



A data-centric perspective on the information needed for hydrological uncertainty predictions

Andreas Auer¹, Martin Gauch², Frederik Kratzert³, Grey Nearing⁴, Sepp Hochreiter¹, and Daniel Klotz⁵

¹ELLIS Unit Linz and LIT AI Lab, Institute for Machine Learning, Johannes Kepler University Linz, Austria

²Google Research, Zurich, Switzerland

³Google Research, Vienna, Austria

⁴Google Research, Mountain View, California, USA

⁵Department of Compound Environmental Risks, Helmholtz Centre for Environmental Research — UFZ, Leipzig, Germany

Correspondence: Andreas Auer (auer@ml.jku.at)

Abstract. Uncertainty estimates are fundamental to assess the reliability of predictive models in hydrology. We use the framework of Conformal Prediction to investigate the impact of temporal and spatial information on uncertainty estimates within hydrological predictions. Integrating recent information significantly enhances overall uncertainty predictions, even with substantial gaps between updates. While local information yields good results on average, it proves insufficient for peak flow predictions. Incorporating global information improves the accuracy of peak flow bounds, corroborating findings from related studies. Overall, the study underscores the importance of continuous data updates and the integration of global information for robust and efficient uncertainty estimation.

1 Introduction

10 Uncertainty estimates are the basis for actionable predictions (e.g., Krzysztofowicz, 2001; Beven, 2016a). In this contribution we take a data-centric perspective to investigate how temporal (recency) and spatial (local vs. global) information impacts the quality of uncertainty estimates. With machine learning, we investigate on the level of information content with respect to different tasks, rather than on the level of process representations. For our study, we utilize Conformal Prediction (CP, Vovk et al. (2005)), a model agnostic framework that adds uncertainty intervals to existing predictions. In contrast to many ad-hoc
15 approaches for uncertainty estimation, CP is motivated in a rigorous, probabilistic manner. Under the right conditions, CP intervals will always achieve finite-sample marginal coverage (Vovk et al., 2005). Further, the model agnostic nature of CP enables the separation of a model's point prediction quality from its uncertainty prediction quality. This is, for example, not possible with the Deep Learning baselines from Klotz et al. (2022), which serve as a reference wherever adequate. Specifically, we apply HopCPT (Auer et al., 2023), the state-of-the-art CP model for time series. HopCPT memory-based architecture
20 determines which data points drive the uncertainty prediction during inference. The identification of data points that are relevant



for uncertainty, can be exploited to analyze the impact of selected information on the predictive outcomes in an indirect yet flexible way.

To our knowledge, we are the first to introduce CP to hydrology. Our goal is to use CP as a tool to empirically study the following research questions (RQ):

25 – **RQI:** *Does up to date information improve uncertainty estimation?*

We show that continual updates greatly benefit the general uncertainty prediction, even if the updating happens in batches (i.e. discontinuously).

– **RQII:** *Is the data from a given basin required to get good uncertainty estimates for that basin?*

30 Our results suggest that in the general case local information is indeed beneficial for uncertainty estimation, since it will generally result in tighter prediction intervals. In a PUB setting, while larger intervals are necessary, it's still possible to achieve quite good performance in terms of uncertainty prediction.

– **RQIII:** *Is data from a single basin enough to get good uncertainty predictions for peak flows?*

Our results indicate that it is necessary to use data from different basins to provide good uncertainty predictions for peak flows at a given basin.

35 **Trading space for time in hydrological modeling.** Multiscale Parameter Regionalization (MPR; Samaniego et al., 2017; Scheppe et al., 2022) is a technique that allows to calibrate hydrological models using a scale independent, global parametrization. MPR thus uses all available data to provide an estimation for the parameters of a given basin (the possible parametrization is, however, still strongly restricted by a priori knowledge in the form of the model structure and the functions that link the spatial information to the model parameters). Similarly, Kratzert et al. (2019b) introduced a Long Short-Term Memory (LSTM)
40 rainfall-runoff model that is globally parametrized. In terms of predictive fidelity, this LSTM based approach outperformed many classical rainfall-runoff models (e.g., Kratzert et al., 2019a, 2021; Mai et al., 2022). Klotz et al. (2022) showed how this LSTM based approach can directly provide predictive uncertainty estimations. An inspection of the importance of data for the LSTM based approach can be found in Gauch et al. (2021). They concluded that adding multiple basins (i.e., the spatial part of the data) is key for reaching good model performances. In time series prediction in general, Montero-Manso and Hyndman
45 (2021) found that global modelling approaches — such as the one discussed here — tend to outperform local ones. A different research direction with similar implication is the contribution by Bertola et al. (2023), who analyzed how floods from different regions are informative of each other. They show that many observed floods fall within the envelope values estimated from previous floods in other basins. This suggests that local flood predictions can benefit from information from different places. We are not aware of any publications that explicitly examine the space and time relationship for uncertainty estimations.

50

Uncertainty estimation in hydrology. Uncertainty has long been recognized as a crucial part of hydrological modelling (e.g., Krzysztofowicz, 2001; Beven, 2016b). Thus, there already exists a wide range of approaches for uncertainty estimation in hydrological modeling. As of today, approaches include — but are not limited to — ensemble-based methods, that define and



sample probability distributions around different model inputs, structures, or outputs (e.g., Li et al., 2017; Demargne et al.,
55 2014; Clark et al., 2016); Bayesian and Bayesian inspired techniques, which weight different parameters, models, or out-
comes (e.g., Kavetski et al., 2006; Beven and Binley, 2014); neural network based methods which estimate the parameters (of
mixtures) of probability distributions (Klotz et al., 2022, e.g.); and even explicit post-processing methods (e.g., Shrestha and
Solomatine, 2008; Montanari and Koutsoyiannis, 2012; Koutsoyiannis and Montanari, 2022). Good overviews can for example
be found in Nearing et al. (2016) or Gupta and Govindaraju (2023).

60

CP for time series. The current state of the art in conformal time series prediction is HopCPT, a CP approach based on deep
learning (Auer et al., 2023). To provide a prediction interval for a given basin and time step, HopCPT learns to retrieve his-
torical time steps that belong to similar regimes — i.e., time steps that had similar error patterns. Intuitively, the CP model
performs a soft nearest-neighbor search with a learned similarity measure. This leads to tighter and more accurate uncertainty
65 estimates than existing approaches, as it incorporates knowledge not just about the marginal distribution but also about the
current system state. HopCPT’s regime definition, which considers regime changes within a time series, is loosely related to
regimes in the hydrological modeling sense (Haines et al., 1988; Harris et al., 2000), which classify rivers according to the
overall flow behavior. The definitions of Quandt (1958) and Hamilton (1990), which model time series with multiple regimes
where the distribution parameters are conditional on the active regime, are closer to our actual use of the regime term.

70

The remainder of this paper is structured as follows: First, Section 2.1 provides an introduction to CP geared towards hydrolog-
ogists and time series prediction. Section 2.2 describes the methods relevant to this study (HopCPT, CMAL), and Section 2.3
describes the metrics used for comparing the different approaches. In Section 3 we present our experiments and corresponding
data. Section 4 details and discusses the results. Finally, we conclude in Section 5.

75 2 Methods

2.1 Conformal Prediction

This section provides a brief overview of CP. For a thorough introduction to Conformal Prediction, we refer the reader to
Angelopoulos and Bates (2021).

A CP procedure consists of two steps: First, CP estimates the “unusualness”, here called non-conformity, of data points
80 within a calibration set — which contains previous hold-out data that are not used for training the prediction model. Then, CP
uses this information directly to construct an uncertainty region that consists of the most “usual” values of the calibration set.

There does not exist a definitive function to measure non-conformity. As a matter of fact, there are infinitely many non-
conformity measures. Even a function that randomly assigns a value of an arbitrary distribution would allow the coverage
guarantee of CP to hold — as long as the distribution does not change between calibration and test samples. However, choosing
85 a non-conformity measure that yields good prediction intervals — in the sense that they are not too broad — is part of the
challenge when applying CP. Vovk et al. (2005), for instance, point out that whether any particular approach is an appropriate



way to measure non-conformity depends greatly on contextual factors. In a regression setting a straight-forward example for a non-conformity measure is the absolute error of the prediction.

Even the most basic CP approaches do not pose specific assumptions to the underlying data distribution (for comparison, CMAL, as proposed by Klotz et al. (2022), assumes a mixture of asymmetric Laplacians) except that the data are exchangeable (i.e., the joint distribution is invariant to permutations of the data). Hydrological tasks — such as streamflow prediction — typically violate this exchangeability assumption, since errors are highly correlate in time and exhibit different behavior for different situations — e.g., a model might have much larger errors in a flood situation than in a low-flow situation. On top of that, environmental processes, especially when considering long time periods, likely exhibit shifts in the data distribution. These can arise from a spectrum of factors, ranging from gradual changes, such as those induced by climate change, to more accelerated transformations, like those stemming from infrastructure projects. Since the prediction interval of CP is based on the calibration set, i.e., based on past observations, such shifts can lead to unreliable prediction intervals. Formally, we can also view this as a break from the exchangeability assumption. Besides that, standard CP generates prediction intervals that provide marginal coverage, which produces unnecessarily wide or too small prediction intervals when different error patterns exist. For example, in time series regression, given the absolute error as a simple non-conformity score, the prediction interval would have the same width over the whole time sequence. Recent advances in CP methods have tackled these inherent problems and thus made it possible to use CP for time series uncertainty estimation. Typically, they either adapt the calibration distribution based on the temporal proximity and/or the time series covariates (Auer et al., 2023; Foygel Barber et al., 2022; Xu and Xie, 2022a, b). Other variants propose an adaptive prediction interval (Gibbs and Candes, 2021; Zaffran et al., 2022; Bhatnagar et al., 2023) that operates in an online fashion and therefore needs access to the label measurements after prediction.

To be more explicit, we consider a setting where we have a calibration dataset $D = \{(x_1; y_2)(x_2, y_2), \dots, (x_n, y_n)\}$. Applying a prediction model $\mu : X \mapsto Y$ gives us the (absolute) errors of the calibration data $E = \{e_1, e_2, ..e_n\}$ — these errors represent the non-conformity scores as a higher error refers to a more “unusual” sample (Figure 1-left). Our goal is to create a prediction interval for a new sample $(x_{n+1}; y_{n+1})$ which covers the unknown error e_{n+1} with probability $1 - \alpha$, where α represents the specified miscoverage rate. The standard CP procedure is as follows: (1) We use the $1 - \alpha$ empirical quantile of E to estimate the score for which a $1 - \alpha$ ratio of the samples have a lower score — i.e., are less “unusual” (Figure 1-middle). (2) Because we assume the data is exchangeable, the non-conformity score of the new sample is represented by the same distribution. Therefore we can simply define the prediction by adding and subtracting the quantile score to the model prediction to arrive at the lower and upper bound respectively (Figure 1-right).

In a classical CP setting it is important that the calibration set is not part of the training data of the model μ . This is because the fit on the training data is biased: models can, for example, overfit. Therefore, the non-conformity score distribution of the training data will likely not generalize to new data. Also, many CP approaches assume that the model itself is already capable of providing uncertainty estimates. This is however not a necessity. Conformal prediction can be used in classification or regression settings for point, interval, and distributional predictions.

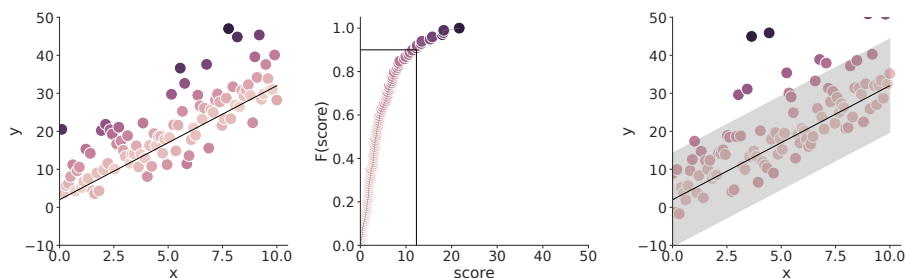


Figure 1. Illustration of standard conformal prediction applied to a regression setting. Left: The black line shows the prediction of a model. The colored points show the real observation — the darker the color the bigger the prediction residual, i.e., the non-conformity score. Middle: CDF of the non-conformity score distribution and the respective cut-off quantile at $\alpha = 0.1$. Right: The prediction of μ , with the CP interval defined by the value on the cut-off quantile, on the test data.

120 2.2 Uncertainty Estimation Models

This section introduces the models and metrics used to answer our research questions. We introduce HopCPT and its “global” variant, HopCPT-G, in Section 2.2.1 and Section 2.2.2 respectively. Section 2.2.3 presents CMAL.

2.2.1 HopCPT

Why are modeling systems often successfully used for decision making even if they do not provide uncertainty predictions?
125 We believe that this is because decision makers implicitly consider the model behavior over time and compare a given forecast to the recent performance and the model behavior in similar situations. We thus design HopCPT (Auer et al., 2023) to explicitly capture this notion in a quantitative way.

