

A crash-testing framework for predictive uncertainty assessment when forecasting high flows in an extrapolation context

Lionel Berthet, François Bourgin, Charles Perrin, Julie Viatgé, Renaud Marty and Olivier Piotte

We would like to thank the reviewers very much for their careful and positive reviews and their final corrections.

Main comment

The term “non-parametric” has a specific meaning in the statistical literature, which implies that an underlying statistical distribution is not used. Examples of these include the NQT, as referred to in this paper. It is incorrect to refer to the Log and Box-Cox transformations used in this work as non-parametric transformations (even if they do not have any calibrated parameters). The authors should resolve this issue in order to avoid reader confusion.

We agree that using “non-parametric” for both statistical approaches and scalar functions (such as the variable transformations) may be rather confusing for the reader. Therefore, we now keep this word only for the statistical approaches. The transformations are now described as with / without calibrated parameters, or as “calibrated” / “uncalibrated”.

Specific comments

Page 2, lines 23-24:

*Change “In order to achieve an efficient crisis management and decision making”
to “In order to achieve efficient crisis management and decision making”*

We corrected this sentence accordingly.

Page 3, line 4:

*Change “The first step consists in identifying the different sources of uncertainty”
to “The first step consists of identifying the different sources of uncertainty”*

We corrected this sentence, as suggested.

Page 3, lines 14-15:

*Change “(e.g., in their study on hydrological prediction, Renard et al. (2010) had not to consider the uncertainty in meteorological forecasts)”
to “(e.g., in their study on hydrological prediction, which did not consider uncertainty in meteorological forecasts)”*

The previous sentence was slightly changed, because the methodology itself was not presented (for the first time by Renard *et al.*, 2010). Therefore the sentence was not change. Indeed the fact that the uncertainty in meteorological forecasts was not taken into account in this study is due to the fact that this study focuses on hydrological prediction (and not on forecasting).

Page 14, line 2:

Change “Since there are only one parameter for the Box-Cox transformation”
to “Since there is only one parameter for the Box-Cox transformation”

This sentence was not change because the full sentence is: “*there are only one parameter [...] and two parameters [...]*”. Both plural and singular seem correct (based on a short search on the Internet). However the singular option suggests an ellipsis interpretation (*there is one parameter [...] and **there are** two parameters [...]*) which may not be obvious for non native-English readers.

Page 21, lines 6-7:

Change “In addition to reliability, we looked at other qualities of the probabilistic forecasts, namely the overall performance (measured by the CRPSS) and accuracy. We also checked their sharpness.”

to “In addition to reliability, we looked at other qualities of the probabilistic forecasts, namely the overall performance (measured by the CRPSS), accuracy (measured by NSE) and sharpness (relative sharpness metric).”

We changed the sentences to add the criteria. We did not merged the two sentences in order to “separate” the forecasts qualities from the sharpness which is only a distribution property: a sharp forecasted distribution is not necessarily a good forecast *per se*.

Additional changes

Two minor changes were made in the sections “*Author contributions*” and “*Acknowledgements*”.

A crash-testing framework for predictive uncertainty assessment when forecasting high flows in an extrapolation context

Lionel Berthet¹, François Bourgin^{2,3}, Charles Perrin³, Julie Viatgé³, Renaud Marty¹, and Olivier Piotte⁴

¹DREAL Centre-Val de Loire, Loire Cher & Indre Flood Forecasting Service, Orléans, France

²GERS-LEE, Univ Gustave Eiffel, IFSTTAR, F-44344 Bouguenais, France

³Université Paris-Saclay, INRAE, UR HYCAR, 92160, Antony, France

⁴Ministry for the Ecological and Inclusive Transition, SCHAPI, Toulouse, France

Correspondence: lionel.berthet@developpement-durable.gouv.fr

Abstract. An increasing number of flood forecasting services assess and communicate the uncertainty associated with their forecasts. While obtaining reliable forecasts is a key issue, it is a challenging task, especially when forecasting high flows in an extrapolation context, i.e. when the event magnitude is larger than what was observed before. In this study, we present a crash-testing framework that evaluates the quality of hydrological forecasts in an extrapolation context. The experiment setup is based on i) a large set of catchments in France, ii) the GRP rainfall-runoff model designed for flood forecasting and used by the French operational services and iii) an empirical hydrologic uncertainty processor designed to estimate conditional predictive uncertainty from the hydrological model residuals. The variants of the uncertainty processor used in this study differ in the data transformation they use (log, Box-Cox and log-sinh) to account for heteroscedasticity and the evolution of the other properties of the predictive distribution with the discharge magnitude. Different data subsets were selected based on a preliminary event selection. Various aspects of the probabilistic performance of the variants of the hydrologic uncertainty processor, reliability, sharpness and overall quality, were evaluated. Overall, the results highlight the challenge of uncertainty quantification when forecasting high flows. They show a significant drop in reliability when forecasting high flows in an extrapolation context and considerable variability among catchments and across lead times. The increase in statistical treatment complexity did not result in significant improvement, which suggests that a parsimonious and easily understandable data transformation such as the log transformation or the Box-Cox transformation can be a reasonable choice for flood forecasting.

Copyright statement.

1 Introduction

1.1 The big one: dream or nightmare for the forecaster?

In many countries, operational flood forecasting services (FFS) issue forecasts routinely throughout the year and during rare or critical events. End-users are mostly concerned by the largest and most damaging floods, when critical decisions have to be

made. For such events, operational flood forecasters must get prepared to deal with extrapolation, i.e. to work on events of a magnitude that they and their models have seldom or never met before.

The relevance of simulation models and their calibration in evolving conditions, such as contrasted climate conditions and climate change has been studied by several authors. For example, Wilby (2005), Vaze et al. (2010), Merz et al. (2011) and Brigode et al. (2013) explored the transferability of hydrological model parameters from one period to another and assessed the uncertainty associated with this parameterization, while Coron et al. (2012) proposed a generalisation of the differential split-sample test (Klemeš, 1986). In spite of its importance in operational contexts, only a few studies have addressed the extrapolation issue for flow forecasting, to the best of our knowledge, with the notable exception of data-driven approaches (e.g., Todini, 2007). Imrie et al. (2000), Cigizoglu (2003) and Giustolisi and Laucelli (2005) evaluated the ability of trained artificial neural networks (ANNs) to extrapolate beyond the calibration data and showed that ANNs used for hydrological modelling may have poor generalisation properties. Singh et al. (2013) studied the impact of extrapolation on hydrological prediction with a conceptual model, and Barbetta et al. (2017) expressed concerns for the extrapolation context defined as floods of a magnitude not encountered during the calibration phase.

Addressing the extrapolation issue involves a number of methodological difficulties. Some data issues are specific to the data used for hydrological modelling, such as the rating curve reliability (Lang et al., 2010). Other well-known issues are related to the calibration process: are the parameters, which are calibrated on a limited set of data, representative or at least somewhat adapted to other contexts? A robust modelling approach for operational flood forecasting, i.e. a method able to provide relevant forecasts in conditions not met during the calibration phase, requires paying special attention to the behaviour of hydrological models and the assessment of predictive uncertainty in an extrapolation context.

1.2 Obtaining reliable forecasts remains a challenging task

Even if significant progress has been made and implemented in operational flood forecasting systems (e.g., Bennett et al., 2014; Demargne et al., 2014; Pagano et al., 2014), some uncertainty remains. In order to achieve an efficient crisis management and decision making, communication of reliable predictive uncertainty information is therefore a prerequisite (Todini, 2004; Pappenberger and Beven, 2006; Demeritt et al., 2007; Verkade and Werner, 2011). Hereafter, reliability is defined as the statistical consistency between the observations and the predictive distributions (Gneiting et al., 2007).

The uncertainty associated with operational forecasts is most often described by a predictive uncertainty distribution. Assessing a reliable predictive uncertainty distribution is challenging because hydrological forecasts yield residuals that show heteroscedasticity, i.e. an increase in the uncertainty variance with discharge, time auto-correlation, skewness, etc. Some studies (e.g., Yang et al., 2007b; Schoups and Vrugt, 2010) account for these properties for the calibration of hydrological models within a Bayesian framework, using specific formulations of likelihood. In an extrapolation context, it is of utter importance that the predictive uncertainty assessment provides a correct description of the evolution of the predictive distribution properties with the discharge magnitude. Bremnes (2019) showed that the skewness of wind speed distribution depends on the forecasted wind. Modelled residuals of discharge forecasts often exhibit high heteroscedasticity (Yang et al., 2007a). McInerney et al.

(2017) focused their study on representing error heteroscedasticity of discharge forecasts with respect to simulated streamflow. To achieve reliable forecasts, a correct description of the heteroscedasticity, either explicitly or implicitly, is necessary.

Various approaches for uncertainty assessment have been developed to assess the uncertainty in hydrological predictions (see e.g., Montanari, 2011). The first step consists in identifying the different sources of uncertainty or at least the most important ones that have to be taken into account given a specific context. In the context of flood forecasting, decomposing the total uncertainty into its two main components is now common: the input uncertainty (mainly the meteorological forecast uncertainty) and the modelling uncertainty, as proposed by Krzysztofowicz (1999). More generally, the predictive uncertainty due to various sources may be explicitly modelled and propagated through the modelling chain, while the “remaining” uncertainty (from the other sources) may then be assessed by statistical post-processing.

1.2.1 Modelling each source of uncertainty

A first approach ~~that~~ intends to model each source of uncertainty separately and to propagate these uncertainties through the modelling chain ~~is presented by Renard et al. (2010)~~ (e.g., Renard et al., 2010). Following this approach, the predictive uncertainty distribution results from the separate modelling of each relevant source of uncertainty (e.g., in their study on hydrological prediction, Renard et al. (2010) had not to consider the uncertainty in meteorological forecasts) and from the statistical model specification. While this approach is promising, operational application can be hindered by the challenge of making the hydrological modelling uncertainty explicit, as pointed out by Salamon and Feyen (2009).

In practice, it is necessary to consider the meteorological forecast uncertainty to issue hydrological forecasts. The ensemble approaches intend to account for this source of uncertainty. They are increasingly popular in the research and the operational forecasting communities. An increasing number of hydrological ensemble forecasting systems are in operational use and have proved their usefulness, e.g. the European Flood Awareness System (EFAS: Ramos et al., 2007; Thielen et al., 2009; Pappenberger et al., 2011, 2016) and the Hydrologic Ensemble Forecast Service (HEFS: e.g. Demargne et al., 2014).

Multi-model approaches can be used to assess modelling uncertainty (Velazquez et al., 2010; Seiller et al., 2017). While promising, this approach requires the implementation and the maintenance of a large number of models, which can be burdensome in operational conditions. There is no evidence that such an approach ensures that the heteroscedasticity of the predictive uncertainty distribution would be correctly assessed.

In forecasting mode, data assimilation schemes based on statistical modelling are of common use to reduce and assess the predictive uncertainty. Some algorithms such as particle filters (Moradkhani et al., 2005a; Salamon and Feyen, 2009; Abbaszadeh et al., 2018) or the ensemble Kalman filter (Moradkhani et al., 2005b) provide an assessment of the predictive uncertainty as a direct result of data assimilation (“in the loop”). Some of these approaches can explicitly account for the desired properties of the predictive uncertainty distribution, such as heteroscedasticity, through the likelihood formulation.

1.2.2 Post-processing approaches

Alternatively, numerous post-processors of deterministic or probabilistic models have been developed to account for the uncertainty from sources that are not modelled explicitly. They differ in several aspects (see a recent review by Li et al., 2017). Most

approaches are conditional: the predictive uncertainty is modelled with respect to a predictor, which most often is the forecasted value (Todini, 2007, 2009). Some methods are based on predictive distribution modelling, while others can be described as "distribution-free", as mentioned by Breiman (2001). Among the former, many approaches are built in a statistical regression framework to assess the total or remaining predictive uncertainty. Examples are the Hydrologic Uncertainty Processor (HUP) in a Bayesian forecasting system (BFS) framework (Krzysztofowicz, 1999; Krzysztofowicz and Maranzano, 2004), the Model Conditional Processor (MCP: Todini, 2008; Coccia and Todini, 2011; Barbetta et al., 2017), the meta-Gaussian model of Montanari and Grossi (2008) or the Bayesian joint probability (BJP) method (Wang et al., 2009), among others. The latter approaches build a description of the predictive residuals from past error series, such as data learning-algorithms (Solomatine and Shrestha, 2009). Some related methods are the non-parametric approach of Van Steenbergen et al. (2012), the empirical hydrological uncertainty processor of Bourgin et al. (2014) or the k-nearest neighbours method of Wani et al. (2017). The Quantile Regression (QR) framework (Weerts et al., 2011; Dogulu et al., 2015; Verkade et al., 2017) lies in between in that it introduces an assumption of a linear relationship between the forecasted discharge and the quantiles of interest.

1.2.3 Combining different approaches

The approaches presented in subsections 1.2.1 and 1.2.2 are not exclusive of each other. Even when future precipitation is the main source of uncertainty, post-processing is often required to produce reliable hydrological ensembles (Zalachori et al., 2012; Hemri et al., 2015; Abaza et al., 2017; Sharma et al., 2018). Thus, many operational flood forecasting services use post-processing techniques to assess hydrological modelling uncertainty, while meteorological uncertainty is taken into account separately (Berthet and Potté, 2014). Post-processors are then trained with "perfect" future rainfall (i.e., equal to the observations). Moreover, even for assessing modelling uncertainty, combining the strengths of several methodologies may improve the results.

Note that many of these approaches use a variable transformation to handle the heteroscedasticity and more generally the evolution of the predictive distribution properties with respect to the forecasted discharge. Some ~~are non-parametric~~have no calibrated parameter, while others ~~use a few~~encompass a few calibrated parameters, allowing more flexibility in the predictive distribution assessment. More details on commonly used variable transformations are presented in Section 2.1.4.

1.3 Scope

In this article, we focus on uncertainty assessment with a post-processing approach based on residuals modelling. Whilst the operational goal is to improve the hydrological forecasting, this study do not consider the meteorological forecast uncertainty: it only focuses on the hydrological modelling uncertainty, as these two main sources of uncertainty can be considered independently (*e.g.*, Bourgin et al., 2014). Del Giudice et al. (2013) and McInerney et al. (2017) presented interesting comparisons of different variable transformations used for residuals modelling. Yet, their studies do not focus on the extrapolation context. Since achieving a reliable predictive uncertainty assessment in an extrapolation context is a challenging task likely to remain imperfect if the stability of the characteristics of the predictive distributions is not properly ensured, it requires a specific crash-testing framework (Andréassian et al., 2009). The objectives of this article are:

- to present a framework aimed at testing the hydrological modelling and uncertainty assessment in the extrapolation context;
 - to assess the ability and the robustness of a post-processor to provide reliable predictive uncertainty assessment for large floods when different variable transformations are used;
- 5 – to provide guidance for operational flood forecasting system development.

We attempt to answer three questions: (a) Can we improve residuals modelling with an adequate variable transformation in an extrapolation context? (b) Do more flexible transformations, such as the log-sinh transformation, help in obtaining more reliable predictive uncertainty assessment? (c) If the performance decreases when extrapolating, is there any driver that can help the operational forecasters to predict this performance loss and question the quality of the forecasts ?

- 10 Section 2 describes the data, the forecast model, the post-processor and the testing methodology chosen to address these questions. Section 3 presents the results of the numerical experiments that are then discussed in Section 4. Finally, a number of conclusions and perspectives are proposed.

2 Data and methods

2.1 Data and forecasting model

2.1.1 Catchments and hydrological data

We used a set of 154 unregulated catchments spread throughout France (Fig. 1) over various hydrological regimes and forecast-
ing contexts to provide robust answers to our research questions (Andréassian et al., 2006; Gupta et al., 2014). They represent
a large variability in climate, topography and geology in France (Table 1), although their hydrological regimes are little or
not at all influenced by snow accumulation. Hourly rainfall, potential evapotranspiration (PE) and streamflow data series were
available over the 1997 – 2006 period. PE was estimated using a temperature-based formula (Oudin et al., 2005). Rainfall and
temperature data come from a radar-based reanalysis produced by Météo-France (Tabary et al., 2012). Discharge data were
extracted from the national streamflow HYDRO archive (www.hydro.eaufrance.fr). When a hydrological model is used to is-
sue forecasts, it is often necessary to compare the lead-time to a characteristic time of the catchment (section 2.1.2). For each
catchment, the lag time LT is estimated as the lag time maximising the cross-correlation between rainfall and discharge time
series.

