

This is the accepted manuscript version of the contribution published as:

Romero-Cuellar, J., Arabzadeh, R., Craig, J.R., Tolson, B.A., **Mai, J.** (2024):
A multi-model evaluation of probabilistic streamflow predictions via residual error modelling
J. Hydrol. **635**, art. 131152

The publisher's version is available at:

<https://doi.org/10.1016/j.jhydrol.2024.131152>

Journal Pre-proofs

Research papers

A multi-model evaluation of probabilistic streamflow predictions via residual error modelling

Jonathan Romero-Cuellar, Rezgar Arabzadeh, James R. Craig, Bryan A. Tolson, Juliane Mai

PII: S0022-1694(24)00547-X
DOI: <https://doi.org/10.1016/j.jhydrol.2024.131152>
Reference: HYDROL 131152

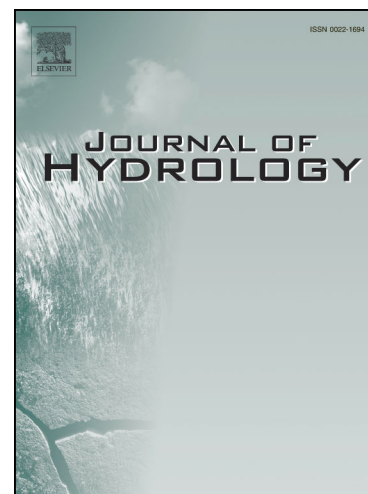
To appear in: *Journal of Hydrology*

Received Date: 13 July 2023
Revised Date: 8 December 2023
Accepted Date: 27 December 2023

Please cite this article as: Romero-Cuellar, J., Arabzadeh, R., Craig, J.R., Tolson, B.A., Mai, J., A multi-model evaluation of probabilistic streamflow predictions via residual error modelling, *Journal of Hydrology* (2024), doi: <https://doi.org/10.1016/j.jhydrol.2024.131152>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2024 Published by Elsevier B.V.



1 **A Multi-model Evaluation of Probabilistic Streamflow Predictions via**
2 **Residual Error Modelling**

3 **Jonathan Romero-Cuellar¹, Rezgar Arabzadeh¹, James R. Craig¹, Bryan A. Tolson¹, and Juliane Mai^{2,3}**

4 ¹ Department of Civil and Environmental Engineering, University of Waterloo, Waterloo, ON,
5 Canada.

6 ² Computational Hydrosystems, Helmholtz Centre for Environmental Research - UFZ, Leipzig,
7 Saxony, Germany.

8 ³ Center for Scalable Data Analytics and Artificial Intelligence - ScaDS.AI, Leipzig, Saxony,
9 Germany.

10 Corresponding author: Jonathan Romero-Cuellar (jromeroc@uwaterloo.ca)

- 11 • **Keywords:** Uncertainty analysis; Probabilistic prediction; Residual error; Streamflow;
12 Hydrological modeling; Postprocessing method

13

14

15

16

17 Abstract

18 Probabilistic streamflow predictions are valuable tools for predictive uncertainty estimation,
19 hydrologic risk management, and support for decision-making in water resources. Usually,
20 predictive uncertainty quantification is developed and assessed using only a single hydrological
21 model, making it difficult to generalize to other model configurations. To tackle this issue, we
22 assess changes in the model performance ranking of diverse streamflow models by applying a
23 residual error model post-processing approach to multiple basins and multiple models. This
24 assessment employed 141 basins from the Great Lakes watershed covering the USA and Canada,
25 and 13 diverse streamflow models, which are evaluated using deterministic and probabilistic
26 performance metrics. As the first study to implement probabilistic methods to diverse streamflow
27 models applied to a multitude of basins, the analysis here examines the dependence of
28 probabilistic streamflow estimation quality on model quality. Our findings show that streamflow
29 model choice influences the robustness of probabilistic predictions. It was found that moving
30 from deterministic to probabilistic predictions using a post-processing approach does not change
31 the streamflow model performance ranking for the best and worst deterministic models, but
32 models of intermediate rank in deterministic evaluation do not have consistent ranking when
33 evaluated in probabilistic mode. Post-processing residual errors of long short-term memory
34 (LSTM) network models are consistently the best-performing model in terms of deterministic
35 and probabilistic metrics. This study highlights the significance of combining deterministic
36 streamflow model predictions with residual error models for improving the quality and
37 increasing the value of hydrological predictions, quantifying uncertainty, and facilitating
38 decision-making in operational water management. It also clarifies the degree to which
39 probabilistic predictions depend upon good model performance and can compensate for poor
40 model performance.

41 1 Introduction

42 Although deterministic model predictions offer information to support hydrologic
43 planning, design, and decision-making in environmental and water resource applications, they do
44 not provide uncertainty estimates needed for hydrologic risk management. Deterministic models
45 omit the model uncertainty associated with simulated responses. Throughout this paper, the term
46 “deterministic” model prediction will refer to models that provide only one value of streamflow
47 at any given time while probabilistic refer to models that provide estimates of streamflow in the
48 form of a probability density. These model predictions suffer from many sources of uncertainty,
49 including input data uncertainty, parameter uncertainty, and model structural uncertainty (Moges
50 et al., 2021; Gupta & Govindaraju, 2023). Uncertainty quantification has become an integral part
51 of risk assessment, and hydrological models have been used to provide information for
52 policymakers who are increasingly demanding that model results include estimates of
53 uncertainty (Bastola et al., 2012; Dias et al., 2020). Unfortunately, many of the most widely used
54 models in hydrology (e.g., SWAT, GR4J, HBV, or VIC) are deterministic and provide only a
55 single trace of streamflow output (point predictions). In contrast, probabilistic predictions
56 provide the uncertainty estimation associated with a dynamic basin system by generating an
57 ensemble comprised of multiple sets of streamflow predictions (Farmer & Vogel, 2016). Given
58 the benefits of uncertainty characterization, it is important to convert deterministic predictions to
59 probabilistic ones that can support hydrological risk management (e.g., Vogel, 2017; Reggiani et
60 al., 2022; Shabestanipour et al., 2023).

61 Deterministic predictions from a given hydrological model can be converted to
62 probabilistic predictions using a post-processing strategy to add a residual error model (REM) to
63 the deterministic predictions (e.g., Evin et al., 2014). In general, the post-processing strategy is
64 implemented in two steps. First, hydrological model parameters are estimated using an objective
65 function and calibration algorithm. Second, a REM form is selected, and the parameters of REM
66 are inferred separately. REMs are statistical models to represent the relationship between model
67 outputs (e.g., streamflow) and observations (Ye et al., 2014). Because they characterize the
68 expected difference between model output and observation, REMs include the merged impact of
69 several sources of uncertainties, such as input, parameter, and model structure (Schoups &
70 Vrugt, 2010; Sorooshian et al., 1983). In addition, REMs correct biases in hydrological models
71 (W. Li et al., 2017). The post-processing strategy differs from the classical joint inference
72 approach, which estimates all parameters (hydrologic and error model) simultaneously using a
73 single likelihood function (Bates & Campbell, 2001; Evin et al., 2014; Smith et al., 2015). The
74 residual error model post-processor (REM-PP) is an attractive approach for converting
75 deterministic into probabilistic predictions for two reasons. First, REM-PP avoids the need to re-
76 calibrate the hydrological model. Second, REM-PP eliminates challenging interactions between
77 water balance parameters and REM parameters (Evin et al., 2014).

78 Several studies have improved hydrological predictions using the REM strategy (Evin et
79 al., 2014; Koutsoyiannis & Montanari, 2022; Kuczera et al., 2006; M. Li et al., 2016; McInerney
80 et al., 2017; Reichert & Schuwirth, 2012; Todini, 2008; Wang et al., 2012; Zhao et al., 2015). In
81 general, two approaches exist to describe residual model errors: (i) a lumped approach that
82 aggregates all sources of uncertainty into an individual model error term (Montanari &
83 Koutsoyiannis, 2012; Vrugt et al., 2022) and (ii) decomposition approaches, which attempt to
84 disentangle the sources of uncertainty (Ajami et al., 2007; Kuczera et al., 2006). From a
85 theoretical perspective, decomposition approaches are more versatile as they provide a
86 comprehensive assessment of uncertainty estimation and identify specific sources of
87 uncertainties. However, these approaches are complex to implement, demand more experience,
88 and are typically more time-consuming in terms of both computational cost and user time. For
89 these reasons, decomposition approaches tend not to be applied in practice. In contrast, lumped
90 approaches evaluate the total uncertainty as estimated from the model-to-observation misfit.
91 These approaches are more adequate for operational hydrology because they are less time-
92 consuming, computationally efficient, and require little experience.

93 The post-processor method has been found to be robust in several previous studies (e.g.,
94 Evin et al., 2014; M. Li et al., 2016). On a daily scale, REMs showed a very good performance
95 using power transformations (McInerney et al., 2017; Todini, 2008; Wang et al., 2012). In the
96 post-processing strategy, the quality of uncertainty estimation has been shown (for single
97 models) to depend on the residual error model and the objective function used to calibrate the
98 hydrological model (Hunter et al., 2021), and characteristics of data (Jiang et al., 2019). The
99 proportion of parametric uncertainty to total predictive uncertainty in streamflow is mostly
100 modest (Kuczera et al., 2006; Sun et al., 2017; Yang et al., 2007) when parsimonious
101 hydrological models are calibrated to long observed time series using a residual error model. In
102 this context, predictions in operational hydrology habitually omit parameter uncertainty and
103 focus on residual errors (Engeland & Steinsland, 2014; McInerney et al., 2017). This was the
104 method employed in this research, where deterministic calibration was assumed solely via

105 optimization. In other words, the parameter estimation method was applied without trying to
106 compute parametric uncertainty.

107 Large-sample hydrologic model intercomparison studies provide a unique opportunity to
108 rank models, especially using standardized experiments. The Great Lakes region is a
109 transboundary domain between the USA and Canada and has been the subject of several model
110 intercomparison studies (Fry et al., 2014; Gaborit et al., 2017; Mai et al., 2021, 2022).
111 Particularly, the Great Lakes Runoff Intercomparison Project phase 4: the Great Lakes (GRIP-
112 GL) delivered a standardized experimental setup using similar geophysical datasets, forcings, a
113 common routing product, and identical locations of performance assessment across the 10^6 km²
114 study domain. This project evaluated 13 diverse streamflow models involving a large range of
115 model types from machine-learning-based, lumped, subbasin-based, and gridded models that are
116 either locally, regionally or globally calibrated (Mai et al., 2022). Like most past hydrological
117 modelling intercomparison studies (e.g., Lake Michigan (GRIP-M; Fry et al., 2014), Lake
118 Ontario (GRIP-O; Gaborit et al., 2017), and Lake Erie (GRIP-E; Mai et al., 2021)), GRIP-GL
119 focused on deterministic model performance and did not consider the uncertainty of model
120 streamflow predictions. Therefore, the GRIP-GL project dataset provides a unique chance to
121 assess the relationship between model deterministic performance (assessed in GRIP-GL) and
122 probabilistic performance (assessed here). This is the first study to the authors knowledge to
123 apply uncertainty analysis to such a large variety of both models (13) and watersheds (141).

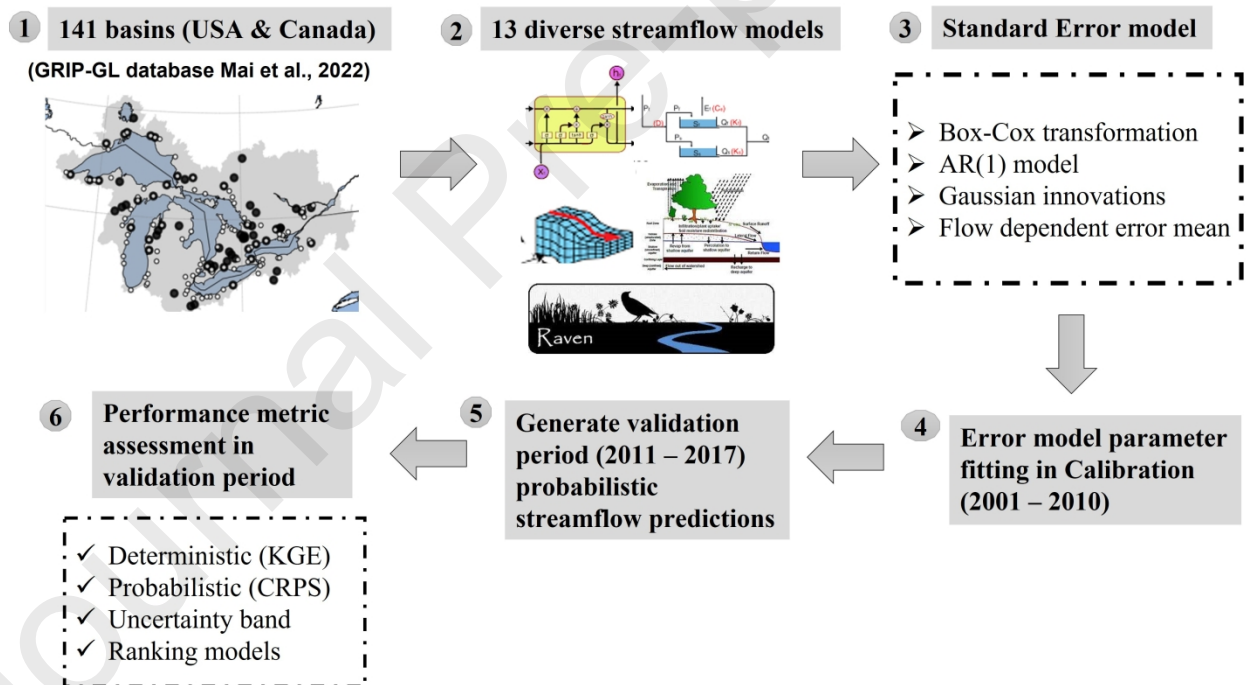
124 In contrast to GRIP-GL project, much of the research on hydrologic model calibration
125 under uncertainty has been developed and assessed using only a single hydrological model,
126 making it difficult to generalize. Certainly, the majority of previous studies was limited to using
127 only one rainfall-runoff model (Evin et al., 2014; Jiang et al., 2019; D. Li et al., 2021; M. Li et
128 al., 2016; McInerney et al., 2017, 2020, 2021; Renard et al., 2010; Salamon & Feyen, 2009;
129 Thyer et al., 2009; Todini, 2008; Weerts et al., 2011; Woldemeskel et al., 2018). To our best
130 knowledge, no studies have evaluated the impact of streamflow model choice on probabilistic
131 predictions using multiple models and basins. This is the first essential research gap addressed in
132 this investigation. Few studies have investigated the predictive uncertainty quantification using
133 different hydrological models (e.g., Ye et al., (2014) and McInerney et al., (2017)). Ye et al.,
134 (2014) compare the effects of post-processing and model calibration on improving water
135 forecasts under different hydroclimatic conditions and models. However, their study reduced
136 comparisons to a deterministic perspective (comparing ensemble means for example). In
137 addition, although they looked at results under multiple hydroclimatic conditions and seven
138 hydrologic models, they did not evaluate the influence of the hydrologic model on the
139 probabilistic prediction (e.g., no comparison of hydrologic model quality with and without post-
140 processing). McInerney et al., (2017) and Morawietz et al., (2011) assessed several residual error
141 models but used only two hydrological models, mainly HBV and GR4J. Neither study focused
142 their analysis on a comparative assessment of relative model performance under deterministic
143 versus probabilistic predictions, nor did they include machine learning models for streamflow
144 like long short-term memory (LSTM) networks in their experimental design.