HopCPT is the current state-of-the-art approach for conformal time series prediction. The uncertainty of time series data often varies heavily between certain periods. One reason are seasonal patterns. For example, for many Alpine rivers, long-
130 lasting periods of low-flows in winter exhibit lower predictive uncertainty than large events in late spring where glacier melt and convectional precipitation events interact. Another reason is the occurrence of irregular but recurring events such as those induced by torrential rain. Note that shifts play an important role also within events, as their frequency and intensity can change over time. HopCPT addresses these challenges by viewing the time series as a soft partition of time periods, where each partition element exhibits individual uncertainty properties. We refer to such a set of time points as a regime (Quandt,
135 1958; Hamilton, 1990) and assume that one can identify them by the covariates and the lagged target of the time series. HopCPT learns to weight which past observations — i.e., calibration points — are likely from the same regime as the current point¹. Based on this information, the different calibration points are weighted differently when constructing the prediction interval. Since the considered time steps in the interval calculation are from the same regime, one can assume that exchangeability is

¹We want to emphasize that calibration data refers here — and throughout the manuscript — to data that is not utilized in training the underlying point prediction model. This aligns with the conventional terminology in probabilistic applications in general and the CP literature in specific.



given and therefore the validity of the prediction interval holds. Yet, since time points from unrelated regimes are disregarded,
140 HopCPT results in tighter intervals than, for example, standard CP.

More formally, HopCPT assigns a weight $a_{t,i}$ to a past time step $i \in M$ in the memory M given the current time step t . The weight reflects the similarity between the given time step i and the current time step t . Intuitively we can say that the weight should be high when i is from the same regime as t , and otherwise low. These weights are constructed with the Modern Hopfield Network component of HopCPT as follows

$$145 \quad a_{t,i} = \beta m(\mathbf{z}_t) \mathbf{W}_q \mathbf{W}_k m(\mathbf{z}_i), \quad (1)$$

$$\bar{a}_{t,i} = \frac{e^{a_{t,i}}}{\sum_{j \in M} e^{a_{t,j}}} \quad \forall i \in M, \quad (2)$$

where \mathbf{z}_i is the representation input for time step i , m represent an encoding module, and \mathbf{W}_q and \mathbf{W}_k are learned weight matrices; β is a hyperparameter that represents the inverse of the so-called softmax temperature, which adjusts the sharpness of softmax-generated probability distributions (low temperatures push all the attention to a single value; while high temperatures
150 distribute the attention uniformly).

The weights form the basis for the weighted conformal prediction interval (Foygel Barber et al., 2022). Besides the weighting, HopCPT deviates from standard CP as it does not consider the one-sided quantile of the absolute errors (non-conformity scores), but instead follows Xu and Xie (2022a) and excludes the lower and upper quantiles of the relative errors. Formally,² the prediction interval of time step t is calculated by

$$155 \quad \hat{C}^\alpha(\mathbf{z}_t, M) = \left[\hat{\mu}(\mathbf{X}_t) + Q_{\frac{\alpha}{2}} \left(\sum_{i \in M} \bar{a}_{t,i} \delta_{\epsilon_i} \right), \hat{\mu}(\mathbf{X}_t) + Q_{1-\frac{\alpha}{2}} \left(\sum_{i \in M} \bar{a}_{t,i} \delta_{\epsilon_i} \right) \right], \quad (3)$$

where Q_τ is the τ -quantile and δ_{ϵ_i} is a point mass distribution (i.e., a degenerate distribution where all the mass is concentrated at a single point) at the prediction error at time step i . This essentially means we compute the τ -quantile across a mixture of M distributions, each corresponding to a distinct time step i in the memory — characterized by a point mass at its prediction error ϵ_i . Thus, each of these distributions is concentrated at ϵ_i with density 1 and 0 elsewhere.

160 HopCPT retrieves the calibration data points during the prediction from the Hopfield Memory (Ramsauer et al., 2021) of the model. Thus, one can simply add every new and available observation to the memory. This corresponds to an automatic recalibration, which accounts for shifts in the data distribution. In addition, HopCPT allows to add a so-called temporal encoding to the time steps. This encoding adds information about the time difference between the predicted time step and the previous time steps. Given this information, HopCPT can learn to weight recent points higher, which further helps to address distribution
165 shifts.

HopCPT has already exhibited good performance in streamflow uncertainty prediction (Auer et al., 2023). In this work we use HopCPT approach as a tool to explore the effects of data availability on uncertainty prediction.

Memory update. In Auer et al. (2023), HopCPT updates the Hopfield Memory after each prediction with the — then —
170 previous prediction error. However, this requires access to the target label of this time step to calculate the error of the prediction

²Note that we use the alternative proposal presented in Appendix E of the original work as it is computationally more efficient

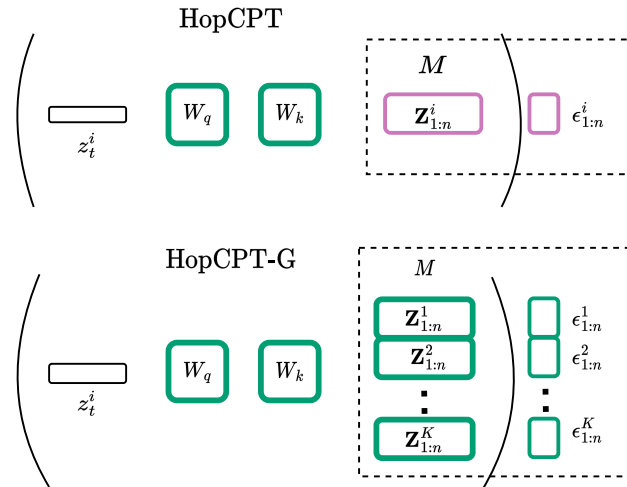


Figure 2. Illustration of the difference between HopCPT and HopCPT-G. The shared weights and vectors are outlined in green, the vectors individual to a specific time series in magenta. HopCPT-G also shares the memory vectors and incorporates all available time steps.

model. In streamflow modeling, we often do not have direct access to the target label (i.e., the streamflow measurement) before a new prediction has to be issued. Therefore, we adapt HopCPT to use a fixed memory that is only based on the calibration data by default. We refer to this adaptation as the “offline mode”, in contrast to the original “online mode”. Given the notation from Equation 3 this means that in the online mode the memory M contains all time steps up to time step $t - 1$, while in the
 175 offline mode, the memory is fixed per time series and only includes the calibration data.

In practice, one could expect that the memory may be updated intermittently, whenever new data becomes available. Hence, we explore the influence of different updating schemes (Experiment I).

Input features. Auer et al. (2023) concatenate (static and dynamic) time series covariates, model predictions, and lagged
 180 targets as input features to generate the representation z_i of a time step i . However, the use of the lagged target requires access to the target label right after the prediction, which is typically not the case in streamflow prediction. A straightforward solution would be to simply exclude the lagged target as an input feature. However, since most current rainfall–runoff models follow a state-space approach, we also explore how lack of labels can be compensated by using the model state. We hypothesize that this
 185 state should include most of the relevant information that would have been provided by the streamflow observations, as well as additional information that is not available to HopCPT in the original publication. Note that (a) the model states potentially contain more information than just the model prediction because the prediction is some projection of its state³ and (2) the state likely encodes important historical information beyond the current time step.

³Here we assume that the model state at time t already considers the input features of time step t .



2.2.2 HopCPT-G

We hypothesize that the union of the existing calibration data represents not only the included streamflow time series, but rather
 190 a general set of existing streamflow regimes. In this scenario one could utilize the error information from similar situations
 in other basins to best model the current situation in a specific basin. Further, the application to ungauged basins should be
 possible. In a simple case, where the error information is only relevant within a certain time series, HopCPT can learn to fall
 back to the local setting and only up-weight the calibration data of its own time series. To examine the potential of global
 memory, we modify HopCPT so that it operates on a global memory for both training and inference. We refer to this new
 195 model as HopCPT-G.

During inference, the association weights and prediction interval calculation of HopCPT-G is very similar to HopCPT.
 Equation 1 and 3 are still applicable, however, the memory M consists of all available time steps of *all* time series — in
 contrast to HopCPT, where only the time series of the predicted basin is used.

Figure 2 illustrates the difference. While in HopCPT only the learned weights are shared between different time steps, in
 200 HopCPT-G both the learned weights as well as the memory vectors are shared between different time series.

This change is also reflected in the training loss. N time steps are drawn without replacement from K randomly selected
 time series. The loss for this batch is then calculated as

$$\mathcal{L} = N^{-1} * \|(|\epsilon_{1:N}| - \mathbf{A}|\epsilon_{1:N}|)^2\|_1. \quad (4)$$

$$\mathbf{A}_{ij} \in \mathbb{R}^{N \times N} = \begin{cases} a_{ij} & \text{for } i \neq j \\ 0 & \text{else.} \end{cases} \quad (5)$$

205 To ensure that the samples of the own time series (which are likely the most relevant) are in the batch, we choose N and K
 in a way that all samples of the K series are in the batch.

PUB Training. It is crucial for HopCPT-G to learn a time step representation which captures rich information about the error
 regime. As shown in multiple works, machine learning models often suffer from shortcut learning (Lapuschkin et al., 2019;
 210 Geirhos et al., 2020). For HopCPT-G a potential shortcut might be to just learn to distinguish the different time series. This
 result would especially harm the PUB prediction performance, as it relies on representations which generalize over individual
 time series. To facilitate more robust representations that avoid this learning shortcut, we propose a training loss that masks out
 all observations of the time series to which the predicted time step belongs. I.e., HopCPT-G can only use the error observations
 from the other time series to form its prediction in the training phase. This way, the shortcut of only learning to consider the
 215 “own” time series is impossible. Formally, this changes the association matrix \mathbf{A} to

$$\mathbf{A}_{ij} \in \mathbb{R}^{N \times N} = \begin{cases} a_{ij} & \text{for } i \neq j \wedge id(i) \neq id(j), \\ 0 & \text{else,} \end{cases} \quad (6)$$

where id maps a time step to an identifier of its corresponding time series.



2.2.3 CMAL

220 CMAL (Klotz et al., 2022) is an LSTM-based mixture density network (Bishop, 1994). The model predicts the parameters of
asymmetric Laplacian distributions. This choice of distribution allows modeling asymmetric uncertainties that are typical for
many hydrological variables. The direct comparison to HopCPT is slightly problematic, as CMAL does not operate on top of
an existing prediction model. However, as one of the best models for uncertainty estimation in streamflow prediction, CMAL
is a good performance yardstick in our evaluation.

2.3 Metrics

225 The evaluation focuses on the validity and efficiency of the prediction intervals of the models. Validity means that, when a
certain coverage, e.g., 90% is specified, also 90% of the test data is actually covered by the prediction interval. This criterion
is measured by the Δ Cov metric which represents the difference between the specified and the empirical coverage. Δ Cov is
at best zero, while a notably negative Δ Cov diminishes the utility of the model. A positive Δ Cov is less problematic, but a
sign that the model could provide more efficient intervals. Here, efficiency refers to the width of the prediction interval: a more
230 narrow prediction interval is more efficient than a wider one. We evaluate this property with the PI-Width metric, which simply
corresponds to the average width of the prediction interval over the evaluation period. A smaller PI-Width value is therefore
better. Additionally, we evaluate the Winkler score, which jointly elicits both criteria and thus allows an easy comparison
between different models. The Winkler score is calculated as

$$\text{WS}_\alpha(\mathbf{z}_t, M, y_t) = \begin{cases} \text{IW}^\alpha(\mathbf{z}_t, M) + \frac{2}{\alpha}(y_t - \widehat{C}^{\alpha,u}(\mathbf{z}_t, M)) & \text{if } y_t > \widehat{C}^{\alpha,u}(\mathbf{z}_t, M), \\ \text{IW}^\alpha(\mathbf{z}_t, M) + \frac{2}{\alpha}(\widehat{C}^{\alpha,l}(\mathbf{z}_t, M) - y_t) & \text{if } y_t < \widehat{C}^{\alpha,l}(\mathbf{z}_t, M), \\ \text{IW}^\alpha(\mathbf{z}_t, M) & \text{else.} \end{cases} \quad (7)$$

235 The score corresponds to the interval width $\text{IW}_t^\alpha = \widehat{C}^{\alpha,u} - \widehat{C}^{\alpha,l}$ whenever the observed value y_t is within the interval of
 $\widehat{C}_t^\alpha(\mathbf{Z}_{t+1})$ — else, the score gives penalty that is weighted by the warranted coverage level α .⁴ The Winkler score as such
calculated for each individual time step t , but we report the average over all time steps and basins (as is common in time series
literature).

3 Data and Experiments

240 We evaluate all approaches based on the predictions of an LSTM based rainfall–runoff model. We evaluate all models with three
coverage levels: $\alpha = \{0.05, 0.10, 0.15\}$. Appendix D provides the technical details of our experimental setup; and Appendix E
describes our hyperparameter search.