Table 1. Characteristics of the 154 catchments, computed over the 1997 – 2006 data series.

| | Quantiles | | | | | | |
|---|-----------|------|------|-------|-------|-------|-------|
| | 0 | 0.05 | 0.25 | 0.50 | 0.75 | 0.95 | 1 |
| Catchment area (km ²) | 9 | 27 | 79 | 184 | 399 | 942 | 3,260 |
| Average altitude (m above sea level) | 64 | 92 | 188 | 376 | 589 | 897 | 1,050 |
| Average slope (%) | 2 | 3 | 6 | 9 | 18 | 32 | 39 |
| Lag time (h) | 3 | 5 | 9 | 12 | 19 | 29 | 33 |
| Mean annual rainfall (mm/yr) | 639 | 727 | 876 | 1,003 | 1,230 | 1,501 | 1,841 |
| Mean annual potential evapotranspiration (mm/yr) | 549 | 549 | 631 | 659 | 700 | 772 | 722 |
| Specific mean annual discharge (mm/yr) | 53 | 142 | 262 | 394 | 583 | 1,114 | 1,663 |
| Mean annual discharge (m ³ /s) | 1 | 1 | 1 | 2 | 5 | 17 | 53 |
| Maximum hourly rainfall (mm/h) | 10 | 12 | 16 | 20 | 26 | 41 | 61 |
| Quantile 0.99 of the hourly discharge (m ³ /s) | 1 | 2 | 6 | 15 | 33 | 115 | 296 |

2.1.2 Hydrological model

We used discharge forecasts computed by the GRP rainfall-runoff model. The GRP model is designed for flood forecasting and
is currently used by the Flood Forecasting Services (FFS) in France in operational conditions (Furusho et al., 2016; Viatgé et al.,
2018). It is a deterministic lumped storage-type model that uses catchment areal rainfall and PE as inputs. In forecasting mode

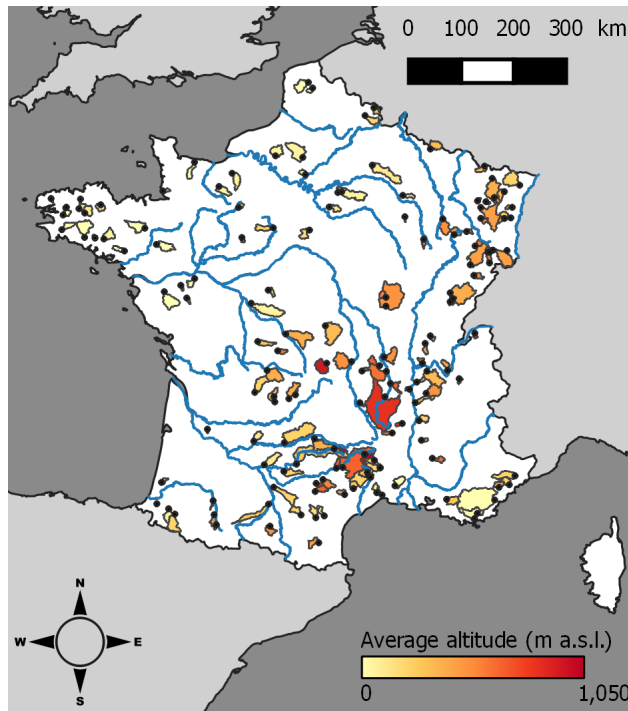


Figure 1. The set of 154 unregulated catchments used in this study. Average altitude is given in meters above sea level (*m a.s.l.*).

(Appendix A), the model also assimilates discharge observations available when issuing a forecast to update the main state variable of the routing function and to update the output discharge. In this study, it is run at the hourly time step and forecasts are issued for several lead times ranging from 1 to 72 h. More details about the GRP model can be found in Appendix A.

Since herein only the ability of the post-processor to extrapolate uncertainty quantification is studied, the model is fed only with observed rainfall (no forecast of precipitation), in order to reduce the impact of the input uncertainty. For the same reason, the model is calibrated in forecasting mode over the 10-year series by minimising the sum of squared errors for a lead time taken as the lag time LT. The results will be presented for four lead times, LT / 2, LT, 2 LT and 3 LT, to cover the different behaviours that can be seen when data assimilation is used to reduce errors in an operational flood forecasting context (Berthet et al., 2009).

2.1.3 Empirical hydrological uncertainty processor (EHUP)

We used the empirical hydrological uncertainty processor (EHUP) presented in Bourgin et al. (2014). It is a data-based and non-parametric approach to estimate the conditional predictive uncertainty distribution. This post-processor was compared to other post-processors in earlier studies and proved to provide relevant results (Bourgin, 2014). It is now used by operational FFS in France under the operational tool called OTAMIN (Viatgé et al., 2019). The main difference with many other post-processors (such as the MCP, the meta-Gaussian processor, the BJP or the BFS) is that no assumption is made about the shape

of the uncertainty distribution, which brings more flexibility to represent the various forecast error characteristics encountered in large sample modelling. We will discuss the impact of this choice in Sect. 4.2.

The basic idea of the empirical hydrological uncertainty processor (EHUP) is to estimate empirical quantiles of errors stratified by different flow groups to account for the variation of the forecast error characteristics with forecast magnitude.

- 5 Since forecast error characteristics also vary with the lead time when data assimilation is used, the EHUP is trained separately for each lead time.

For each lead time separately, the following steps are used:

1. Training:

- 10 – The flow groups are obtained by first ordering the forecast-observation pairs according to the forecasted values and then stratifying the pairs into a chosen number of groups (in this study, we used 20 groups), so that each group contains the same number of pairs.
- Within each flow group, errors are calculated as the difference between the two values of each forecast-observation pairs and several empirical quantiles (we used 99 percentiles) are calculated in order to characterize the distribution of the error values.

15 2. Application:

- The predictive uncertainty distribution that is associated to a given (deterministic) forecasted value is defined by adding this forecasted value to the empirical quantiles that belong to the same flow group as the forecasted value.

Since this study focuses on the extrapolation case, the validation is achieved on deterministic forecasts higher than the highest one used for the calibration. Therefore, only the highest flow group of the calibration data is used to estimate the uncertainty assessment (to be used on the control data). This highest flow group contains the top 5% pairs of the whole training data, ranked by forecasted values. This threshold is chosen as a compromise between focusing at the highest values and using a sufficiently large number of forecast-observation pairs when estimating empirical quantiles of errors. In extrapolation, when the forecast discharge is higher than the highest value of the training period, the predictive distribution of the error is kept constant, i.e., the same values of the empirical quantiles of errors are used, as illustrated in Fig. 5.

25 The EHUP can be applied after a preliminary data transformation, and by adding a final step to back-transform the predictive distributions obtained in a transformed space. In previous work, we used the log transformation because it ensures that no negative values are obtained when estimating the predictive uncertainty for low flows (Bourgin et al., 2014). When estimating the predictive uncertainty for high flows, the data transformation has a strong impact in extrapolation, because the variation of the extrapolated predictive distribution, which is constant in the transformed space, is controlled in the real space by the

30 behaviour of the inverse transformation, as explained below.

2.1.4 The different transformation families

Many uncertainty assessment methods mentioned in the introduction use a variable transformation to handle the heteroscedasticity of the residuals and account for the variation of the prediction distributions with the magnitude of the predicted variable. Here, we briefly recall a number of variable transformations commonly used in hydrological modelling. Let y and \tilde{y} be the observed and forecasted variables (here, the discharge) and $\varepsilon = y - \tilde{y}$ the residuals. When using a transformation g , we consider the residuals $\varepsilon' = g(y) - g(\tilde{y})$.

Three analytical transformations are often met in hydrological studies: the log, Box-Cox and log-sinh transformations. The log transformation is commonly used (e.g., Morawietz et al., 2011):

$$g_1 : y \mapsto \log(y + a)$$

- where a is a small positive constant to deal with y values close to 0. It can be taken equal to 0 when focusing on large discharge values. This transformation ~~is then non-parametric (no parameter has~~ has no parameter to be calibrated). Applying a statistical model on residuals computed on log-transformed variables may be interpreted as using a corresponding model on multiplicative error (e.g., assuming a Gaussian model for residuals of log-transformed discharges is equivalent to a log-Normal model on the multiplicative errors y/\tilde{y}). Therefore, it may be adapted to strongly heteroscedastic behaviours. It has been used successfully to assess hydrological uncertainty (Yang et al., 2007b; Schoups and Vrugt, 2010; Bourgin et al., 2014).

The Box-Cox transformation (Box and Cox, 1964) is a classic transformation that is quite popular in the hydrological community (e.g., Yang et al., 2007b; Wang et al., 2009; Hemri et al., 2015; Reichert and Mieleitner, 2009; Singh et al., 2013; Del Giudice et al., 2013):

$$g_{2[\lambda]} : y \mapsto \frac{(y + a)^\lambda - 1}{\lambda} \quad \text{if } \lambda \neq 0 \quad \text{and} \quad g_{2[\lambda]} = g_1 \quad \text{if } \lambda = 0$$

- Here a is taken equal to 0 (as for the log transformation); the Box-Cox transformation is then a one-parameter transformation. It makes it possible to cover very different behaviours. The log transformation is a special case of the Box-Cox transformation when the calibration results in $\lambda = 0$. In contrast, applying the Box-Cox transformation with $\lambda = 1$ on the variable y to model the distribution of their residuals is equivalent to applying no transformation (Fig. 2). McInerney et al. (2017) obtained their most reliable and sharpest results with $\lambda = 0.2$ over 17 perennial catchments.
- More recently, the log-sinh transformation has been proposed (Wang et al., 2012; Pagano et al., 2013). It is a two-parameter transformation:

$$g_{3[\alpha, \beta]} : y \mapsto \beta \cdot \log \left(\sinh \left(\frac{\alpha + y}{\beta} \right) \right)$$

This transformation provides more flexibility. Indeed, for $y \gg \alpha$ and $y \gg \beta$, the log-sinh transformation reduces to no transformation, while for $\alpha \ll y \ll \beta$, it is equivalent to the log transformation (Fig. 2). Thus, with the same parameterization, it can result in very different behaviours depending on the magnitude of the discharge. Applying no transformation may be intuitively attractive to model the residuals distribution for large discharge values, when the variance is not longer expected to increase

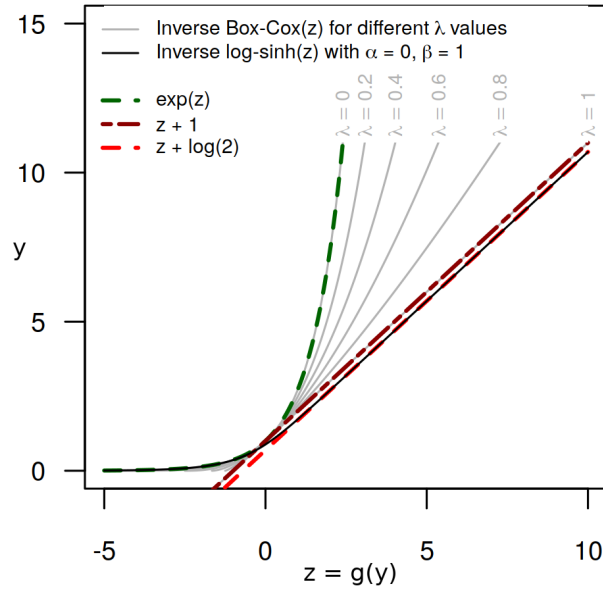


Figure 2. The inverse transformation explains the final effect on the uncertainty assessment: the constant probability distribution in the transformed space (provided by the EHUP) will result in a distribution in the untransformed space, whose evolution depends on the behaviour of the inverse data transformation. Here, the Box-Cox transformation provides different behaviours, depending on its parameter value (λ). It ranges from an affine transformation (equivalent to no transformation; red dashed line) to the log transformation (thick green dashed line). With a single parameterization, the log-sinh transformation can be equivalent to the log transformation for values of y much smaller than the value of its parameter β and equivalent to an affine transformation for large values of y (much higher than β ; see Appendix B).

(homoscedastic behaviour). It is then particularly attractive when modelling predictive uncertainty in an extrapolation context, in order to avoid an excessively "explosive" assessment of the predictive uncertainty for large discharge values.

In addition to the log transformation used by Bourgin et al. (2014), in this study we tested the Box-Cox and the log-sinh transformations to explore more flexible ways to deal with the challenge of extrapolating prediction uncertainty distributions (Fig. 2). The impacts of the data transformations used in this study are illustrated in Fig. 3.

Another common variable transformation is the normal quantile transformation (NQT; *e.g.*, Kelly and Krzysztofowicz, 1997). It is a ~~non-parametric transformation~~ transformation without calibrated parameter linking a given distribution and the Gaussian distribution, quantile by quantile. While several hydrological processors such as the HUP, MCP and the Quantile Regression encompass the NQT-transformed variables, Bogner et al. (2012) warn against the drawbacks of this transformation, which is by construction not suited for the extrapolation context and requires additional assumptions to model the tails of the distribution (Weerts et al., 2011; Coccia and Todini, 2011). This is why we did not test this transformation in this study focused on the extrapolation context.

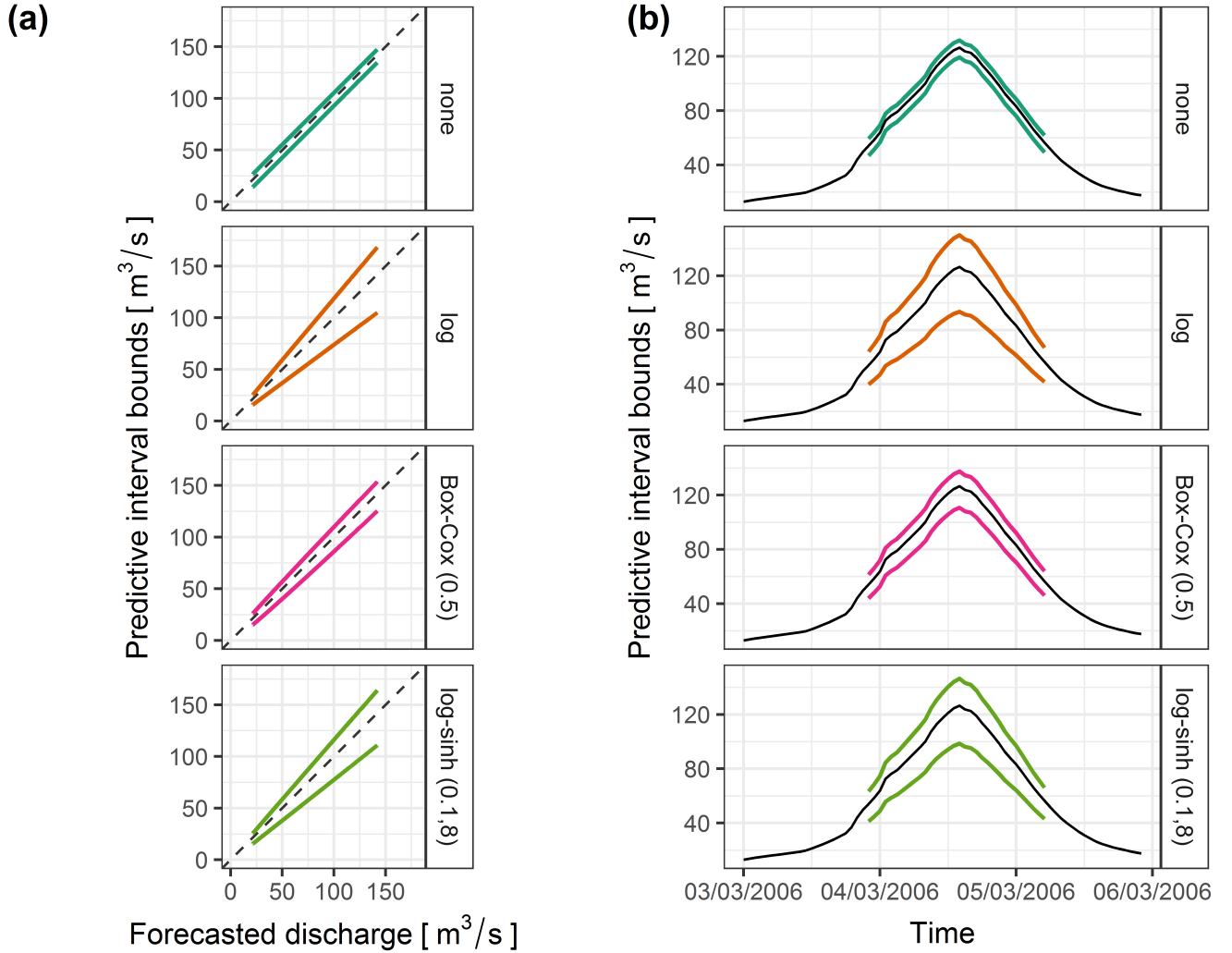


Figure 3. Predictive 0.1 and 0.9 quantiles when assessed with no transformation, the log-transformation, the Box-Cox transformation with its λ parameter equal to 0.5 and the log-sinh with its α and β parameters equal to 0.1 and 8, on the Ill River at Didenheim (668 km^2) for lead time LT: (a) as a function of the deterministic forecasted discharge and (b) for the flood event on 3 March 2006. The heteroscedasticity strongly differs from one variable transformation to another.

2.2 Methodology: a testing framework designed for extrapolation context assessment

2.2.1 Testing framework

The empirical hydrological uncertainty processor (EHUP) is a non-parametric approach based on the characteristics of the residuals distribution over a training data set. Moreover, the Box-Cox and the log-sinh transformations are parametric and

Table 2. Characteristics of the events selected for the lead time LT over the 1997 – 2006 data series.

| | Quantiles | | | | | | |
|--|-----------|-------|-------|-------|-------|-------|-------|
| | 0 | 0.05 | 0.25 | 0.50 | 0.75 | 0.95 | 1 |
| Total length of events (days) | 663 | 1178 | 1299 | 1505 | 1808 | 2061 | 2762 |
| Number of events G1 | 28 | 65 | 114 | 141 | 193 | 276 | 434 |
| Number of events G2 | 8 | 16 | 23 | 32 | 43 | 58 | 170 |
| Number of events G3 | 7 | 12 | 19 | 26 | 35 | 46 | 162 |
| Specific median value of the peak discharges G1 (mm/h) | 0.009 | 0.022 | 0.041 | 0.062 | 0.098 | 0.184 | 0.266 |
| Specific median value of the peak discharges G2 (mm/h) | 0.026 | 0.064 | 0.140 | 0.217 | 0.342 | 0.706 | 1.042 |
| Specific median value of the peak discharges G3 (mm/h) | 0.050 | 0.131 | 0.295 | 0.439 | 0.701 | 1.896 | 3.071 |

require a calibration step. Therefore, the methodology adopted for this study is a split-sample scheme test inspired by the differential split-sample scheme of Klemeš (1986) and based on three data subsets: a data set for training the EHUP, a data set for calibrating the parameters of the variable transformation, and a control data set for evaluating the predictive distributions when extrapolating high flows.

