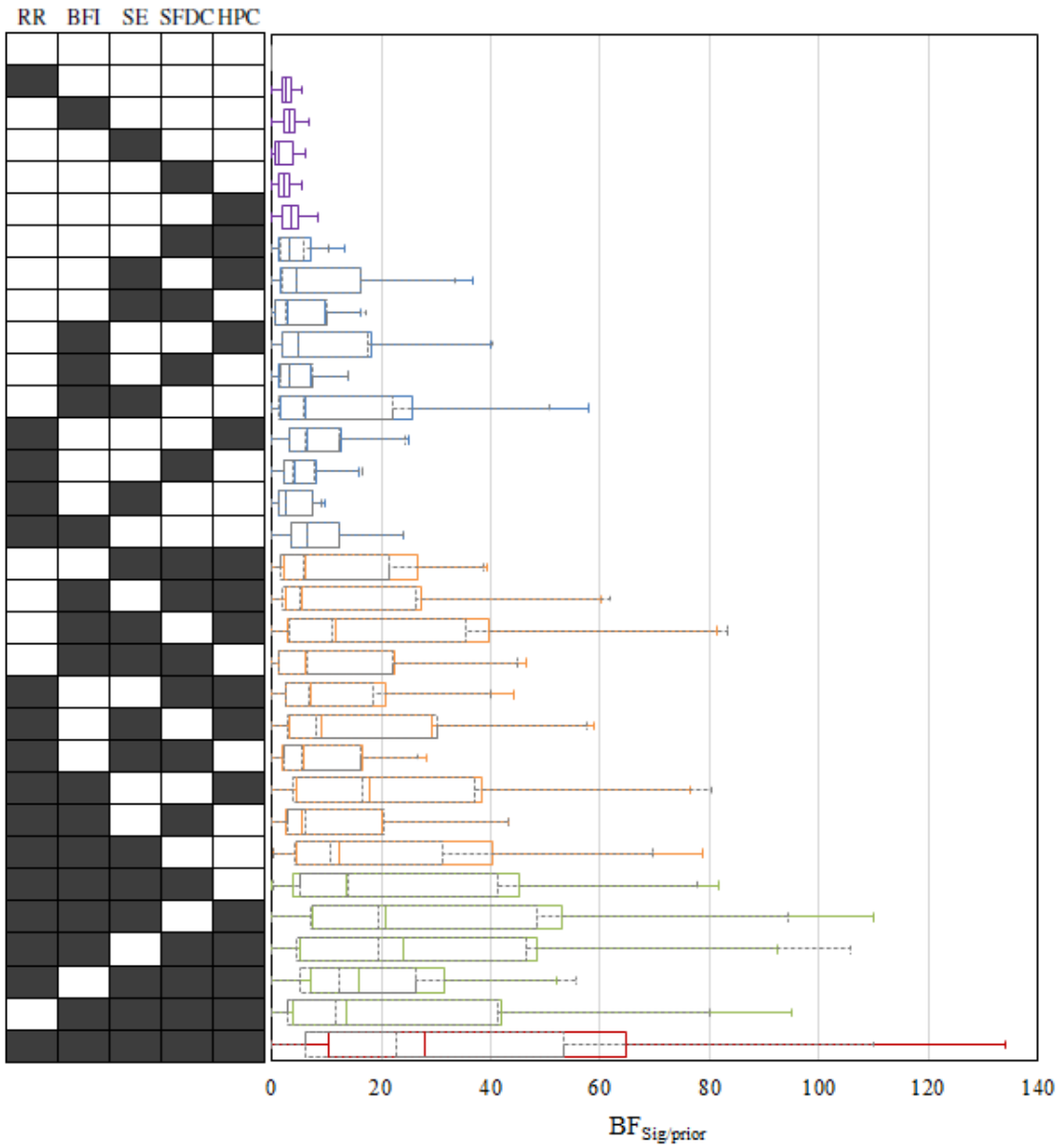


Figure 3. Boxplots representing the distribution of the Bayes factor for each combination of signatures for synthetic streamflow data. The colored boxplots correspond to the results obtained when inter-signature error correlations are considered in the likelihood function, whereas the grey dashed boxplots correspond to the results obtained assuming that the inter-signature errors are independent.