⁴ $\widehat{C}^{\alpha,u}$ and $\widehat{C}^{\alpha,l}$ refer to the upper and lower bound of the interval respectively



3.1 Data

All experiments are based on the public Catchment Attributes and Meteorology for Large-Sample Studies (CAMELS) data set (Newman et al., 2015; Addor et al., 2017). CAMELS comprises basis-averaged daily meteorological forcings derived from three different gridded data products across the United States of America. We used the same 531 basins which were used in the original benchmark (Newman et al., 2017) and in related follow-up work (Kratzert et al., 2019a; Klotz et al., 2022). This subset contains basins ranging in size from 4 to 2,000 km². The dataset contains daily meteorological forcings (precipitation, temperature, short-wave radiation, humidity) from three different data sources (NLDAS, Maurer, DayMet), daily streamflow discharge data from the US Geological Survey, as well as basin-averaged catchment attributes related to soil, geology, vegetation, and climate. The dataset is split into three parts across the time axis. These parts represent (1) the training data for the prediction model, (2) the calibration data used by the uncertainty model, and (3) the test data for evaluation. For CMAL, which does not need any calibration data, both the training and calibration split are used for training.

3.2 Experiment I

This experiment assesses how recent measurements affect the performance of HopCPT. We split it into two parts:

Experiment I-a compares identical HopCPT models once with “offline” memory and once with “online” memory (see Section 2.2.1). We examine two different input feature configurations — once using the time series covariates and the model prediction and once using the model state and model prediction (see Section 2.2.1). In addition to the quantitative comparison, we qualitatively analyze the individual time series for which the performance gap between “online” and “offline” is the greatest.

This yields insights into potential shifts in the prediction error in such cases.

Experiment I-b analyzes intermediate memory update strategies that fall between fully “offline” and fully “online”. In real-world scenarios, gathering the labels with some delay might sometimes be feasible, and at least partially help to mitigate the impact of distribution shifts. Therefore, we evaluate HopCPT with a memory update frequency of 1 week, 1 month, 3 months, 6 months, 1 year, and 2 years (note that the sampling frequency of the series is one day and the overall test period is 9 years) and compare the results to the “offline” and fully “online” (i.e., frequency of 1 day).

3.3 Experiment II

In Experiment II we investigate how much the data from a given basin contributes towards the uncertainty predictions for said basin itself. We do so by first comparing HopCPT with HopCPT-G in a gauged setting (Experiment II-a), and then comparing HopCPT-G with CMAL in an ungauged setting (Experiment II-b). The details for both comparisons are explained in the following.

Experiment II-a compares HopCPT-G to the originally proposed HopCPT in the gauged setting. We focus on the HopCPT feature configuration which includes the model state and prediction, as we hypothesize that this configuration is less likely to simply down-weight samples from time series that are different from the predicted one. We additionally evaluate HopCPT-G with PUB training, i.e., the adapted training loss (see Equation 6). Although we do not evaluate on ungauged basins in this



275 setting, PUB training can increase the tendency of HopCPT to incorporate data from other basins — and that could potentially lead to more robust representations.

Experiment II-b investigates if uncertainty estimates are also possible without any local information from the predicted basin. To examine this scenario we loosely follow the PUB setting from Kratzert et al. (2019a): The set of time series is split into 11 mutually exclusive subsets of equal size. The principle is similar to k-fold cross validation (but not on a per-sample basis): We define that for the k-th fold, the gauged basins are the union of all but the k-th subset, while the k-th subset represents the ungauged basins. We reserve one of the folds for hyperparameter tuning and exclude it from the evaluation. For the other 10 folds, we individually train the prediction and uncertainty model on the “gauged” basins (9 out of 10 subsets) and evaluate them on the ungauged basins from the remaining subset. Note that within each subset, each time series is (as in the standard case) split into training, calibration, and test data. CMAL’s training data encompasses both the training and calibration split to ensure that in total each model has the same amount of data available. We evaluate only the test period of the ungauged basins (this avoids information leakage from the train and calibration period of the gauged basins). We evaluate two variants of HopCPT: (1) HopCPT-G with normal training and (2) HopCPT-G with PUB training. For both variants, we use the model state and prediction as input features.

3.4 Experiment III

290 Experiment III examines the uncertainty estimation for peak flows specifically. Peak flow uncertainties are especially hard to capture. Firstly, because prediction models tend to make larger errors (i.e., high aleatoric uncertainty); and secondly because the occurrence of peak flow events is limited (i.e., high epistemic uncertainty). To measure the respective performance we calculate the metrics only using time steps where the streamflow observations are in the top x % of the corresponding basin. Specifically, we evaluate for $x \in \{2, 5, 10, 20, 30, 50, 100\}$ (and 100 comprises all data and hence corresponds to the “standard” evaluation) within the gauged basin setting of Experiment II. For the sake of completeness, Appendix G3 also presents the peak flow evaluation for the other experiments.

4 Results and Discussion

4.1 Experiment I

Experiment I-a compares the offline and online mode of HopCPT. Table 1 shows that the main performance advantage of the online setting is its better coverage. This holds especially true for the HopCPT variant which uses the model states as features. The width of the prediction interval stays almost constant between offline and online approaches. This suggests that the change in Winkler score is a direct effect of the lower/higher coverage. We argue that the slight loss in coverage for the offline setting is due to a distribution shift in some basins (Fig. 3): While the coverage distribution of the basins with highest coverage is very similar between the offline and online cases, the biggest change happens in the 10 – 20% of basins with the lowest coverage.



Table 1. HopCPT performance of the two evaluated input combinations in offline and online mode for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score) — significance tested with a Mann–Whitney U test at $p < 0.005$. For PI-Width and Winkler score lower values are better — for Δ Cov non-negative values close to 0 are best. The values in parenthesis represent the standard deviation over the different seeds.

α		0.05			0.10			0.15		
		Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
X / YHat	offline	−.005 (0.002)	2.88 (0.09)	1.21 (0.02)	−.011 (0.003)	2.15 (0.06)	0.94 (0.01)	−.017 (0.003)	1.75 (0.05)	0.80 (0.01)
	online	.005 (0.002)	2.87 (0.09)	1.16 (0.02)	.004 (0.002)	2.14 (0.06)	0.91 (0.01)	.001 (0.002)	1.74 (0.05)	0.78 (0.01)
Model States / YHat	offline	−.017 (0.003)	1.92 (0.06)	1.02 (0.02)	−.029 (0.004)	1.44 (0.05)	0.78 (0.01)	−.039 (0.005)	1.17 (0.04)	0.66 (0.01)
	online	.005 (0.002)	1.96 (0.06)	0.90 (0.01)	.002 (0.003)	1.46 (0.04)	0.71 (0.01)	−.004 (0.004)	1.19 (0.04)	0.61 (0.00)

305 A particularly striking examples of this shift as shown in Figure 4: In the prediction year 2000 (approx. 1 year after the
 memory end), the offline setting provides reasonable intervals. However, in the year 2004, the real streamflow seems to be
 shifted upwards for low flows and downwards for some higher flows. Since the output of the prediction model remains roughly
 constant, we suggest that the phenomenon at hand is a shift in the runoff that not visible in the input patterns. The online model
 can accommodate for this shift since the new information (in the form of the shifted observations) are incorporated into the
 310 memory. The offline model simply has no mechanism to account for that and construct invalid intervals. Figures G1 and G2 in
 the appendix show additional examples of this error shift behavior.

Experiment I-b investigates the effect of HopCPT with infrequent memory updates, as intermediate settings between the
 edge cases of fully online and offline settings. Table 2 presents the results. The slight coverage loss increases gradually with
 a higher update delay — as illustrated in Figure 5. However, an update frequency of 1 year already halves the coverage loss
 315 compared to the offline setting. The changes in the PI-Width are less clear and non-monotonic, however, the variation is any-
 way negligible.

Answer to RQI: Does up to date information improve the general uncertainty estimation?

Our results indeed suggest that a continuous incorporation of new data improves uncertainty predictions. For the center of the
 predictive distribution (associated with large α values; Tab. 1) our results are less pronounced than for the tails (associated
 with small α values; Tab. 1). Further, continuously updating the uncertainty estimates as new data comes along is very useful,
 given that environmental processes are associated with all kinds of distribution shifts. In our experiments, a memory update
 mechanism was advantageous even in cases where real-world constraints only allow for very infrequent updates (Tab. 2).

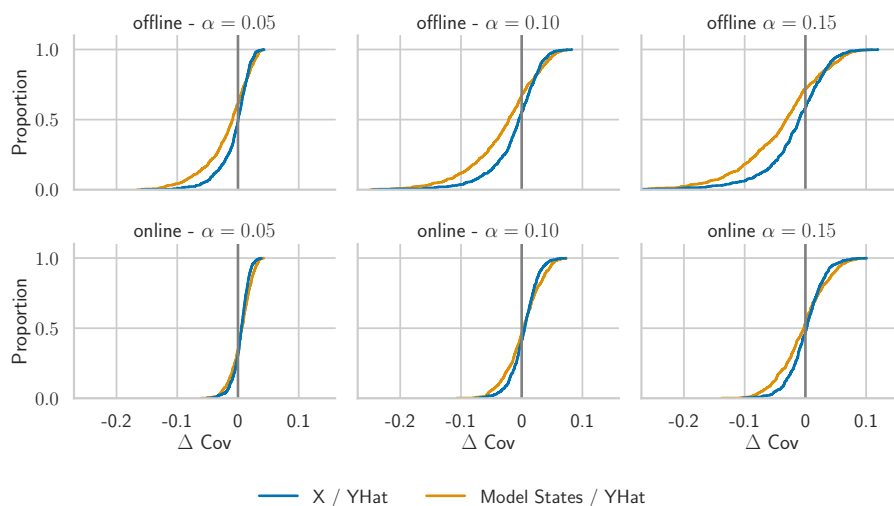


Figure 3. CDF of ΔCov over individual basins for the models evaluated in Experiment I-a. The bottom row shows the results for the “online” mode, i.e., the memory gets updated after each prediction, and the top row for the “offline” mode, i.e., the memory does not get updated during the test period. For ΔCov non-negative values close to 0 are best.

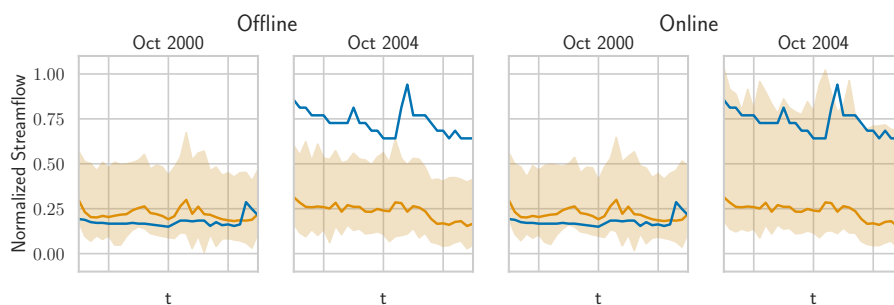


Figure 4. The real streamflow (blue), the prediction (gold), and HopCPT’s prediction interval (light gold) for a basin (id: 12447390) in October 2000 and October 2004. The two plots one the left show HopCPT in the offline mode — the right two plots in the online mode. The online mode allows HopCPT to account for the distribution shift in October 2004.



Table 2. Performance of HopCPT with different memory update behaviour for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The “Delay” column indicates the update frequency of the memory. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score). For PI-Width and Winkler score lower values are better — for Δ Cov non-negative values close to 0 are best. The values in parenthesis represent the standard deviation over the different seeds.

Delay	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
Online	0.005 (0.002)	1.96 (0.06)	0.900 (0.010)	0.002 (0.003)	1.46 (0.04)	0.710 (0.010)	−0.004 (0.004)	1.19 (0.04)	0.610 (0.000)
1 Week	0.002 (0.002)	1.98 (0.05)	0.922 (0.009)	−0.002 (0.002)	1.47 (0.04)	0.719 (0.006)	−0.009 (0.003)	1.20 (0.03)	0.617 (0.004)
1 Month	−0.002 (0.002)	1.97 (0.05)	0.943 (0.010)	−0.007 (0.003)	1.47 (0.04)	0.731 (0.006)	−0.014 (0.003)	1.20 (0.03)	0.625 (0.004)
3 Months	−0.005 (0.002)	1.97 (0.05)	0.956 (0.010)	−0.011 (0.003)	1.47 (0.04)	0.739 (0.006)	−0.019 (0.003)	1.20 (0.03)	0.631 (0.004)
6 Months	−0.006 (0.002)	1.97 (0.05)	0.961 (0.010)	−0.013 (0.003)	1.47 (0.04)	0.742 (0.006)	−0.021 (0.003)	1.20 (0.03)	0.633 (0.004)
1 Year	−0.007 (0.002)	1.97 (0.05)	0.964 (0.010)	−0.014 (0.003)	1.47 (0.04)	0.743 (0.006)	−0.023 (0.003)	1.20 (0.03)	0.634 (0.004)
2 Years	−0.008 (0.002)	1.97 (0.05)	0.969 (0.011)	−0.016 (0.003)	1.47 (0.04)	0.747 (0.007)	−0.025 (0.003)	1.20 (0.03)	0.637 (0.005)
Offline	−0.017 (0.003)	1.92 (0.06)	1.020 (0.020)	−0.029 (0.004)	1.44 (0.05)	0.780 (0.010)	−0.039 (0.005)	1.17 (0.04)	0.660 (0.010)

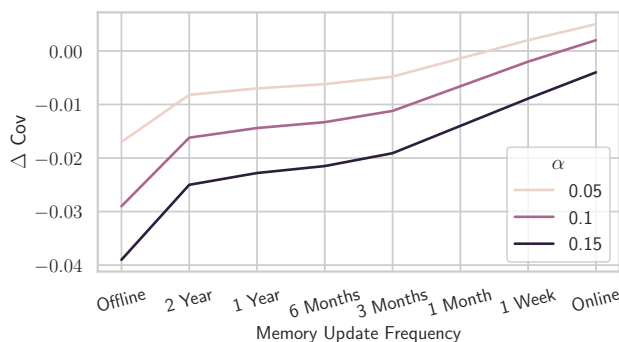


Figure 5. Δ Cov of HopCPT for different memory update frequency settings. Increasing update frequency monotonically improves the coverage for each evaluated coverage level. For Δ Cov non-negative values close to 0 are best.