5 **2.2.2 Events selection**

To populate the three data subsets with independent data, separate flood events were first selected by an iterative procedure similar to those detailed by Lobligeois et al. (2014) and Ficchi et al. (2016): (1) the maximum forecasted discharge of the whole time series was selected, (2) within a 20-day period before (after) the peak flow, the beginning (end) of the event was placed at the preceding (following) time step closest to the peak at which the streamflow is lower than 20% (25%) of the peak flow value, (3) the event was kept if there was less than 10% missing values, if the beginning and end of the event were lower than 66% of the peak flow value and if the peak value was higher than the median value of the time series. The process is then iterated over the remaining data to select all events. A minimum time lapse of 24 h was enforced between two events, ensuring that consecutive events are not overlapping and that the autocorrelation between the time steps of two separate events remains limited.

15 The number of events and their characteristics vary greatly among catchments, as summarised in Table 2. Note that the events selected for one catchment can slightly differ for the four different lead times considered in this study, because the selection was made using the forecasted discharge and not the observed discharge.

2.2.3 Selection of the data subsets

20 The selected events were then gathered into three events sets, G1, G2 and G3, based on the magnitude of their peaks and the number of useful time steps for each test phase (training of the EHUP post-processor, calibration of the variable transformations, evaluation of the predictive distributions): G1 contains the lowest events, while the highest events are in G3.

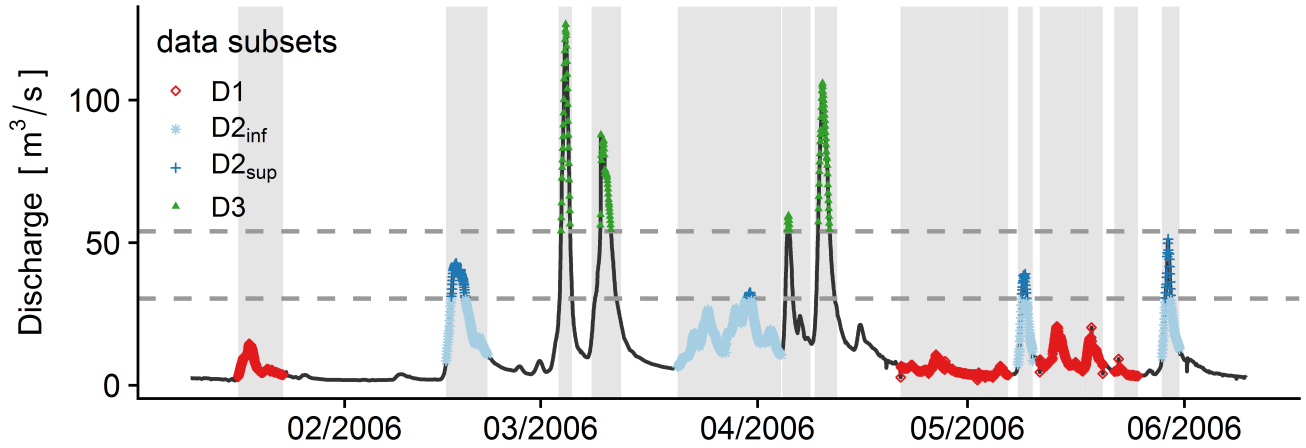


Figure 4. Illustration of the selection of the data subsets for the Ill River at Didenheim (668 km²). First, the events are selected (grey highlighting). Then, the four data subsets are populated according to the thresholds (horizontal dashed lines). See Sect. 2.1.3 for more details.

The selection of the data subsets was tailored to study the behaviour of the post-processing approach in an extrapolation context. The control data subset had to encompass only time steps with simulated discharge values higher than those met during the training and calibration steps. Similarly, the calibration data subset had to encompass time steps with simulated discharge values higher than those of the training subset.

- 5 To achieve these goals, only the time steps within flood events were used. We distinguished four data subsets, as illustrated in Fig. 4. The subset D1 gathered all the time steps of the events of the G1 group. Then, the set D2 of the time steps of the events of the G2 group was split into two subsets: D2_{sup} gathered all the time steps with forecasted discharge values higher than the maximum met on D1, and D2_{inf} was filled with the other time steps. Finally, D3 was similarly filled with all the time steps of the G3 events with forecasted discharge values higher than the maximum met on D2.
- 10 The discharge thresholds used to populate the D1, D2_{sup} and D3 subsets from the events belonging to the G1, G2 and G3 groups were chosen to ensure a sufficient number of time steps in every subset. We chose to set the minimum number of time steps in D3 and D2_{sup} to 720 as a compromise between having enough data to evaluate the methods and keeping the extrapolation range sufficiently large. We lowered this limit to 500 for the top 5% pairs of D1, since this subset was only used to build the empirical distribution by estimating the percentiles during the training step and not used for evaluating the quality
- 15 of the uncertainty assessment.

2.2.4 Calibration and evaluation steps

Since there are only one parameter for the Box-Cox transformation and two parameters for the log-sinh transformation, a simple calibration approach of the transformation parameters was chosen: the parameter space was explored by testing several parameter set values. For the Box-Cox transformation, 17 values for the λ parameter were tested: from 0 to 1 with a step of 0.1 and with a refined mesh for the sub-intervals 0 – 0.1 and 0.9 – 1. For the log-sinh transformation, a grid of 200 (α, β) values was designed for each catchment based on the maximum value of the forecasted discharge on the D2 subset, as explained in greater detail in the Appendix B.

Note that the hydrological model was calibrated over the whole set of data (1997 – 2006) to make the best use of the data set, since this study focuses on the effect of extrapolation on the predictive uncertainty assessment only.

We used a two-step procedure, as illustrated in Fig. 5. The first step was to calibrate the parametric transformations. For each transformation and for each parameter set, the empirical residuals were computed over D1 and D2_{sup} where the EHUP was trained on the highest flow group of D1 (see section 2.1.3) and the calibration criterion was computed on D2_{sup}. Indeed, the data transformations have almost no impact on the uncertainty estimation by EHUP for events of the same magnitude as those of the training subset. Therefore, the calibration subset has to encompass events of a larger magnitude (D2_{sup}). The parameter set obtaining the best criterion value was selected. In the second step, the EHUP was trained on the highest flow group of D1, D2_{inf} and D2_{sup} combined and using the parameter set obtained during the calibration step. Then, the predictive uncertainty distribution was evaluated on the control data set D3. Training the EHUP on the highest flow group of the union of D1, D2_{inf} and D2_{sup} allows to control the uncertainty assessment from small to large degrees of extrapolation (on D3). Indeed if we had kept the training on D1 only, we would have not been able to test small degrees of extrapolation on independent data for every catchment (see the discussion in Sect. 3.3).

2.3 Performance criteria and calibration

2.3.1 Probabilistic evaluation framework

Reliability was first assessed by a visual inspection of the Probability integral transform (PIT) diagrams (Laio and Tamea, 2007; Renard et al., 2010). Since this study was carried out over a large sample of catchments, two standard numerical criteria were used to summarise the results: the α -index, which is directly related to the PIT diagram (Renard et al., 2010), and the coverage rate of the 80% predictive intervals (bounded by the 0.1 and 0.9 quantiles of the predictive distributions), used by the French operational FFS (a perfect uncertainty assessment would reach the 0.8 value). The α -index is equal to $1 - 2 \cdot A$, where A is the area between the PIT curve and the bisector, and its value ranges from 0 to 1 (perfect reliability).

The overall quality of the probabilistic forecasts was evaluated with the Continuous Rank Probability Score (CRPS, Hersbach, 2000), which compares the predictive distribution to the observation:

$$\text{CRPS} = \frac{1}{N} \sum_{k=1}^N \int_0^{+\infty} [F_k(Q) - H(Q - Q_{k,obs})]^2 dQ$$

where N is the number of time steps, $F(Q)$ is the predictive cumulative distribution, H the Heaviside function and Q_{obs} is the observed value. We used a skill score (CRPSS) to compare the mean CRPS to a reference, here the mean CRPS obtained from the unconditional climatology, i.e. from the distribution of the observed discharges over the same data subset. Values range from $-\infty$ to 1 (perfect forecasts): a positive value indicates that the modelling is better than the climatology.

- 5 For operational purposes, the sharpness of the probabilistic forecasts was checked by measuring the mean width of the 80% predictive intervals. A dimensionless relative sharpness index was obtained by dividing the mean width by the mean runoff:

$$1 - \frac{\sum_{k=1}^N [q_{0.9}(Q_k) - q_{0.1}(Q_k)]}{\sum_{k=1}^N Q_{k,obs}}$$

where $q_{90}(Q)$ and $q_{10}(Q)$ are the upper and the lower bounds of the 80% predictive interval for each forecast. The sharper the forecasts, the closer this index is to 1.

- 10 In addition to the probabilistic criteria presented above, the accuracy of the forecasts was assessed using the Nash-Sutcliffe Efficiency (NSE) calculated with the mean values of the predictive distributions (best value: 1).

2.3.2 The calibration criterion

- Since the calibration step aims at selecting the most reliable description of the residuals in extrapolation, the α -index was used to select the parameter set that yields the highest reliability for each catchment, each lead time and each transformation. While
 15 other choices were possible, we followed the paradigm presented by Gneiting et al. (2007): reliability has to be ensured before sharpness. Note that the CRPS could have been chosen, since it can be decomposed as the sum of two terms: reliability and sharpness (Hersbach, 2000). However, in the authors' experience, the latter is the controlling factor (Bourgin et al., 2014). Moreover, the CRPS values were often quite insensitive to the values of the log-sinh transformation parameters.

- In cases where an equal value of the α -index was obtained, we selected the parameter set that gave the best sharpness
 20 index. For the log-sinh transformation, there were still a few cases where an equal value of the sharpness index was obtained, revealing the lack of sensitivity of the transformation in some areas of the parameter space. For those cases, we chose to keep the parameter set that had the lowest α value and the β value closest to $max_{D2}(\tilde{Q})$.

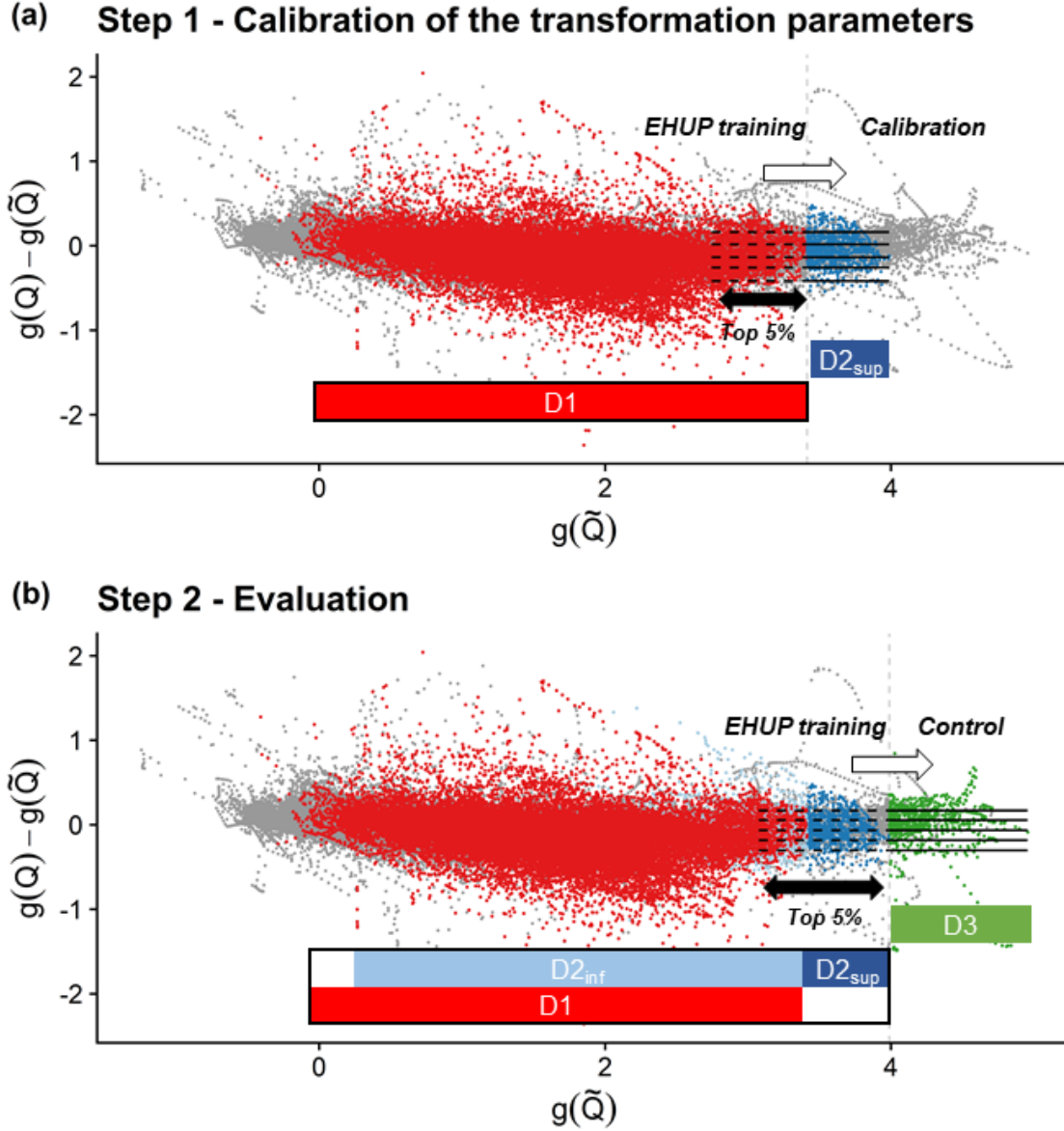


Figure 5. Residuals as a function of the forecast discharges in the transformed space. The horizontal dashed lines represent the 0.1, 0.25, 0.5, 0.75 and 0.9 quantiles of the residuals computed during the training phase of the EHUP post-processor, for the highest flow group (top 5% pairs of the training data ranked by forecasted values). The straight horizontal lines represent their use to assess the predictive uncertainty in extrapolation during (a) the calibration step of the variable transformation parameters and (b) the evaluation step of the predictive uncertainty. During the calibration step (a), many parameters set values are tested, while only the calibrated set of transformation parameters is used during the evaluation step (b). The vertical dashed lines show the beginning of the extrapolation range. The grey dots are the data pairs which are not selected in D1, nor in D2_{inf}, D2_{sup} and D3: they are not used during the calibration step and then the evaluation step. Illustration from the Ill River at Didenheim, 668 km². Data used for the EHUP training at each step is sketched by a thick rectangle. In this study, only the top 5% pairs of the training data ranked by forecasted values is used to estimate the residual distribution, since we focus on the extrapolation behaviour of the EHUP.

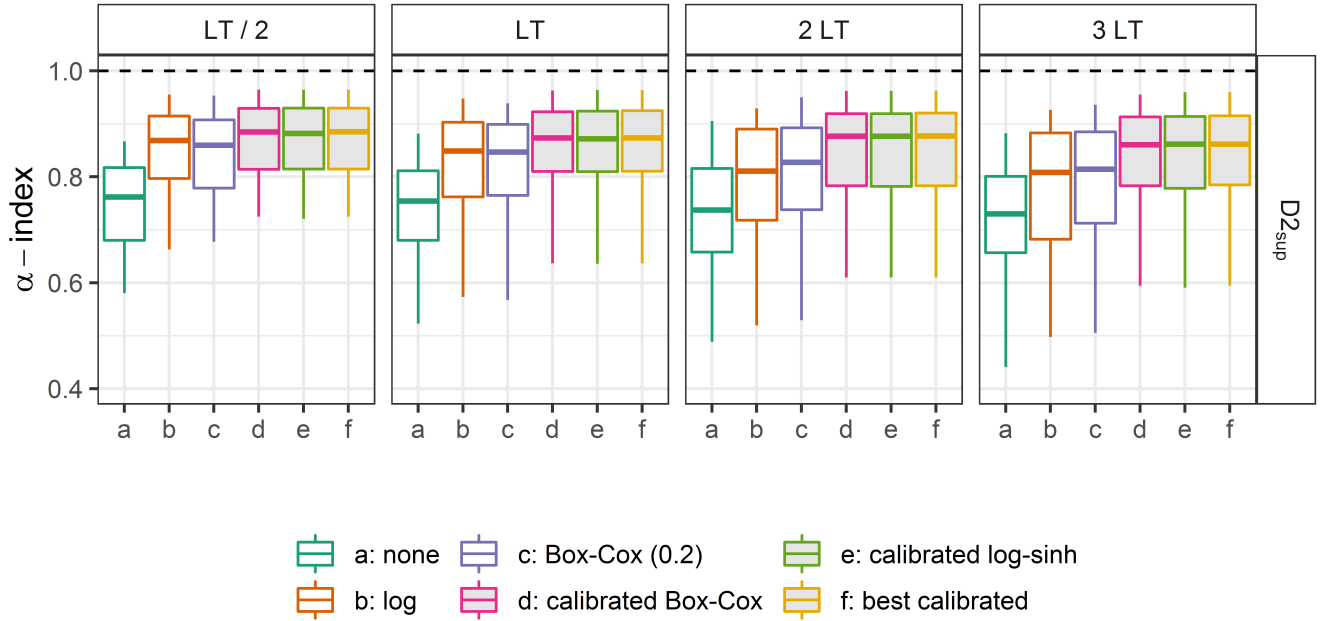


Figure 6. Distributions of the α -index values on the calibration data set $D2_{sup}$, obtained with different transformations for four lead times (the filled boxplots represent the calibrated distributions). Boxplots (5th, 25th, 50th, 75th and 95th percentiles) synthesise the variety of scores over the catchments of the data set. The optimal value is represented by the horizontal dashed lines.

