1  **Evaluating residual error approaches for post-processing monthly**

2  **and seasonal streamflow forecasts**

3  Fitsum Woldemeskel[1], David McInerney[2],   Julien Lerat[3], Mark Thyer[2], Dmitri Kavetski[2,4],

4  Daehyok Shin[1], Narendra Tuteja[3] and George Kuczera[4]

5  (1) Bureau of Meteorology, VIC, Australia

6  (2) School of Civil, Environmental and Mining Engineering, University of Adelaide, SA, Australia

7  (3) Bureau of Meteorology, ACT, Australia

8  (4) School of Engineering, University of Newcastle, Callaghan, NSW, Australia

9

10  Correspondence email: fitsum.woldemeskel@bom.gov.au

11

12

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

Hydrology and
Earth System
Sciences

Open Access

Discussions

30  **Abstract**

31  Streamflow forecasting is prone to substantial uncertainty due to errors in meteorological forecasts,

32  hydrological model structure and parameterization, as well as in the observed rainfall and streamflow

33  data used to calibrate the models. Statistical streamflow post-processing is an important technique

34  available to improve the probabilistic properties of the forecasts. This study evaluates three residual error

35  models based on the logarithmic (Log), log-sinh (Log-Sinh) and Box-Cox with $\lambda = 0.2$ (BC0.2)

36  transformation schemes and identifies the best performing scheme for post-processing monthly and

37  seasonal (3-months) streamflow forecasts, such as those produced by the Australian Bureau of

38  Meteorology. Using the Bureau's operational dynamic streamflow forecasting system, we carry out

39  comprehensive analysis of the three post-processing schemes across 300 Australian catchments with a

40  wide range of hydro-climatic conditions. Forecast verification is assessed using reliability and sharpness

41  metrics, as well as the Continuous Ranked Probability Skill Score (CRPSS). Results show that the

42  uncorrected forecasts (i.e. without post-processing) are unreliable at half of the catchments. Post-

43  processing using the three residual error models substantially improves reliability, with more than 90%

44  of forecasts classified as reliable. In terms of sharpness, the BC0.2 scheme significantly outperforms the

45  Log and Log-Sinh schemes. Overall, the BC0.2 scheme achieves reliable and sharper-than-climatology

46  forecasts at a larger number of catchments than the Log and Log-Sinh error models. This study is

47  significant because the reliable and sharper forecasts obtained using the BC0.2 post-processing scheme

48  will help water managers and users of the forecasting service to make better-informed decisions in

49  planning and management of water resources.

50  **Keywords**: seasonal streamflow forecasts, residual error models, post-processing, Box-Cox

51  transformation

52

53

54

55

56

57

Hydrology and
Earth System
Sciences

Open Access

EGU

Discussions

58

## **Key points**

60    1. Uncorrected and post-processed streamflow forecasts (using three residual error models, based
61       on the Log, Log-Sinh and BC0.2 transformations respectively) are evaluated over 300 diverse
62       Australian catchments.
63    2. Post-processing enhances streamflow forecast reliability, increasing the percentage of sites with
64       reliable predictions from 50% to over 90%.
65    3. The BC0.2 transformation achieves significantly better forecast sharpness than the Log-sinh and
66       Log transformations, particularly in dry catchments.

67

## 1 Introduction

68

69 Hydrological forecasts provide crucial supporting information on a range of water resource management

70 decisions, including (depending on the forecast lead-time) flood emergency response, water allocation

71 for various uses, and drought risk management (Li et al., 2016; Turner et al., 2017). The forecasts,

72 however, should be thoroughly verified and proved to be of sufficient quality to support decision-making

73 and to meaningfully benefit the economy, environment and society.

74 Sub-seasonal and seasonal streamflow forecasting systems can be broadly classified into two types

75 (Crochemore et al., 2016):

76 *i. Dynamic modelling systems.* Here, a hydrological model is commonly developed at a daily time-step

77 to capture key hydrological processes. The model is calibrated against observed streamflow using

78 historical rainfall and potential evaporation data. Once the model is calibrated, rainfall forecasts from a

79 numerical climate model are used as an input to produce daily streamflow forecasts, which are then

80 aggregated to the time scale of interest and post-processed using statistical models. Examples of

81 operational services based on the dynamic approach include the Australian Bureau of Meteorology's

82 dynamic modelling system (Laugesen et al., 2011; Tuteja et al., 2011; Lerat et al., 2015); the

83 Hydrological Ensemble Forecast Service (HEFS) of the US National Weather Service (NWS) (Brown

84 et al., 2014; Demargne et al., 2014); the Hydrological Outlook UK (HOUK) (Prudhomme et al., 2017);

85 and the short-term forecasting European Flood Alert System (EFAS) (Cloke et al., 2013).