145 This study examines the effect of streamflow model choice on probabilistic predictions
146 using a residual error model post-processing (REM-PP) approach to multiple streamflow models
147 and basins. We addressed the following question: How does streamflow model choice impact
148 probabilistic predictions? Finally, this study provides new insights into the influence of

149 hydrologic model choice on estimations of probabilistic streamflow predictions and changes in a
 150 model performance ranking.

151 2 Material and methods

152 We apply REM-PP to time series of streamflow in 141 calibration basins simulated using
 153 13 diverse streamflow models. We adopt the Box-Cox (BC0.2) transformation residual error
 154 model as the REM-PP because it is an established practice in the hydrology literature (e.g.,
 155 McInerney et al., (2017), McInerney et al., (2018), and Hunter et al., (2021)). McInerney et al.,
 156 (2017) showed that this error model is effective for probabilistic predictions and uncertainty
 157 estimation in that it helps to remove heteroscedasticity and to deal with time series
 158 autocorrelation. In fact, McInerney et al., (2017) evaluated various error schemes, including
 159 standard and weighted least square, the Box-Cox transformations (with fixed and calibrated
 160 power parameters), and the log-sinh transformation across 23 basins. They recommended the
 161 Box-Cox approach due to it having the highest average performance and achieving high-quality
 162 daily probabilistic predictions. Figure 1 shows a graphical summary of the processes and
 163 methods followed in this study, and a detailed description of the REM-PP is provided in the
 164 following section.



165
 166 **Figure 1.** Flowchart for methods employed in this study.

167 2.1 Streamflow post-processing model

168 Consider a deterministic hydrological model (h) that represent streamflow predictions
 169 time series ($q_t^{\theta_h}$) with hydrological parameters (θ_h) and forcings (x).

$$q_t^{\theta_h} = h(\theta_h; x) \quad (1)$$

170 Then, define a probabilistic model of streamflow (Q_t), which is formulated from a
 171 deterministic hydrological model (h) by adding a random residual error term η_t intended to
 172 characterize the predictive uncertainty due to the combined effect of uncertainties in the input
 173 data, model parameters, and model structure.

$$Q_t(\theta, x) = f(q_t^{\theta_h}, \eta_t^{\theta_\eta}) \quad (2)$$

174 where θ is comprised of hydrological model parameters (θ_h) and error model parameters (θ_η).
 175 The streamflow post-processing procedure used in this study consists of fitting a statistical model
 176 to the streamflow residual errors, defined by the differences between the observed and predicted
 177 daily streamflow time series. Residual errors of hydrological models are commonly
 178 heteroscedastic (i.e., have larger errors in large flows), asymmetrical, non-Gaussian, and
 179 persistent (i.e., often possess several days of consecutive errors with the same sign and similar
 180 magnitude) (Bates & Campbell, 2001; Smith et al., 2015; Sorooshian & Dracup, 1980). In many
 181 cases, residual error models perform well at the daily scale using the Box-Cox transformation
 182 (McInerney et al., 2017) and a first-order autoregressive AR(1) model (Evin et al., 2014). The
 183 Box-Cox transformation is intended to reduce the degree of heteroscedasticity and asymmetry in
 184 the error distribution, while the AR model is used to minimize the presence of persistence.

185 The probabilistic model Q_t is constructed by linking a deterministic component $q_t^{\theta_h}$ and a
 186 random residual error component η_t . Here, we formulate the residual error model as additive in
 187 transformed space,

$$z(Q_t; \theta_z) = z(q_t^{\theta_h}; \theta_z) + \eta_t \quad (3)$$

188 where z is a transformation function with parameters θ_z . The deterministic discharge time series
 189 $q_t^{\theta_h}$ is generated by a hydrological model, which is a function of hydrological model parameters
 190 θ_h and forcings x_t . We use the Box-Cox transformation to reduce the heteroscedasticity and
 191 skewness in the residuals,

$$z(q_t; \theta_z) = \begin{cases} \frac{(q_t + A)^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \log(q_t + A), & \text{otherwise} \end{cases} \quad (4)$$

192 with parameters $\theta_z = \{A, \lambda\}$, where A is a shift or offset parameter and λ is a power parameter
 193 (Box & Cox, 1964). Following the recommendation of McInerney et al., (2017), the parameters

194 are fixed a priori at $\lambda = 0.2$ and $A = 0$. The temporal persistence of residual errors is modelled as
 195 a first-order autoregressive (AR1) model,

$$\eta_t = \mu_t + \phi_\eta(\eta_{t-1} - \mu_{t-1}) + y_t \quad (5)$$

196 where ϕ_η is the lag-1 autoregressive parameter, μ_t is the mean of residuals at time t , and y_t
 197 represents the innovation, which is a random component or noise term, at time t . The innovations
 198 were assumed to follow a Gaussian distribution with a mean μ_t and standard deviation σ_y^2 ,

$$y_t \sim \mathcal{N}(\mu_t, \sigma_y^2) \quad (6)$$

199 In this study, we consider that the residual mean μ_t is flow-dependent. We assumed a
 200 linear dependence with respect to the transformed streamflow (Jiang et al., 2019; M. Li et al.,
 201 2017; Romero-Cuellar et al., 2019),

$$\mu_t = \alpha + \beta z(\theta_z, q_t^{\theta_h}) \quad (7)$$

202 where parameters α and β symbolize the intercept and slope respectively. The parameters of the
 203 residual error model are indicated as $\theta_\eta = \{\theta_\mu, \sigma_y, \phi_\eta\}$, where θ_μ represents parameters of the
 204 model for the mean μ_t . These parameters are grouped as $\theta_\mu = \{\alpha, \beta\}$. The full parameter set of
 205 the probabilistic model Q_t is $\theta = \{\theta_h, \theta_z, \theta_\eta\}$. The REM-PP was implemented using the
 206 ‘‘ProbPred’’ R package (Hunter et al., 2021).

207 2.1.1 Parameter estimation

208 The hydrologic model parameters θ were inferred from observed model forcings $\tilde{x} =$
 209 $\{\tilde{x}_t; t = 1, \dots, N_t\}$ and observed streamflow time series $\tilde{q} = \{\tilde{q}_t; t = 1, \dots, N_t\}$, where N_t was the total
 210 number of time steps. We used a two-stage post-processing approach for parameter estimation of
 211 hydrological models and then the residual error model. In stage 1, the hydrological model
 212 parameters θ_h were estimated previously by the modelling teams in GRIP-GL, typically by using
 213 optimization algorithms (see Supplements of Mai et al., (2022) for full details for each model).
 214 Some of the models were calibrated individually at each site, such that parameters were basin-
 215 specific (local optimization). Other models were calibrated simultaneously across all sites (global
 216 optimization). In stage 2, we estimated the residual error model parameters θ_η using the method-
 217 of-moments (Hazelton, 2011) (see the appendix A of Hunter et al., (2021) for details on the
 218 method-of-moments used in this study). Stage 2 computations are extremely fast because it
 219 works only with observed data and optimal streamflow predictions from stage 1. Therefore, it
 220 does not require additional hydrological model runs. The calibration period used was from
 221 January 2001 to December 2010, while the validation period was from January 2011 to
 222 December 2017.

223 2.1.2 Probabilistic predictions

224 The predictive distribution of streamflow is delivered by the probabilistic model Q_t
 225 combining the deterministic hydrological model prediction and residual error probability model.

226 At time step t , assuming the deterministic streamflow prediction $q_t^{\hat{\theta}_h}$ has been calculated, a
 227 sample of size 1 from the predictive streamflow distribution, $q_t^{pred(s)}$, is produced as follows:

228 1) Sample random components (innovations) from a Gaussian distribution

$$y_t^{(s)} \leftarrow \mathcal{N}(\hat{\mu}_t, \hat{\sigma}_y^2) \quad (8)$$

229 2) Estimate residuals using Equation (5)

$$\eta_t^{(s)} = \hat{\mu}_t + \hat{\phi}_\eta(\eta_{t-1}^{(s)} - \hat{\mu}_{t-1}) + y_t^{(s)} \quad (9)$$

230 Note that for $t = 1$, we sample $\eta_1^{(s)} \leftarrow \mathcal{N}(0, \hat{\sigma}_y^2)$

231 3) Employ the inverse transformation

$$q_t^{pred(s)} = z^{-1} \left[z(q_t^{\hat{\theta}_h}; \theta_z) + \eta_t^{(s)}; \theta_z \right] \quad (10)$$

232 The predictive distribution on day t can be characterized by repeating the above sampling
 233 steps s times and then the complete sample of size s characterizing the predictive streamflow
 234 distribution on day t is

$$q^{pred} = \{q_t^{pred(s)}; t = 1, \dots, T; s = 1, \dots, N_s\} \quad (11)$$

235 where N_s is the number of samples.

236 In general, this procedure describes a post-processor approach for predictive uncertainty
 237 estimation. It provides practitioners with a straightforward method to convert deterministic
 238 predictions to probabilistic ones using previously calibrated hydrological models (McInerney et
 239 al., 2018).

240 2.2 Data, streamflow models, and study domain

241 Data used in this study were simulated and observed daily streamflows from the Great
 242 Lakes Runoff Intercomparison Project phase 4: the Great Lakes (GRIP-GL) database (Mai et al.,

243 2022). The case study considers a subset of 141 calibration gauges of the GRIP-GL database.
 244 The study period consists of a warmup period (2000), a calibration period from January 2001 to
 245 December 2010, and a validation period from January 2011 to December 2017. Note that the
 246 number of gauges used in this study was different from the GRIP-GL database because we did
 247 not use spatial validation gauges. The GRIP-GL database contains 13 different streamflow
 248 models (Table 1).

249 The 13 streamflow models were grouped according to their main calibration strategy. The
 250 first group is the machine-learning-based model, which used a global calibration strategy. The
 251 global calibration means that the model was trained for all calibration stations at the same time,
 252 resulting in a single trained model for the entire study domain. The second group is comprised of
 253 the seven models that are locally calibrated, which means that the models were trained for each
 254 of the calibration stations individually. This leads to one calibrated model setup per gauging
 255 station basin. The third group is the five models that followed a regional calibration strategy,
 256 which means simultaneously calibrating model parameters to all calibration stations within a
 257 sub-region of the entire study region (rather than estimating parameters for each gauged basin
 258 individually).

259 The streamflow model parameters were generally estimated in Mai et al. (2022) (in stage
 260 1) by optimization algorithms which minimize the aggregated differences between observed and
 261 simulated streamflow values. The optimization algorithms varied between models. Following
 262 this GRIP-GL split sample, all residual error models here are calibrated on the same calibration
 263 period and all results are analysed for the validation period. Detailed model descriptions can be
 264 found in Mai et al., (2022).

265 **Table 1.** Streamflow model structures investigated in the study. Models were set up, calibrated,
 266 and outputs derived in GRIP-GL (Mai et al., 2022).

Model ID	Model name	Calibration Strategy	Model Type	Spatial resolution
LSTM.lumped	Long Short-Term Memory network	Global	Machine learning	Basins
LBRM.CC.lumped	The Large Basin Runoff Model	Local	Conceptual model	Basins
HYMOD2.lumped	HYMOD2	Local	Conceptual model	Basins
GR4J.lumped	mode'le du Ge'nie Rural a' 4 parametres Journalier	Local	Conceptual model	Basins

HMETs.lumped	Hydrological Model of École de technologie supérieure	Local	Conceptual model	Basins
Blended.lumped	Blended aggregated	Local	Conceptual model with weights	Basins
Blended.Raven	Blended semi-distributed	Local	Conceptual model with weights	Subbasins
VIC.Raven	Variable Infiltration Capacity	Local	Water and energy balance	Grid
SWAT.Raven	The Soil and Water Assessment Tool	Regional	Physically based model	Subbasins
WATFLOOD.Raven	WATFLOOD	Regional	Physically based model	Grid
MESH.CLASS.Raven	Modélisation Environnementale communautaire-Surface and Hydrology	Regional	Conceptual Hydrology - Land Surface Model	Grid
MESH.SVS.Raven	Modélisation Environnementale communautaire-Surface and Hydrology	Regional	Conceptual Hydrology - Land Surface Model	Grid
GEM.Hydro.Watroute	GEM-Hydro	Regional	Physically based model	Grid

267

268 2.3 Performance metrics

269 Various deterministic and probabilistic performance metrics are evaluated for the
 270 validation period (1 January 2011 – 31 December 2017). Deterministic predictions were
 271 evaluated using the Kling-Gupta efficiency (KGE) (Gupta et al., 2009), while probabilistic
 272 predictions were evaluated using the mean continuous ranked probability score (CRPS)
 273 (Hersbach, 2000).