3 Results

3.1 Results on the calibration data set $D2_{sup}$

Figure 6 shows the distributions of the α -index values obtained with different transformations on the calibration data set ($D2_{sup}$) for lead times LT / 2, LT, 2 LT, 3 LT. The distributions are summarised with boxplots. Clearly, not using any transformation leads to poorer reliability than any tested transformation. In addition, we note that the calibrated transformations provide better results (although not perfect) than the ~~non-parametric-non calibrated~~ ones on the calibration data set, as expected, and that no noticeable difference can be seen in Fig. 6 between the calibrated Box-Cox transformation (d), the calibrated log-sinh transformation (e) and the best calibrated transformation (f). Nevertheless, the uncalibrated log transformation and Box-Cox transformation with parameter λ set at 0.2 ($BC_{\lambda=0.2}$) reach quite reliable forecasts. Comparing the results obtained for the different lead times reveals that less reliable predictive distributions are obtained for longer lead times, in particular for the ~~non-parametric-transformations-transformations without calibrated parameter~~.

Figures 7 and 8 show the distribution of parameter values obtained for the Box-Cox and the log-sinh transformation during the calibration step. The distributions vary with lead time. While the log-transformation behaviour is frequently chosen for LT/2 and LT, the additive behaviour (corresponding to the use of no transformation, see Sect. 2.1.4) becomes more frequent

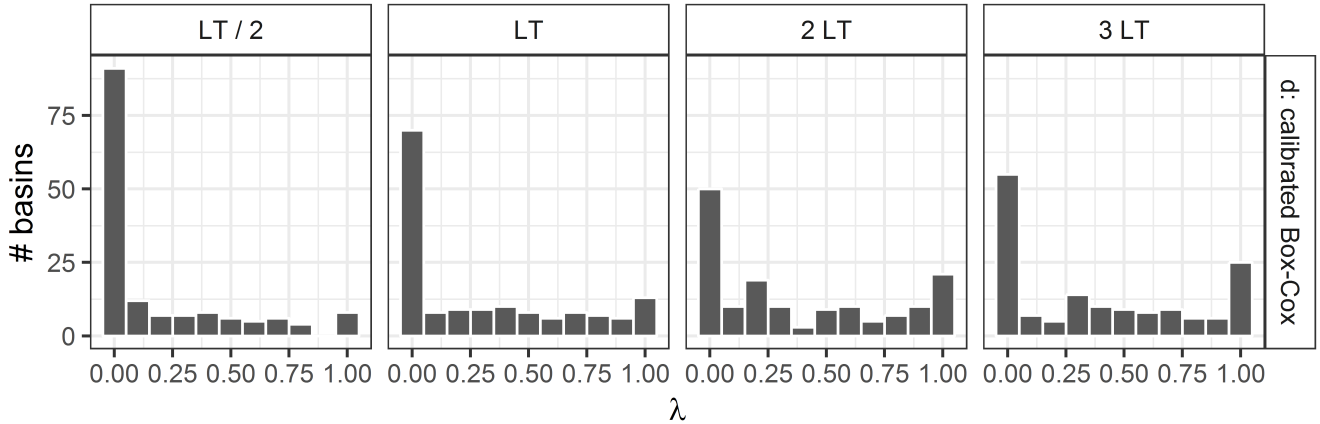


Figure 7. Distribution over the basins of the values of the Box-Cox transformation parameter obtained during the calibration step for the four different lead times.

for 2 LT and 3 LT. A similar conclusion can be drawn for the log-sinh transformation: a low value of α and a high value of β yield a multiplicative behaviour that is frequently chosen, for all lead times, but less for 2 LT and 3 LT than for LT / 2 and LT. This explains in particular the loss of reliability that can be seen for the log transformation for LT 3 in Fig. 6. These results reveal that the extrapolation behaviour of the residuals distributions is complex. It varies among catchments and with lead time because of the strong impact of data assimilation.

3.2 Results on the D3 control data set

3.2.1 Reliability

First, we conducted a visual inspection of the PIT diagrams, which convey an evaluation of the overall reliability of the probabilistic forecasts (examples in Fig. 9). In some cases, the forecasts are reliable (*e.g.*, the Thérain River at Beauvais, 755 km²; except if no transformation is used). Alternatively, these diagrams may show quite different patterns, highlighting bias (*e.g.*, the Meu River at Montfort-sur-Meu, 477 km²) or under-dispersion (*e.g.*, the Aa River at Wizernes, 392 km², or the Sioulet River at Miremont, 473 km²; for the latter, the calibration on D2_{sup} leads to log-sinh and Box-Cox transformations equivalent to no transformation, which turns out not to be relevant on the control data set where the log and the Box-Cox transformations are more reliable).

Then the distribution of the α -index values in Fig. 10 reveals a significant loss of reliability compared to the values obtained on the calibration data set (Fig. 6). We note that the log transformation is the most reliable approach for LT / 2 and is comparable to the Box-Cox transformation ($BC_{\lambda=0.2}$) for LT. With increasing lead time, the $BC_{\lambda=0.2}$ transformation becomes slightly better than the other transformations, including the calibrated ones. In addition, comparing the results obtained for the different lead times confirms that it is challenging to produce reliable predictive distributions when extrapolating at longer lead times.

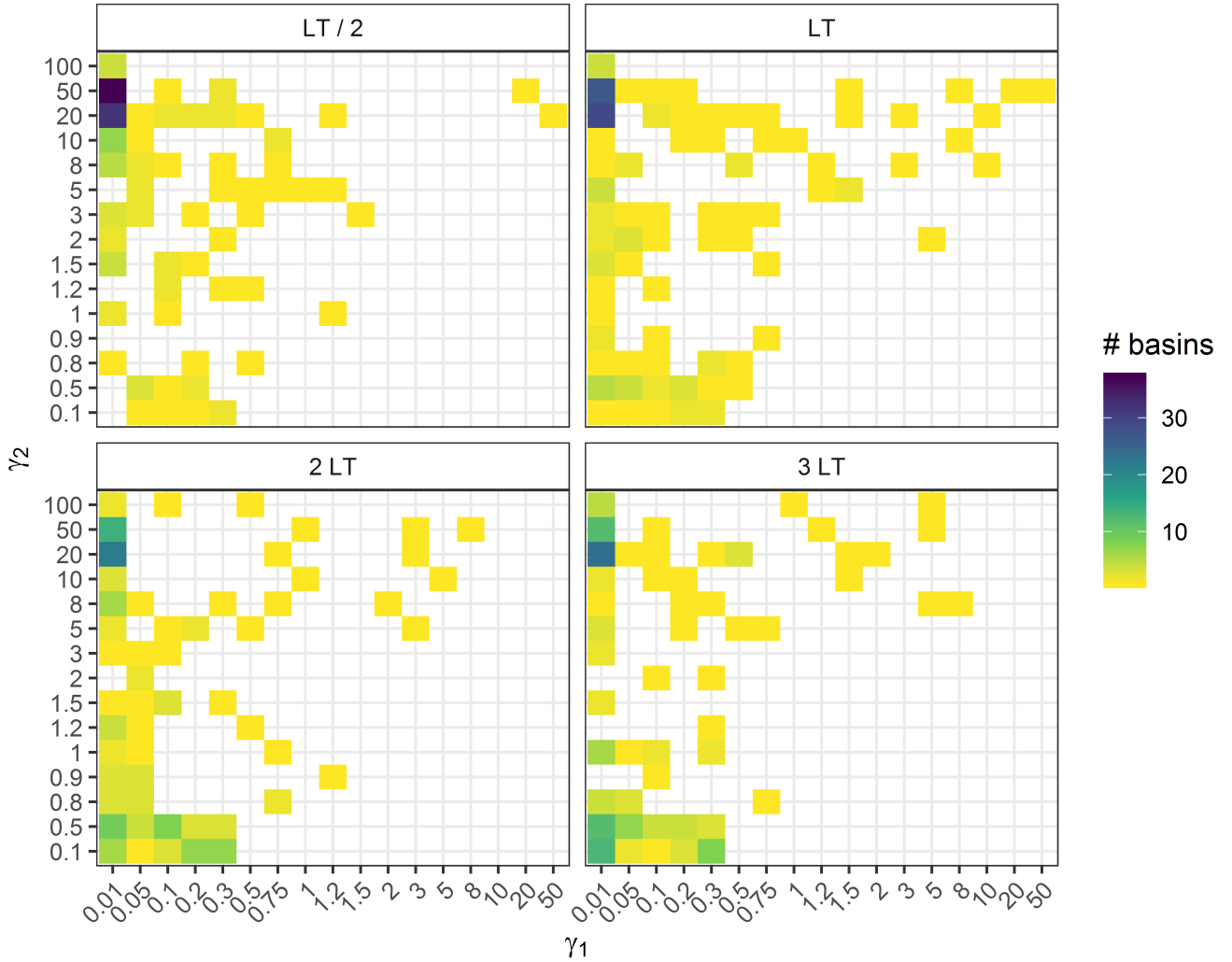


Figure 8. Distribution over the basins of the values of the log-sinh transformation parameters obtained during the calibration step for the four different lead times. $\gamma_1 = \alpha / \max_{D^2}(\tilde{Q})$ and $\gamma_2 = \beta / \max_{D^2}(\tilde{Q})$ (see Appendix B).

Overall, it means that the added value of the flexibility brought by the calibrated transformations is not transferable in an independent extrapolation situation, as illustrated in Fig. 11.

In operational settings, non-exceedance frequencies of the quantiles of the predictive distribution which are the lower and upper bounds of the predictive interval communicated to the authorities are of particular interest. The 80%-predictive interval (bounded by the 0.1 and 0.9 quantiles) is mostly used in France. It is expected that the non-exceedance frequency of the lower bound and the exceedance frequency of the upper bound remain close to 10% for a reliable predictive distribution. Deviations from these frequencies indicates biases in the estimated quantiles. Figure 12 reveals that on average the 0.1 quantile

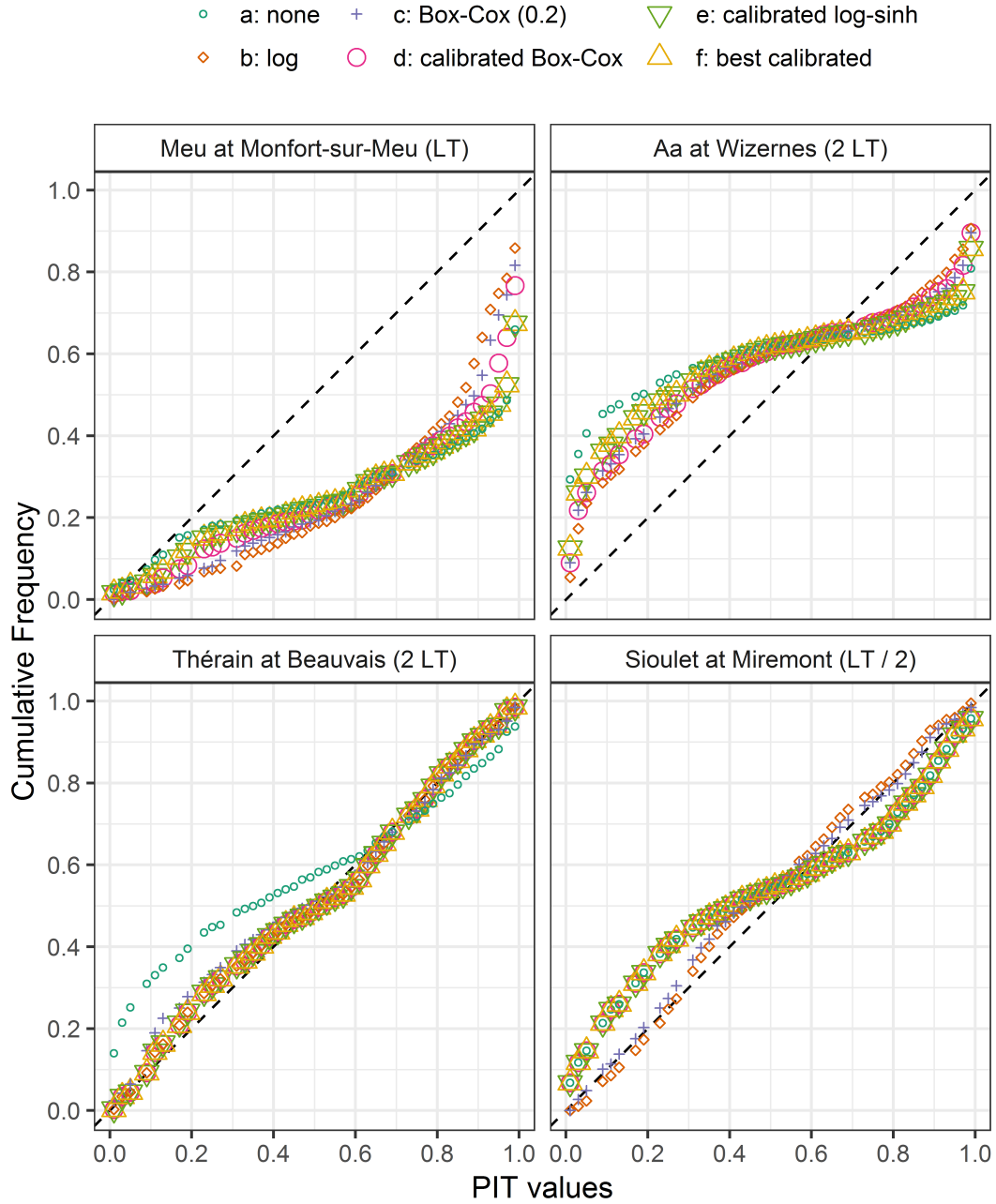


Figure 9. Examples of PIT diagrams obtained on the control data set D3, with different transformations at four locations.

is generally better assessed than the 0.9 quantile on average, though the latter is generally more sought for operational purposes. More importantly, the lack of reliability of the log transformation for the 3-LT lead time seen in Fig. 10 appears to be related

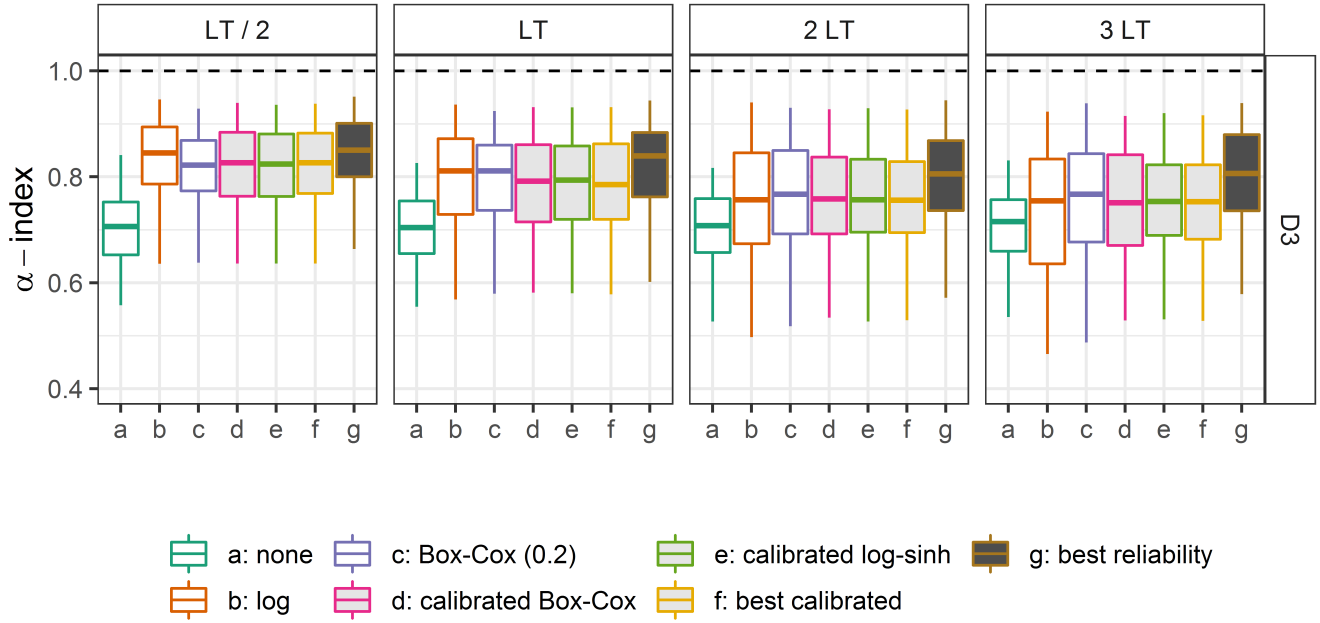


Figure 10. Distributions of the α -index values on the control data set D3, obtained with different transformations for four lead times (the filled boxplots represent the calibrated distributions). Boxplots (5th, 25th, 50th, 75th and 95th percentiles) synthesise the variety of scores over the catchments of the data set. The optimal values are represented by the horizontal dashed lines. Option "g" gives the best performance that could be achieved with this model and this post-processor for these catchments (see the discussion in Sect. 4.1).

to an underestimation of the 0.1 quantile, which is higher than for the other tested transformations, while the 0.9 quantile is less underestimated than for the other transformations. These results highlight that reliability can have different complementary facets and that some parts of the predictive distributions can be more or less reliable. In a context of flood forecasting, particular attention should be given to the upper part of the predictive distribution.