86 *ii. Statistical modelling systems.* Here, a statistical model based on relevant predictors is applied directly

87 at the time scale of interest. A number of predictors have been considered in the literature, including

88 antecedent rainfall and streamflow, soil moisture, depth and extent of snow cover, and climate indices

89 derived from sea surface temperature (Robertson and Wang, 2009, 2011; Wang et al., 2009; Tang and

90 Lettenmaier, 2010; Lü et al., 2016; Zhao et al., 2016). The Bureau of Meteorology's Bayesian Joint

91 Probability (BJP) forecasting system is an example of an operational service based on a statistical

92 approach (Senlin et al., 2017).

93 Hybrid systems that share some characteristics of dynamic and statistical approaches have also been

94 investigated. For example, Robertson et al. (2013) and Humphrey et al. (2016) used dynamic model

95 simulations as predictors for statistical models.

96 Dynamic and statistical approaches have distinct advantages and limitations. Dynamic systems can

97 potentially provide realistic responses in unfamiliar climate situations as it is possible to impose physical

98 constraints in such situations (Wood and Schaake, 2008). In comparison, statistical models have the

99 flexibility to include features that may lead to more reliable predictions. For example, the BJP model

100    uses climate indices (e.g. NINO3.4), which are typically not used in dynamic approaches. That said, the

101    suitability of statistical models for the analysis of non-stationary catchment and climate conditions is

102    questionable (Wood and Schaake, 2008).

103    Streamflow forecasts built on hydrological models are affected by uncertainty in a number of factors,

104    including rainfall forecasts, observed rainfall and streamflow data, as well as the parametric and

105    structural uncertainty of the hydrological model. Progress has been made towards reducing biases and

106    characterizing the sources of uncertainty in streamflow forecasting. These advances include improving

107    rainfall forecasts through post-processing ( Robertson et al., 2013b; Crochemore et al., 2016), accounting

108    for input, parametric and/or structural uncertainty (Kavetski et al., 2006; Kuczera et al., 2006; Renard et

109    al., 2011; Tyralla and Schumann, 2016) and using data assimilation techniques (Dechant and

110    Moradkhani, 2011). Although these steps may improve some aspects of the forecasting system, a

111    residual bias may nonetheless remain. Such bias can only be reduced via post-processing, which, if

112    successful, will improve forecast accuracy and reliability (Madadgar et al., 2014; Lerat et al., 2015).

113    This study focuses on improving streamflow forecasting using dynamic approaches, by identifying

114    residual error models suitable for post-processing hydrological forecasts at monthly and seasonal time-

115    scales. A number of post-processing approaches have been investigated in the literature, including

116    quantile mapping (Hashino et al., 2007), Bayesian frameworks (Pokhrel et al., 2013; Robertson et al.,

117    2013a), as well as methods based on state-space models and wavelet transformations (Bogner and Kalas,

118    2008). Wood and Schaake (2008) used the correlation between forecast ensemble means and

119    observations to generate a conditional forecast. Compared with the traditional approach of correcting

120    individual forecast ensembles, the correlation approach improved forecast skill and reliability. In another

121    study, Pokhrel et al. (2013) implemented a Bayesian Joint Probability (BJP) method to correct biases,

122    update predictions and quantify uncertainty in monthly hydrological model predictions in 18 Australian

123    catchments. The study found that the accuracy and reliability of forecasts improved. More recently,

124    Mendoza et al. (2017) evaluated a number of seasonal streamflow forecasting approaches, including

125    purely statistical, purely dynamical, and hybrid approaches. Based on analysis of catchments

126    contributing to five reservoirs, the study concluded that incorporating catchment and climate information

127    into post-processing improves forecast skill. While the above review mainly focused on post-processing

128    at sub-seasonal and seasonal forecasts (as it is the main focus of the current study), post-processing is

129    also commonly applied to short-range forecasts (e.g. Li et al., 2016; Seo et al., 2006) and to long-range

130    forecasts up to 12 months ahead (Bennett et al., 2016).

131    In most streamflow post-processing approaches, a residual error model is applied to quantify forecast

132    uncertainty. Most residual error models are based on least squares techniques with weights and/or data

133    transformations (e.g. Carpenter and Georgakakos, 2001; Li et al., 2016; Seo et al., 2006). In order to

134    produce post-processed streamflow forecasts, a daily-scale residual error model is used in the calibration

135    of hydrological model parameters, and a monthly/seasonal-scale residual error model used as part of

136    streamflow post-processing to quantify the forecast uncertainty. In a recent study, McInerney et al.