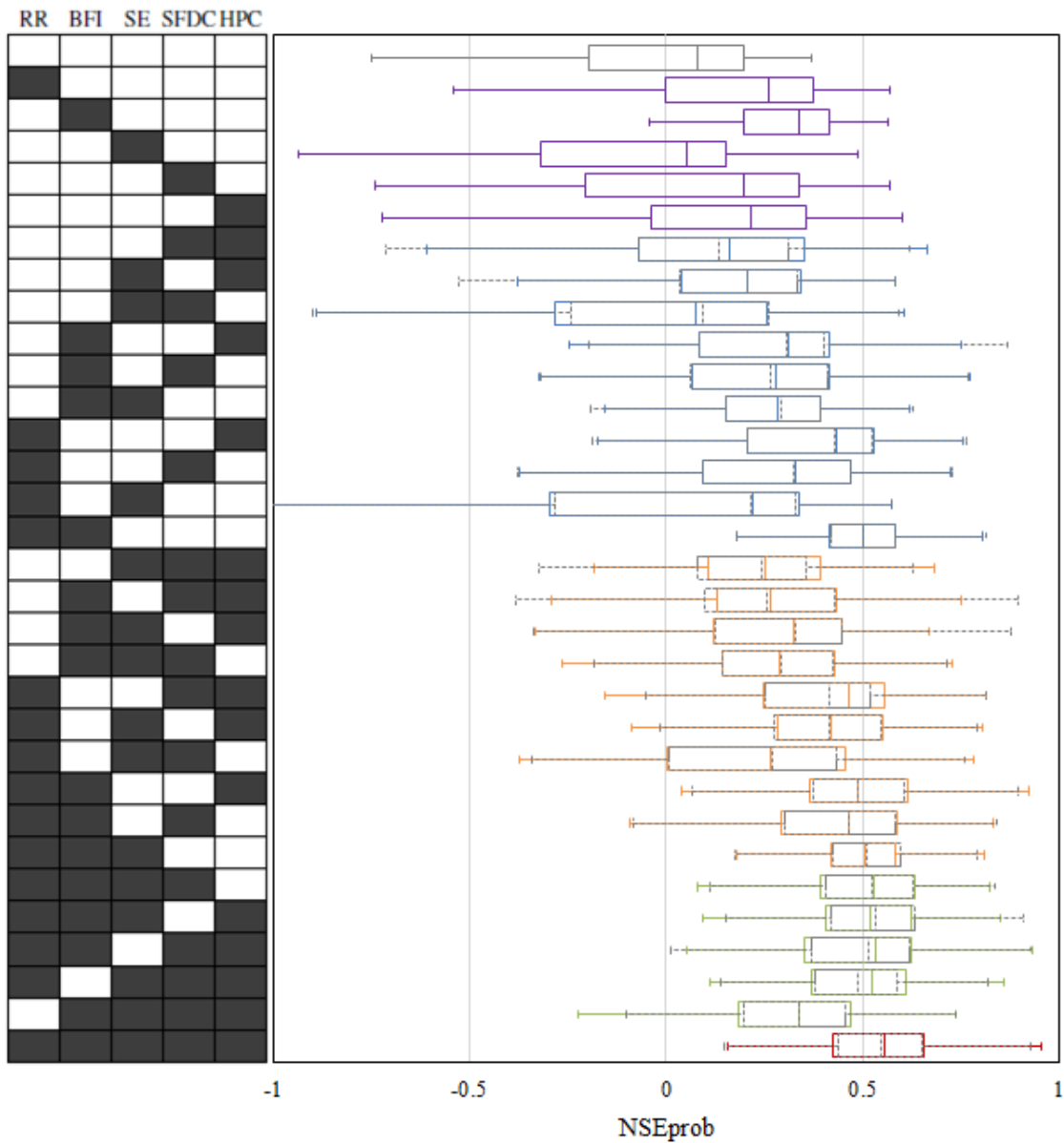


Figure 4. Boxplots representing the distribution of NSEprob values for each combination of signatures for synthetic streamflow data. The colored boxplots correspond to the results obtained when inter-signature error correlations are considered in the likelihood function, whereas the grey dashed boxplots correspond to the results obtained assuming that the inter-signature errors are independent.