5 3.2.2 Other performance metrics

In addition to reliability, we looked at other qualities of the probabilistic forecasts, namely the overall performance (measured by the CRPSS) and accuracy ([measured by NSE](#)). We also checked their sharpness ([relative sharpness metric](#)). The distributions of four performance criteria are showed for lead time LT in Fig. 13. We note that the log transformation has the closest median value for the coverage ratio, at the expense of a lower median relative sharpness value, because of larger predictive interval widths caused by the multiplicative behaviour of the log transformation. In addition, the CRPSS and the NSE distributions have limited sensitivity to the variable transformation (also shown by Woldemeskel et al. (2018) for the CRPS), even if we can see that not using any transformation yields slightly better results. This confirms that the CRPSS itself is not sufficient to evaluate the adequacy of uncertainty estimation. Similar results were obtained for the other lead times (Supplementary materials).

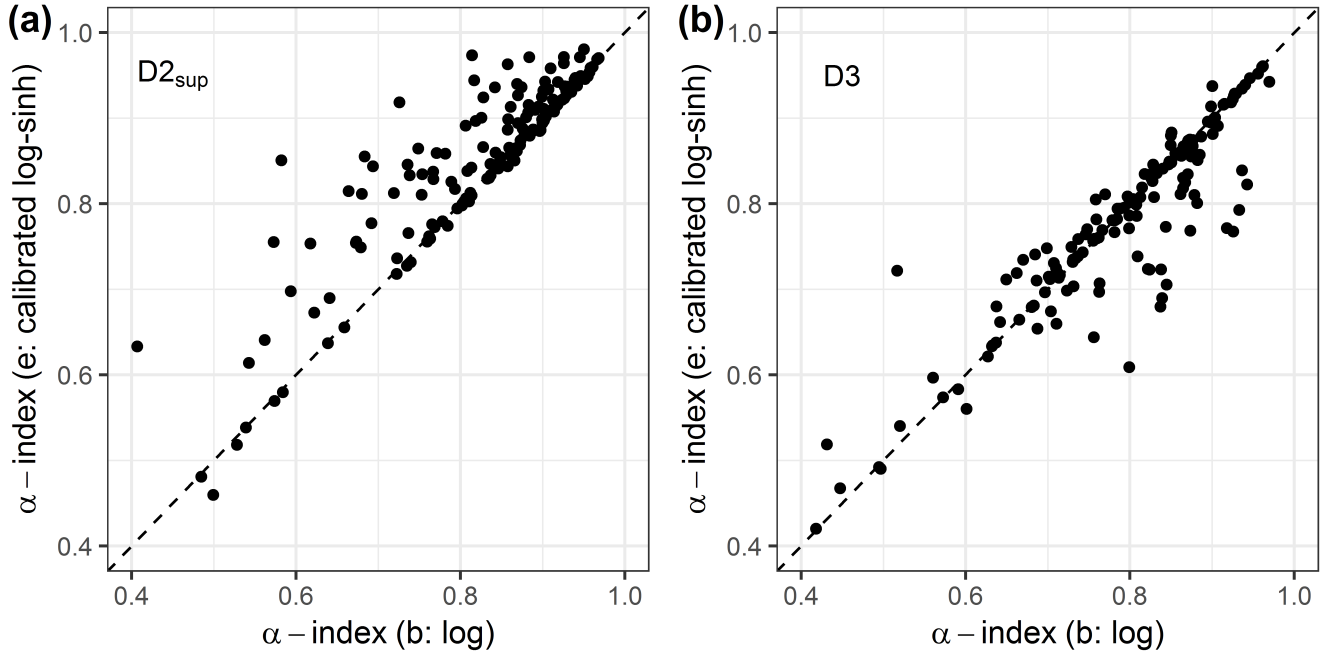


Figure 11. Scatter plots of the reliability α -index obtained with the log transformation and the log-sinh transformation: (a) on $D2_{sup}$ in the calibration step and (b) on $D3$ in the control step.

3.3 Investigating the performance loss in an extrapolation context

For operational forecasters, it is important to be able to predict when they can trust the forecasts issued by their models and when their quality becomes questionable. Therefore we investigated whether the reliability and reliability loss observed in an extrapolation context were correlated with some properties of the forecasts. First, Fig. 14 shows the relationship between the α -index values obtained in $D2_{sup}$ and those obtained in $D3$ for three representative transformations. The results indicate that it is not possible to anticipate the α -index values when extrapolating high flows in $D3$ based on the α -index values obtained when extrapolating high flows in $D2_{sup}$.

In addition, two indices were chosen to describe the degree of extrapolation: the ratio of the median of the forecasted discharges on $D3$ over the median of the forecasted discharge on $D2_{sup}$, and the ratio of the median of the forecasted discharges on $D3$ over the discharge for a return period of 20 years (for catchments where the assessment of the vicennial discharge was available in the national database: <http://www.hydro.eaufrance.fr>). In both cases, no trend appears, regardless of the variable transformation used, with Spearman coefficients values (much) lower than 0.33. The reliability can remain high for some catchments even when the magnitude of the events of the control data set are much higher than those of the training data set (see Supplementary materials for figures).

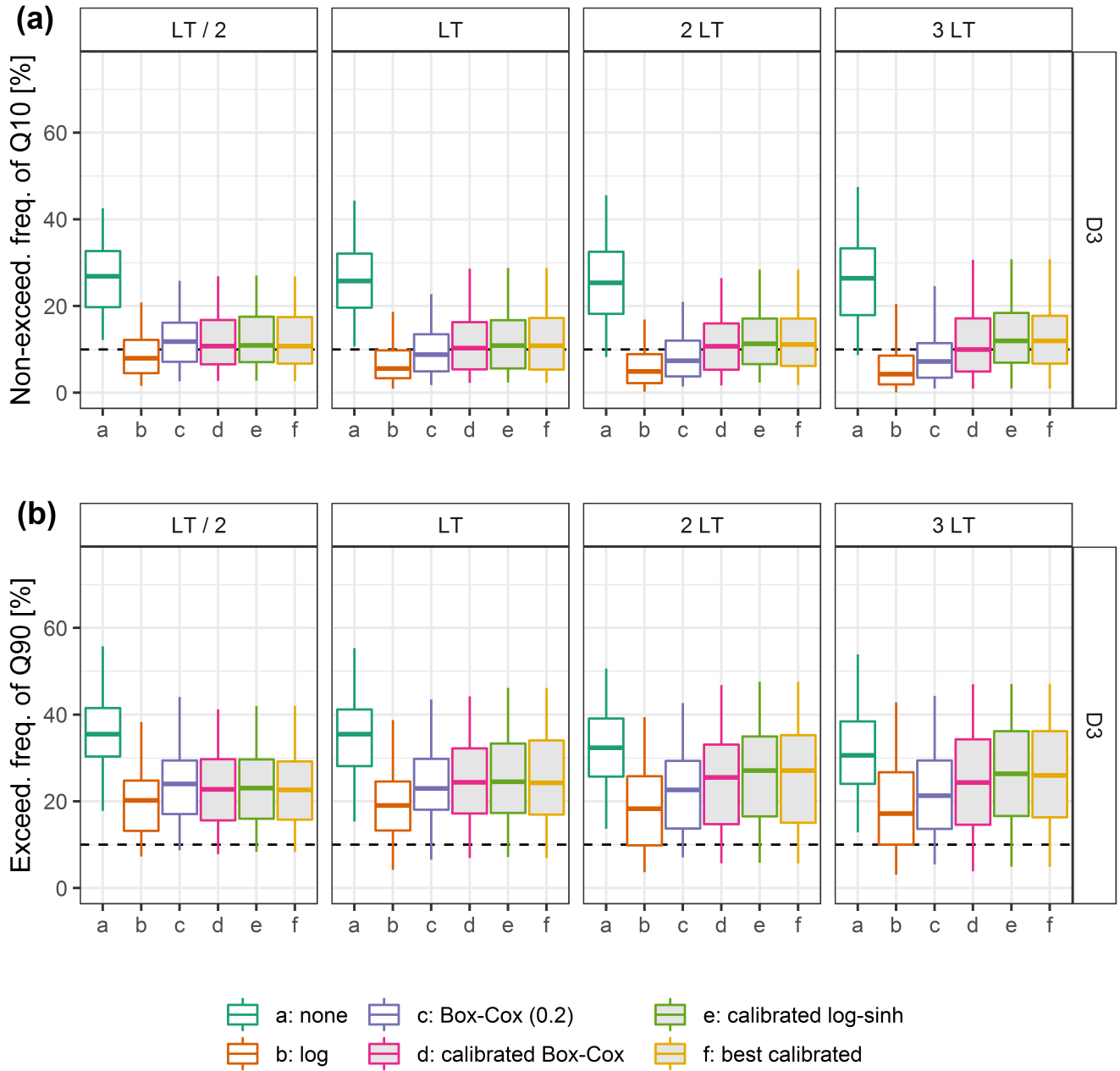


Figure 12. Distributions over the catchment set of a) the non-exceedance frequency of the 0.1 quantile and b) the exceedance frequency of the 0.9 quantile on the control data set D3, obtained with the different transformations tested (the filled boxplots are related to calibrated transformations). The optimal values are represented by the horizontal dashed lines.

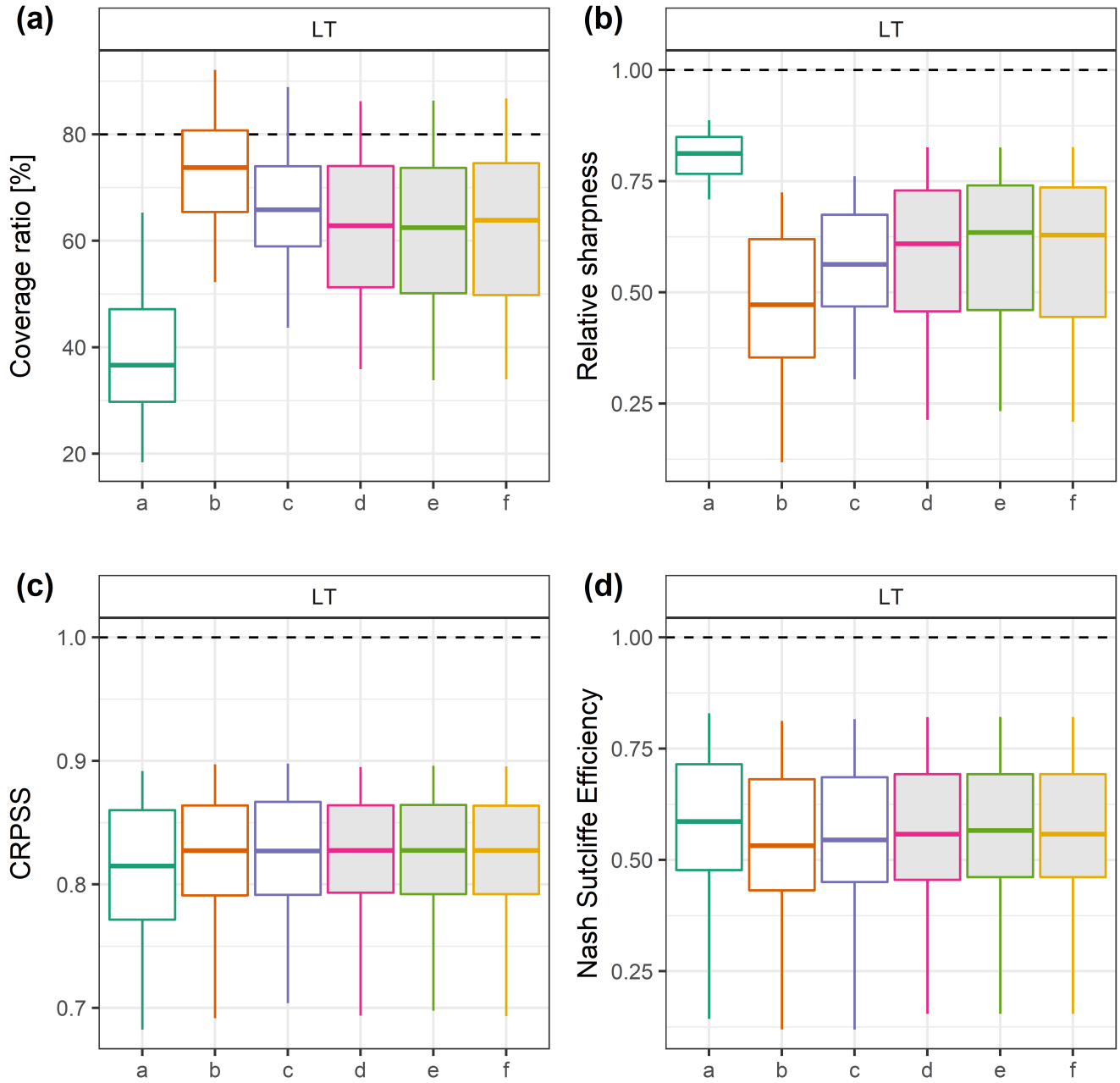


Figure 13. Distributions of coverage rate, relative sharpness, CRPSS and NSE values over the catchment set on the control data set D3, obtained with the different transformations tested (the filled boxplots are related to calibrated transformations). The optimal values are represented by the horizontal dashed lines.

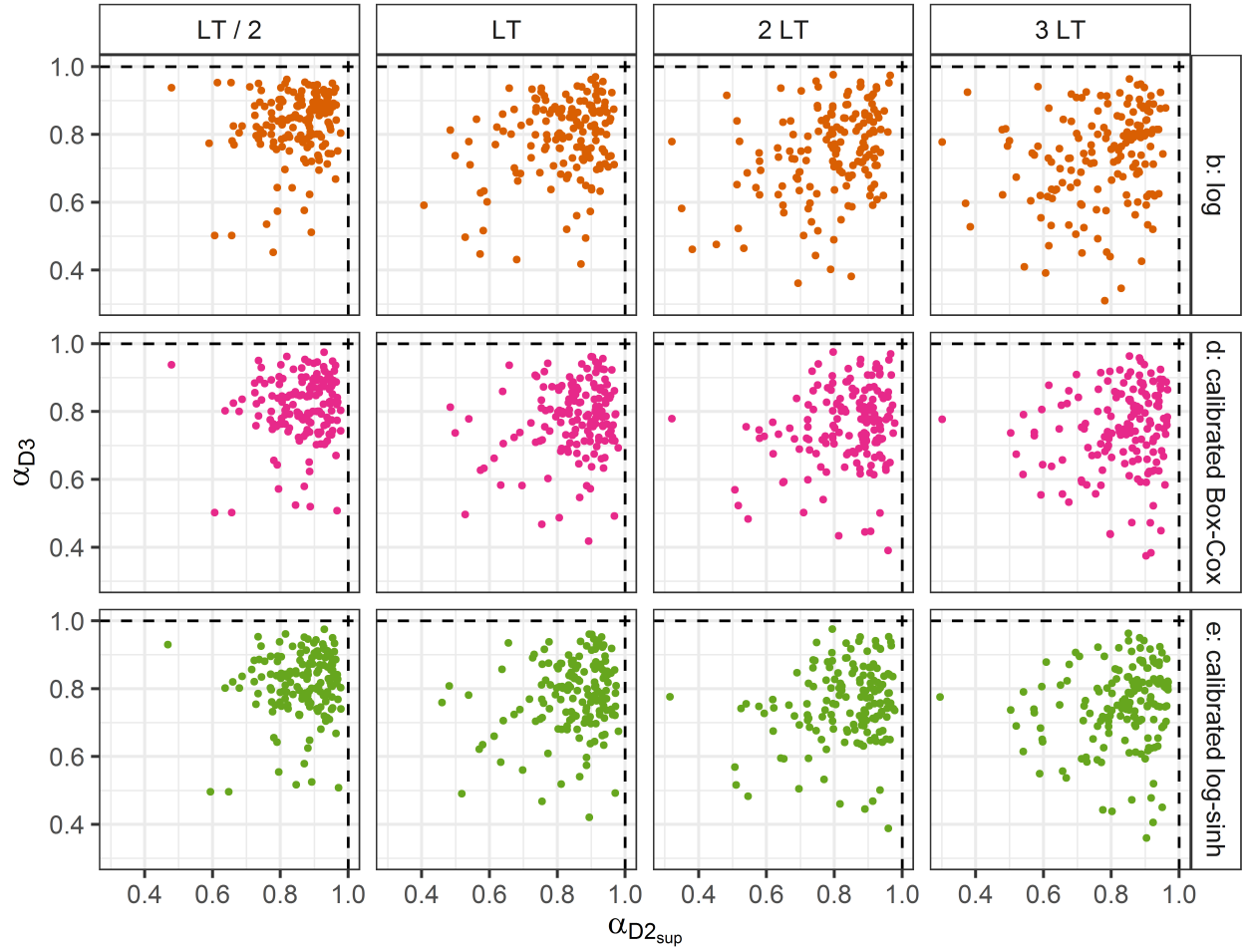


Figure 14. Comparison of the α -index values obtained in $D2_{sup}$ and $D3$. One point for each catchment. The optimal values are represented by the dashed lines. Similar results were obtained for the three other transformations (not shown).