4.2 Experiment II

320 **Experiment II-a** probes the effect of global data in the “gauged” basin setting. Table 3 shows the result of HopCPT-G (with
and without PUB training) and compares it to the non-global HopCPT. Surprisingly, the shared memory of HopCPT-G (no PUB
training) does not improve the results. However, it also does not do any notable harm. We hypothesize that the calibration data
of each basin is already so comprehensive that all relevant error regimes for the respective basin are covered with sufficient
325 all non-relevant basins, or (ii) the missing temporal encoding in HopCPT-G. The PUB training procedure, on the other hand,
notably improves the coverage. However, this improvement comes at the cost of efficiency. Hence, the resulting Winkler Score
is similar to the other approaches.

Figure 6 shows the coverage distribution over the individual basins. The distributions of HopCPT and HopCPT-G are very
similar. This indicates that HopCPT-G cannot profit from the cross-time-series information, but successfully learns to exclude
330 irrelevant basins in order to arrive at the “right” error distribution.

Experiment II-b investigates the PUB setting. Table 4 shows the result of the evaluation. CMAL, which has very good cov-
erage in a gauged basin setting (see Appendix A), exhibits a noticeable under-coverage. HopCPT-G, on the other hand, only
slightly loses coverage compared to the non-PUB setting. HopCPT-G with PUB training reduces the coverage loss notably. Its
efficiency is comparable to CMAL and HopCPT-G. CMAL achieves the best Winkler score despite its high under-coverage.
335 This suggests that the uncovered samples are relatively close to the border of the prediction interval. Figure 7 shows the distri-
bution of Δ Cov over the individual time series. Since the PUB setting does still yield acceptable performances but does not
allow to use information from a given basin, one can conclude that the global setting is able to transfer uncertainty information
about the uncertainty from other basins to the unseen basins. Quantitative metrics for the individual folds can be found in the
appendix in Table G1.

340

Answer to RQII: Is the data from a given basin required to get good uncertainty estimates?

Our result suggests that local information is beneficial for efficient uncertainty intervals. Given that enough local information
is available and the average quality of the estimate is the focus (in contrast to RQIII), it is also sufficient to provide reasonable
uncertainty estimates (Experiment II-a). However, local information is not strictly necessary to produce sensible uncertainty
estimations, and information transfer via global information is possible (Experiment II-b).



Table 3. Performance of the different HopCPT variants with for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score) — significance tested with a Mann–Whitney U test at $p < 0.005$. For PI-Width and Winkler score lower values are better — for Δ Cov non-negative values close to 0 are best. The values in parenthesis represent the standard deviation over the different seeds.

α	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
HopCPT	−.017 (0.003)	1.92 (0.06)	1.02 (0.02)	−.029 (0.004)	1.44 (0.05)	0.78 (0.01)	−.039 (0.005)	1.17 (0.04)	0.66 (0.01)
HopCPT-G	−.023 (0.003)	2.10 (0.07)	1.06 (0.01)	−.031 (0.004)	1.61 (0.05)	0.81 (0.00)	−.036 (0.004)	1.35 (0.05)	0.68 (0.00)
HopCPT-G (PUB Train)	.018 (0.006)	3.16 (0.18)	1.04 (0.07)	.021 (0.008)	2.36 (0.13)	0.81 (0.04)	.019 (0.011)	1.93 (0.11)	0.69 (0.02)

Table 4. Performance of different models on the PUB experiment for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all test basins of all splits. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score) — significance tested with a Mann–Whitney U test at $p < 0.005$. For PI-Width and Winkler score lower values are better — for Δ Cov non-negative values close to 0 are best. The values in parenthesis represent the standard deviation over the different seeds.

α	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
HopCPT-G	−.052 (0.001)	2.15 (0.01)	1.78 (0.02)	−.070 (0.001)	1.57 (0.01)	1.29 (0.01)	−.080 (0.002)	1.27 (0.01)	1.05 (0.01)
HopCPT-G (PUB Train)	−.007 (0.003)	3.00 (0.08)	1.34 (0.04)	−.017 (0.005)	2.27 (0.07)	1.04 (0.03)	−.026 (0.006)	1.87 (0.06)	0.89 (0.02)
CMAL	−.119 (0.008)	2.44 (0.08)	1.18 (0.03)	−.155 (0.009)	1.91 (0.06)	0.93 (0.02)	−.172 (0.009)	1.61 (0.05)	0.80 (0.02)

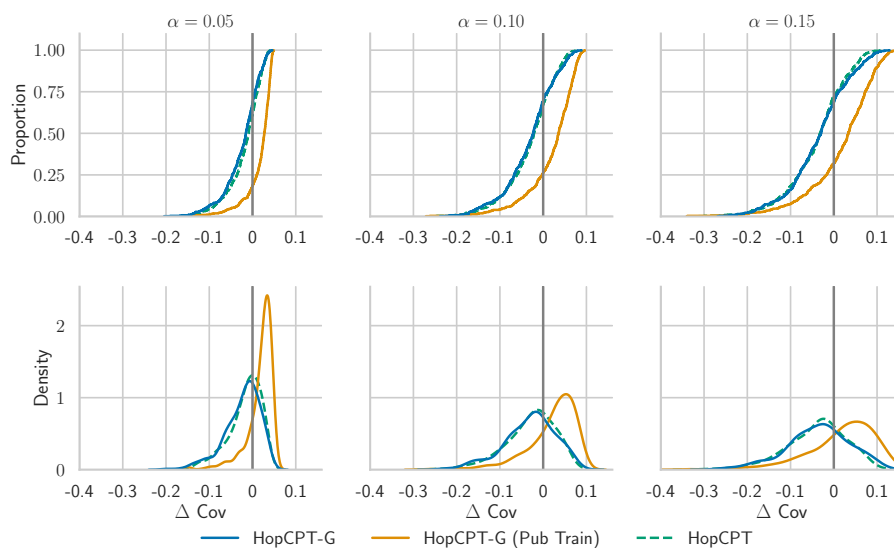


Figure 6. CDF and PDF of Δ Cov over individual basins for models evaluated in Experiment II-a. For Δ Cov non-negative values close to 0 are best.

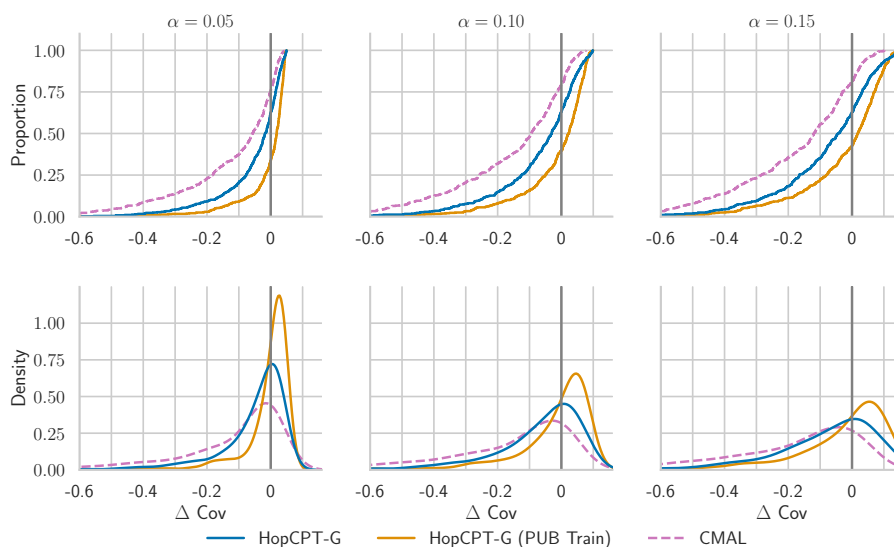


Figure 7. CDF and PDF of Δ Cov over all individual test basins of all PUB folds for models evaluated in Experiment II-b. For Δ Cov non-negative values close to 0 are best.



4.3 Experiment III

In contrast to our results from Experiment II, the peak flow performance exhibits a notable difference between the HopCPT variants (Fig. 8). HopCPT-G provides better coverage than HopCPT, PUB training boosts this further. The coverage difference between the HopCPT variants also increases with increasing α . CMAL, which is also a “fully global” model, achieves the best Winkler scores in the peak flow setting. This indicates that the information from other basins is beneficial for peak flow regimes — even when plenty of past information for the basin is available. We argue that this is because the data is more scarce in these regimes and certain situations might not be available in the observed past of the individual basin. Thus, in this situation it becomes useful to leverage information from the other basins in order to obtain good uncertainty estimates (we note that this is in contrast to the results from Experiment II, which showed that for the average runoff predictions no improvement is obtained by considering global information).

Answer to RQIII: Is data from a single basin enough to get good uncertainty predictions for peak flows?

Our results indicate that using only data from a single basin — that is, only considering the local information — leads to worse uncertainty estimates for peak flow settings. The more restricted the “peak flow categorization“ is chosen, the more pronounced this effect gets (Figure 8). Using global information — i.e., information from other basins — improves the uncertainty estimates. In particular, this holds in terms of reliability of the estimate (i.e., coverage).

5 Conclusions and Outlook

This contribution investigates how the temporal (recency) and spatial (local vs. global) dimensions of information impact the quality of uncertainty estimates. To conduct our study we apply the Conformal Prediction (CP) framework in the form of the Hopfield Conformal Prediction for Timeseries (HopCPT) approach, which extends the CP framework for time series predictions. In short, we find that (a) the inclusion of the most recent information has high benefits for the general uncertainty predictions; and (b) that global information is not important for general uncertainty predictions, but pivotal to provide good bounds for peak flows.

Regarding (a), our analysis suggests that incorporating recent information helps to improve the uncertainty estimates. This is even true if the information can only be provided after longer periods (e.g., when for technical reasons only yearly updates are possible). We could qualitatively link this phenomenon to distribution shifts that appear over time. We conclude that continuous monitoring and incorporation of the newly obtained data in the prediction process are vital components of prediction systems that strive to provide reliable and efficient uncertainty estimations.

Regarding (b), our results indicate that local information is sufficient to provide good uncertainty estimates on average (assuming that a reasonable history over multiple years is provided). However, for peak flows it is not. We argue that this is because the estimation problem is particularly hard since high-flow situations are diverse, have a high measurement variance, and happen very infrequently. Global information is able to improve the estimates. These observations are in accordance with the results from Frame et al. (2022) and Bertola et al. (2023) — which indicate that the signals from different basins can be

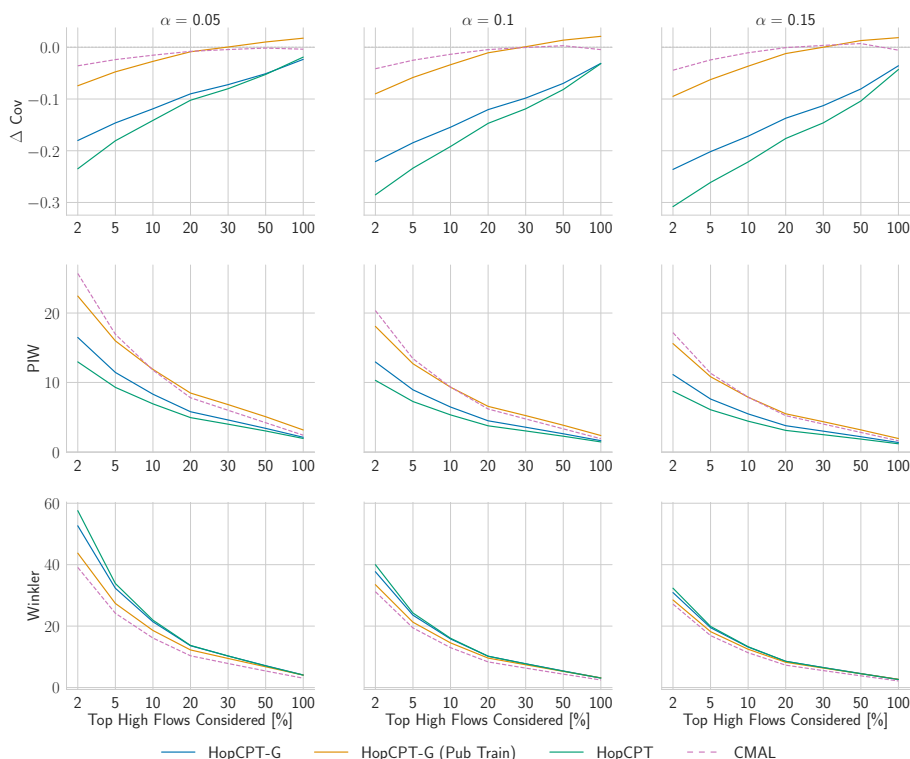


Figure 8. Evaluation metrics for Experiment III. Only peak flow time steps — defined via varying share of the overall steps (x-axis of the plots) — are considered for the respective metric. For PI-Width and Winkler score lower values are better — for Δ Cov non-negative values close to 0 are best.

leveraged to make better predictions for peak flows. This can be taken to the extreme for predictions in ungauged basins. Indeed, our experiments in this regard suggest that reasonable uncertainty estimates can be also provided without any local information. However, whenever global information is used for uncertainty estimation it is especially important that the respective model can selectively consider the relevant information — i.e. can decide which parts of the overall information should be used for the “current” estimation. HopCPT-G and CMAL are both able to do that.