274 2.3.1 Kling-Gupta efficiency (KGE)

275 KGE was developed by decomposing the Nash-Sutcliffe efficiency (NSE) (Nash &
 276 Sutcliffe, 1970) into three parts, that is, correlation, mean bias, and variability bias:

$$KGE = 1 - \sqrt{(r - 1)^2 + (\beta - 1)^2 + (\gamma - 1)^2} \quad (12)$$

277 where r symbolises the correlation between predicted and observed datasets, β represents the
 278 ratio of prediction mean to observation mean, and γ represents the ratio of standard deviation of
 279 predictions to observations. KGE is a positively oriented score with an ideal value of 1.

280 2.3.2 Continuous Ranked Probability Score (CRPS)

281 The CRPS is a familiar metric to assess the complete performance of probabilistic
 282 predictions. It measures the quadratic difference between the predicted and the cumulative
 283 distribution function (CDF) of probabilistic predictions. The CRPS has the dimension of the
 284 predicted variable, i.e., CRPS' units are the same as area-normalized streamflow (mm d^{-1}). In
 285 addition, the CRPS can be decomposed into a reliability part and a sharpness part (Hersbach,
 286 2000). For deterministic predictions, the CRPS is equivalent to the mean absolute error (MAE).
 287 In practice, the mean CRPS is defined as the average CRPS over the entire evaluation period
 288 because the CRPS measures the prediction for a single time step. We compute the CRPS from
 289 the empirical CDF of probabilistic predictions.

$$CRPS_i = \int_{-\infty}^{\infty} (F(q_t) - H_i(q_t \geq \bar{q}_t))^2 dq_t \quad (13)$$

290 where $F(q_t)$ symbolises the cdf of the simulated flow predictive distribution at time t and H_i is
 291 the Heaviside step function, which equal 1 when the predictions are greater than the observations
 292 and equals 0 otherwise. q_t is the simulated flow and \bar{q}_t is the observed flow at time t . Note that
 293 we evaluate and report the mean $CRPS_i$ across all time steps in the validation period. In general,
 294 a high-quality prediction should have a CRPS as low as possible. The CRPS was calculated
 295 using the “scoringRules” R package (Jordan et al., 2019).

296 2.4 Evaluation of performance differences between models

297 Statistical significance testing was used to test if the streamflow models had better or
 298 worse performance metrics over the range of basins. Practical significance testing assesses
 299 whether the difference between performance metric values is large enough to represent
 300 substantial differences in a practical and necessarily subjective sense (Vaske et al., 2010). We
 301 used practical significance testing to evaluate differences in performance metrics. We decided
 302 that a difference by more than 10% of the change metric value for the hydrological models was
 303 practically significant. The practical significance testing is implemented using the method of
 304 McInerney et al., (2019). The pairwise Wilcoxon signed rank test (Bauer, 1972) was used to
 305 compare between distributions of a performance metric over all basins.

306 2.5 Normalised performance metrics

307 To facilitate a comparison between deterministic and probabilistic performance metrics in
 308 Section 3.3, KGE and CRPS were normalised for each basin according to this equation:

$$z_{KGE_i} = \frac{KGE_i - \min(KGE)}{\max(KGE) - \min(KGE)}; \quad (14)$$

$$z_{CRPS_i} = 1 - \left(\frac{CRPS_i - \min(CRPS)}{\max(CRPS) - \min(CRPS)} \right)$$

309 where KGE_i , $\min(KGE)$, and $\max(KGE)$ is the KGE metric value for model i , the minimum
 310 KGE across all models applied to that basin, and the maximum KGE across all models applied to
 311 that basin, respectively. Notation for CRPS is the same except note that CRPS refers to the
 312 average CRPS across all time steps. After we normalised the KGE and CRPS by basin, we
 313 compute the median for each model. Performance metrics range are $[0, 1]$, where 1 is the best
 314 and 0 is the worst.

315 **3 Results**

316 3.1. Representative example (single station)

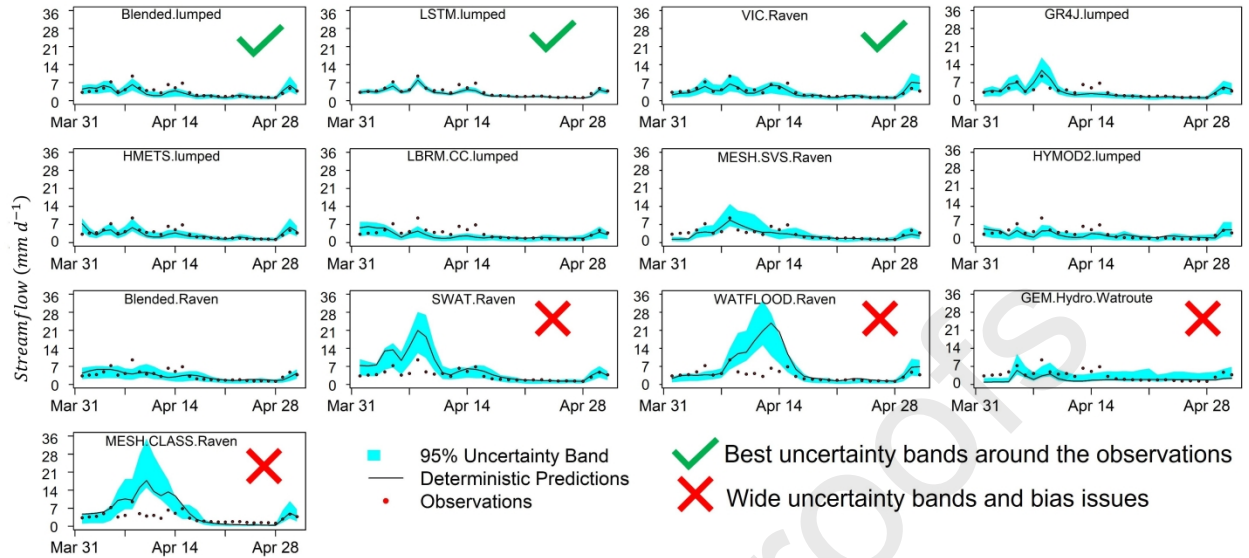
317 This section presents deterministic and probabilistic performance metrics of the 13 streamflow
 318 models for the Ganaraska River above Dale (Gauge 02HD012) during the validation period (1
 319 January 2011 – 31 December 2017). Results indicate that the variation between model structures
 320 was wide in terms of deterministic and probabilistic performance metrics (Table 2). Specifically,
 321 KGE ranged from 0.816 (Blended.lumped) to 0.153 (MESH.CLASS.Raven) with 81% change
 322 and CRPS ranged from 0.180 (LSTM.lumped) to 0.299 (MESH.SVS.Raven) with 40% change.
 323 As seen in Table 2, deterministic model skill (KGE) is correlated to probabilistic prediction skill
 324 (CRPS). In the representative station (02HD012) where the deterministic model performs well,
 325 the performance of the probabilistic predictions in terms of the CRPS will generally be good
 326 (e.g., LSTM.lumped). In contrast, when the deterministic model performs poorly, the flow
 327 uncertainty estimation is also poor (e.g., MESH.CLASS.Raven). This result shows that the
 328 residual error model post-processor could not “remedy” a poor deterministic model prediction.

329 **Table 2.** Deterministic and probabilistic performance metrics of streamflow model
 330 predictions based on 13 model structures for the single station (02HD012) during validation
 331 period (1 January 2011 – 31 December 2017).

Model	Metric	
	KGE	CRPS

Blended.lumped	0.816	0.216
LSTM.lumped	0.803	0.180
VIC.Raven	0.802	0.240
GR4J.lumped	0.802	0.223
HMETS.lumped	0.774	0.234
LBRM.CC.lumped	0.730	0.280
MESH.SVS.Raven	0.667	0.299
HYMOD2.lumped	0.664	0.277
Blended.Raven	0.597	0.284
SWAT.Raven	0.241	0.265
WATFLOOD.Raven	0.175	0.270
GEM.Hydro.Watroute	0.168	0.289
MESH.CLASS.Raven	0.153	0.260

332 Figure 2 illustrates sample flow time series with 95% uncertainty band for 13 streamflow
333 models during a freshet event in 2014. In this paper, 95% uncertainty bands were calculated to be
334 within 2.5 and 97.5 percentiles, as calculated from the Gaussian error model in transform space.
335 Figure 2 shows that the Blended.lumped, LSTM.lumped, and VIC.Raven models produce the
336 best predictions, with a narrow 95% uncertainty band around the observed data. In contrast,
337 MESH.CLASS.Raven, SWAT.Raven, and WATFLOOD.Raven models present wide 95%
338 uncertainty bands. These models are over-confident, with 44%, 41%, and 29%, respectively, of
339 observed flows outside the 95% uncertainty band. These models also show a bias issue, which
340 overestimate the high flow and underestimate the low flow.



341

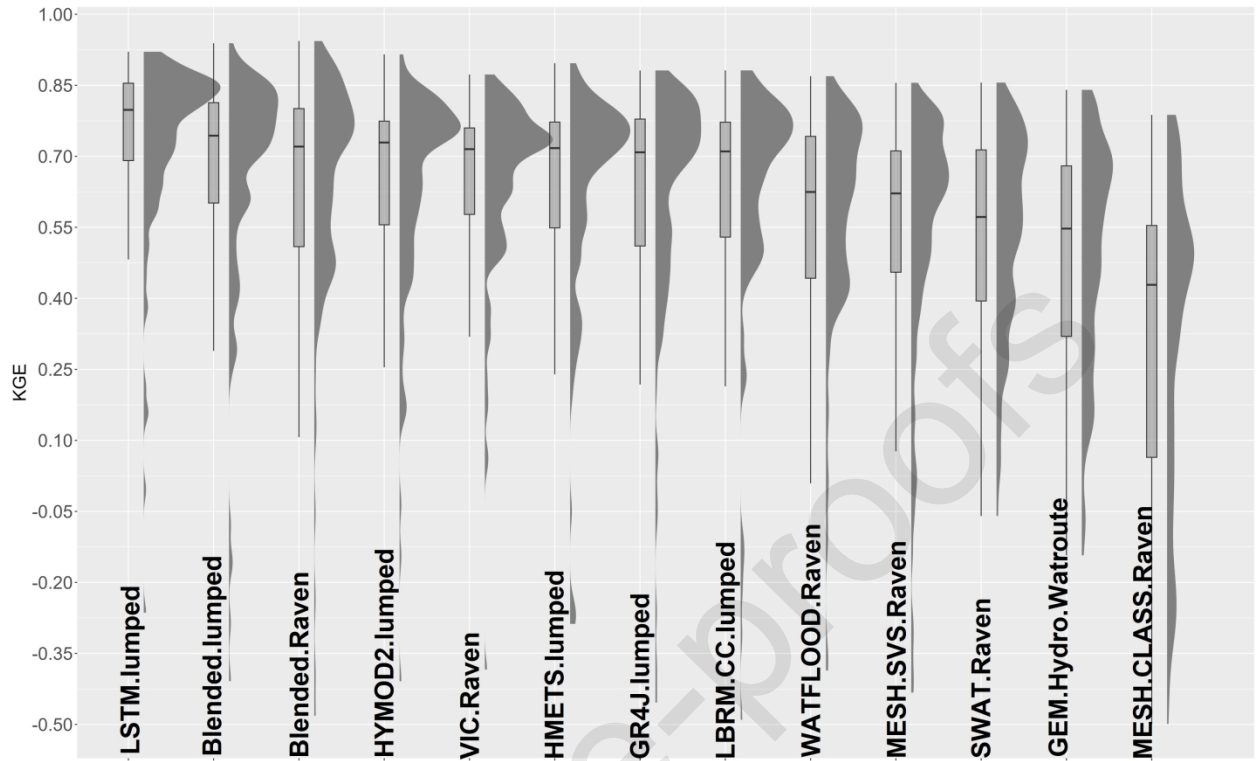
342 **Figure 2.** Sample time series with coloured uncertainty bands of streamflow based on 13 model
 343 structures for the single station (02HD012) during validation period (1 January 2011 – 31
 344 December 2017).

345 These results show the influence of streamflow model structures on streamflow
 346 predictive uncertainty quantification (Figure 2). For the same station and using the same residual
 347 error model post-processor, streamflow uncertainty quantification results were dissimilar for
 348 different streamflow models.

349 3.2 Model performance evaluation using full set of basins

350 3.2.1 Deterministic streamflow predictions (KGE)

351 Figure 3 compares the KGE of 13 streamflow models across all 141 case study basins.
 352 The LSTM.lumped model delivers the best KGE, indicated by the high median KGE
 353 performance metric of 0.8. The LSTM.lumped model also provides practically and statistically
 354 significant improvements in KGE. This outcome is consistent with the results of Mai et al.,
 355 (2022), Kratzert et al., (2018), and Kratzert et al., (2019). These studies have shown that
 356 globally-optimized LSTM models generally outperform conventional physically-based
 357 hydrological models used in water resource modelling. The Blended.lumped, Blended.raven, and
 358 HYMOD2.lumped models achieve a similar performance with median values of KGE 0.72, 0.72,
 359 and 0.71 respectively. In contrast to the LSTM.lumped model, the MESH.CLASS.Raven model
 360 provides the lowest KGE with median value of 0.42. It is interesting that the newer structurally-
 361 flexible models such as LSTM and Blended provide top performance in terms of the
 362 deterministic performance metric (KGE).