Finally, we found no correlation with the relative accuracy of the deterministic forecasts either. The goodness-of-fit during the calibration phase cannot be used as an indicator of the robustness of the uncertainty estimation in an extrapolation context (see Supplementary materials for figures).

4 Discussion

4.1 Do more complex parametric transformations yield better results in an extrapolation context?

Overall, the results obtained on the control data set suggest that the log transformation and the fixed Box-Cox transformation ($BC_{\lambda=0.2}$) can yield relatively satisfactory α -index and coverage ratio values given their multiplicative or near-multiplicative behaviour in extrapolation. More tapered behaviours that can be obtained with the calibrated Box-Cox or log-sinh transformations do not show advantages when extrapolating high flows on an independent data set. In other words, what is learnt during the calibration of the more complex parametric transformations does not yield better results in an extrapolation context.

These results could be explained by the fact that the calibration did not result in the optimally relevant parameter set. To investigate whether another calibration strategy could yield better results, we compared the performance on the D3 data set, when the calibration is achieved on $D2_{sup}$ ("f: best calibrated") or on D3 ("g: best reliability"). The results shown in Fig. 10 reveal that even when the best parameter set is chosen among the 217 possibilities tested in this study (17 for the Box-Cox and 200 for the log-sinh), the α -index distributions are far from perfect and reliability decreases with increasing lead time. It suggests that the stability of the residuals distributions when extrapolating high flows might be a greater issue than the choice of the variable transformation. Nonetheless, the gap between the distributions of the ~~non-parametric transformations~~ transformations without calibrated parameter ("b" and "c"), the best calibrated transformation ("f") and the best performance that could be achieved ("g") highlights that it might be possible to obtain better results with a more advanced calibration strategy. This is, however, beyond the scope of this study and is therefore left for further investigations.

4.2 Empirical-based versus distribution-based approaches : does the distribution shape choice impact the uncertainty assessment in an extrapolation context?

Besides the reduction of heteroscedasticity, many studies use post-processors which are explicitly based on the assumption of a Gaussian distribution and use data transformations to fulfil this hypothesis (Li et al., 2017). Examples are the MCP or the meta-Gaussian model; the NQT was designed to precisely achieve it. In their study on autoregressive error models used as post-processors, Morawietz et al. (2011) showed that error models with an empirical distribution for the description of the standardized residuals perform better than those with a normal distribution. We first checked whether the variable transformation helped to reach a Gaussian distribution of the residuals computed with the transformed variables. Then we investigated whether better performance can be achieved using empirical transformed residuals distributions or using Gaussian distributions calibrated on these empirical distributions.

We used the Shapiro-Francia test where the null hypothesis is that the data is normally distributed. For each parametric transformation, we selected the parameter set of the calibration grid which obtains the highest p-value. For more than 98% of the catchments, the p-value is lower than 0.018 (respectively, 0.023) when the Box-Cox transformation (respectively, the log-sinh transformation) is used. This indicates that there are only a few catchments for which the normality assumption is

not to be rejected. In a nutshell, the variable transformations can stabilise the variance, but they do not necessarily ensure the normality of the residuals. It is important not to overlook this frequently encountered issue in hydrological studies.

Even if there is no theoretical advantage to using the Gaussian distribution calibrated on the transformed-variable residuals rather than the empirical distribution to assess the predictive uncertainty, we tested the impact of this choice. For each transformation, the predictive uncertainty assessment obtained with the empirical transformed-variable residuals distribution is compared to the assessment based on the Gaussian distribution whose mean and variance are those of the empirical distribution. Figure 15 shows the α -index distributions obtained over the catchments for both options on the control data set D3. We note that no clear conclusion can be drawn. No transformation (or identity transformation), which does not reduce the heteroscedasticity at all, benefits from the use of the Gaussian distributions for all lead times. In contrast, the predictive uncertainty assessment based on the empirical distribution with the log transformation is more reliable than the one based on the Gaussian distribution. For short lead times, it is slightly better to use the empirical distributions for the calibrated transformations (Box-Cox and log-sinh), but we observe a different behaviour for longer lead times. For these longer lead times, assessing the predictive uncertainty by the Gaussian distribution fitted on the empirical distributions of transformed residuals obtained with the calibrated log-sinh or Box-Cox transformations is the most reliable option. It is better than using the log transformation with the empirical distribution, but not very different from using the $BC_{\lambda=0.2}$ transformation.

Investigations on the impact of the choice between the empirical or the Gaussian distributions on the post processor performance are shown in Supplementary materials. They show that the choice of the distribution is not the dominant factor.

4.3 A need for a focus change?

In most modelling studies, several methodological steps depend on the range of the observations. First, calibration is designed to limit the residual errors in the available historical data. However the largest residuals are often associated with the highest discharge values. It is well known that removing the largest flood events from a data set can significantly modify the resulting calibrated parameter set. This is particularly true with the use of some common criteria such as quadratic criteria which strongly emphasise the largest errors (Legates and McCabe Jr., 1999; Berthet et al., 2010; Wright et al., 2015). Conversely, it is likely that “unavailable data” such as a physically realistic but (so far) unseen flood would significantly change the calibration results if it could be included in the calibration data set. Moreover, model conceptualisation (building) itself is often based on the understanding of how a catchment behaves “on average”. In some studies, outliers may even be considered as disturbing and be discarded (Liano, 1996).

However, to provide robust models for operational purposes, we also need to focus on rare (rarely observed) events, still keeping in mind all the well-known issues associated with working with (too) few data (Anctil et al., 2004; Perrin et al., 2007). For predictive uncertainty assessment, this issue is exacerbated by the seasonality of hydrological extremes (Allamano et al., 2011; Li et al., 2013) for most approaches, which rely heavily on data (beyond data-learning approaches, all models which need to be calibrated). Therefore, there is an urgent need to gather and compile data on extreme events (Gaume et al., 2009). Nevertheless, operational forecasters must still prepare themselves to work in an extrapolation context, as pointed out by Andréassian et al. (2012).

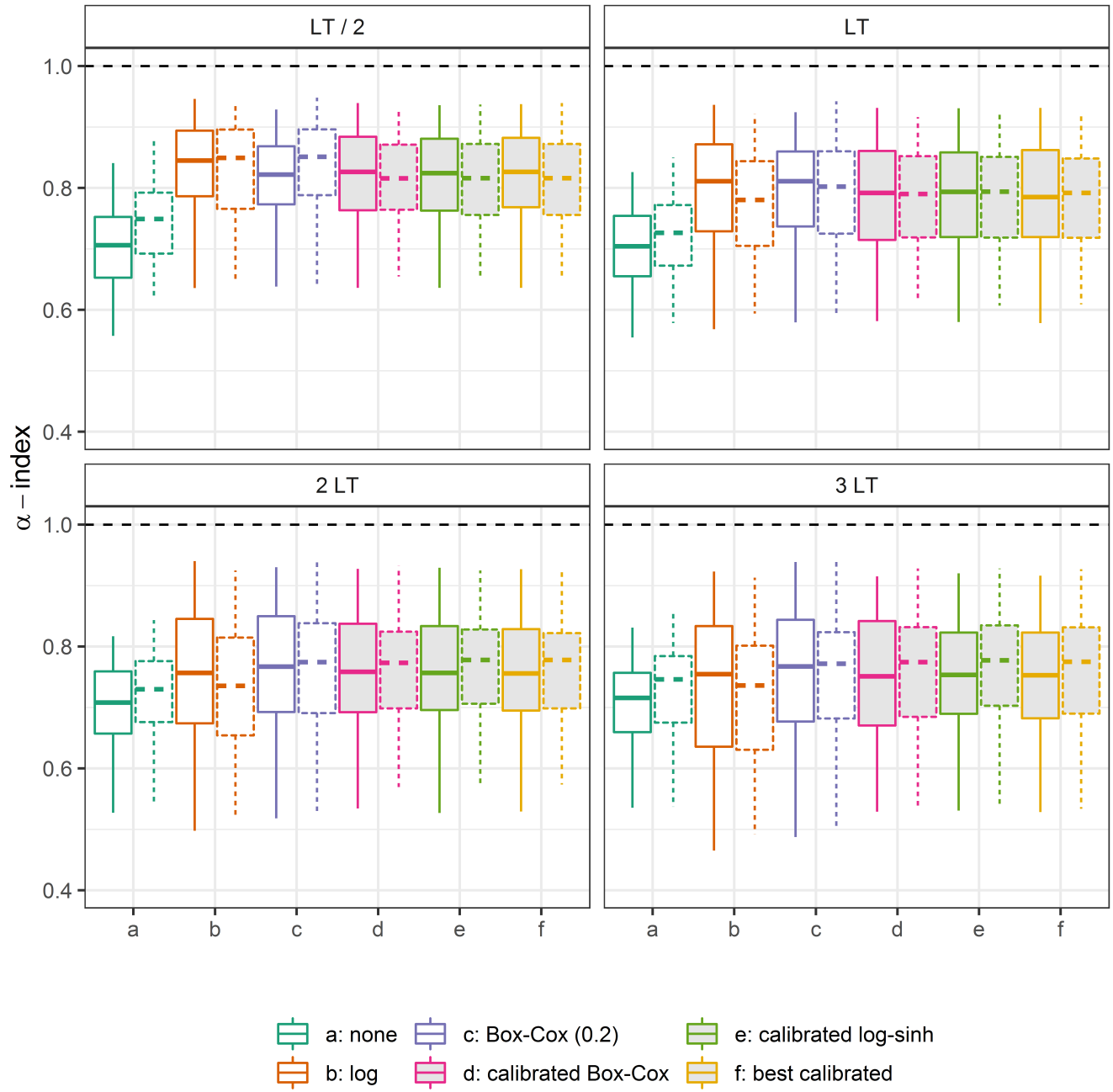


Figure 15. Distributions of the α -index values on the control data set D3, obtained with different transformations for four lead times, when using the empirical residuals distributions (straight boxplots) and the Gaussian distributions (dashed boxplots). The optimal values are represented by the horizontal dashed lines.

5 Conclusions

Even if major floods are rare, it is of the utmost importance that the forecasts issued during such events are reliable to facilitate an efficient crisis management. Like Lieutenant Drogo in the Tartar Steppe who spent his entire life fulfilling his day-to-day duties, but waiting in his fortress for the invasion by foes (Buzzati, 1940), many forecasters expect and are preparing for a major event, even if their routine involves only minor events. That is why a strong concern for the extrapolation context should be encouraged in all modelling and calibration steps. This article proposes a control framework focusing on the forecasting performance in an extrapolation context.

We use this framework to test the predictive uncertainty assessment using a statistical post-processing of a rainfall-runoff model, based on a variable transformation. The latter has to handle the heteroscedasticity and the evolution of the other predictive uncertainty distribution properties with the discharge magnitude, which is very problematic in an extrapolation context to issue reliable uncertainty assessment. As pointed out by McInerney et al. (2017), the choice of the heteroscedastic error modelling approach makes a significant impact on the predictive performance. This is true as well in an extrapolation context.

5.1 Main findings

Using the proposed framework for an evaluation in an extrapolation context, we showed that:

- (a) Using an appropriate variable transformation can significantly improve the predictive distribution and its reliability. However, a performance loss still remains in an extrapolation context with any of the three transformations we tested.
- (b) The ~~more-parametric-transformations~~ transformations with more calibrated parameters do not achieve significantly better results than the ~~non-parametric-transformations~~ transformations with no calibrated parameter:
 - while it allows a flexibility which can theoretically be very attractive in an extrapolation context, the log-sinh transformation is not more reliable in such a context;
 - the ~~non-parametric-not calibrated~~ log transformation and Box-Cox transformation with the λ parameter set at 0.2 are robust and compare favourably.
- (c) We did not find any variable significantly correlated with the performance loss in an extrapolation context.

The findings reported herein corroborate the results of McInerney et al. (2017) within the context of flood forecasting and extrapolation: calibrating the Box-Cox or log-sinh transformation can be counter-productive. We therefore suggest that operational flood forecasting services could consider the less flexible but more robust options: using the log transformation or the Box-Cox transformation with its λ parameter set closed to 0.2.

Importantly, these results reveal significant performance losses on some catchments when it comes to extrapolation, whatever variable transformation is used. Even if the scheme tested yields satisfying results in terms of reliability for the majority of catchments, it fails on a significant number of catchments and further investigations are needed to gain a deeper understanding of when and why failures occur.

5.2 Limitations and perspectives

We used the framework designed by Krzysztofowicz et al. and already applied in various studies, which separates the input uncertainty (mainly the observed and forecasted rainfall) and the hydrological uncertainty. Our study focuses only on the ‘effect’ of the extrapolation degree in the hydrological uncertainty when using the best available rainfall product. Future works should combine both input uncertainty (rainfall) and hydrological uncertainty (Bourgin et al., 2014) to evaluate the impact of using uncertain (forecasted) rainfall in a forecasting context.

Though no variable was found to be correlated to the performance loss, the investigations should be continued using a wider set of variables. First, it may open new perspectives to explain these losses and improve our understanding of the flaws of the hydrological model and of the EHUP. Furthermore, it would be very useful to help operational forecasters to detect the hydrological situations for which their forecasts have to be questioned (in particular during major events when forecasts are made in an extrapolation context).

Furthermore, improving the calibration strategy and using a regionalisation of the predictive distribution assessment, as proposed in Bourgin et al. (2015) and Bock et al. (2018), could help build more robust assessment of uncertainty quantification when forecasting high flows.

Finally, more studies focusing on the extrapolation context may help to better understand the limitations of the modelling (hydrological model structure, calibration, post-processing, etc.) and their consequences for practical matters. It is to be encouraged as a key for better and more reliable flood forecasting.

Appendix A: GRP model

The GRP model belongs to the suite of GR models (Michel, 1983). These models are designed to be as simple as possible but efficient for hydrological studies and for various operational needs, resulting in parsimonious model structures (<https://webgr.irstea.fr/en/models/a-brief-history/>, last access: 1 April 2019). The GRP model is designed for forecasting purposes (Berthet, 2010). It is a deterministic continuous lumped storage-type model (Fig. A1). The inputs are limited to areal rainfall and (interannual) potential evapotranspiration (both data may be available in real time). It can be understood as the sequence of two hydrological functions:

- a production function which is the same as in the well-known GR4J model developed by Perrin et al. (2003);
- a routing function which is a simplified version of the GR4J's routing function, since it only counts one flow branch composed of a unit hydrograph and quadratic routing store. The tests showed that the performance of the GRP and GR4J structures was similar in a forecasting mode.

A snow module (Valéry et al., 2014) may be implemented on top of the model if necessary.

Like any GR model, it is parsimonious. It has only three parameters: (a) an adjustment factor of effective rainfall which contributes to finding a good water balance, (b) the unit hydrograph time base used to account for the time lag between rainfall and streamflow and (c) the capacity of the routing store, which temporally smooths the output signal.

Its main difference with the other GR models is the implementation “in the loop” of two data assimilation schemes:

- a state updating procedure which modifies the main state of the routing function as a function of the last discharge values;
- an output updating based on an autoregressive model of the multiplicative error or an artificial neural network (multi-layer perception) whose inputs are the last discharge value and the two last forecasting errors. In this study, the autoregressive model was used.

The parameters are calibrated in forecasting mode, i.e., with the application of the updating procedures. This model is used by the French Flood Forecasting Services, some hydroelectricity suppliers and canal managers at the hourly time step in order to issue real-time forecasts for lead times ranging from a few hours to a few days at several hundred sites. Recently, Ficchi et al. (2016) pave the way to a multi-time step GRP version.