In future work, experiments with more scarce local data could reveal further advantages of approaches that incorporate global information. More broadly, CP offers a fresh perspective to uncertainty-aware streamflow prediction. Appendix B shows a first exploration of how the principles from CP can improve existing hydrological uncertainty estimations with minimal interventions. In the future, we want to explore how ideas from CP are able to refine other (perhaps less formal) approaches that are currently used in practice.



Code and data availability. The code repository is available at <https://doi.org/10.5281/zenodo.10674231>. The trained base model (LSTM) and the utilized model states as well as the global HopCPT models are available at <https://doi.org/10.5281/zenodo.10653863>. The trained CMAL models for the non-PUB experiments is available at <https://doi.org/10.5281/zenodo.10654345>, the CMAL models for the PUB experiments are available at <https://doi.org/10.5281/zenodo.10654399>. The data for CAMELS can be accessed for free on the NCAR's official website (<https://gdex.ucar.edu/dataset/camels/file.html>). The expanded Maurer forcings, which include data on daily minimum and maximum temperatures, are available for download at <https://doi.org/10.4211/hs.17c896843cf940339c3c3496d0c1c077> (Kratzert, 2019).



Appendix: Overview

The appendix is structured as follows: Appendix A compares the overall performance of different uncertainty methods, including HopCPT and CMAL. Appendix B investigates the impact of applying the principles of CP to the existing uncertainty estimation method Bluecat. Appendix C investigates how different input features influence the uncertainty estimation performance — i.e. provides an orthogonal dimension to the spatial and temporal consideration of the main paper. Appendix D and Appendix E provide details about the experimental setup and the hyperparameter tuning respectively. Appendix F presents the point prediction metrics of the experiments. Finally, Appendix G provides additional results to the experiments in the main paper.

Appendix A: Model intercomparison

This section compares the performance of the different (non-global methods⁵). For HopCPT, we evaluate the model offline (i.e., without updating the memory as outlined in Section 2.2.1). Two variants for the input features are examined: (1) The time series features and the model prediction, similar to the original work but without the lagged target; (2) the model states and prediction as proposed in Section 2.2.1. For Bluecat, we selected the best model considering both the original Bluecat and the adapted version suggested in Sec. B. We follow Auer et al. (2023) and also evaluate against kNN as a naive similarity-based baseline (see Appendix A1). While Auer et al. (2023) show that kNN on the time series features is not sufficient, we analyze the performance of kNN with (1) the model states and prediction (corresponds to HopCPT) and (2) with only the model prediction (corresponds to Bluecat) as input. Additionally, we evaluate CMAL as the current state-of-the-art approach.

A1 kNN as Navive Similarity Approach

HopCPT uses the similarity representation that is learned by the MHN network. A more naive way to consider such a similarity would be to simply use a kNN model based on the model inputs. Following the original work, we use such a kNN model as a baseline in the experiments.

A2 Model intercomparison - Results

Table A1 shows the evaluation of the different models. CMAL, closely followed by HopCPT, provides the best coverage at all coverage levels. While HopCPT with the model states and input reaches the most efficient — i.e., smallest — prediction intervals, CMAL performs best in terms of the Winkler score. This is slightly surprising, given that the Winkler score encompasses the coverage and PI-Width. However, the fact that Winkler score additionally considers the distance of the uncovered test samples leads us to the hypothesis that CMAL, as it is optimized to maximize the likelihood of all samples, results in smaller distances in that regard. In contrast, HopCPT and Bluecat do not consider the uncovered samples in optimization and calibration. The comparison between the HopCPT evaluations with different input features shows that adding the model states

⁵CMAL can be considered a global model



Table A1. Performance of the evaluated models for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score), given that $\Delta \text{Cov} \leq -\alpha$, i.e., the specific algorithm reached reasonable coverage (the result is grayed otherwise) — significance tested with a Mann–Whitney U test at $p < 0.005$. The values in parenthesis represent the standard deviation over the different seeds (results without these are from deterministic models).

α	0.05			0.10			0.15		
	ΔCov	PI-Width	Winkler	ΔCov	PI-Width	Winkler	ΔCov	PI-Width	Winkler
HopCPT	−.005 (0.002)	2.88 (0.09)	1.21 (0.02)	−.011 (0.003)	2.15 (0.06)	0.94 (0.01)	−.017 (0.003)	1.75 (0.05)	0.80 (0.01)
HopCPT (model states)	−.017 (0.003)	1.92 (0.06)	1.02 (0.02)	−.029 (0.004)	1.44 (0.05)	0.78 (0.01)	−.039 (0.005)	1.17 (0.04)	0.66 (0.01)
Bluecat' (M=100)	−.050	2.72	1.19	−.057	2.00	0.90	−.061	1.63	0.76
kNN (yHat)	−.176	1.95	1.84	−.202	1.47	1.26	−.215	1.19	1.01
kNN (model states)	−.193	1.37	1.74	−.219	1.04	1.15	−.231	0.85	0.90
CMAL	−.004 (0.004)	2.40 (0.06)	0.78 (0.01)	−.004 (0.008)	1.89 (0.05)	0.63 (0.01)	−.006 (0.011)	1.60 (0.04)	0.55 (0.01)

as inputs enhances the efficiency of HopCPT while keeping an approximate coverage, which results in an overall lower Winkler score. A more detailed analysis of the effect of different input modalities is given in Appendix C. The best variant of Bluecat falls behind HopCPT and CMAL in all three metrics. The kNN model, as a naive similarity measure baseline, results in high under-coverage also when considering the model states. This reinforces the motivation for a learned similarity measure, as already mentioned in the original HopCPT paper.

Individual basin analysis Figure A1 gives more detailed insights into how the metrics are distributed around the individual basins. CMAL, which archives the best overall coverage, does not provide better coverage for all basins but the high mean is driven by the higher over-coverage of the upper half of basins. This aligns with the analysis in (Klotz et al., 2022). Overall, the coverage distribution of HopCPT is more centered around zero and better bounds the lowest coverage. As the cumulative distribution plot in Figure A2 shows, for the 50 (approx 10%) basins with worst coverage, HopCPT without model states performs best, and the worst coverage is best for HopCPT (no matter if with or without model states). The distribution of the PI-Width (Figure A3) is similar for all models, however, the efficiency advantage of HopCPT is most pronounced for the basins with larger prediction intervals.

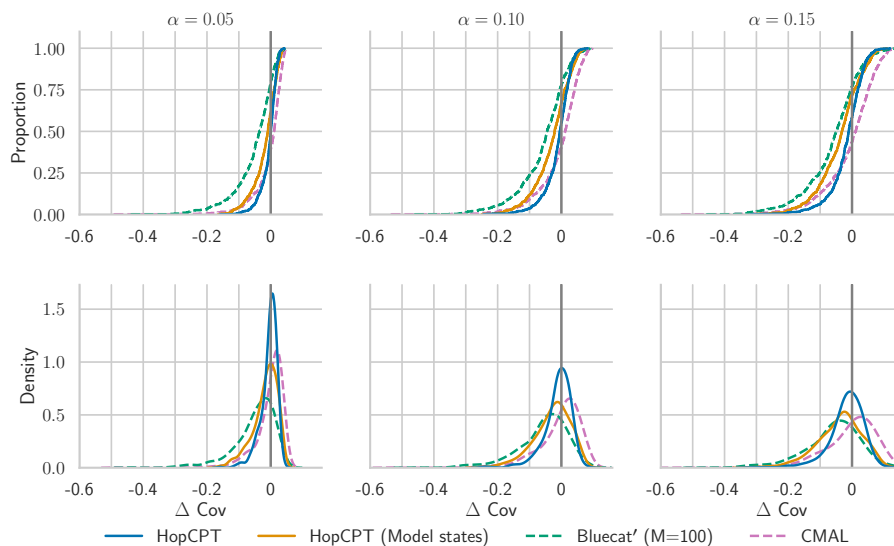


Figure A1. CDF and PDF of ΔCov over individual basins for models with approximate average coverage ($\Delta \text{Cov} \leq -\alpha$) evaluated in Appendix A.

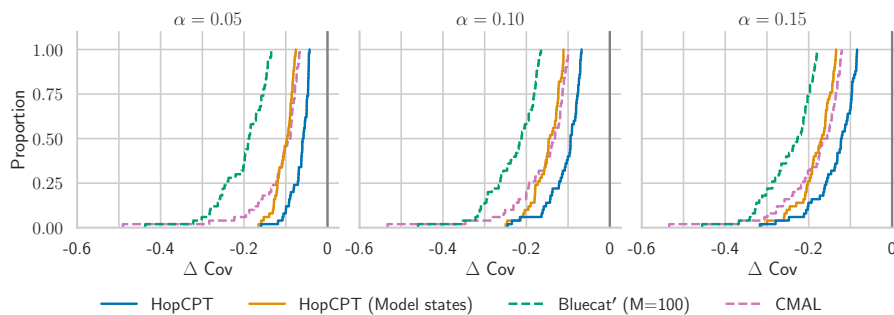


Figure A2. CDF of ΔCov over the 50 basins with the highest misscoverage for models with approximate average coverage ($\Delta \text{Cov} \leq -\alpha$) evaluated in Appendix A.

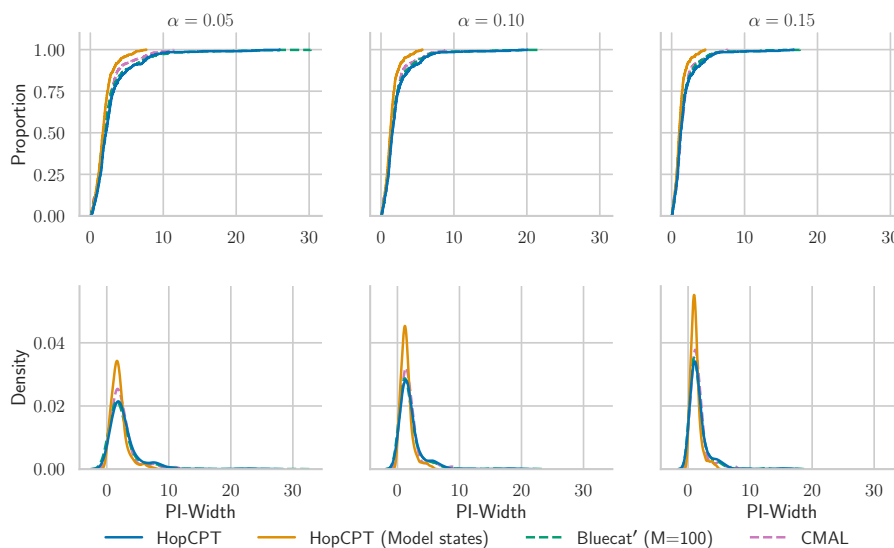


Figure A3. CDF and PDF of PI-Width over individual basins for models with approximate average coverage ($\Delta \text{Cov} \leq -\alpha$) evaluated in Appendix A.

430 Appendix B: Bluecat

Similar to HopCPT, the hydrological uncertainty estimation approach Bluecat (Koutsoyiannis and Montanari, 2022) aims to provide a prediction interval given a point prediction model and a set of data points (calibration data). Bluecat first evaluates the calibration data. Then, for a new prediction point, the calibration point with the most similar model prediction and all points within a certain distance to this point are considered. This can be seen as a special form of kNN where (1) the distance metric between two points is the distance of the model predictions and (2) not k points but all points within a certain distance
 435 threshold are considered (Rozos et al., 2022). Given this set of selected calibration points, similar to the CP framework, the specified quantiles are calculated and provide the bounds of the confidence band.

The CP perspective on Bluecat (Bluecat')

Viewing Bluecat through a CP lens suggests a relatively low effort improvement in the form of the introduction of a cali-
 440 bration set for the interval prediction. Bluecat only uses the training data of the prediction model to determine the prediction interval, i.e., prediction training data = calibration data. However, the CP literature suggests that an independent calibration set can give a more unbiased estimate, since prediction models can easily overfit training data (and, as a result, the prediction intervals might be overly optimistic and do not fulfill the specified coverage criterion). We refer to this adaptation as Bluecat' throughout the paper.