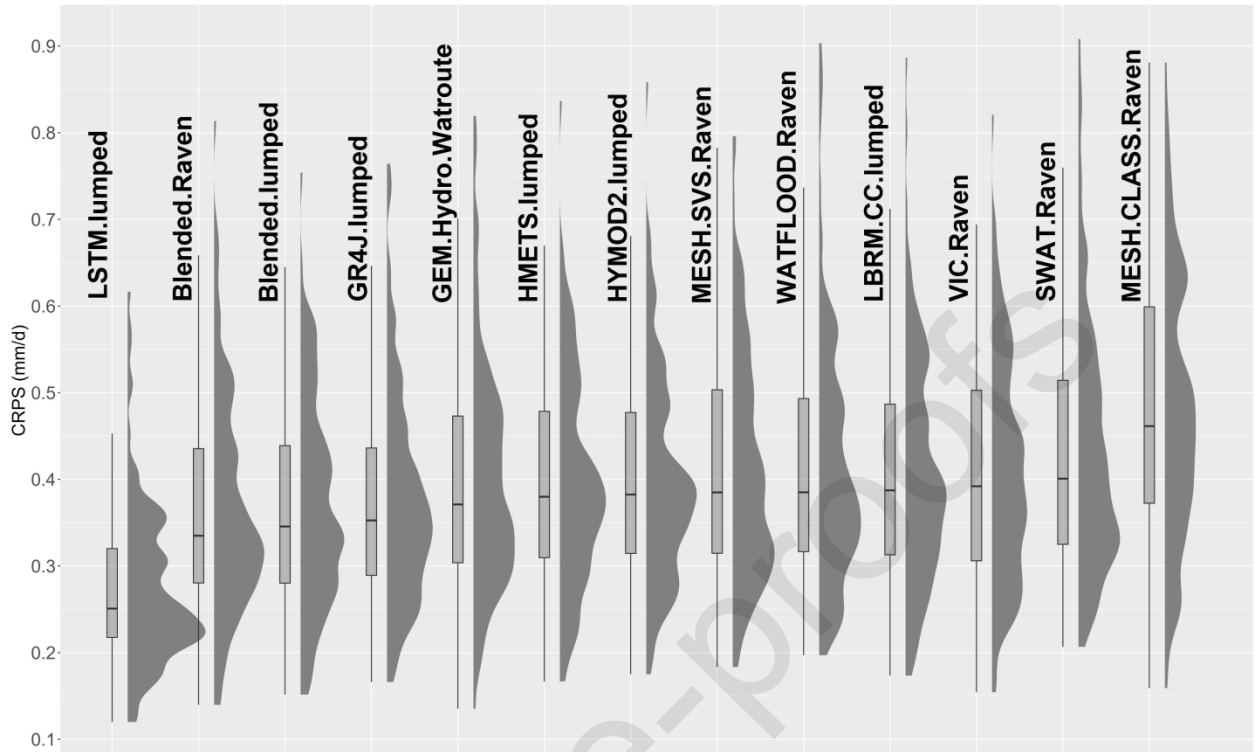
363

364 **Figure 3.** Density curves and boxplots of mean KGE over 13 streamflow model structures for all
 365 case study basins during validation period (1 January 2011 – 31 December 2017). For each
 366 model, the distribution of mean KGE scores across basins is summarized by a box and whisker
 367 plot highlighting the median and the 25 and 75 quantiles as well as the entire empirical
 368 distribution.

369

370 3.2.2 Probabilistic streamflow predictions (CRPS)

371 Figure 4 reports the CRPS of 13 streamflow models overall case study basins. The
 372 LSTM.lumped model is, once again, the best-performing model, with a median CRPS metric of
 373 0.25, which is practically and statistically significant better than the other numerical models. On
 374 the contrary, the MESH.CLASS.Raven model performed generally worse, with a median CRPS
 375 metric of 0.46.



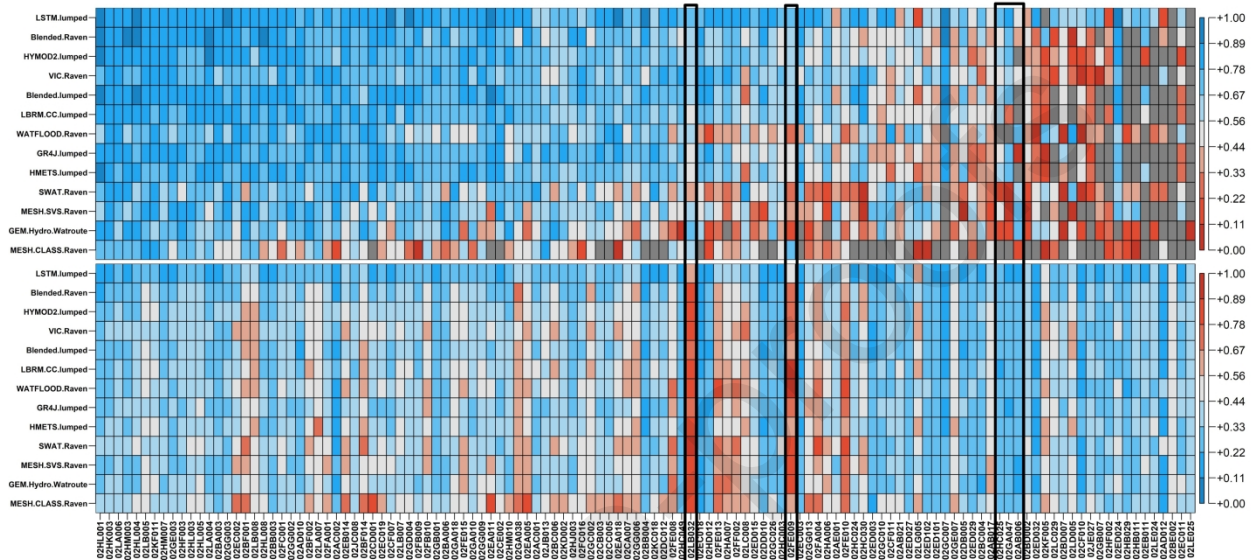
376

377 **Figure 4.** Density curves and boxplots of mean CRPS over 13 streamflow model structures for
 378 all case study basins during validation period (1 January 2011 – 31 December 2017). For each
 379 model, the distribution of mean CRPS scores across basins is summarized by a box and whisker
 380 plot highlighting the median and the 25 and 75 quantiles as well as the entire empirical
 381 distribution.

382

383 Figure 5 shows a heatmap of KGE (top) and CRPS (bottom) across the basins and models
 384 during the validation period. This figure helps to explore the distributions of figure 3 and 4 in
 385 detail at the individual basin level. Results confirm that LSTM.lumped delivers better KGE and
 386 CRPS values in many basins (blue tones), while MESH.CLASS.Raven provides the lowest KGE
 387 and CRPS values in many basins (red tones). Results also suggest that after post-processing there
 388 is less variation by models because there are vertical strips of the same color on the bottom (after
 389 post-processing) versus the colour variations on the top (no post-processing). This pattern seems
 390 especially true for basins with poor KGE (basins located in the right part of figure 5) and for
 391 basins with diverse KGE (e.g., 02HC025, 02AB006, and 02GA047 which are highlighted on the
 392 x-axis with black squares). A possible explanation for this pattern might be the rationality behind
 393 the REM-PP. In most cases, REM-PP struggles to match model outputs (e.g., streamflow) with
 394 observations. These matches employ statistical models that reduce model biases from all sources
 395 of uncertainty and characterize the associated uncertainty (Ye et al., 2014). Therefore, if REM-
 396 PP works reasonably well, it is sensible that REM-PP's outputs show similar performance. In
 397 other words, for some basins, the post-processing can equalise all the model's performance in
 398 terms of CRPS. This explanation is in line with the work of Ehlers et al., (2019) who found that
 399 the k-NN post-processor was able to compensate a poor deterministic simulation. Results also

400 indicate that in some basins it is difficult to obtain satisfactory performance in terms of KGE
 401 (dark grey cells) due to these basins being highly flow regulated. Similarly, in terms of CRPS,
 402 some basins seem difficult to obtain satisfactory performance; these basins show spiky
 403 hydrographs (e.g., 02FE009 and 02LB032 and these are highlighted on the x-axis with red
 404 squares).



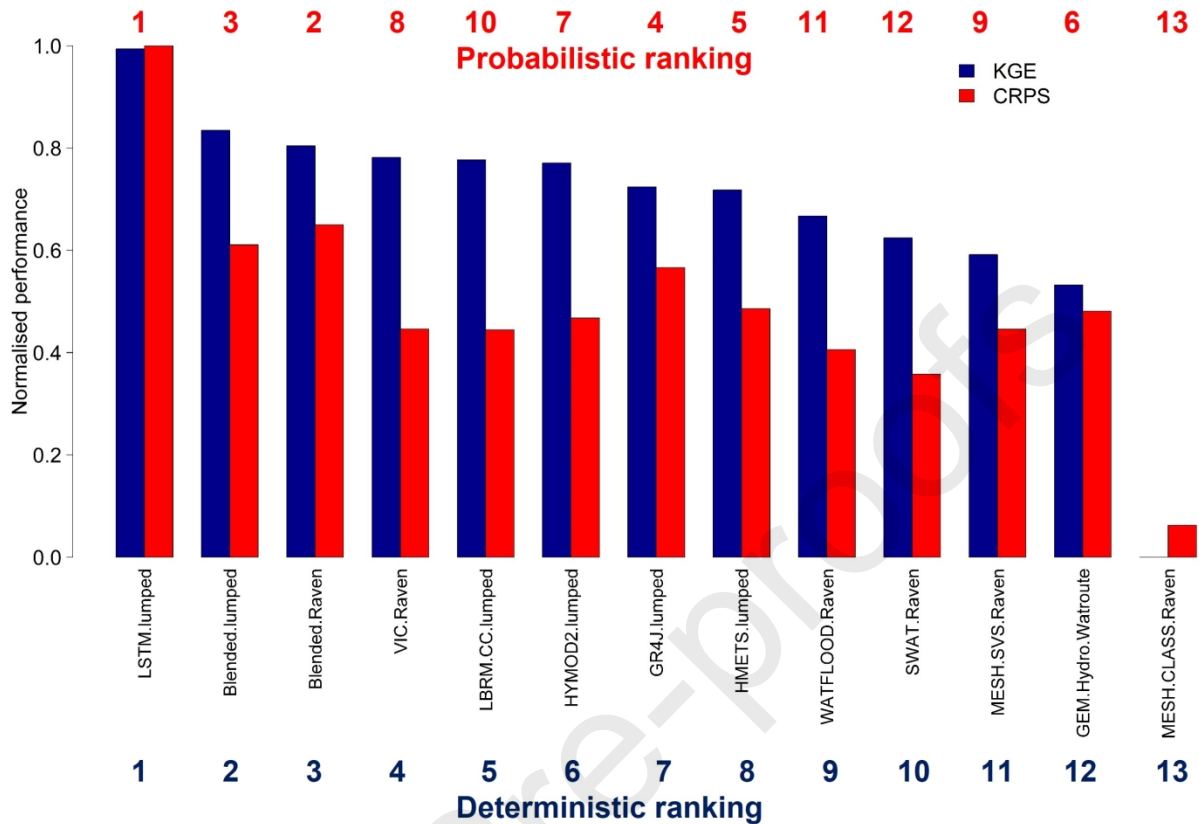
405
 406 **Figure 5.** Heatmap of KGE (top) and CRPS (bottom) over 13 streamflow model structures for all
 407 case study basins during the validation period (1 January 2011 – 31 December 2017). Horizontal
 408 axes denote different basins, which are ordered according to the mean KGE performance metric
 409 (best on the left), while vertical axes denote different models. Dark grey cells symbolize values
 410 of KGE below 0 meaning undesirable values. Blue tones indicate good (preferred) scores, and
 411 red tones signify bad scores. Some basin IDs are highlighted with black squares because they
 412 present particular behavior, which is examined in the result section.

413

414 3.3 Ranking models

415 Figure 6 represents the normalised performance metric of the KGE and CRPS of 13
 416 streamflow models across all case study basins during the validation period (1 January 2011 – 31
 417 December 2017). Results show that LSTM.lumped model provides the best performance across
 418 deterministic (KGE) and probabilistic (CRPS) metrics. In contrast to the LSTM.lumped model,
 419 the MESH.CLASS.Raven model delivers the poorest performance in terms of KGE and CRPS.
 420 The LSTM.lumped model is consistently the best-performing model. In addition, Figure 6 shows
 421 the ranking based on the median normalised KGE and CRPS. The order of the numerical models
 422 was determined based on the ranking of the median normalised KGE.

423



424

425 **Figure 6.** Comparing the model ranking and normalised deterministic and probabilistic
 426 performance of 13 streamflow model structures during the validation period (1 January 2011 –
 427 31 December 2017) using the median of the normalised KGE and CRPS for all case study. Red
 428 numbers in the top represent the probabilistic ranking (CRPS) and blue numbers in the bottom
 429 represent the deterministic ranking (KGE). Smaller number in the ranking indicates better
 430 performance.

431 A key finding is that moving from deterministic to probabilistic predictions using a post-
 432 processor strategy does not change the streamflow model performance ranking for the best and
 433 worst models in GRIP-GL, but it can change the ranks considerably for models of intermediate
 434 rank. In particular, post-processing greatly improves GR4J and GEM.Hydro.Watroute
 435 performance rank (relative to deterministic) while post-processing seems to degrade VIC and
 436 LBRM.CC.lumped rankings (relative to deterministic) (Figure 6). In part, these ranking changes
 437 reflect the comparable performance of probabilistic models in terms of CRPS – many of the
 438 scores fall within the range 0.4 to 0.5, where a swap in rank may simply be an artefact of the
 439 subtle variations in scores. However, results show a general pattern. When ranked by both KGE
 440 and CRPS, the highest ranked models are LSTM and Blended while MESH.CLASS, SWAT, and
 441 WATFLOOD are the lowest-ranked models (Figure 6). While not absolutely true, the
 442 probabilistic prediction quality tends to improve with deterministic performance, with both
 443 higher sharpness and robustness. As might be expected, poorly performing deterministic models
 444 tend to have much wider error models to compensate for their poorer performance, and therefore
 445 probabilistic prediction ranges are much wider. Inconsistency or performance likely leads to

446 reduced robustness of the models, as it becomes more difficult for a single error model to capture
447 model uncertainty which varies seasonally or with varying flow regime.