137    (2017) concluded that residual error models based on Box-Cox transformations with fixed parameter

138    values are particularly effective for daily scale predictions, yielding significant improvements in dry

139    catchments. While McInerney et al. (2017) used observed rainfall to force the hydrological model, and

140    evaluated daily streamflow predictions, this study investigates whether these findings generalize to

141    monthly and seasonal forecasts using forecast rainfall.

142    An important aspect of this work is its focus on general findings applicable over diverse hydro-

143    climatological conditions. Most of the studies in the published literature use a limited number of

144    catchments and case studies to test prospective methods. Dry catchments, characterised by intermittent

145    flows and frequent low flows, pose the greatest challenge to hydrological models (Ye et al., 1997;

146    Knoche et al., 2014). Yet the provision of good quality forecasts across a large number of sites is an

147    essential attribute of national scale operational forecasting service, especially in large countries with

148    diverse climatic and catchment conditions, such as Australia.

149    This paper aims to develop streamflow post-processing approaches suitable for use in an operational

150    streamflow forecasting service. More specifically, our aims are:

151    **Aim 1**: Evaluate the value of streamflow forecast post-processing by comparing forecasts with no post-

152    processing (hereafter called 'uncorrected' forecasts) against post-processed forecasts.

153    **Aim 2**: Evaluate three residual error models proposed in recent publications, namely the Log, Box-Cox

154    (McInerney et al., 2017) and Log-Sinh (Wang et al., 2012) schemes, for monthly and seasonal

155    streamflow post-processing.

156    **Aim 3**: Evaluate the generality of results over a diverse range of hydro-climatic conditions, in order to

157    ensure the recommendations are robust in the context of an operational streamflow forecasting service.

158    To achieve these aims, we use the operational monthly and seasonal (3-months) dynamic streamflow

159    forecasting system of the Australian Bureau of Meteorology (Lerat et al., 2015). We evaluate the residual

160    error models across 300 catchments across Australia, with detailed analysis of dry and wet catchments.

161    Forecast verification is carried out using Continuous Ranked Probability Skill Score (CRPSS) as well

162    as metrics measuring reliability and sharpness, which are important aspects of a probabilistic forecast

163    (Wilks, 2011). These metrics are used by the Bureau of Meteorology to describe streamflow forecast

164    performance of the operational service.

165   The rest of the paper is organised as follows. The forecasting methodology is described in Section 2 and

166   application studies are described in Section 3. Results are presented in Section 4, followed by discussions

167   and conclusions in Sections 5 and 6 respectively.

## 2   Seasonal Streamflow forecasting methodology

### 2.1   Overview

170   The streamflow forecasting system adopted in this study is based on the Bureau of Meteorology's

171   dynamic modelling system (Figure 1). This dynamic modelling system uses daily rainfall forecasts as

172   inputs into a daily rainfall-runoff model to produce daily streamflow forecasts. These streamflow

173   forecasts are then aggregated in time and post-processed to produce monthly and seasonal streamflow

174   forecasts, which are issued each month. In general, two steps are involved: simulation and forecasting.

### 2.2   Simulation Step

176   In the simulation step, the daily rainfall-runoff model is calibrated to observed daily streamflow using

177   observed rainfall (Jeffrey et al., 2001) as forcing.

178   The rainfall-runoff model GR4J (Perrin et al., 2003) is used as it has been proven to provide (on average)

179   good performance across a large number of catchments ranging from semi-arid to temperate and tropical

180   humid (Perrin et al., 2003; Tuteja et al., 2011). The calibration of the hydrological model is based on the

181   weighted least squares likelihood function, similar to that outlined in Evin et al. (2014). Markov Chain

182   Monte Carlo (MCMC) analysis is used to estimate posterior parametric uncertainty (Tuteja et al., 2011).

183   Following MCMC analysis, 40 random sets of GR4J parameters are retained and used in the forecast

184   step.

### 2.3   Forecast Step

186   Once the hydrological model is calibrated, daily downscaled rainfall forecast from the Bureau of

187   Meteorology's global climate model, namely the Predictive Ocean Atmosphere Model for Australia

188   POAMA-2 (Hudson et al., 2013), are routed through the hydrological model to produce daily

189   uncorrected streamflow forecasts. The atmospheric component of POAMA-2 uses a spatial scale of

190   approximately $250 \times 250$ km (Charles et al., 2013). To estimate catchment-scale rainfall, a statistical

191   downscaling model based on an analogue approach (Timbal and McAvaney, 2001) was applied. In the

192   analogue approach, local climate information is obtained by matching analogous previous situations to

193   the predicted climate. To this end, an ensemble of 166 rainfall forecast time series (33 POAMA

194   ensembles $\times$ 5 replicates from downscaling + 1 ensemble mean) were generated. These forecasts are

195   then input into GR4J and propagated using the 40 GR4J parameter sets to obtain 6640 ($166 \times 40$) daily

196   streamflow forecasts. The daily streamflow forecasts generated using GR4J are then aggregated to

Hydrology and
Earth System
Sciences
Discussions

197 monthly and seasonal time scales to produce ensembles of 6640 uncorrected monthly and seasonal
198 forecasts.