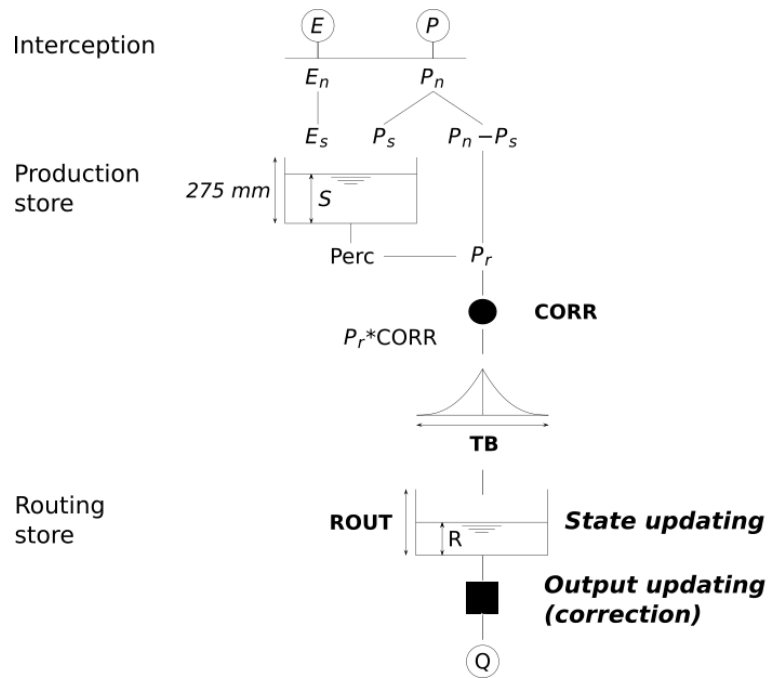


Figure A1. The GRP model flow-chart. After an interception step, the production function splits the net rainfall (P_n), according to the level of the production store. The effective rainfall (P_r) is the sum of the direct flow and the percolation from this store. A corrective multiplicative coefficient (CORR) is then applied. Then the flow runs through a unit hydrograph (time base: TB) and the routing store (capacity: ROUT).

Table B1. Summary of the different behaviours of the log-sinh transformation.

| Cases | Behaviours |
|-------------------------------------|--|
| $\alpha > 3 \cdot \beta$ | Additive error model |
| $y > 3 \cdot \beta$ | Additive error model |
| $\alpha \ll y \ll \beta$ | Multiplicative error model |
| $y \ll \alpha \ll \beta$ | Additive error model |
| Otherwise if $\alpha + y \ll \beta$ | Multiplicative error model (with an additive constant) |

Appendix B: The log-sinh transformation behaviours

B1 Log-sinh transformation formulations

In this study, we used the formulations of the log-sinh transformation chosen by Del Giudice et al. (2013):

$$g_{3 [\alpha, \beta]} : y \mapsto \beta \cdot \log \left(\sinh \left(\frac{\alpha + y}{\beta} \right) \right) \quad (\text{B1})$$

5 It is strictly equivalent to the formulation used by McInerney et al. (2017):

$$g'_{3 [a, b]} : y \mapsto \frac{1}{b} \log (\sinh (a + b \cdot y)) \quad (\text{B2})$$

B2 Different behaviours

Depending on the relative values of α and β and on the value range of y , compared to α and β , the log-sinh transformation can be reduced to an affine transformation (i.e. $g_{3 [\alpha, \beta]}(y + \varepsilon) - g_{3 [\alpha, \beta]}(y) = \delta \cdot \varepsilon$) or to the log transformation. The former case is
10 equivalent to no transformation (identity function; additive error model), whereas the latter one is equivalent to a multiplicative error model (Table B1).

When $y \ll \beta$,

$$\begin{aligned} \sinh \left(\frac{\alpha + y}{\beta} \right) &= \sinh \left(\frac{\alpha}{\beta} \right) \cdot \cosh \left(\frac{y}{\beta} \right) + \\ &\quad \sinh \left(\frac{y}{\beta} \right) \cdot \cosh \left(\frac{\alpha}{\beta} \right) \\ 15 \quad &\simeq \sinh \left(\frac{\alpha}{\beta} \right) \cdot \left(1 + \frac{y^2}{2 \cdot \beta^2} + \dots \right) + \\ &\quad \cosh \left(\frac{\alpha}{\beta} \right) \cdot \left(\frac{y}{\beta} + \dots \right) \end{aligned}$$

Thus,

$$g_{3 [\alpha, \beta]}(y) \simeq \beta \cdot \log \left(\sinh \left(\frac{\alpha}{\beta} \right) + \cosh \left(\frac{\alpha}{\beta} \right) \cdot \frac{y}{\beta} \right)$$

When $\alpha \ll \beta$, the latter results in:

$$g_{3 [\alpha, \beta]}(y) \simeq \beta \cdot \log \left(\frac{\alpha + y}{\beta} \right) \quad (\text{B3})$$

B2.1 Cases where the log-sinh transformation is equivalent to an affine transformation

If $x > 3$, $\sinh(x) \approx e^x/2$, then $\log(\sinh(x)) = x - \log(2)$. Therefore when $z = (\alpha + y)/\beta > 3$, the log-sinh transformation is
 5 equivalent to an affine transformation. In such cases, $g_{3 [\alpha, \beta]}^{-1}(g_{3 [\alpha, \beta]}(y) + \varepsilon) = y + \varepsilon$.

This happens when

- $\alpha > 3 \cdot \beta$;
- the β value is chosen large enough so that for any y value in the discharge range, $y > 3 \cdot \beta$.

Moreover, when $y \ll \alpha \ll \beta$, (B3) gives:

$$\begin{aligned} 10 \quad g_{3 [\alpha, \beta]} &\simeq \beta \cdot \left[\log \left(\frac{\alpha + y}{\alpha} \right) + \log(\alpha) - \log(\beta) \right] \\ &\simeq \frac{\beta \cdot y}{\alpha} + \beta \cdot (\log(\alpha) - \log(\beta)) \end{aligned}$$

The log-sinh transformation is then equivalent to an affine transformation.

B2.2 Cases where the log-sinh transformation is equivalent to a log transformation

As pointed out by McInerney et al. (2017), when $\alpha \ll y \ll \beta$, (B3) gives:

$$g_{3 [\alpha, \beta]}(y) \simeq \beta \cdot \log \left(\frac{y}{\beta} \right)$$

The log-sinh transformation is then equivalent to a mere log transformation.

15 B3 Calibration

The α and β parameters are in the same physical dimension as the y variable. Since this study is dedicated to the extrapolation context, we used the following adimensional parameterization to calibrate the variable transformation on various catchments. α and β are compared to the maximum forecasted discharge on the $D2_{\text{sup}}$ data subset:

$$\alpha = \gamma_1 \cdot \max_{D2}(\tilde{Q}) \quad \text{and} \quad \beta = \gamma_2 \cdot \max_{D2}(\tilde{Q}) \quad (\text{B4})$$

20 In the calibration step, the parameter space is explored on a (γ_1, γ_2) grid: 18 values of γ_1 from 0.01 to 100 and 15 values of γ_2 from 0.1 to 100 were tested, excluding combined values leading to the very same behaviours, such as $\gamma_1 \gg 3 \cdot \gamma_2$, equivalent to no transformation (additive error model). A total of 200 (γ_1, γ_2) combinations were tested for each calibration.

Author contributions. All co-authors contributed to the study and the article. Lionel Berthet and François Bourgin contributed equally and share co-first authorship.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. The authors thank [the editor Prof. Dimitri Solomatine](#), Dr. Kolbjorn Engeland and one anonymous reviewer. Both
5 [reviewers](#) provided insightful comments which helped to greatly improve this text. Thanks are also due to Météo-France for providing the meteorological data and Banque HYDRO data base for the hydrological data. The contribution of the authors from Univ Gustave Eiffel and INRAE was financially supported by SCHAPI (Ministry of the Environment).

References

- Abaza, M., Anctil, F., Fortin, V., and Perreault, L.: On the incidence of meteorological and hydrological processors: Effect of resolution, sharpness and reliability of hydrological ensemble forecasts, *Journal of Hydrology*, 555, 371–384, <https://doi.org/10.1016/j.jhydrol.2017.10.038>, 2017.
- 5 Abbaszadeh, P., Moradkhani, H., and Yan, H.: Enhancing hydrologic data assimilation by evolutionary Particle Filter and Markov Chain Monte Carlo, *Advances in Water Resources*, 111, 192–204, <https://doi.org/10.1016/j.advwatres.2017.11.011>, 2018.
- Allamano, P., Laio, F., and Claps, P.: Effects of disregarding seasonality on the distribution of hydrological extremes, *Hydrology and Earth System Sciences*, 15, 3207–3215, <https://doi.org/10.5194/hess-15-3207-2011>, 2011.
- Anctil, F., Perrin, C., and Andréassian, V.: Impact of the length of observed records on the performance of ANN and of conceptual parsimonious rainfall-runoff forecasting models, *Environmental Modelling & Software*, 19, 357 – 368, 2004.
- 10 Andréassian, V., Hall, A., Chahinian, N., and Schaake, J.: Introduction and Synthesis: Why should hydrologists work on a large number of basin data sets?, *IAHS-AISH Publication*, n° 307, 1–5, 2006.
- Andréassian, V., Perrin, C., Berthet, L., Le Moine, N., Lerat, J., Loumagne, C., Oudin, L., Mathevet, T., Ramos, M. H., and Valery, A.: HESS Opinions 'Crash tests for a standardized evaluation of hydrological models', *Hydrology and Earth System Sciences*, 13, 1757–1764, 2009.
- 15 Andréassian, V., Le Moine, N., Perrin, C., Ramos, M.-H., Oudin, L., Mathevet, T., Lerat, J., and Berthet, L.: All that glitters is not gold: the case of calibrating hydrological models, *Hydrological Processes*, 26, 2206–2210, <https://doi.org/10.1002/hyp.9264>, 2012.
- Barbetta, S., Coccia, G., Moramarco, T., Brocca, L., and Todini, E.: The multi temporal/multi-model approach to predictive uncertainty assessment in real-time flood forecasting, *Journal of Hydrology*, 551, 555–576, <https://doi.org/10.1016/j.jhydrol.2017.06.030>, 2017.
- Bennett, J. C., Robertson, D. E., Shrestha, D. L., Wang, Q., Enever, D., Hapuarachchi, P., and Tuteja, N. K.: A System for Continuous Hydrological Ensemble Forecasting (SCHEF) to lead times of 9 days, *Journal of Hydrology*, 519, Part D, 2832–2846, <https://doi.org/10.1016/j.jhydrol.2014.08.010>, 2014.
- 20 Berthet, L.: Flood forecasting at the hourly time-step: for a better assimilation of flow information in hydrological modelling (in French), Ph.D. thesis, AgroParisTech (Paris), Irstea (Antony), Doctoral School GRN, 2010.
- Berthet, L. and Potté, O.: International survey for good practices in forecasting uncertainty assessment and communication, in: *EGU General Assembly*, vol. 16, pp. EGU2014–8579, 2014.
- 25 Berthet, L., Andreassian, V., Perrin, C., and Javelle, P.: How crucial is it to account for the antecedent moisture conditions in flood forecasting? Comparison of event-based and continuous approaches on 178 catchments, *Hydrology and Earth System Sciences*, 13, 819–831, 2009.
- Berthet, L., Andréassian, V., Perrin, C., and Loumagne, C.: How significant are quadratic criteria? Part 2. On the relative contribution of large flood events to the value of a quadratic criterion, *Hydrological Sciences Journal*, 55, 1063–1073, <https://doi.org/10.1080/02626667.2010.505891>, 2010.
- 30 Bock, A. R., Farmer, W. H., and Hay, L. E.: Quantifying uncertainty in simulated streamflow and runoff from a continental-scale monthly water balance model, *Advances in Water Resources*, 122, 166 – 175, <https://doi.org/10.1016/j.advwatres.2018.10.005>, 2018.
- Bogner, K., Pappenberger, F., and Cloke, H. L.: Technical Note: The normal quantile transformation and its application in a flood forecasting system, *Hydrology and Earth System Sciences*, 16, 1085–1094, <https://doi.org/10.5194/hess-16-1085-2012>, 2012.
- 35 Bourgin, F.: How to quantify predictive uncertainty in hydrological modelling? Exploratory work on a large sample of catchments (in French), Ph.D. thesis, AgroParisTech (Paris), Irstea (Antony), Doctoral School GRNE, 2014.

- Bourgin, F., Ramos, M., Thirel, G., and Andréassian, V.: Investigating the interactions between data assimilation and post-processing in hydrological ensemble forecasting, *Journal of Hydrology*, 519, 2775–2784, <https://doi.org/10.1016/j.jhydrol.2014.07.054>, 2014.
- Bourgin, F., Andréassian, V., Perrin, C., and Oudin, L.: Transferring global uncertainty estimates from gauged to ungauged catchments, *Hydrology and Earth System Sciences*, 19, 2535–2546, <https://doi.org/10.5194/hess-19-2535-2015>, 2015.
- 5 Box, G. E. P. and Cox, D. R.: An Analysis of Transformations, *Journal of the Royal Statistical Society. Series B (Methodological)*, 26, 211–252, 1964.
- Breiman, L.: Statistical modeling: The two cultures, *Statistical Sciences*, 16, 199–215, <https://doi.org/10.1214/ss/1009213726>, 2001.
- Bremnes, J. B.: Constrained Quantile Regression Splines for Ensemble Postprocessing, *Monthly Weather Review*, 147, 1769–1780, <https://doi.org/10.1175/MWR-D-18-0420.1>, <https://doi.org/10.1175/MWR-D-18-0420.1>, 2019.
- 10 Brigode, P., Oudin, L., and Perrin, C.: Hydrological model parameter instability: A source of additional uncertainty in estimating the hydrological impacts of climate change?, *Journal of Hydrology*, 476, 410–425, <https://doi.org/10.1016/j.jhydrol.2012.11.012>, 2013.
- Buzzati, D.: *Il deserto dei Tartari (The Tartar Steppe)*, Rizzoli, Milano, 1940.
- Cigizoglu, H.: Estimation, forecasting and extrapolation of river flows by artificial neural networks, *Hydrological Sciences Journal*, 48, 349–362, 2003.
- 15 Coccia, G. and Todini, E.: Recent developments in predictive uncertainty assessment based on the model conditional processor approach, *Hydrology and Earth System Sciences*, 15, 3253–3274, <https://doi.org/10.5194/hess-15-3253-2011>, 2011.
- Coron, L., Andreassian, V., Perrin, C., Lerat, J., Vaze, J., Bourqui, M., and Hendrickx, F.: Crash testing hydrological models in contrasted climate conditions: An experiment on 216 Australian catchments, *Water Resources Research*, 48, W05552, <https://doi.org/10.1029/2011wr011721>, 2012.
- 20 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P., and Rieckermann, J.: Improving uncertainty estimation in urban hydrological modeling by statistically describing bias, *Hydrology and Earth System Sciences*, 17, 4209–4225, <https://doi.org/10.5194/hess-17-4209-2013>, 2013.
- Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D.-J., Hartman, R., Herr, H. D., Fresch, M., Schaake, J., and Zhu, Y.: The Science of NOAA’s Operational Hydrologic Ensemble Forecast Service, *Bulletin of the American Meteorological Society*, 95, 79–98, <https://doi.org/10.1175/BAMS-D-12-00081.1>, 2014.
- 25 Demeritt, D., Cloke, H., Pappenberger, F., Thielen, J., Bartholmes, J., and Ramos, M.-H.: Ensemble predictions and perceptions of risk, uncertainty, and error in flood forecasting, *Environmental Hazards*, 7, 115–127, 2007.
- Dogulu, N., López López, P., Solomatine, D. P., Weerts, A. H., and Shrestha, D. L.: Estimation of predictive hydrologic uncertainty using the quantile regression and UNEEC methods and their comparison on contrasting catchments, *Hydrology and Earth System Sciences*, 19, 3181–3201, <https://doi.org/10.5194/hess-19-3181-2015>, 2015.
- 30 Ficchi, A., Perrin, C., and Andréassian, V.: Impact of temporal resolution of inputs on hydrological model performance: An analysis based on 2400 flood events, *Journal of Hydrology*, 538, 454–470, <https://doi.org/10.1016/j.jhydrol.2016.04.016>, 2016.
- Furusho, C., Perrin, C., Viatgé, J., R., L., and Andréassian, V.: Collaborative work between operational forecasters and scientists for better flood forecasts, *La Houille Blanche*, pp. 5–10 (in French), <https://doi.org/10.1051/lhb/2016033>, 2016.
- 35 Gaume, E., Bain, V., Bernardara, P., Newinger, O., Barbuc, M., Bateman, A., Blaskovicová, L., Blöschl, G., Borga, M., Dumitrescu, A., Daliakopoulos, I., Garcia, J., Irimescu, A., Kohnova, S., Koutroulis, A., Marchi, L., Matreata, S., Medina, V., Preciso, E., Sempere-Torres, D., Stancalie, G., Szolgay, J., Tsanis, I., Velasco, D., and Viglione, A.: A compilation of data on European flash floods, *Journal of Hydrology*, 367, 70 – 78, <https://doi.org/10.1016/j.jhydrol.2008.12.028>, 2009.