448 These results illustrate that streamflow model choice strongly influences probabilistic
449 performance. Different streamflow model structures produce varied deterministic performance
450 metrics (KGE), and these deterministic performance differences produce distinct probabilistic
451 performance metrics (CRPS) using the same residual error model post-processor. This variance
452 of CRPS values indicates that the quality of probabilistic predictions depends strongly on the
453 quality of the deterministic prediction. This finding is in line with Morawietz et al., (2011), who
454 evaluated different versions of the autoregressive error model as post-processors for a
455 probabilistic streamflow forecast system in Norway.

456 **4 Discussion and conclusions**

457 This research examines the effects of streamflow model choice on probabilistic
458 performance (CRPS) using a residual error model post-processing (REM-PP) approach to
459 multiple streamflow models and basins. We assess 141 calibration basins from the Great Lakes
460 region in the USA and Canada. This assessment employed 13 diverse streamflow predictions
461 based on machine-learning, lumped, subbasin, and gridded models. Using daily-scale streamflow
462 model predictions, we apply an additive REM-PP in transformed space using the Box-Cox
463 transformation, a first-order autoregressive AR(1) model, Gaussian innovations, and a flow-
464 dependent error mean. The main findings are as follows:

- 465 (1) Streamflow model choice strongly influences the robustness of probabilistic
466 predictions in terms of CRPS.
- 467 (2) Moving from deterministic to probabilistic predictions using a post-processor
468 approach does not change the streamflow model performance ranking for the best and
469 worst deterministic models, but models of intermediate rank in deterministic
470 evaluation do not have consistent ranking when evaluated in probabilistic mode
471 (Figure 6). This provides some evidence that probabilistic forecasts can, in some
472 sense, compensate for differences in model performance.
- 473 (3) Post-processing residual errors of long short-term memory (LSTM) network models
474 are consistently the best-performing model in terms of KGE (Figure 3) and CRPS
475 (Figure 4, Figure 5, and Figure 6). Notably, the normalized performance gap between
476 the LSTM and the second-best model is roughly doubled when the model comparison
477 moves from a deterministic to a probabilistic perspective. That is, the LSTM model
478 looks even more dominant over the other models when we use a post-processing
479 approach to consider model prediction uncertainty (Figure 6).

480 Based on empirical case studies, these findings suggest that merging deterministic
481 predictions of LSTM with the statistical modelling of the residual errors (post-processing)
482 provides the best approach to convert deterministic predictions to probabilistic ones. This
483 approach would be a promising route to characterize and reduce uncertainty in hydrological
484 predictions. These results are in line with the benchmarking study of Mai et al., (2022) who
485 showed that LSTM models outperformed the conventional hydrological models used to predict

486 daily streamflow. The superior performance of LSTM can be attributed to its capacity to digest
487 more training data (flow locations, basin characteristics, and forcings) and its global setup. Note
488 that LSTM was the only globally calibrated model. In addition, the range of error model
489 parameters are similar across all streamflow models except for the LSTM, which appear to be
490 unique. Boxplots of error model parameter values are showed in Appendix A.

491 In a general sense, model intercomparisons are sensitive to the performance evaluation
492 framework. If the models are evaluated from a deterministic framework, the same performance
493 of the models cannot be expected from a probabilistic framework. In other words, assuming that
494 the deterministic model ranking will be the same when moving to a probabilistic performance
495 evaluation is an inappropriate assumption. When probabilistic performance is implemented,
496 models are ranked differently. The residual error model could be an integral component in model
497 intercomparison. In past model intercomparisons, a few studies focussed on deterministic model
498 performance ranking (e.g., Lake Michigan (GRIP-M; Fry et al., 2014), Lake Ontario (GRIP-O;
499 Gaborit et al., 2017), and Lake Erie (GRIP-E; Mai et al., 2021)), but no previous study has
500 investigated the change between deterministic and probabilistic model performance rankings.
501 The novelty of this study lies in assessing the influence of multiple and diverse streamflow
502 predictions (i.e., multiple hydrologic models) on the quality of uncertainty estimates relative to
503 deterministic model performance. In other words, this study shows that streamflow model choice
504 affects the uncertainty estimation using large-sample catchments and a standardized
505 experimental setup.

506 In this work, we implemented a straightforward method to transform deterministic into
507 probabilistic predictions. The method is based on the post-processor strategy for the residual
508 error model, which is empirically more flexible and robust than joint strategies because it avoids
509 challenging interactions between water balance parameters and error model parameters (Evin et
510 al., 2014). The REM-PP used the Method-of-Moments inference method, which does not need a
511 likelihood function or optimisation, and therefore demands less computational costs. Although
512 the Method-of-Moments is simple, it can deliver predictive inference like Maximum Likelihood
513 or Bayesian methods (McInerney et al., 2018). The REM-PP quantifies the predictive uncertainty
514 without the need to change the objective function or to re-calibrate the hydrological model
515 (Hunter et al., 2021).

516 This study has potential limitations. We only addressed one simple post-processing
517 algorithm configuration. The results may not directly extend to other methods of uncertainty
518 analysis, such as Monte Carlo, informal or formal Bayesian methods, etc. The results may not
519 extend to different regions where models may rank differently in terms of accuracy.

520 Given that this research has focused on the effects of the hydrological model choice on
521 streamflow predictive uncertainty quantification, future work should focus on a more
522 comprehensive assessment using several hydroclimatological regimes. Especially, for example,
523 is whether predictive uncertainty estimation depends more on the hydrological model or on the
524 hydroclimatic regime.

525 Mainly in future work, we will include the uncertainty of model inputs (precipitation and
526 temperature) using a stochastic modelling generator. Further development is needed to improve
527 the predictive uncertainty estimation of low flows, including the treatment of zero flows

528 (McInerney et al., 2019). They can occur, for instance, in the extremely cold winter low flow
529 conditions in the Arctic and Prairie basins, in Canada. Remarkably, future research also should
530 focus on predictive uncertainty quantification at ungauged catchments which was not considered
531 here.

532 In general, this study helps to provide insight into the influence of streamflow model
533 choice on the estimation of prediction uncertainty and connects the gap between deterministic
534 and probabilistic prediction methods in practical hydrological applications. Moreover, this
535 research uses several case studies and streamflow models, including machine-learning-based,
536 lumped, subbasin-based, and gridded models, and applies a straightforward method to
537 researchers and practitioners who need to transform deterministic predictions to probabilistic
538 ones to characterize predictive uncertainty and support the management of hydrological risk.

539 Finally, this study is the first to perform a large-sample model intercomparison (13
540 models and 141 basins) from both a deterministic and probabilistic perspective and then compare
541 the differences in model rankings. The results of this study emphasise how crucial it is to
542 combine residual error models with the best deterministic streamflow model predictions (like
543 LSTM) to improve the quality and value of hydrological predictions, quantify uncertainty, and
544 better inform decision-making in operational water management.

545 **5 Data available**

546 The streamflow data used for this analysis are available at [https://hydrohub.org/grip-](https://hydrohub.org/grip-gl/maps_streamflow.html)
547 [gl/maps_streamflow.html](https://hydrohub.org/grip-gl/maps_streamflow.html) (accessed on July 10, 2023). The code to implement the residual error
548 model is available via the “ProbPred” R package (<https://github.com/AdWater/ProbApp>).

549

550 **CRedit authorship contribution statement**

551 **Jonathan Romero-Cuellar:** Conceptualisation, Investigation, Formal analysis, Visualisation,
552 Software, Methodology, Writing - original draft. **Rezgar Arabzadeh:** Software, Data curation,
553 Visualisation, Formal analysis. **James R. Craig:** Conceptualisation, Methodology, Writing -
554 review & editing, Supervision, Resources. **Bryan A. Tolson:** Conceptualisation, Methodology,
555 Writing - review & editing, Supervision. **Juliane Mai:** Writing - review & editing.

556

557 **Declaration of Competing Interest**

558 The authors declare that they have no known competing financial interests or personal
559 relationships that could have appeared to influence the work reported in this paper.

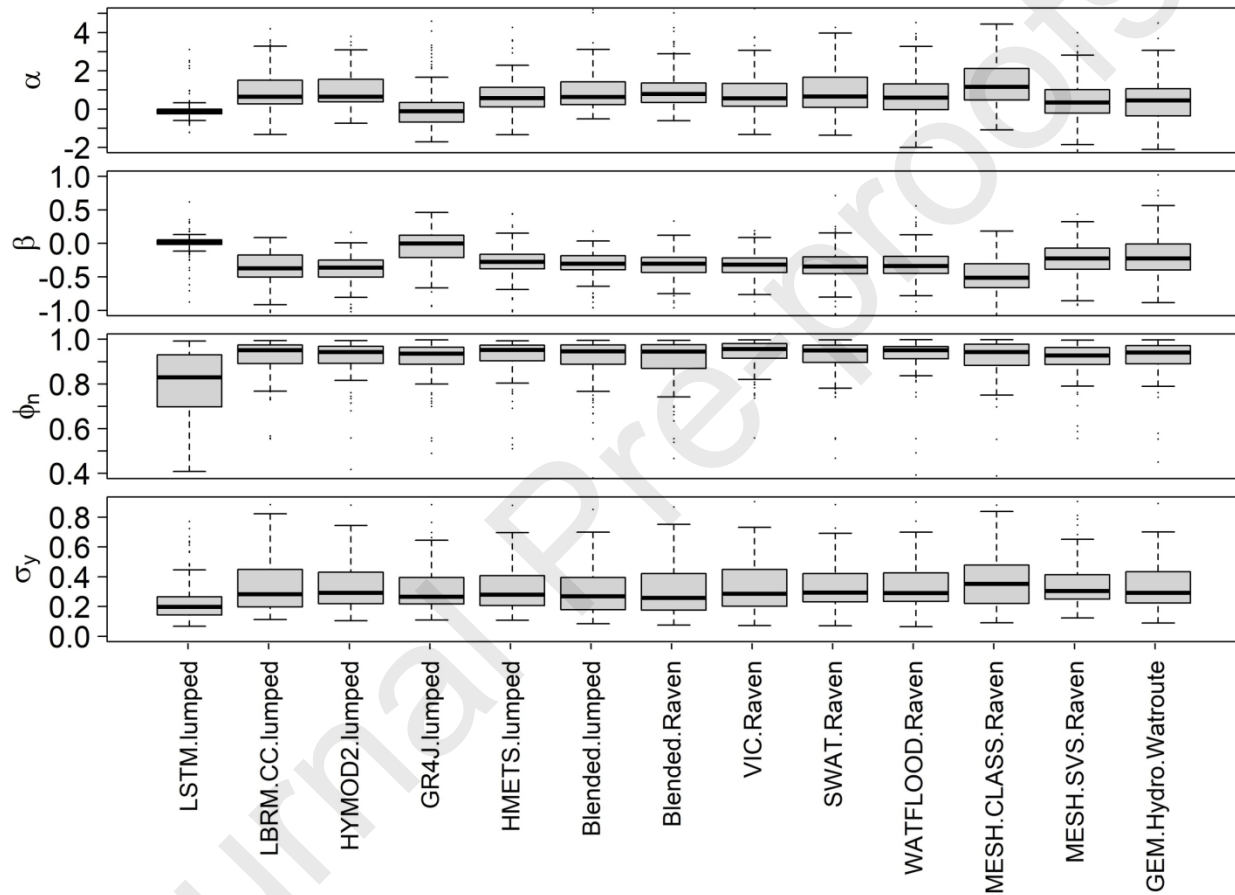
560

561 **Acknowledgements**

562 This research is funded by the Natural Resources Canada Emergency Management program, the
 563 NSERC Canada Research Chairs program, and NSERC Discovery Individual grants held by Drs.
 564 Craig and Tolson. Special thanks to David McInerney for his support in ProbPred R package.
 565 We sincerely thank the associate editor and two anonymous reviewers for their thoughtful and
 566 helpful feedback, which significantly improved the manuscript.

567

568 Appendix A: REMs' parameter ranges



569

570 Figure A. Boxplots of error model parameter values across streamflow models. α is the intercept
 571 and β is the slope of flow linear regression function respectively. ϕ_η is the lag-1 autoregressive
 572 parameter, and σ_y standard deviation of Gaussian distribution.