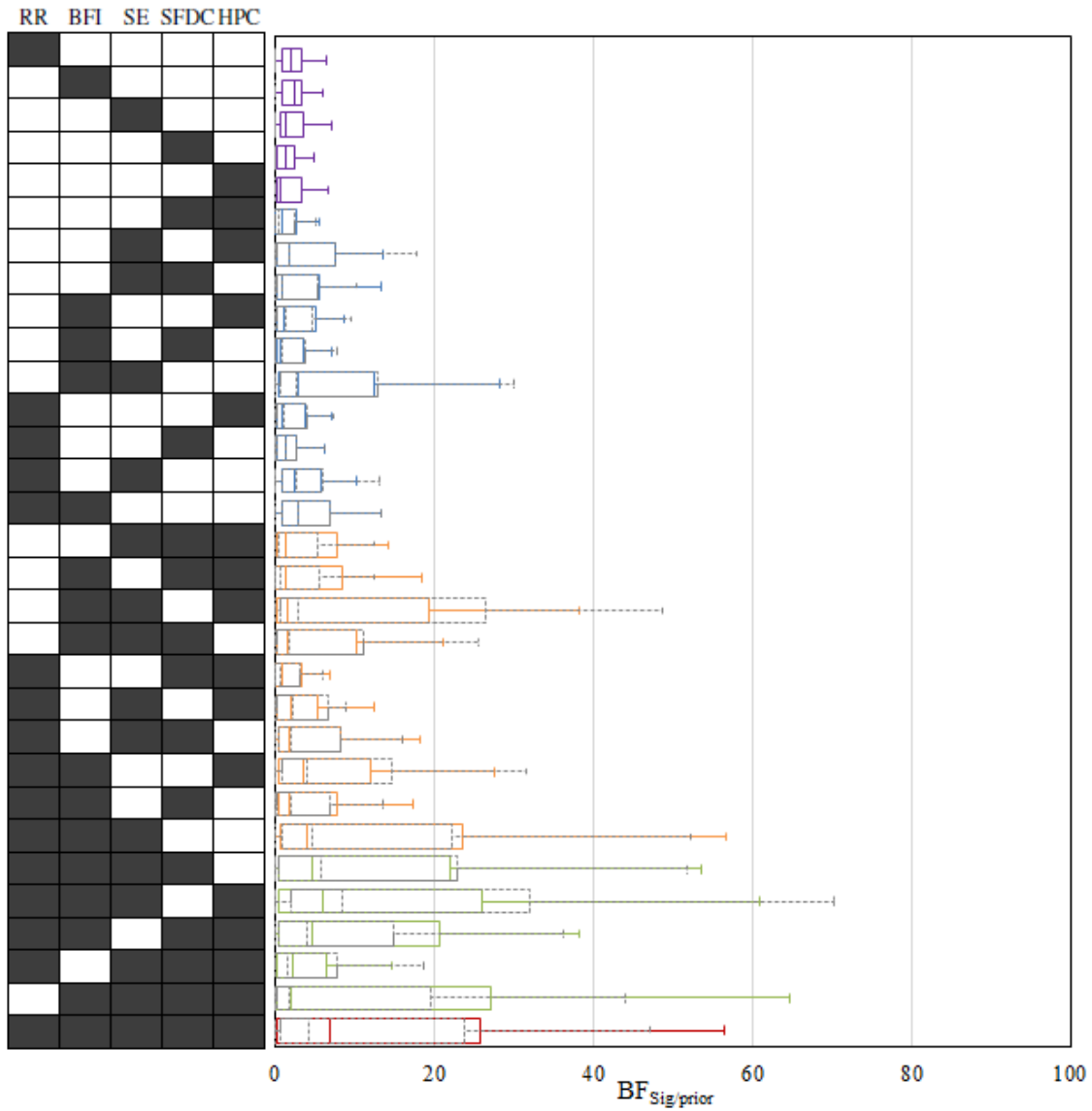


Figure 5. Boxplots representing the distribution of the Bayes factor for each combination of signatures for observed streamflow data. The colored boxplots correspond to the results obtained when inter-signature error correlations are considered in the likelihood function, whereas the grey dashed boxplots correspond to the results obtained assuming that the inter-signature errors are independent.