445 We evaluate Bluecat once as originally proposed (Bluecat) and once with the adaption (Bluecat'). Bluecat requires setting the hyperparameter M , which controls how many close data points are considered. In the original publication, there is no clear guidance on how to select the parameter for new datasets, which is why we evaluate two variations: (1) tuning the



Table B1. Performance of the evaluated Bluecat variants for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score). As Bluecat is deterministic no standard deviation is given.

α	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
Bluecat (M=100)	-.092	2.09	1.51	-.112	1.55	1.06	-.122	1.26	0.86
Bluecat (M=225)	-.071	2.22	1.55	-.091	1.60	1.10	-.102	1.28	0.89
Bluecat' (M=100)	-.050	2.72	1.19	-.057	2.00	0.90	-.061	1.63	0.76
Bluecat' (M=225)	-.028	2.88	1.24	-.035	2.06	0.93	-.039	1.64	0.78

hyperparameters with the same model selection criteria as in HopCPT (resulting in $M = 225$), and (2) using the hyperparameter value from the original publication ($M = 100$).

450 B1 Bluecat - Results

Table B1 shows the result of the comparison. The adaptation Bluecat' achieves considerably improved coverage. The lower Winkler scores further indicate the overall better performance of Bluecat'. These results support the hypothesis that Bluecat intervals are overly optimistic since they are based on the — most likely overfitted — training data predictions of the prediction model. Further, the hyperparameter-tuned versions ($M = 225$) of both Bluecat variants enhance the coverage of the models at
 455 the cost of lower interval efficiency. This makes sense, as the model selection favors models with better coverage as long as no model fully achieves the specified coverage.

Appendix C: Input information

This section assesses how different HopCPT inputs affect its performance. The experiment compares all combinations of the three potential feature components: time series covariates (X), model prediction ($Y\hat{H}$), and model state. We use an LSTM
 460 prediction model. Hence, the model state refers to the two internal state vectors of the LSTM, which are updated in each new step and are — together with the current input — the basis for the output of the model. Specifically, we consider the state vectors right after the prediction, since these states already consider the current input features. We tuned hyperparameters individually for each feature set combination. Additionally, we analyze results where we include the lagged target (Y) as a feature, as in the original HopCPT paper.



Table C1. HopCPT performance of the evaluated input combinations for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score) — significance tested with a Mann–Whitney U test at $p < 0.005$. The values in parenthesis represent the standard deviation over the different seeds.

α	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
X	−.003 (0.002)	3.50 (0.15)	1.61 (0.02)	−.005 (0.003)	2.35 (0.10)	1.17 (0.01)	−.007 (0.005)	1.80 (0.08)	0.96 (0.01)
X / YHat	−.005 (0.002)	2.88 (0.09)	1.21 (0.02)	−.011 (0.003)	2.15 (0.06)	0.94 (0.01)	−.017 (0.003)	1.75 (0.05)	0.80 (0.01)
Model States	−.017 (0.004)	1.93 (0.06)	1.01 (0.01)	−.029 (0.005)	1.45 (0.04)	0.77 (0.01)	−.041 (0.006)	1.18 (0.03)	0.66 (0.00)
Model States / YHat	−.017 (0.003)	1.92 (0.06)	1.02 (0.02)	−.029 (0.004)	1.44 (0.05)	0.78 (0.01)	−.039 (0.005)	1.17 (0.04)	0.66 (0.01)
X / Model States	−.019 (0.003)	1.90 (0.04)	1.02 (0.02)	−.032 (0.005)	1.43 (0.03)	0.78 (0.01)	−.043 (0.006)	1.17 (0.03)	0.66 (0.01)
X / Model States / YHat	−.020 (0.004)	1.90 (0.05)	1.01 (0.01)	−.032 (0.005)	1.43 (0.04)	0.77 (0.01)	−.043 (0.005)	1.18 (0.03)	0.65 (0.01)

465 Table C1 shows the performance of HopCPT with the different input feature configurations. Interestingly, only using the
 time series covariates results in the best coverage. However, the prediction intervals in these settings are rather large, which
 is also reflected in notably higher Winkler scores. Therefore, we argue that using the model states as input features is indeed
 beneficial. This argument is further strengthened by the results in Experiment I, where we see that the coverage loss is likely
 due to distribution shifts. Accounting for this shift by updating the memory leads to almost perfect coverage for the model
 470 state setting and unnecessary over-coverage in the covariate-only setting. Both adding the model prediction and the model
 state provide additional information. The model state does not only include information about the current covariates, but also
 includes history information due to the recursive nature of the LSTM. The model prediction, on the other hand, can be seen
 as a projection of the model state where some relevant information could be lost. Extending the model state with additional
 inputs hardly changes the performance of HopCPT, which supports the hypothesis that the model state already includes the
 475 vast majority of the required information.

Table C2, shows the results when the lagged target variable is included in the feature set, as done in the original work Auer
 et al. (2023). Note that this is typically not feasible in streamflow prediction. However, to how this affects the performance we
 also evaluated these input combinations as it is possible in the lab setting.



Table C2. HoCPT performance of the input combinations which include the lagged target y for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all basins. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score). The values in parenthesis represent the standard deviation over the different seeds.

α	0.05			0.10			0.15		
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
X / Y	-.004 (0.001)	2.69 (0.05)	1.27 (0.01)	-.007 (0.003)	1.92 (0.03)	0.96 (0.01)	-.010 (0.004)	1.54 (0.03)	0.81 (0.01)
X / Y / YHat	.011 (0.008)	2.75 (0.19)	0.96 (0.03)	.012 (0.015)	2.16 (0.14)	0.77 (0.02)	.009 (0.019)	1.82 (0.12)	0.68 (0.01)
X / Y / Model States	-.021 (0.003)	1.90 (0.04)	1.00 (0.01)	-.035 (0.004)	1.43 (0.03)	0.77 (0.01)	-.046 (0.005)	1.18 (0.02)	0.65 (0.00)
X / Y / YHat / Model States	-.018 (0.003)	1.91 (0.06)	1.00 (0.01)	-.030 (0.004)	1.44 (0.04)	0.76 (0.01)	-.041 (0.004)	1.19 (0.03)	0.65 (0.00)

Appendix D: Experiment Details

480 Each non-deterministic experiments, apart from experiments with HopCPT-G, are repeated with 12 different seeds. HopCPT-G experiments are repeated with 8 different seeds as the global models are computationally more demanding. For the latter experiments, we remove outlier runs where the training does not converge.

Appendix E: Hyperparameter Search

We conducted a hyperparameter grid search for each model. In the case of HopCPT, we did this individually for each input
 485 feature set. Each hyperparameter search was repeated with 3 seeds — the best average validation score was used as selection criteria. The validation score for the LSTM and CMAL is the NSE metric, for HopCPT, Bluecat, and kNN we followed (Auer et al., 2023) and used the smallest PI-Width at an epoch with Δ Cov ≤ 0 . To limit the number of grid search combinations for the LSTM and CMAL model, we split the hyperparameter into two sets which were tuned sequentially. The second set was trained given the result from the first set. Table E1 shows the parameters used in the hyperparameter search for the LSTM,
 490 CMAL, HopCPT, and HopCPT-G models. For Bluecat we evaluated for $M = \{25, 50, 75, 100, 125, 150, 200\}$. For kNN we varied the share k_s of samples which defines the k parameter (i.e. $k = k_s * \text{number of memory samples}$) and evaluated for $k_s = \{0.025, 0.05, 0.1, 0.15, 0.2, 0.25, 0.3, 0.35\}$.



Table E1. Parameters used in the hyperparameter search.

Model	Parameter	Values		
CMAL	Set1	hidden size	60, 125, 250, 500	
		output dropout	0.4, 0.5	
		target noise	0.05, 0.1, 0.2	
		# distrubtions	1,3,5,10	
	Set2	batch size	256, 512	
		learning rate	0.0005, 0.0001, 0.001	
	LSTM	Set1	hidden size	60, 125, 250, 500
			output dropout	0.4, 0.5
target noise			0.05, 0.1, 0.2	
Set2		batch size	256, 512	
		learning rate	0.0005, 0.0001, 0.001	
HopCPT		learning rate	0.001,0.001	
	encode hidden layer	1,2,3		
	encode dropout	0, 0.1		
	temporal encoding	yes,no		
HopCPT-G	learning rate	0.001,0.001		
	encode hidden layer	0,1,2,3		
	encode dropout	0, 0.1		

Appendix F: Point Prediction Metrics

Table F1 shows the point prediction performance of LSTM and CMAL for Experiments d. The performance for the PUB
 495 prediction setting (Experiment II-b) is presented in Table F2. As in Klotz et al. (2022), CMAL outperforms the LSTM slightly.
 This could imply an advantage for the interval predicting. However, this is justifiable, considering CMAL's greater volume of
 training data as it does not need any calibration data.

Appendix G: Extended Results

G1 Experiment I-b - Additional Shift Examples

500 In Experiment I-b we show an example where a shift in the data highlights the advantages of the online mode for HopCPT.
 Figure G1 and G2 show two additional examples, both from the same basin, at different months.



Table F1. Point prediction metrics for the experiments 1-5.1. NSE: Nash–Sutcliffe efficiency ($-\infty, 1$]; high values are better. MSE: Mean Squared Error; low values are better. KGE: Kling–Gupta efficiency ($-\infty, 1$]; high values are better. The variability of the metrics over the different basins and seeds is provided in the form of the standard deviation for the mean aggregation and in the form of the IQR for the median aggregation.

		LSTM	CMAL
NSE	Mean	0.717 (0.297)	0.716 (0.385)
	Median	0.754 (0.136)	0.790 (0.157)
MSE	Mean	2.363 (3.63)	2.293 (5.025)
	Median	1.257 (1.955)	0.996 (1.769)
KGE	Mean	0.754 (0.254)	0.733 (0.215)
	Median	0.817 (0.140)	0.786 (0.168)

Table F2. Point prediction metrics Experiment 5.2 (PUB). The values represent the average over all test basins of all splits. NSE: Nash–Sutcliffe efficiency ($-\infty, 1$]; high values are better. MSE: Mean Squared Error; low values are better. KGE: Kling–Gupta efficiency ($-\infty, 1$]; high values are better. The variability of the metrics over the different basins and seeds is provided in the form of the standard deviation for the mean aggregation and in the form of the IQR for the median aggregation.

		LSTM	CMAL
NSE	Mean	0.444 (2.111)	0.472 (2.394)
	Median	0.703 (0.243)	0.690 (0.253)
MSE	Mean	3.396 (6.603)	3.362 (6.718)
	Median	1.365 (2.560)	1.369 (2.545)
KGE	Mean	0.496 (7.245)	0.524 (0.591)
	Median	0.654 (0.326)	0.638 (0.280)

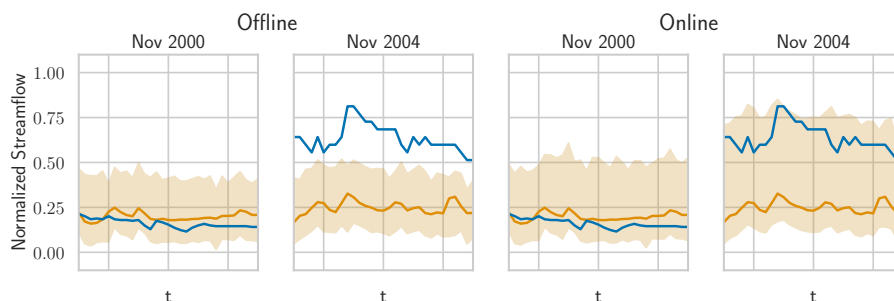


Figure G1. The real streamflow (gold), the prediction (blue), and HopCPT’s prediction interval (light gold) for a basin (id: 12447390) in November 2000 and October 2004. The two plots one the left show HopCPT in the offline mode — the right two plots in the online mode. The online mode allows HopCPT to account for the distribution shift in November 2004.

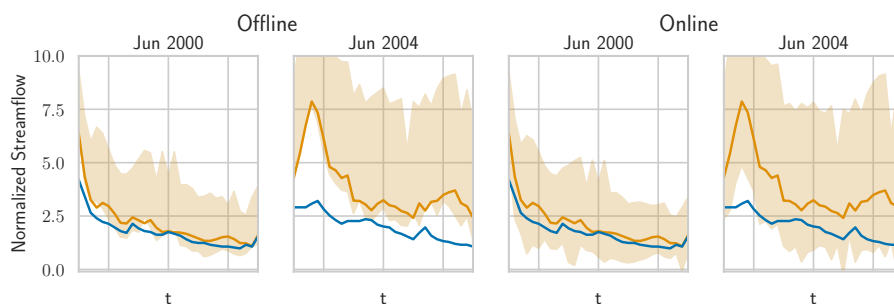


Figure G2. The real streamflow (gold), the prediction (blue), and HopCPT’s prediction interval (light gold) for a basin (id: 12447390) in June 2000 and October 2004. The two plots one the left show HopCPT in the offline mode — the right two plots in the online mode. The online mode allows HopCPT to account for the distribution shift in June 2004.