573 References

- 574 Ajami, N. K., Duan, Q., & Sorooshian, S. (2007). An integrated hydrologic Bayesian multimodel combination
 575 framework: Confronting input, parameter, and model structural uncertainty in hydrologic prediction. *Water*
 576 *Resources Research*, 43(1), W01403. <https://doi.org/10.1029/2005WR004745>
- 577 Bastola, S., Murphy, C., & Fealy, R. (2012). Generating probabilistic estimates of hydrological response for Irish
 578 catchments using a weather generator and probabilistic climate change scenarios. *Hydrological Processes*,
 579 26(15), 2307–2321. <https://doi.org/10.1002/hyp.8349>

- 580 Bates, B. C., & Campbell, E. P. (2001). A Markov Chain Monte Carlo Scheme for parameter estimation and
 581 inference in conceptual rainfall-runoff modeling. *Water Resources Research*, 37(4), 937–947.
 582 <https://doi.org/10.1029/2000WR900363>
- 583 Bauer, D. F. (1972). Constructing Confidence Sets Using Rank Statistics. *Journal of the American Statistical*
 584 *Association*, 67(339), 690. <https://doi.org/10.2307/2284469>
- 585 Box, G.E.P., & Cox, D.R. (1964). An Analysis of Transformations. *Journal of the Royal Statistical Society, Series*
 586 *B*, 26 (2), 211–252. <http://www.jstor.org/stable/2984418>
- 587 Dias, L. F., Aparício, B. A., Nunes, J. P., Morais, I., Fonseca, A. L., Pastor, A. V., & Santos, F. D. (2020).
 588 Integrating a hydrological model into regional water policies: Co-creation of climate change dynamic adaptive
 589 policy pathways for water resources in southern Portugal. *Environmental Science and Policy*, 114, 519–532.
 590 <https://doi.org/10.1016/j.envsci.2020.09.020>
- 591 Ehlers, L.B., Wani, O., Koch, J., Sonnenborg, T.O., & Refsgaard, J.C. (2019). Using a simple post-processor to
 592 predict residual uncertainty for multiple hydrological model outputs. *Advances in Water Resources*, 129, 16-
 593 30. <https://doi.org/10.1016/j.advwatres.2019.05.003>
- 594 Engeland, K., & Steinsland, I. (2014). Probabilistic postprocessing models for flow forecasts for a system of
 595 catchments and several lead times. *Water Resources Research*, 50(1), 182–197.
 596 <https://doi.org/10.1002/2012WR012757>
- 597 Evin, G., Thyer, M., Kavetski, D., McInerney, D., & Kuczera, G. (2014). Comparison of joint versus postprocessor
 598 approaches for hydrological uncertainty estimation accounting for error autocorrelation and heteroscedasticity.
 599 *Water Resources Research*, 50(3), 2350–2375. <https://doi.org/10.1002/2013WR014185>
- 600 Farmer, W. H., & Vogel, R. M. (2016). On the deterministic and stochastic use of hydrologic models. *Water*
 601 *Resources Research*, 52(7), 5619–5633. <https://doi.org/10.1002/2016WR019129>
- 602 Fry, L. M., Gronewold, A. D., Fortin, V., Buan, S., Clites, A. H., Luukkonen, C., Holtschlag, D., Diamond, L.,
 603 Hunter, T., Seglenieks, F., Durnford, D., Dimitrijevic, M., Subich, C., Klyszejko, E., Kea, K., & Restrepo, P.
 604 (2014). The Great Lakes Runoff Intercomparison Project Phase 1: Lake Michigan (GRIP-M). *Journal of*
 605 *Hydrology*, 519, 3448–3465. <https://doi.org/10.1016/J.JHYDROL.2014.07.021>
- 606 Gaborit, É., Fortin, V., Tolson, B., Fry, L., Hunter, T., & Gronewold, A. D. (2017). Great Lakes Runoff Inter-
 607 comparison Project, phase 2: Lake Ontario (GRIP-O). *Journal of Great Lakes Research*, 43(2), 217–227.
 608 <https://doi.org/10.1016/J.JGLR.2016.10.004>
- 609 Gupta, A., & Govindaraju, R. S. (2023). Uncertainty quantification in watershed hydrology: Which method to use?
 610 *Journal of Hydrology*, 616, 128749. <https://doi.org/10.1016/j.jhydro.2022.128749>
- 611 Gupta, H. V., Kling, H., Yilmaz, K. K., & Martinez, G. F. (2009). Decomposition of the mean squared error and
 612 NSE performance criteria: Implications for improving hydrological modelling. *Journal of Hydrology*, 377(1-
 613 2), 80–91. <https://doi.org/10.1016/j.jhydro.2009.08.003>
- 614 Hazelton, M. L. (2011). Methods of Moments Estimation BT. In M. Lovric (Ed.), *International Encyclopedia of*
 615 *Statistical Science* (pp. 816–817). Springer Berlin Heidelberg. <https://doi.org/10.1007/978-3-642-04898-2>
- 616 Hersbach, H. (2000). Decomposition of the continuous ranked probability score for ensemble prediction systems.
 617 *Weather and Forecasting*, 15(5), 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)

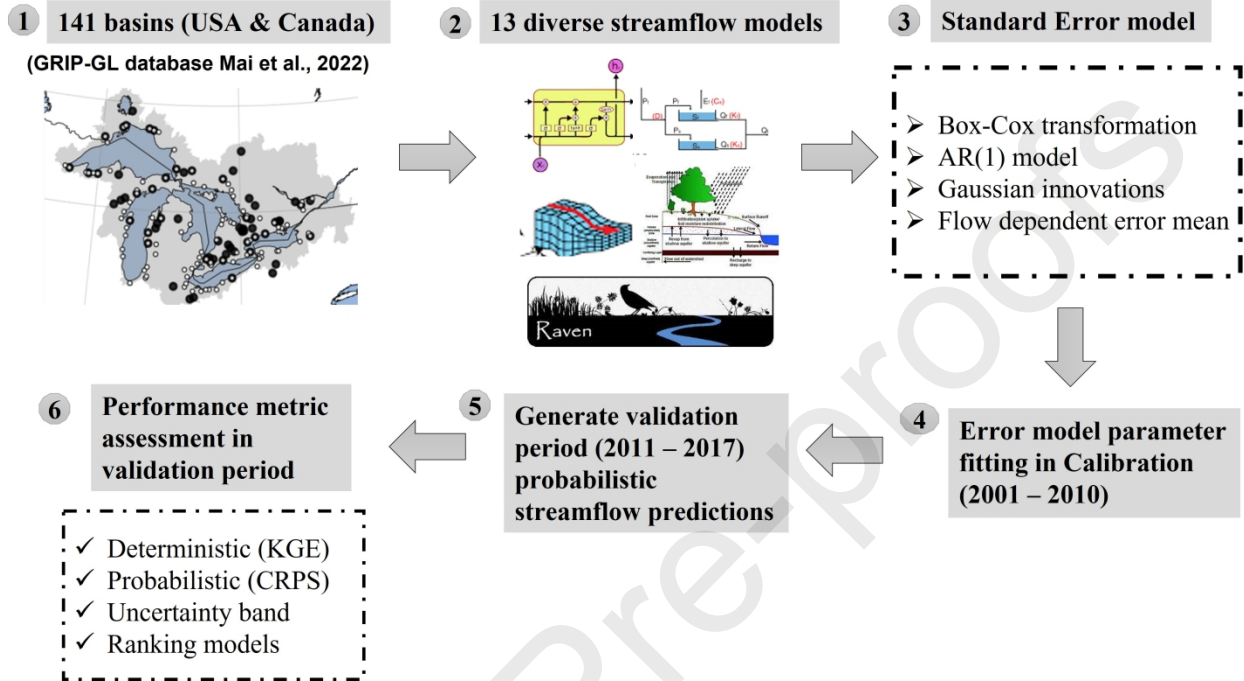
- 619 Hunter, J., Thyer, M., McInerney, D., & Kavetski, D. (2021). Achieving high-quality probabilistic predictions from
620 hydrological models calibrated with a wide range of objective functions. *Journal of Hydrology*, 603, 126578.
621 <https://doi.org/10.1016/J.JHYDROL.2021.126578>
- 622 Jiang, X., Gupta, H. V., Liang, Z., & Li, B. (2019). Toward Improved Probabilistic Predictions for Flood Forecasts
623 Generated Using Deterministic Models. *Water Resources Research*, 55(11), 9519–9543.
624 <https://doi.org/10.1029/2019WR025477>
- 625 Jordan, A., Krüger, F., & Lerch, S. (2019). Evaluating Probabilistic Forecasts with scoringRules. *Journal of*
626 *Statistical Software*, 90, 1–37. <https://doi.org/10.18637/JSS.V090.I12>
- 627 Koutsoyiannis, D., & Montanari, A. (2022). Bluecat: A Local Uncertainty Estimator for Deterministic Simulations
628 and Predictions. *Water Resources Research*, 58(1), e2021WR031215. <https://doi.org/10.1029/2021WR031215>
- 629 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall-runoff modelling using Long
630 Short-Term Memory (LSTM) networks. *Hydrol. Earth Syst. Sci*, 22, 6005–6022. [https://doi.org/10.5194/hess-](https://doi.org/10.5194/hess-22-6005-2018)
631 [22-6005-2018](https://doi.org/10.5194/hess-22-6005-2018)
- 632 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019). Towards learning universal,
633 regional, and local hydrological behaviors via machine learning applied to large-sample datasets. *Hydrology*
634 *and Earth System Sciences*, 23(12), 5089–5110. <https://doi.org/10.5194/hess-23-5089-2019>
- 635 Kuczera, G., Kavetski, D., Franks, S., & Thyer, M. (2006). Towards a Bayesian total error analysis of conceptual
636 rainfall-runoff models: Characterising model error using storm-dependent parameters. *Journal of Hydrology*,
637 331(1–2), 161–177. <https://doi.org/10.1016/j.jhydrol.2006.05.010>
- 638 Li, D., Marshall, L., Liang, Z., Sharma, A., & Zhou, Y. (2021). Characterizing distributed hydrological model
639 residual errors using a probabilistic long short-term memory network. *Journal of Hydrology*, 603.
640 <https://doi.org/10.1016/j.jhydrol.2021.126888>
- 641 Li, M., Wang, Q. J., Bennett, J. C., & Robertson, D. E. (2016). Error reduction and representation in stages (ERRIS)
642 in hydrological modelling for ensemble streamflow forecasting. *Hydrology and Earth System Sciences*, 20(9),
643 3561–3579. <https://doi.org/10.5194/hess-20-3561-2016>
- 644 Li, M., Wang, Q. J., Robertson, D. E., & Bennett, J. C. (2017). Improved error modelling for streamflow forecasting
645 at hourly time steps by splitting hydrographs into rising and falling limbs. *Journal of Hydrology*, 555, 586–
646 599. <https://doi.org/10.1016/j.jhydrol.2017.10.057>
- 647 Li, W., Duan, Q., Miao, C., Ye, A., Gong, W., & Di, Z. (2017). A review on statistical postprocessing methods for
648 hydrometeorological ensemble forecasting. *Wiley Interdisciplinary Reviews: Water*, 4(6), e1246.
649 <https://doi.org/10.1002/wat2.1246>
- 650 Mai, J., Shen, H., Tolson, B. A., Gaborit, É., Arsenault, R., Craig, J. R., Fortin, V., Fry, L. M., Gauch, M., Klotz, D.,
651 Kratzert, F., O'Brien, N., Princz, D. G., Rasiya Koya, S., Roy, T., Seglenieks, F., Shrestha, N. K., Temgoua,
652 A. G. T., Vionnet, V., & Waddell, J. W. (2022). The Great Lakes Runoff Intercomparison Project Phase 4:
653 The Great Lakes (GRIP-GL). *Hydrology and Earth System Sciences*, 26(13), 3537–3572.
654 <https://doi.org/10.5194/HESS-26-3537-2022>
- 655 Mai, J., Tolson, B. A., Shen, H., Gaborit, É., Fortin, V., Gasset, N., Awoye, H., Stadnyk, T. A., Fry, L. M., Bradley,
656 E. A., Seglenieks, F., Temgoua, A. G. T., Princz, D. G., Gharari, S., Haghnegahdar, A., Elshamy, M. E.,
657 Razavi, S., Gauch, M., Lin, J., ... Pietroniro, A. (2021). Great Lakes Runoff Intercomparison Project Phase 3:
658 Lake Erie (GRIP-E). *Journal of Hydrologic Engineering*, 26(9), 05021020.
659 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0002097](https://doi.org/10.1061/(ASCE)HE.1943-5584.0002097)

- 660 McInerney, D., Kavetski, D., Thyer, M., Lerat, J., & Kuczera, G. (2019). Benefits of Explicit Treatment of Zero
661 Flows in Probabilistic Hydrological Modeling of Ephemeral Catchments. *Water Resources Research*, 55(12),
662 11035–11060. <https://doi.org/10.1029/2018WR024148>
- 663 McInerney, D., Thyer, M., Kavetski, D., Bennett, B., Lerat, J., Gibbs, M., & Kuczera, G. (2018). A simplified
664 approach to produce probabilistic hydrological model predictions. *Environmental Modelling & Software*, 109,
665 306–314. <https://doi.org/10.1016/j.envsoft.2018.07.001>
- 666 McInerney, D., Thyer, M., Kavetski, D., Laugesen, R., Tuteja, N., & Kuczera, G. (2020). Multi-temporal
667 Hydrological Residual Error Modeling for Seamless Subseasonal Streamflow Forecasting. *Water Resources*
668 *Research*, 56, e2019WR026979. <https://doi.org/10.1029/2019WR026979>
- 669 McInerney, D., Thyer, M., Kavetski, D., Laugesen, R., Woldemeskel, F., Tuteja, N., & Kuczera, G. (2021).
670 Improving the Reliability of Sub-Seasonal Forecasts of High and Low Flows by Using a Flow-Dependent
671 Nonparametric Model. *Water Resources Research*, 57, e2020WR029317.
672 <https://doi.org/10.1029/2020WR029317>
- 673 McInerney, D., Thyer, M., Kavetski, D., Lerat, J., & Kuczera, G. (2017). Improving probabilistic prediction of daily
674 streamflow by identifying Pareto optimal approaches for modeling heteroscedastic residual errors. *Water*
675 *Resources Research*, 53(3), 2199–2239. <https://doi.org/10.1002/2016WR019168>
- 676 Moges, E., Demissie, Y., Larsen, L., & Yassin, F. (2021) Review: Sources of Hydrological Model Uncertainties and
677 Advances in Their Analysis. *Water*, 13(1),28. <https://doi.org/10.3390/w13010028>
- 678 Montanari, A., & Koutsoyiannis, D. (2012). A blueprint for process-based modeling of uncertain hydrological
679 systems. *Water Resources Research*, 48, W09555. <https://doi.org/10.1029/2011WR011412>
- 680 Morawietz, M., Xu, C. Y., Gottschalk, L., & Tallaksen, L. M. (2011). Systematic evaluation of autoregressive error
681 models as post-processors for a probabilistic streamflow forecast system. *Journal of Hydrology*, 407(1–4), 58–
682 72. <https://doi.org/10.1016/j.jhydrol.2011.07.007>
- 683 Nash, J. E., & Sutcliffe, J. V. (1970). River flow forecasting through conceptual models part I — A discussion of
684 principles. *Journal of Hydrology*, 10(3), 282–290. [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6)
- 685 Reichert, P., & Schuwirth, N. (2012). Linking statistical bias description to multiobjective model calibration. *Water*
686 *Resources Research*, 48(9), 9543. <https://doi.org/10.1029/2011WR011391>
- 687 Reggiani, P., Talbi, A., & Todini, E. (2022). Towards informed water resources planning and management.
688 *Hydrology*, 9(8), 136. <https://doi.org/10.3390/hydrology9080136>
- 689 Renard, B., Kavetski, D., Kuczera, G., Thyer, M., & Franks, S. W. (2010). Understanding predictive uncertainty in
690 hydrologic modeling: The challenge of identifying input and structural errors. *Water Resources Research*,
691 46(5), W05521. <https://doi.org/10.1029/2009WR008328>
- 692 Romero-Cuellar, J., Abbruzzo, A., Adelfio, G., & Francés, F. (2019). Hydrological post-processing based on
693 approximate Bayesian computation (ABC). *Stoch Environ Res Risk Assess*, 33(7), 1361–1373.
694 <https://doi.org/10.1007/s00477-019-01694-y>
- 695 Salamon, P., & Feyen, L. (2009). Assessing parameter, precipitation, and predictive uncertainty in a distributed
696 hydrological model using sequential data assimilation with the particle filter. *Journal of Hydrology*, 376(3–4),
697 428–442. <https://doi.org/10.1016/j.jhydrol.2009.07.051>
- 698 Schoups, G., & Vrugt, J. A. (2010). A formal likelihood function for parameter and predictive inference of
699 hydrologic models with correlated, heteroscedastic, and non-Gaussian errors. *Water Resources Research*,
700 46(10), W10531. <https://doi.org/10.1029/2009WR008933>