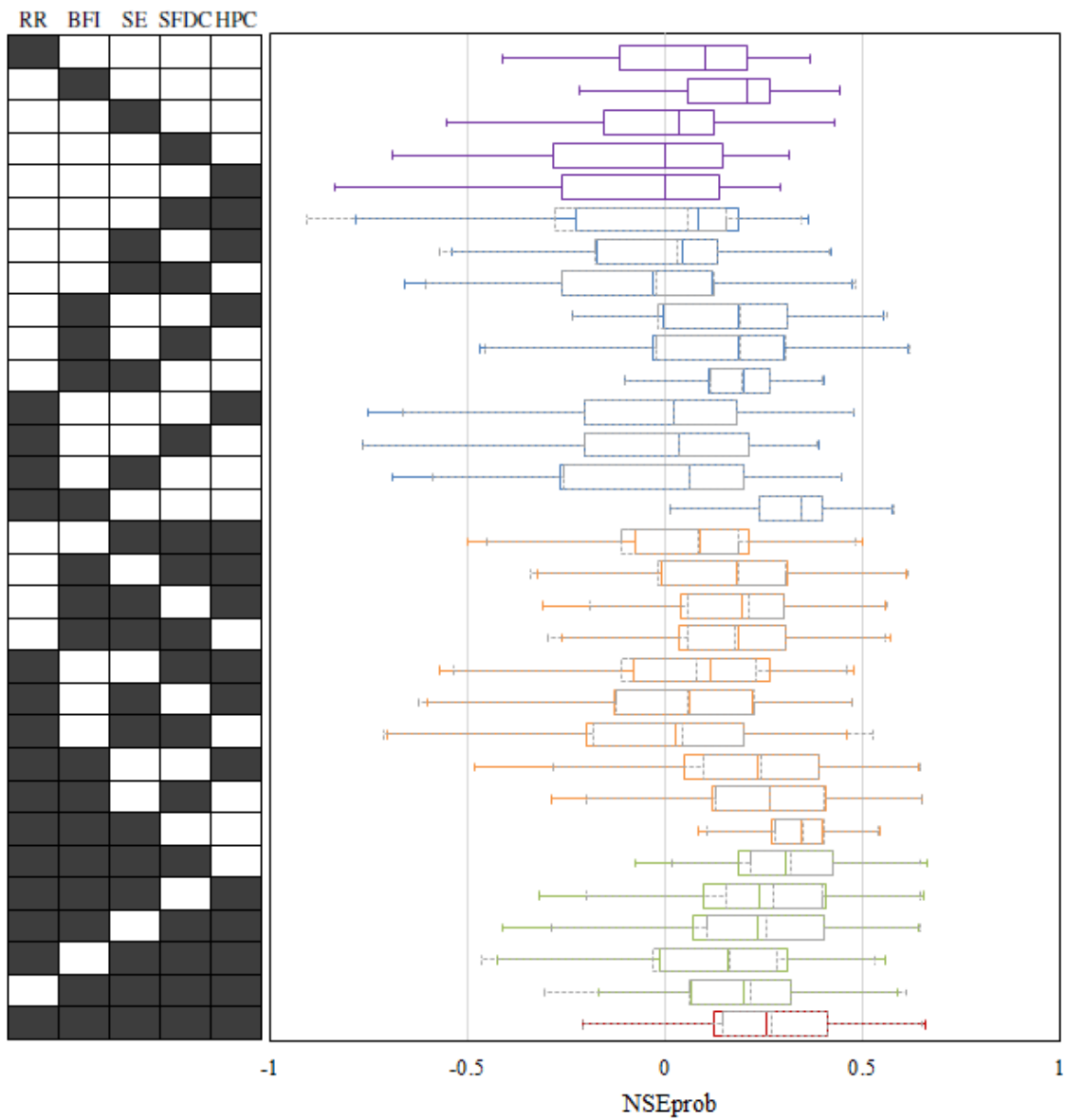


Figure 6. Boxplots representing the distribution of NSEprob values for each combination of signatures for observed streamflow data. The colored boxplots correspond to the results obtained when inter-signature error correlations are considered in the likelihood function, whereas the grey dashed boxplots correspond to the results obtained assuming that the inter-signature errors are independent.

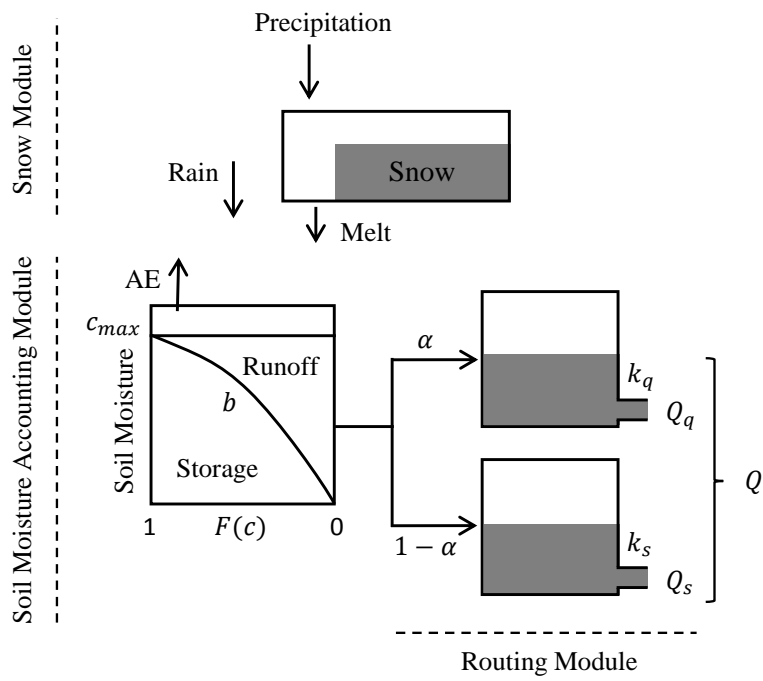


Figure A1. Schematic representation of the rainfall-runoff conceptual model structure used.