G2 Experiment II-b - PUB Results per Fold

Table G1 shows the results of the PUB experiment for the individual folds. The performance varies between the folds, however, the ranking of the different models is consistent over all folds.

505 G3 Experiment III - Additional Peak Flow Evaluations

Figure G3 shows the peak flow results for HopCPT with different memory update frequencies (see Experiment I-b). The coverage difference between different update frequencies does not notably change when only peak flows are considered. Figure G4 presents the peak flow results for the PUB setting (see Experiment II-b). Similarly as in the “gauged” basin setting the PUB training improves the coverage of HopCPT notably. CMAL coverage is rather constant over the different peak flow shares.



Table G1. PUB performance over the individual folds for the miscoverage levels $\alpha = \{0.05, 0.10, 0.15\}$. The values represent the average over all test basins of the respective fold. Bold numbers correspond to the best result for the respective metric in the experiment (PI-Width and Winkler score). The values in parenthesis represent the standard deviation over the different seeds.

α	0.05			0.10			0.15			
	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	
Fold 1	HopCPT-G	-.035 (0.003)	2.21 (0.03)	1.59 (0.04)	-.048 (0.005)	1.62 (0.03)	1.17 (0.02)	-.056 (0.005)	1.32 (0.02)	0.96 (0.02)
	HopCPT-G (PUB Train)	.000 (0.006)	3.13 (0.25)	1.23 (0.03)	-.008 (0.009)	2.34 (0.20)	0.97 (0.02)	-.018 (0.012)	1.92 (0.17)	0.84 (0.01)
	CMAL	-.097 (0.030)	2.37 (0.15)	1.10 (0.04)	-.125 (0.039)	1.86 (0.12)	0.87 (0.03)	-.140 (0.043)	1.57 (0.11)	0.76 (0.02)
Fold 2	HopCPT-G	-.052 (0.005)	2.13 (0.05)	2.06 (0.04)	-.064 (0.006)	1.52 (0.04)	1.43 (0.02)	-.068 (0.008)	1.22 (0.03)	1.14 (0.01)
	HopCPT-G (PUB Train)	-.020 (0.011)	2.75 (0.24)	1.48 (0.10)	-.034 (0.015)	2.08 (0.17)	1.13 (0.05)	-.049 (0.015)	1.71 (0.14)	0.96 (0.03)
	CMAL	-.179 (0.026)	2.36 (0.20)	1.31 (0.12)	-.218 (0.027)	1.85 (0.16)	1.00 (0.07)	-.232 (0.028)	1.55 (0.13)	0.85 (0.05)
Fold 3	HopCPT-G	-.050 (0.004)	1.95 (0.04)	1.39 (0.02)	-.072 (0.006)	1.41 (0.03)	1.02 (0.01)	-.089 (0.008)	1.13 (0.03)	0.84 (0.01)
	HopCPT-G (PUB Train)	-.011 (0.002)	2.53 (0.12)	1.16 (0.02)	-.020 (0.003)	1.92 (0.10)	0.88 (0.01)	-.032 (0.007)	1.58 (0.08)	0.74 (0.01)
	CMAL	-.121 (0.019)	1.98 (0.14)	0.90 (0.06)	-.155 (0.019)	1.55 (0.11)	0.71 (0.04)	-.170 (0.018)	1.30 (0.10)	0.62 (0.03)
Fold 4	HopCPT-G	-.030 (0.005)	2.18 (0.04)	1.47 (0.02)	-.042 (0.007)	1.58 (0.03)	1.08 (0.01)	-.050 (0.008)	1.27 (0.02)	0.90 (0.01)
	HopCPT-G (PUB Train)	-.002 (0.008)	2.74 (0.23)	1.18 (0.06)	-.005 (0.013)	2.11 (0.18)	0.92 (0.03)	-.008 (0.015)	1.76 (0.15)	0.78 (0.02)
	CMAL	-.090 (0.011)	2.47 (0.18)	1.05 (0.05)	-.119 (0.015)	1.94 (0.14)	0.82 (0.03)	-.135 (0.018)	1.63 (0.12)	0.71 (0.03)
Fold 5	HopCPT-G	-.079 (0.002)	2.34 (0.04)	2.59 (0.05)	-.106 (0.004)	1.72 (0.02)	1.82 (0.02)	-.119 (0.005)	1.39 (0.02)	1.46 (0.01)
	HopCPT-G (PUB Train)	-.035 (0.006)	3.32 (0.09)	1.79 (0.08)	-.057 (0.010)	2.51 (0.07)	1.41 (0.05)	-.073 (0.011)	2.07 (0.06)	1.20 (0.03)
	CMAL	-.128 (0.032)	2.92 (0.27)	1.58 (0.13)	-.170 (0.040)	2.30 (0.20)	1.26 (0.10)	-.189 (0.044)	1.94 (0.17)	1.09 (0.09)



α		0.05			0.10			0.15		
		Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler	Δ Cov	PI-Width	Winkler
Fold 6	HopCPT-G	-.050 (0.005)	2.25 (0.03)	1.58 (0.03)	-.078 (0.006)	1.66 (0.03)	1.19 (0.02)	-.093 (0.007)	1.34 (0.03)	1.00 (0.01)
	HopCPT-G (PUB Train)	.004 (0.011)	3.18 (0.32)	1.28 (0.04)	-.007 (0.016)	2.41 (0.24)	1.01 (0.03)	-.020 (0.018)	1.99 (0.19)	0.87 (0.02)
	CMAL	-.103 (0.021)	2.55 (0.22)	1.18 (0.07)	-.143 (0.023)	2.00 (0.18)	0.93 (0.04)	-.161 (0.022)	1.68 (0.15)	0.81 (0.03)
Fold 7	HopCPT-G	-.050 (0.006)	2.15 (0.04)	1.77 (0.04)	-.065 (0.006)	1.55 (0.03)	1.27 (0.02)	-.072 (0.008)	1.25 (0.02)	1.04 (0.01)
	HopCPT-G (PUB Train)	.000 (0.011)	3.14 (0.33)	1.34 (0.08)	-.001 (0.014)	2.34 (0.23)	1.05 (0.04)	-.001 (0.014)	1.91 (0.18)	0.90 (0.03)
	CMAL	-.105 (0.018)	2.46 (0.23)	1.17 (0.07)	-.142 (0.025)	1.92 (0.17)	0.93 (0.05)	-.163 (0.030)	1.62 (0.14)	0.81 (0.05)
Fold 8	HopCPT-G	-.046 (0.004)	2.15 (0.04)	1.66 (0.03)	-.064 (0.005)	1.58 (0.03)	1.22 (0.02)	-.072 (0.005)	1.28 (0.02)	1.01 (0.01)
	HopCPT-G (PUB Train)	-.002 (0.007)	3.13 (0.22)	1.34 (0.03)	-.010 (0.007)	2.38 (0.17)	1.04 (0.01)	-.020 (0.010)	1.97 (0.14)	0.89 (0.01)
	CMAL	-.123 (0.024)	2.38 (0.09)	1.18 (0.03)	-.161 (0.024)	1.87 (0.08)	0.94 (0.03)	-.180 (0.023)	1.57 (0.07)	0.82 (0.03)
Fold 9	HopCPT-G	-.047 (0.006)	2.09 (0.07)	1.71 (0.08)	-.060 (0.006)	1.53 (0.05)	1.23 (0.04)	-.065 (0.006)	1.23 (0.05)	1.00 (0.03)
	HopCPT-G (PUB Train)	.007 (0.008)	2.92 (0.15)	1.25 (0.04)	.004 (0.017)	2.21 (0.12)	0.97 (0.02)	-.002 (0.023)	1.82 (0.10)	0.83 (0.02)
	CMAL	-.122 (0.031)	2.34 (0.17)	1.15 (0.05)	-.162 (0.034)	1.82 (0.13)	0.90 (0.03)	-.178 (0.033)	1.53 (0.11)	0.77 (0.02)
Fold 10	HopCPT-G	-.079 (0.007)	2.10 (0.07)	1.96 (0.08)	-.100 (0.010)	1.55 (0.06)	1.39 (0.04)	-.109 (0.012)	1.26 (0.05)	1.13 (0.03)
	HopCPT-G (PUB Train))	-.020 (0.011)	3.18 (0.18)	1.42 (0.03)	-.035 (0.012)	2.41 (0.13)	1.12 (0.02)	-.047 (0.014)	1.99 (0.11)	0.96 (0.01)
	CMAL	-.127 (0.031)	2.62 (0.20)	1.24 (0.05)	-.157 (0.031)	2.06 (0.16)	0.96 (0.04)	-.168 (0.030)	1.74 (0.13)	0.82 (0.04)

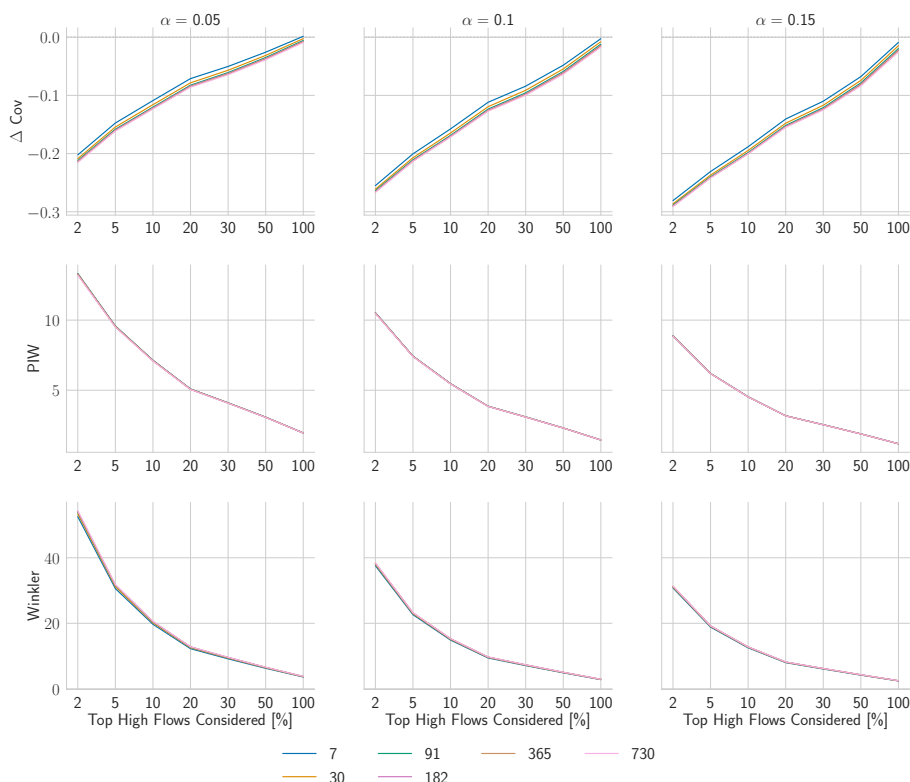


Figure G3. Evaluation metrics for high flows with a varying share of considered time steps. The different runs represent different update frequencies. (setting of Experiment I-b)

510 *Author contributions.* AA, DK, and MG designed all experiments. GN and FK checked the adequacy of the experiments and provided feedback. AA conducted all experiments, the results were analyzed by AA, DK, and MG. GN, DK, AA, and MG designed the structure of the exposition All authors contributed to the writing process with AA, DK, and MG doing the majority of the writing. AA created all figures. DK and SH supervised the manuscript.