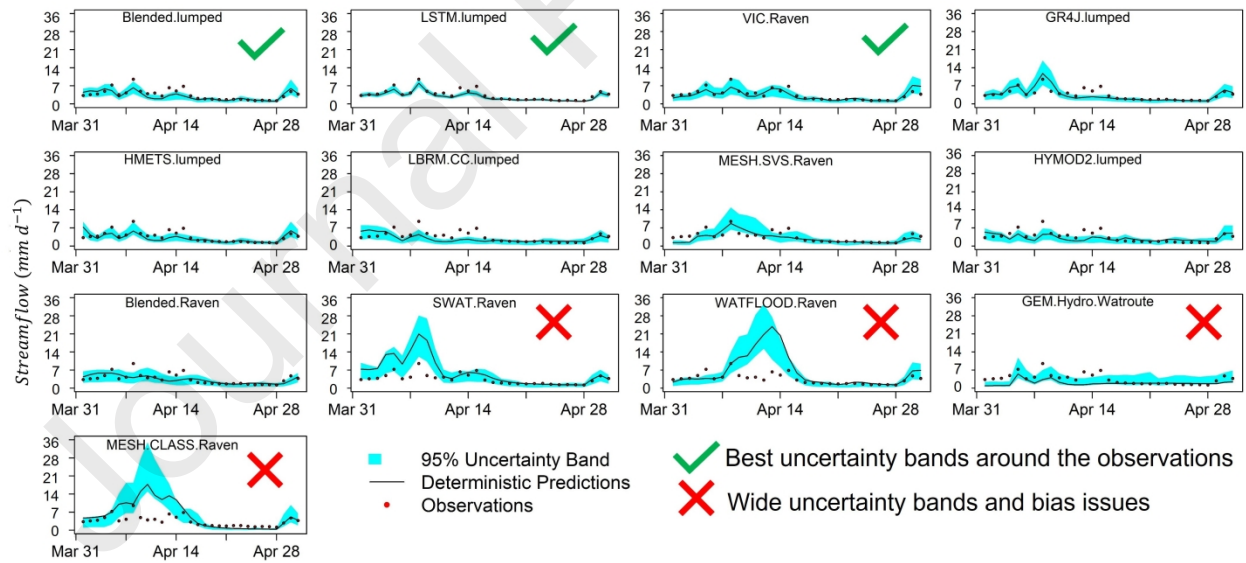
- 701 Shabestanipour, G., Brodeur, Z., Farmer, W. H., Steinschneider, S., Vogel, R. M., & Lamontagne, J. R. (2023).
 702 Stochastic Watershed Model Ensembles for Long-Range Planning: Verification and Validation. *Water*
 703 *Resources Research*, 59(2), e2022WR032201. <https://doi.org/10.1029/2022WR032201>
- 704 Smith, T., Marshall, L., & Sharma, A. (2015). Modeling residual hydrologic errors with Bayesian inference. *Journal*
 705 *of Hydrology*, 528, 29–37. <https://doi.org/10.1016/j.jhydrol.2015.05.051>
- 706 Sorooshian, S., & Dracup, J. A. (1980). Stochastic parameter estimation procedures for hydrologic rainfall-runoff
 707 models: Correlated and heteroscedastic error cases. *Water Resources Research*, 16(2), 430–442.
 708 <http://doi.wiley.com/10.1029/WR016i002p00430>
- 709 Sorooshian, S., Gupta, V. K., & Fulton, J. L. (1983). Evaluation of Maximum Likelihood Parameter estimation
 710 techniques for conceptual rainfall-runoff models: Influence of calibration data variability and length on model
 711 credibility. *Water Resources Research*, 19(1), 251–259. <https://doi.org/10.1029/WR019i001p00251>
- 712 Sun, R., Yuan, H., & Liu, X. (2017). Effect of heteroscedasticity treatment in residual error models on model
 713 calibration and prediction uncertainty estimation. *Journal of Hydrology*, 554, 680–692.
 714 <https://doi.org/10.1016/J.JHYDROL.2017.09.041>
- 715 Thyer, M., Renard, B., Kavetski, D., Kuczera, G., Franks, S. W., & Srikanthan, S. (2009). Critical evaluation of
 716 parameter consistency and predictive uncertainty in hydrological modeling: A case study using Bayesian total
 717 error analysis. *Water Resources Research*, 45(12). <https://doi.org/10.1029/2008WR006825>
- 718 Todini, E. (2008). A model conditional processor to assess predictive uncertainty in flood forecasting. *International*
 719 *Journal of River Basin Management*, 6(2), 123–137. <https://doi.org/10.1080/15715124.2008.9635342>
- 720 Vaske, J. J., Gliner, J. A., & Morgan, G. A. (2010). Communicating Judgments About Practical Significance: Effect
 721 Size, Confidence Intervals and Odds Ratios. *Human Dimensions of Wildlife*, 7(4), 287–300.
 722 <https://doi.org/10.1080/10871200214752>
- 723 Vogel, R. M. (2017). Stochastic watershed models for hydrologic risk management. *Water Security*, 1, 28-35.
 724 <https://doi.org/10.1016/j.wasec.2017.06.001>
- 725 Vrugt, J. A., de Oliveira, D. Y., Schoups, G., & Diks, C. G. H. (2022). On the use of distribution-adaptive likelihood
 726 functions: Generalized and universal likelihood functions, scoring rules and multi-criteria ranking. *Journal of*
 727 *Hydrology*, 615, 128542. <https://doi.org/10.1016/J.JHYDROL.2022.128542>
- 728 Wang, Q. J., Shrestha, D. L., Robertson, D. E., & Pokhrel, P. (2012). A log-sinh transformation for data
 729 normalization and variance stabilization. *Water Resources Research*, 48(5).
 730 <https://doi.org/10.1029/2011WR010973>
- 731 Weerts, A. H., Winsemius, H. C., & Verkade, J. S. (2011). Estimation of predictive hydrological uncertainty using
 732 quantile regression: examples from the National Flood Forecasting System (England and Wales). *Hydrol.*
 733 *Earth Syst. Sci.*, 115, 255–265. <https://doi.org/10.5194/hess-15-255-2011>
- 734 Woldemeskel, F., McInerney, D., Lerat, J., Thyer, M., Kavetski, D., Shin, D., Tuteja, N., & Kuczera, G. (2018).
 735 Evaluating post-processing approaches for monthly and seasonal streamflow forecasts. *Hydrology and Earth*
 736 *System Sciences*, 22(12), 6257–6278. <https://doi.org/10.5194/HESS-22-6257-2018>
- 737 Yang, J., Reichert, P., & Abbaspour, K. C. (2007). Bayesian uncertainty analysis in distributed hydrologic modeling:
 738 A case study in the Thur River basin (Switzerland). *Water Resources Research*, 43(10), W10401.
 739 <https://doi.org/10.1029/2006WR005497>
- 740 Ye, A., Duan, Q., Yuan, X., Wood, E. F., & Schaake, J. (2014). Hydrologic post-processing of MOPEX streamflow
 741 simulations. *Journal of Hydrology*, 508, 147–156. <https://doi.org/10.1016/j.jhydrol.2013.10.055>

742 Zhao, T., Wang, Q. J., Bennett, J. C., Robertson, D. E., Shao, Q., & Zhao, J. (2015). Quantifying predictive
 743 uncertainty of streamflow forecasts based on a Bayesian joint probability model. *Journal of Hydrology*, 528,
 744 329–340. <https://doi.org/10.1016/j.jhydrol.2015.06.043>