- Giustolisi, O. and Laucelli, D.: Improving generalization of artificial neural networks in rainfall-runoff modelling, *Hydrological Sciences Journal*, 50, 439–457, 2005.
- Gneiting, T., Balabdaoui, F., and Raftery, A. E.: Probabilistic forecasts, calibration and sharpness, *Journal of the Royal Statistical Society Series B-Statistical Methodology*, 69, 243–268, <https://doi.org/10.1111/j.1467-9868.2007.00587.x>, 2007.
- 5 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, *Hydrology and Earth System Sciences*, 18, 463–477, <https://doi.org/10.5194/hess-18-463-2014>, 2014.
- Hemri, S., Lisniak, D., and Klein, B.: Multivariate postprocessing techniques for probabilistic hydrological forecasting, *Water Resources Research*, 51, 7436–7451, <https://doi.org/10.1002/2014WR016473>, 2015.
- Hersbach, H.: Decomposition of the continuous ranked probability score for ensemble prediction systems, *Weather and Forecasting*, 15, 559–570, [https://doi.org/10.1175/1520-0434\(2000\)015<0559:dotcrp>2.0.co;2](https://doi.org/10.1175/1520-0434(2000)015<0559:dotcrp>2.0.co;2), 2000.
- 10 Imrie, C., Durucan, S., and Korre, A.: River flow prediction using artificial neural networks: Generalisation beyond the calibration range, *Journal of Hydrology*, 233, 138–153, 2000.
- Kelly, K. S. and Krzysztofowicz, R.: A bivariate meta-Gaussian density for use in hydrology, *Stoch. Hydrol. Hydraul.*, 11, 17–31, 1997.
- Klemeš, V.: Operational testing of hydrological simulation models, *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, 31, 13–24, <https://doi.org/10.1080/02626668609491024>, 1986.
- 15 Krzysztofowicz, R.: Bayesian theory of probabilistic forecasting via deterministic hydrologic model, *Water Resources Research*, 35, 2739–2750, 1999.
- Krzysztofowicz, R. and Maranzano, C. J.: Hydrologic uncertainty processor for probabilistic stage transition forecasting, *Journal of Hydrology*, 293, 57–73, <https://doi.org/10.1016/j.jhydrol.2004.01.003>, 2004.
- 20 Laio, F. and Tamea, S.: Verification tools for probabilistic forecasts of continuous hydrological variables, *Hydrology and Earth System Sciences*, 11, 1267–1277, 2007.
- Lang, M., Pobanz, K., Renard, B., Renouf, E., and Sauquet, E.: Extrapolation of rating curves by hydraulic modelling, with application to flood frequency analysis, *Hydrological Sciences Journal*, 55, 883–898, <https://doi.org/10.1080/02626667.2010.504186>, 2010.
- Legates, D. and McCabe Jr., G.: Evaluating the use of 'goodness-of-fit' measures in hydrologic and hydroclimatic model validation, *Water Resources Research*, 35, 233–241, 1999.
- 25 Li, M., Wang, Q. J., and Bennett, J.: Accounting for seasonal dependence in hydrological model errors and prediction uncertainty, *Water Resources Research*, 49, 5913–5929, <https://doi.org/10.1002/wrcr.20445>, 2013.
- Li, W., Duan, Q., Miao, C., Ye, A., Gong, W., and Di, Z.: A review on statistical postprocessing methods for hydrometeorological ensemble forecasting, *Wiley Interdisciplinary Reviews: Water*, 4, e1246, <https://doi.org/10.1002/wat2.1246>, 2017.
- 30 Liano, K.: Robust error measure for supervised neural network learning with outliers, *IEEE Transactions on Neural Networks*, 7, 246–250, 1996.
- Lobligeois, F., Andréassian, V., Perrin, C., Tabary, P., and Loumagne, C.: When does higher spatial resolution rainfall information improve streamflow simulation? An evaluation using 3620 flood events, *Hydrology and Earth System Sciences*, 18, 575–594, <https://doi.org/10.5194/hess-18-575-2014>, 2014.
- 35 McInerney, D., Thyer, M., Kavetski, D., Lerat, J., and Kuczera, G.: Improving probabilistic prediction of daily streamflow by identifying Pareto optimal approaches for modeling heteroscedastic residual errors, *Water Resources Research*, 53, 2199–2239, <https://doi.org/10.1002/2016WR019168>, 2017.

- Merz, R., Parajka, J., and Blöschl, G.: Time stability of catchment model parameters: Implications for climate impact analyses, *Water Resources Research*, 47, W02 531, <https://doi.org/10.1029/2010WR009505>, 2011.
- Michel, C.: Que peut-on faire en hydrologie avec un modèle conceptuel à un seul paramètre?, *La Houille blanche*, pp. 39–44 (in French), 1983.
- 5 Montanari, A.: Uncertainty of Hydrological Predictions, in: *Treatise on Water Science*, edited by Wilderer, P., pp. 459–478, Elsevier, Oxford, 2011.
- Montanari, A. and Grossi, G.: Estimating the uncertainty of hydrological forecasts: A statistical approach, *Water Resources Research*, 44, W00B08, <https://doi.org/10.1029/2008wr006897>, 2008.
- Moradkhani, H., Hsu, K. L., Gupta, H., and Sorooshian, S.: Uncertainty assessment of hydrologic model states and parameters: Sequential data assimilation using the particle filter, *Water Resources Research*, 41, W05 012, <https://doi.org/10.1029/2004wr003604>, 2005a.
- 10 Moradkhani, H., Sorooshian, S., Gupta, H. V., and Houser, P. R.: Dual state–parameter estimation of hydrological models using ensemble Kalman filter, *Advances in Water Resources*, 28, 135–147, <https://doi.org/10.1016/j.advwatres.2004.09.002>, 2005b.
- Morawietz, M., Xu, C.-Y., Gottschalk, L., and Tallaksen, L. M.: Systematic evaluation of autoregressive error models as post-processors for a probabilistic streamflow forecast system, *Journal of Hydrology*, 407, 58–72, <https://doi.org/10.1016/j.jhydrol.2011.07.007>, 2011.
- 15 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andreassian, V., Anctil, F., and Loumagne, C.: Which potential evapotranspiration input for a lumped rainfall-runoff model? Part 2 - Towards a simple and efficient potential evapotranspiration model for rainfall-runoff modelling, *Journal of Hydrology*, 303, 290–306, <https://doi.org/10.1016/j.jhydrol.2004.08.026>, 2005.
- Pagano, T. C., Shrestha, D. L., Wang, Q. J., Robertson, D., and Hapuarachchi, P.: Ensemble dressing for hydrological applications, *Hydrological Processes*, 27, 106–116, <https://doi.org/10.1002/hyp.9313>, 2013.
- 20 Pagano, T. C., Wood, A. W., Ramos, M.-H., Cloke, H. L., Pappenberger, F., Clark, M. P., Cranston, M., Kavetski, D., Mathevet, T., Sorooshian, S., and Verkade, J. S.: Challenges of Operational River Forecasting, *J. Hydrometeor.*, 15, 1692–1707, <https://doi.org/10.1175/JHM-D-13-0188.1>, 2014.
- Pappenberger, F. and Beven, K. J.: Ignorance is bliss: Or seven reasons not to use uncertainty analysis, *Water Resources Research*, 42, W05 302, <https://doi.org/10.1029/2005wr004820>, 2006.
- 25 Pappenberger, F., Thielen, J., and Del Medico, M.: The impact of weather forecast improvements on large scale hydrology: analysing a decade of forecasts of the European Flood Alert System, *Hydrological Processes*, 25, 1091–1113, <https://doi.org/10.1002/hyp.7772>, 2011.
- Pappenberger, F., Pagano, T. C., Brown, J. D., Alfieri, L., Lavers, D. A., Berthet, L., Bressand, F., Cloke, H. L., Cranston, M., Danhelka, J., Demargne, J., Demuth, N., de Saint-Aubin, C., Feikema, P. M., Fresch, M. A., Garçon, R., Gelfan, A., He, Y., Hu, Y. Z., Janet, B., Jurdy, N., Javelle, P., Kuchment, L., Laborda, Y., Langsholt, E., Le Lay, M., Li, Z. J., Mannesiez, F., Marchandise, A., Marty, R., Meißner, D.,
- 30 Manful, D., Organde, D., Pourret, V., Rademacher, S., Ramos, M. H., Reinbold, D., Tibaldi, S., Silvano, P., Salamon, P., Shin, D., Sorbet, C., Sprokkereef, E., Thiémig, V., Tuteja, N. K., van Andel, S. J., Verkade, J. S., Vehviläinen, B., Vogelbacher, A., Wetterhall, F., Zappa, M., Van der Zwan, R. E., and Thielen-del Pozo, J.: *Hydrological Ensemble Prediction Systems Around the Globe*, pp. 1–35, Springer Berlin Heidelberg, Berlin, Heidelberg, https://doi.org/10.1007/978-3-642-40457-3_47-1, 2016.
- Perrin, C., Oudin, L., Andréassian, V., Rojas-Serna, C., Michel, C., and Mathevet, T.: Impact of limited streamflow data on the efficiency and the parameters of rainfall-runoff models, *Hydrological Sciences Journal*, 52, 131–151, 2007.
- 35 Ramos, M.-H., Bartholmes, J., and Pozo, J. T.-d.: Development of decision support products based on ensemble forecasts in the European flood alert system, *Atmospheric Science Letters*, 8, 113–119, <https://doi.org/10.1002/asl.161>, 2007.

- Reichert, P. and Mieleitner, J.: Analyzing input and structural uncertainty of nonlinear dynamic models with stochastic, time-dependent parameters, *Water Resources Research*, 45, W10 402, <https://doi.org/10.1029/2009wr007814>, 2009.
- Renard, B., Kavetski, D., Kuczera, G., Thyer, M., and Franks, S. W.: Understanding predictive uncertainty in hydrologic modeling: The challenge of identifying input and structural errors, *Water Resources Research*, 46, W05 521, <https://doi.org/10.1029/2009wr008328>, 2010.
- Salamon, P. and Feyen, L.: Assessing parameter, precipitation, and predictive uncertainty in a distributed hydrological model using sequential data assimilation with the particle filter, *Journal of Hydrology*, 376, 428–442, <https://doi.org/10.1016/j.jhydrol.2009.07.051>, 2009.
- Schoups, G. and Vrugt, J. A.: A formal likelihood function for parameter and predictive inference of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors, *Water Resources Research*, 46, W10 531, <https://doi.org/10.1029/2009wr008933>, 2010.
- Seiller, G., Anctil, F., and Roy, R.: Design and experimentation of an empirical multistructure framework for accurate, sharp and reliable hydrological ensembles, *Journal of Hydrology*, 552, 313–340, <https://doi.org/10.1016/j.jhydrol.2017.07.002>, 2017.
- Sharma, S., Siddique, R., Reed, S., Ahnert, P., Mendoza, P., and Mejia, A.: Relative effects of statistical preprocessing and postprocessing on a regional hydrological ensemble prediction system, *Hydrology and Earth System Sciences*, 22, 1831–1849, <https://doi.org/10.5194/hess-22-1831-2018>, 2018.
- Singh, S. K., McMillan, H., and Bardossy, A.: Use of the data depth function to differentiate between case of interpolation and extrapolation in hydrological model prediction, *Journal of Hydrology*, 477, 213–228, <https://doi.org/10.1016/j.jhydrol.2012.11.034>, 2013.
- Solomatine, D. P. and Shrestha, D. L.: A novel method to estimate model uncertainty using machine learning techniques, *Water Resources Research*, 45, W00B11, <https://doi.org/10.1029/2008wr006839>, 2009.
- Tabary, P., Dupuy, P., L'Henaff, G., Gueguen, C., Moulin, L., Laurantin, O., Merlier, C., and Soubeyroux, J.-M.: A 10-year (1997-2006) reanalysis of Quantitative Precipitation Estimation over France: methodology and first results, vol. 351, pp. 255–260, IAHS, 2012.
- Thielen, J., Bartholmes, J., Ramos, M.-H., and de Roo, A.: The European Flood Alert System - Part 1: Concept and development, *Hydrology and Earth System Sciences*, 13, 125 – 140, 2009.
- Todini, E.: Role and treatment of uncertainty in real-time flood forecasting, *Hydrological Processes*, 18, 2743–2746, <https://doi.org/10.1002/hyp.5687>, 2004.
- Todini, E.: Hydrological catchment modelling: past, present and future, *Hydrology and Earth System Sciences*, 11, 468–482, 2007.
- Todini, E.: A model conditional processor to assess predictive uncertainty in flood forecasting, *International Journal of River Basin Management*, 6, 123–137, <https://doi.org/10.1080/15715124.2008.9635342>, 2008.
- Todini, E.: Predictive uncertainty assessment in real time flood forecasting, in: *Uncertainties in Environmental Modelling and Consequences for Policy Making*, edited by Baveye, P. C., Laba, M., and Mysiak, J., pp. 205–228, Springer Netherlands, Dordrecht, 2009.
- Valéry, A., Andréassian, V., and Perrin, C.: 'As simple as possible but not simpler': What is useful in a temperature-based snow-accounting routine? Part 2 – Sensitivity analysis of the Cemaneige snow accounting routine on 380 catchments, *Journal of Hydrology*, 517, 1176 – 1187, <https://doi.org/10.1016/j.jhydrol.2014.04.058>, 2014.
- Van Steenbergen, N., Ronsyn, J., and Willems, P.: A non-parametric data-based approach for probabilistic flood forecasting in support of uncertainty communication, *Environmental Modelling & Software*, 33, 92–105, <https://doi.org/10.1016/j.envsoft.2012.01.013>, 2012.
- Vaze, J., Post, D., Chiew, F., Perraud, J.-M., Viney, N., and Teng, J.: Climate non-stationarity – Validity of calibrated rainfall–runoff models for use in climate change studies, *Journal of Hydrology*, 394, 447–457, <https://doi.org/10.1016/j.jhydrol.2010.09.018>, 2010.

- Velazquez, J. A., Anctil, F., and Perrin, C.: Performance and reliability of multimodel hydrological ensemble simulations based on seventeen lumped models and a thousand catchments, *Hydrology and Earth System Sciences*, 14, 2303–2317, <https://doi.org/10.5194/hess-14-2303-2010>, 2010.
- Verkade, J., Brown, J., Davids, F., Reggiani, P., and Weerts, A.: Estimating predictive hydrological uncertainty by dressing deterministic and ensemble forecasts; a comparison, with application to Meuse and Rhine, *Journal of Hydrology*, 555, 257–277, <https://doi.org/10.1016/j.jhydrol.2017.10.024>, 2017.
- Verkade, J. S. and Werner, M. G. F.: Estimating the benefits of single value and probability forecasting for flood warning, *Hydrology and Earth System Sciences*, 15, 3751–3765, <https://doi.org/10.5194/hess-15-3751-2011>, 2011.
- Viatgé, J., Pinna, T., Perrin, C., Dorchie, D., and Garandeau, L.: Towards an enhanced temporal flexibility of the GRP flood forecasting operational model, in: *Proceedings of the SHF conference De la prévision des crues à la gestion de crise (From flood forecasting to crisis management)*, Avignon (France), November 14th - 16th, Société Hydrotechnique de France, 12 p. (in French), 2018.
- Viatgé, J., Berthet, L., Marty, R., Bourgin, F., Piotte, O., Ramos, M. H., and Perrin, C.: Towards the real-time production of predictive intervals around streamflow forecasts in Vigicrues in France, *LHB*, pp. 63–71, <https://doi.org/10.1051/lhb/2019016>, 2019.
- Wang, Q. J., Robertson, D. E., and Chiew, F. H. S.: A Bayesian joint probability modeling approach for seasonal forecasting of streamflows at multiple sites, *Water Resources Research*, 45, W05 407, <https://doi.org/10.1029/2008wr007355>, 2009.
- Wang, Q. J., Shrestha, D. L., Robertson, D. E., and Pokhrel, P.: A log-sinh transformation for data normalization and variance stabilization, *Water Resources Research*, 48, W05 514, <https://doi.org/10.1029/2011wr010973>, 2012.
- Wani, O., Beckers, J. V. L., Weerts, A. H., and Solomatine, D. P.: Residual uncertainty estimation using instance-based learning with applications to hydrologic forecasting, *Hydrology and Earth System Sciences*, 21, 4021–4036, <https://doi.org/10.5194/hess-21-4021-2017>, 2017.
- Weerts, A. H., Winsemius, H. C., and Verkade, J. S.: Estimation of predictive hydrological uncertainty using quantile regression: examples from the National Flood Forecasting System (England and Wales), *Hydrology and Earth System Sciences*, 15, 255–265, <https://doi.org/10.5194/hess-15-255-2011>, 2011.
- Wilby, R. L.: Uncertainty in water resource model parameters used for climate change impact assessment, *Hydrological Processes*, 19, 3201–3219, <https://doi.org/10.1002/hyp.5819>, 2005.
- Woldemeskel, F., McInerney, D., Lerat, J., Thyer, M., Kavetski, D., D. S., Tuteja, N., and Kuczera, G.: Evaluating post-processing approaches for monthly and seasonal streamflow forecasts, *Hydrology and Earth System Sciences*, 22, 6257–6278, <https://doi.org/10.5194/hess-22-6257-2018>, 2018.
- Wright, D. P., Thyer, M., and Westra, S.: Influential point detection diagnostics in the context of hydrological model calibration, *Journal of Hydrology*, 527, 1161 – 1172, <https://doi.org/10.1016/j.jhydrol.2015.05.047>, 2015.
- Yang, J., Reichert, P., and Abbaspour, K. C.: Bayesian uncertainty analysis in distributed hydrologic modeling: a case study in the Thur River basin (Switzerland), *Water Resources Research*, 43, W10 401, <https://doi.org/10.1029/2006wr005497>, 2007a.
- Yang, J., Reichert, P., Abbaspour, K. C., and Yang, H.: Hydrological modelling of the Chaohe Basin in China: Statistical model formulation and Bayesian inference, *Journal of Hydrology*, 340, 167–182, <https://doi.org/10.1016/j.jhydrol.2007.03.006>, 2007b.
- Zalachori, I., Ramos, M. H., Garçon, R., Mathevet, T., and Gailhard, J.: Statistical processing of forecasts for hydrological ensemble prediction: a comparative study of different bias correction strategies, *Advances in Sciences and Research*, 8, 135–141, <https://doi.org/10.5194/asr-8-135-2012>, 2012.