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

515 *Acknowledgements.* The ELLIS Unit Linz, the LIT AI Lab, the Institute for Machine Learning, are supported by the Federal State Upper Austria. We thank the projects Medical Cognitive Computing Center (MC3), INCONTROL-RL (FFG-881064), PRIMAL (FFG-873979), S3AI (FFG-872172), DL for GranularFlow (FFG-871302), EPILEPSIA (FFG-892171), AIRI FG 9-N (FWF-36284, FWF-36235), AI4GreenHeatingGrids (FFG- 899943), INTEGRATE (FFG-892418), ELISE (H2020-ICT-2019-3 ID: 951847), Stars4Waters (HORIZON-CL6-2021-CLIMATE-01-01). We thank Audi.JKU Deep Learning Center, TGW LOGISTICS GROUP GMBH, Silicon Austria Labs (SAL), FILL Gesellschaft

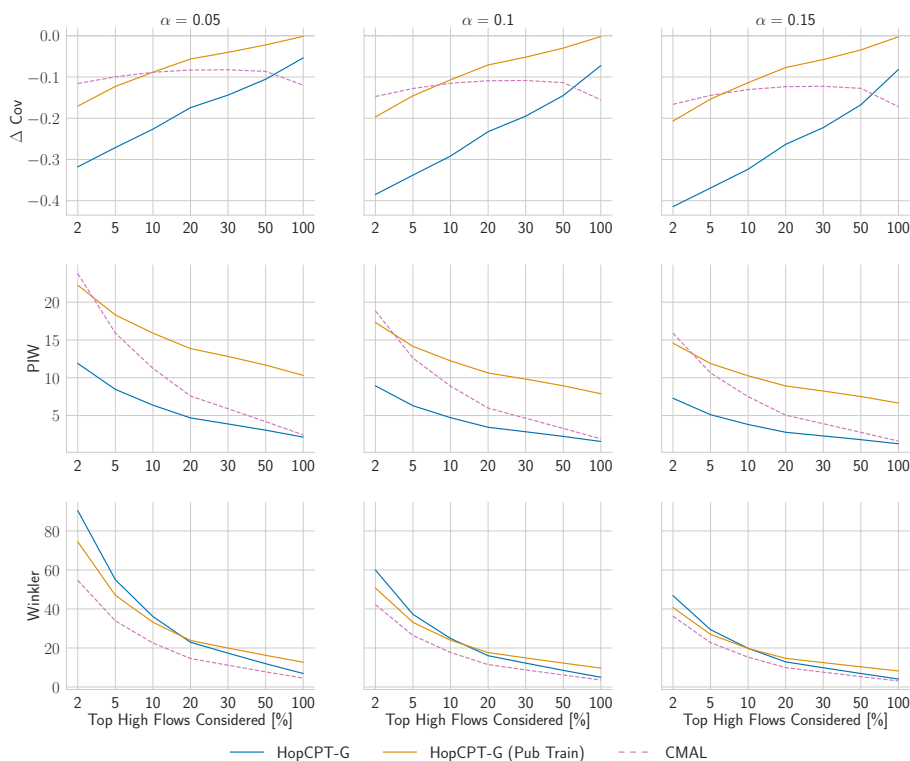


Figure G4. Evaluation metrics for high flows with varying share of considered time steps in the “ungauged basins setting” (setting of Experiment II-b)

520 mbH, Anyline GmbH, Google, ZF Friedrichshafen AG, Robert Bosch GmbH, UCB Biopharma SRL, Merck Healthcare KGaA, Verbund AG, GLS (Univ. Waterloo), Software Competence Center Hagenberg GmbH, Borealis AG, TÜV Austria, Frauscher Sensoric, TRUMPF and the NVIDIA Corporation. We acknowledge EuroHPC Joint Undertaking for awarding us access to LUMI at CSC, Finland. Use as many instances of the pattern LUMI at CSC, Finland as the number of systems awarded via EuroHPC.



References

- 525 Addor, N., Newman, A. J., Mizukami, N., and Clark, M. P.: The CAMELS data set: catchment attributes and meteorology for large-sample studies, *Hydrology and Earth System Sciences (HESS)*, 21, 5293–5313, 2017.
- Angelopoulos, A. N. and Bates, S.: A gentle introduction to conformal prediction and distribution-free uncertainty quantification, arXiv preprint arXiv:2107.07511, 2021.
- Auer, A., Gauch, M., Klotz, D., and Hochreiter, S.: Conformal Prediction for Time Series with Modern Hopfield Networks, in: Thirty-seventh
530 Conference on Neural Information Processing Systems, <https://openreview.net/forum?id=KTRwpWCMsC>, 2023.
- Bertola, M., Blöschl, G., Bohac, M., Borga, M., Castellarin, A., Chirico, G. B., Claps, P., Dallan, E., Danilovich, I., Ganora, D., et al.: Megafloods in Europe can be anticipated from observations in hydrologically similar catchments, *Nature Geoscience*, pp. 1–7, 2023.
- Beven, K.: Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis testing, and communication, *Hydrological Sciences Journal*, 61, 1652–1665, <https://doi.org/10.1080/02626667.2015.1031761>, 2016a.
- 535 Beven, K.: Facets of uncertainty: epistemic uncertainty, non-stationarity, likelihood, hypothesis testing, and communication, *Hydrological Sciences Journal*, 61, 1652–1665, 2016b.
- Beven, K. and Binley, A.: GLUE: 20 years on, *Hydrological processes*, 28, 5897–5918, 2014.
- Bhatnagar, A., Wang, H., Xiong, C., and Bai, Y.: Improved Online Conformal Prediction via Strongly Adaptive Online Learning, in: International Conference on Machine Learning, <https://api.semanticscholar.org/CorpusID:256868761>, 2023.
- 540 Bishop, C. M.: Mixture density networks, Tech. rep., Neural Computing Research Group, 1994.
- Clark, M. P., Wilby, R. L., Gutmann, E. D., Vano, J. A., Gangopadhyay, S., Wood, A. W., Fowler, H. J., Prudhomme, C., Arnold, J. R., and Brekke, L. D.: Characterizing uncertainty of the hydrologic impacts of climate change, *Current climate change reports*, 2, 55–64, 2016.
- Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D.-J., Hartman, R., Herr, H. D., Fresch, M., et al.: The science of NOAA’s operational hydrologic ensemble forecast service, *Bulletin of the American Meteorological Society*, 95, 79–98, 2014.
- 545 Foygel Barber, R., Candes, E. J., Ramdas, A., and Tibshirani, R. J.: Conformal prediction beyond exchangeability, arXiv preprint arXiv:2202.13415, 2022.
- Frame, J. M., Kratzert, F., Klotz, D., Gauch, M., Shalev, G., Gilon, O., Qualls, L. M., Gupta, H. V., and Nearing, G. S.: Deep learning rainfall–runoff predictions of extreme events, *Hydrology and Earth System Sciences*, 26, 3377–3392, 2022.
- Gauch, M., Mai, J., and Lin, J.: The proper care and feeding of CAMELS: How limited training data affects streamflow prediction, *Environmental Modelling & Software*, 135, 104926, 2021.
- 550 Geirhos, R., Jacobsen, J.-H., Michaelis, C., Zemel, R. S., Brendel, W., Bethge, M., and Wichmann, F.: Shortcut learning in deep neural networks, *Nature Machine Intelligence*, 2, 665 – 673, <https://api.semanticscholar.org/CorpusID:215786368>, 2020.
- Gibbs, I. and Candes, E. J.: Adaptive conformal inference under distribution shift, *Advances in Neural Information Processing Systems*, 34, 1660–1672, 2021.
- 555 Gupta, A. and Govindaraju, R. S.: Uncertainty quantification in watershed hydrology: Which method to use?, *Journal of Hydrology*, 616, 128749, 2023.
- Haines, A., Finlayson, B., and McMahon, T.: A Global Classification of River Regimes, *Applied Geography*, 8, 255–272, [https://doi.org/10.1016/0143-6228\(88\)90035-5](https://doi.org/10.1016/0143-6228(88)90035-5), 1988.
- Hamilton, J. D.: Analysis of time series subject to changes in regime, *Journal of econometrics*, 45, 39–70, 1990.



- 560 Harris, N. M., Gurnell, A. M., Hannah, D. M., and Petts, G. E.: Classification of river regimes: a context for hydroecology, *Hydrological Processes*, 14, 2831–2848, [https://doi.org/https://doi.org/10.1002/1099-1085\(200011/12\)14:16/17<2831::AID-HYP122>3.0.CO;2-O](https://doi.org/https://doi.org/10.1002/1099-1085(200011/12)14:16/17<2831::AID-HYP122>3.0.CO;2-O), 2000.
- Kavetski, D., Kuczera, G., and Franks, S. W.: Bayesian analysis of input uncertainty in hydrological modeling: 2. Application, *Water resources research*, 42, 2006.
- 565 Klotz, D., Kratzert, F., Gauch, M., Keefe Sampson, A., Brandstetter, J., Klambauer, G., Hochreiter, S., and Nearing, G.: Uncertainty estimation with deep learning for rainfall–runoff modeling, *Hydrology and Earth System Sciences*, 26, 1673–1693, <https://doi.org/10.5194/hess-26-1673-2022>, 2022.
- Koutsoyiannis, D. and Montanari, A.: Bluecat: A Local Uncertainty Estimator for Deterministic Simulations and Predictions, *Water Resources Research*, 58, e2021WR031 215, <https://doi.org/10.1029/2021WR031215>, 2022.
- 570 Kratzert, F.: CAMELS Extended Maurer Forcing Data, <https://doi.org/10.4211/hs.17c896843cf940339c3c3496d0c1c077>, 2019.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A. K., Hochreiter, S., and Nearing, G. S.: Toward Improved Predictions in Ungauged Basins: Exploiting the Power of Machine Learning, *Water Resources Research*, 55, 11 344–11 354, <https://doi.org/10.1029/2019WR026065>, 2019a.
- Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., and Nearing, G.: Towards learning universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets, *Hydrology and Earth System Sciences*, 23, 5089–5110, <https://doi.org/10.5194/hess-23-5089-2019>, 2019b.
- Kratzert, F., Klotz, D., Hochreiter, S., and Nearing, G. S.: A note on leveraging synergy in multiple meteorological data sets with deep learning for rainfall–runoff modeling, *Hydrology and Earth System Sciences*, 25, 2685–2703, 2021.
- Krzysztofowicz, R.: The case for probabilistic forecasting in hydrology, *Journal of hydrology*, 249, 2–9, 2001.
- 580 Lapuschkin, S., Wäldchen, S., Binder, A., Montavon, G., Samek, W., and Müller, K.-R.: Unmasking Clever Hans predictors and assessing what machines really learn, *Nature communications*, 10, 1096, 2019.
- 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, 2017.
- Mai, J., Shen, H., Tolson, B. A., Gaborit, É., Arsenault, R., Craig, J. R., Fortin, V., Fry, L. M., Gauch, M., Klotz, D., et al.: The great lakes runoff intercomparison project phase 4: the great lakes (GRIP-GL), *Hydrology and Earth System Sciences*, 26, 3537–3572, 2022.
- 585 Montanari, A. and Koutsoyiannis, D.: A blueprint for process-based modeling of uncertain hydrological systems, *Water Resources Research*, 48, 2012.
- Montero-Manso, P. and Hyndman, R. J.: Principles and algorithms for forecasting groups of time series: Locality and globality, *International Journal of Forecasting*, 37, 1632–1653, 2021.
- 590 Nearing, G. S., Tian, Y., Gupta, H. V., Clark, M. P., Harrison, K. W., and Weijs, S. V.: A philosophical basis for hydrological uncertainty, *Hydrological Sciences Journal*, 61, 1666–1678, 2016.
- Newman, A. J., Clark, M. P., Sampson, K., Wood, A., Hay, L. E., Bock, A., Viger, R. J., Blodgett, D., Brekke, L., Arnold, J. R., Hopson, T., and Duan, Q.: Development of a large-sample watershed-scale hydrometeorological data set for the contiguous USA: data set characteristics and assessment of regional variability in hydrologic model performance, *Hydrology and Earth System Sciences*, 19, 209–223, <https://doi.org/10.5194/hess-19-209-2015>, 2015.
- 595 Newman, A. J., Mizukami, N., Clark, M. P., Wood, A. W., Nijssen, B., and Nearing, G.: Benchmarking of a Physically Based Hydrologic Model, *Journal of Hydrometeorology*, 18, 2215–2225, <https://doi.org/10.1175/JHM-D-16-0284.1>, 2017.



- Quandt, R. E.: The estimation of the parameters of a linear regression system obeying two separate regimes, *Journal of the American Statistical Association*, 53, 873–880, 1958.
- 600 Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlović, M., Sandve, G. K., Greiff, V., Kreil, D., Kopp, M., Klambauer, G., Brandstetter, J., and Hochreiter, S.: Hopfield networks is all you need, in: 9th International Conference on Learning Representations (ICLR), 2021.
- Rozos, E., Koutsoyiannis, D., and Montanari, A.: KNN vs. Bluecat — Machine Learning vs. Classical Statistics, *Hydrology*, 9, <https://doi.org/10.3390/hydrology9060101>, 2022.
- 605 Samaniego, L., Kumar, R., Thober, S., Rakovec, O., Zink, M., Wanders, N., Eisner, S., Müller Schmied, H., Sutanudjaja, E. H., Warrach-Sagi, K., et al.: Toward seamless hydrologic predictions across spatial scales, *Hydrology and Earth System Sciences*, 21, 4323–4346, 2017.
- Schwappe, R., Thober, S., Müller, S., Kelbling, M., Kumar, R., Attinger, S., and Samaniego, L.: MPR 1.0: a stand-alone multiscale parameter regionalization tool for improved parameter estimation of land surface models, *Geoscientific Model Development*, 15, 859–882, 2022.
- Shrestha, D. L. and Solomatine, D. P.: Data-driven approaches for estimating uncertainty in rainfall-runoff modelling, *International Journal of River Basin Management*, 6, 109–122, 2008.
- 610 Vovk, V., Gammerman, A., and Shafer, G.: *Algorithmic learning in a random world*, Springer Science & Business Media, 2005.
- Xu, C. and Xie, Y.: Conformal prediction for time series, arXiv preprint arXiv:2010.09107, 2022a.
- Xu, C. and Xie, Y.: Sequential Predictive Conformal Inference for Time Series, arXiv preprint arXiv:2212.03463, 2022b.
- Zaffran, M., Dieuleveut, A., F'eron, O., Goude, Y., and Josse, J.: Adaptive Conformal Predictions for Time Series, in: International Conference on Machine Learning, <https://api.semanticscholar.org/CorpusID:246863519>, 2022.
- 615