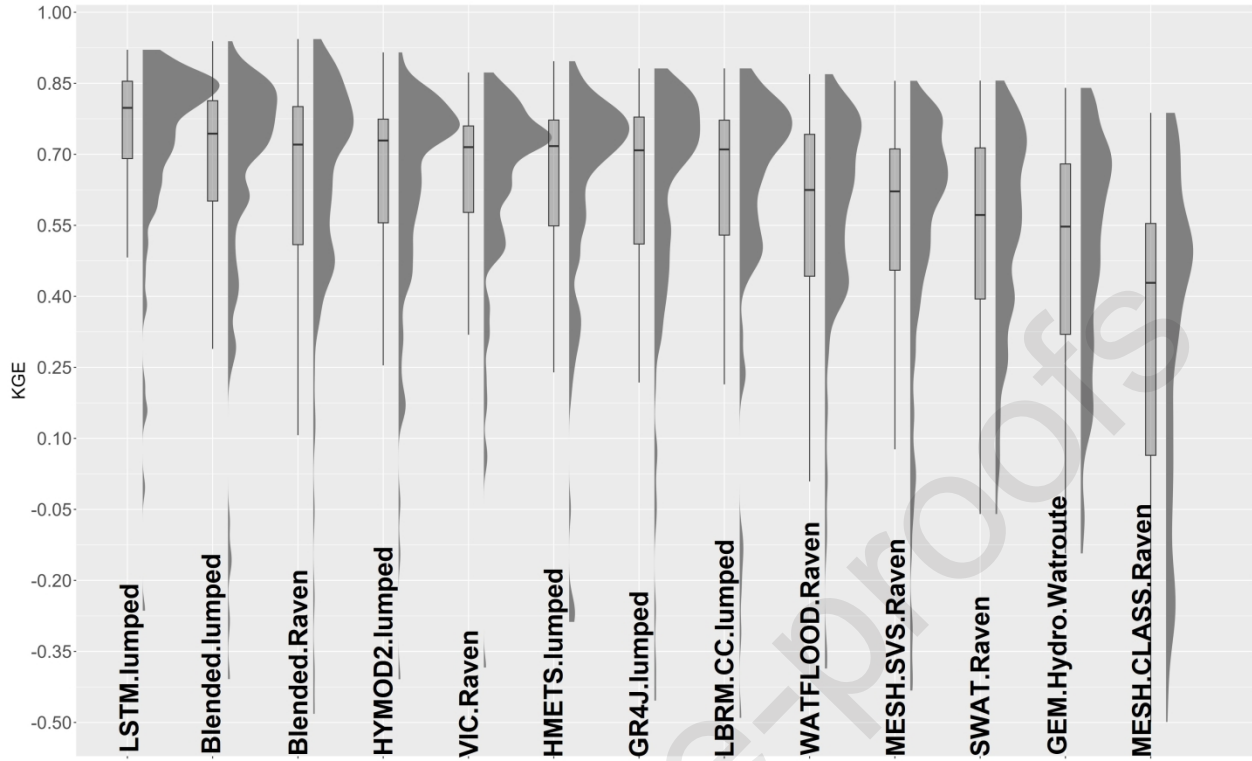
745



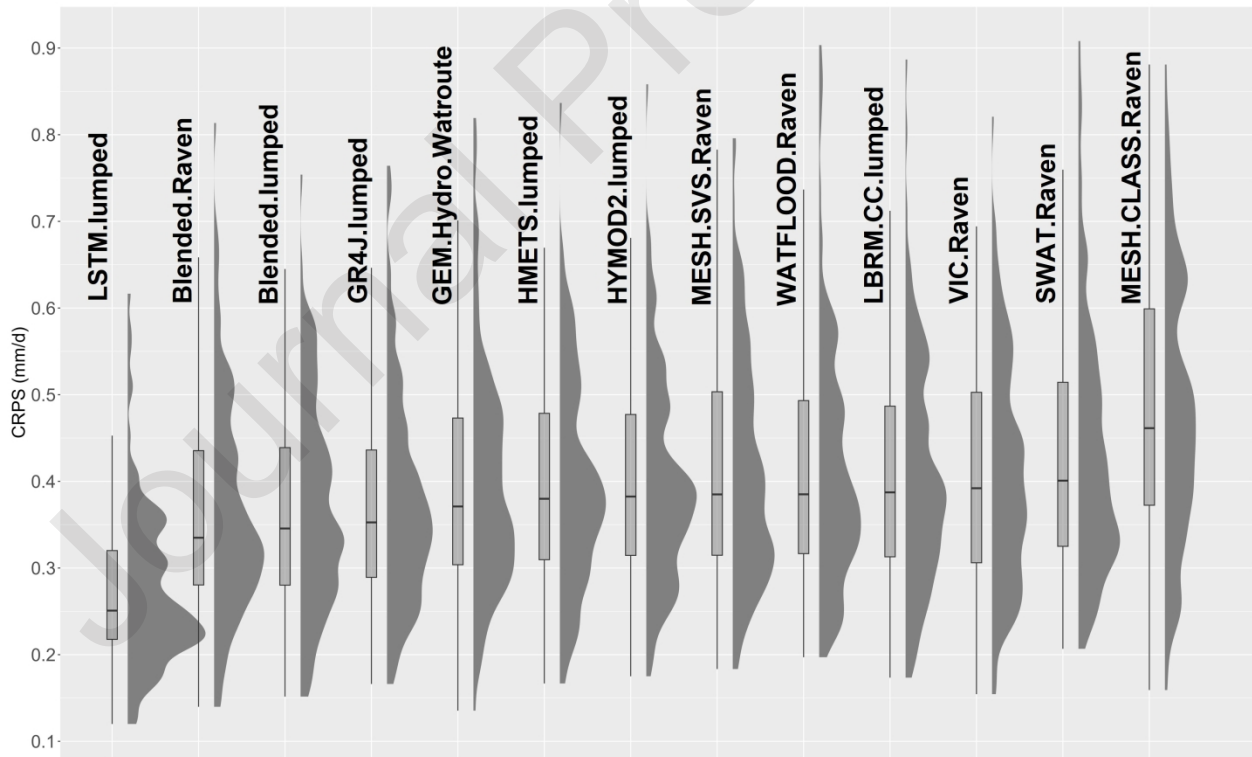
746



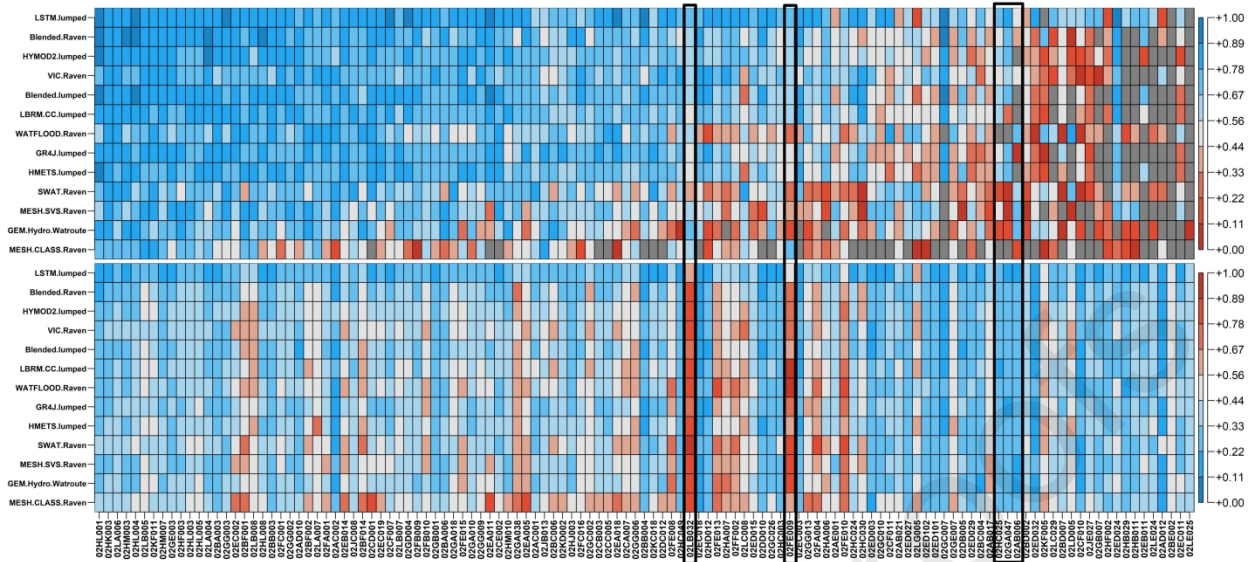
747



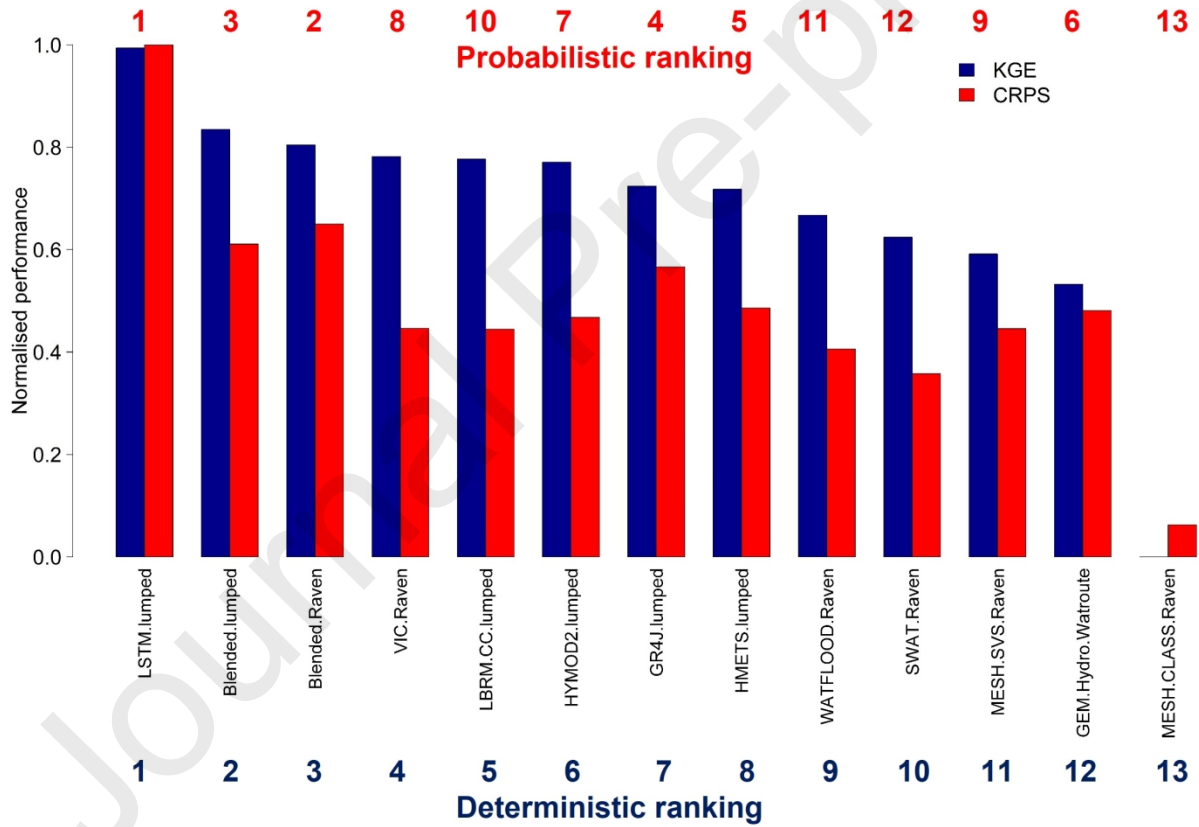
748



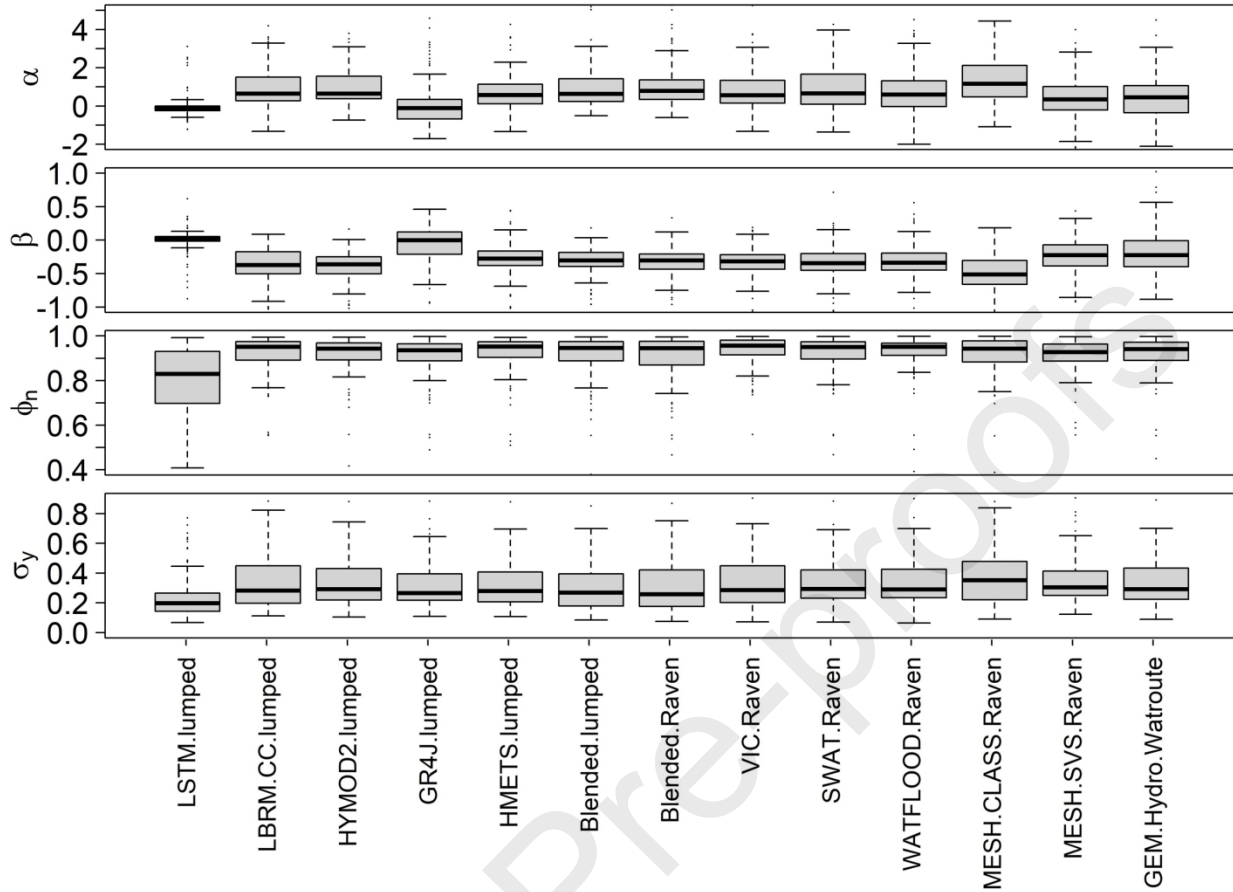
749



750



751

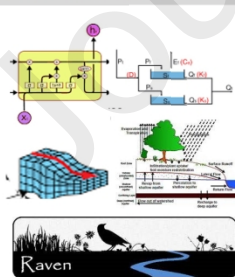


752

1 141 basins (USA & Canada)
(GRIP-GL database Mai et al., 2022)



2 13 diverse streamflow models



3 Standard Error model

- Box-Cox transformation
- AR(1) model
- Gaussian innovations
- Flow dependent error mean

4 Calibration (2001 – 2010)

5 Validation (2011 – 2017)

6 Performance metrics

- ✓ Deterministic (KGE)
- ✓ Probabilistic (CRPS)
- ✓ Uncertainty band
- ✓ Ranking models

Main Results

- I. A comprehensive multi-model evaluation of flow uncertainty quantification via error modeling
- II. The streamflow model choice matters for uncertainty estimation because it can change the rank of intermediate models
- III. Post-processing residual errors of long short-term memory (LSTM) network models is consistently the best-performing approach

753

754 **Highlights:**

- 755 • Uncertainty assessments were estimated in 141 basins for 13 hydrologic model
756 predictions
- 757 • Quality of estimates was loosely correlated to model accuracy
- 758 • LSTM models were consistently best performing in terms of KGE and CRPS

759

760

761 **CRedit authorship contribution statement**

762 **Jonathan Romero-Cuellar:** Conceptualisation, Investigation, Formal analysis, Visualisation,
763 Software, Methodology, Writing - original draft. **Rezgar Arabzadeh:** Software, Data curation,
764 Visualisation, Formal analysis. **James R. Craig:** Conceptualisation, Methodology, Writing -
765 review & editing, Supervision, Resources. **Bryan A. Tolson:** Conceptualisation, Methodology,
766 Writing - review & editing, Supervision. **Juliane Mai:** Writing - review & editing.

767

768

769 **Abstract**

770 Probabilistic streamflow predictions are valuable tools for predictive uncertainty estimation,
771 hydrologic risk management, and support for decision-making in water resources. Usually,
772 predictive uncertainty quantification is developed and assessed using only a single hydrological
773 model, making it difficult to generalize to other model configurations. To tackle this issue, we
774 assess changes in the model performance ranking of diverse streamflow models by applying a
775 residual error model post-processing approach to multiple basins and multiple models. This
776 assessment employed 141 basins from the Great Lakes watershed covering the USA and Canada,
777 and 13 diverse streamflow models, which are evaluated using deterministic and probabilistic
778 performance metrics. As the first study to implement probabilistic methods to diverse streamflow
779 models applied to a multitude of basins, the analysis here examines the dependence of
780 probabilistic streamflow estimation quality on model quality. Our findings show that streamflow
781 model choice influences the robustness of probabilistic predictions. It was found that moving
782 from deterministic to probabilistic predictions using a post-processing approach does not change
783 the streamflow model performance ranking for the best and worst deterministic models, but
784 models of intermediate rank in deterministic evaluation do not have consistent ranking when
785 evaluated in probabilistic mode. Post-processing residual errors of long short-term memory
786 (LSTM) network models are consistently the best-performing model in terms of deterministic
787 and probabilistic metrics. This study highlights the significance of combining deterministic
788 streamflow model predictions with residual error models for improving the quality and
789 increasing the value of hydrological predictions, quantifying uncertainty, and facilitating
790 decision-making in operational water management. It also clarifies the degree to which

791 probabilistic predictions depend upon good model performance and can compensate for poor
792 model performance.

793

794

Journal Pre-proofs