## 2.4 Streamflow post-processing

200 Post-processing of streamflow forecasts is intended to remove systemic biases in the mean, variability
201 and persistence of the uncorrected forecasts, which arise due to inaccuracies in the downscaled rainfall
202 forecasts (e.g. errors in downscaled forecast rainfall from approximately a 250 km grid to the catchment
203 scale) and in the hydrological model (e.g. due to the effects of data errors on the model calibration and
204 due to structural errors in the model itself).

205 The streamflow post-processing method used in this work consists of fitting a statistical model to the
206 streamflow forecast residual errors, defined by the differences between the observed and forecast
207 streamflow time series over a calibration period. Typically these residual errors are heteroscedastic and
208 exhibit persistence. Heteroscedasticity is handled using data transformations (e.g. the Box-Cox
209 transformation), whereas persistence is represented using autoregressive models (e.g., the lag-one
210 autoregressive model, AR(1)). We begin by describing the two major steps of the streamflow post-
211 processing procedure (Sections 2.4.1 and 2.4.2), and then describe the transformations under
212 consideration (Section 2.5).

### 2.4.1 Calibration of residual error model parameters

214 The parameters of the streamflow post-processing model are calibrated in the following three steps:

215 *Step 1*: Compute the transformed forecast residuals for month or season $t$ of the calibration period:

$$\eta_t = Z(\widetilde{Q_t}) - Z(Q_t^F) \tag{1}$$

217 where $\eta_t$ is the normalised residual, $\widetilde{Q_t}$ is the observed streamflow, $Q_t^F$ is the median of the uncorrected
218 streamflow forecast ensemble, and $Z$ is a transformation function used to reduce the heteroscedasticity
219 and skewness of the residuals (Wang et al., 2012; McInerney et al., 2017). The data transformation
220 functions are detailed in Section 2.5.

221 *Step 2*: Compute the standardised residuals according to:

$$\nu_t = (\eta_t - \mu_\eta^{m(t)}) / \sigma_\eta^{m(t)} \tag{2}$$

223 where $\mu_\eta^{m(t)}$ and $\sigma_\eta^{m(t)}$ are the monthly mean and standard deviation of the residuals in the calibration
224 period for the month $m(t)$. The standardisation process in equation (2) aims to account for seasonal
225 variations in the distribution of residuals.

226  The quantities $\mu_\eta^{m(t)}$ and $\sigma_\eta^{m(t)}$ are calculated independently as the sample mean and standard deviation of

227  residuals for each monthly period (for a monthly forecast) or three-monthly period (for seasonal

228  forecasts). The standardised residuals $v_t$ are assumed to have a zero mean and unit standard deviation.

229  *Step 3*: Assume the standardised residuals are described by a first order autoregressive (AR(1)) model:

230  $$v_{t+1} = \rho v_t + y_{t+1} \tag{3}$$

231  where $\rho$ is the AR(1) coefficient and $y_{t+1} \sim N(0, \sigma_y)$ is the innovation assumed to follow a Gaussian

232  distribution.

233  The parameters $\rho$ and $\sigma_y$ are estimated based on the method of moments: $\rho$ is set to the sample auto-

234  correlation of the standardized residuals $\mathbf{v}$, and $\sigma_y$ is set to the sample standard deviation of the

235  observed innovations $\mathbf{y}$, which are calculated from the standardized residuals $\mathbf{v}$ by re-arranging

236  equation (3).

### 2.4.2 Streamflow forecasting

238  Once the streamflow post-processor has been calibrated, the post-processed streamflow forecasts for a

239  given period are computed. For a given ensemble member *j*, the following steps are applied (note the

240  additional subscript *j* for the ensemble number):

241  *Step 1*: Sample the innovation $y_{t+1,j} \leftarrow N(0, \sigma_y)$.

242  *Step 2*: Generate the standardized residuals $v_{t+1,j}$ using equation (3). Here $v_{t,j}$ is determined using

243  equation (2) and $\eta_{t,j}$ using equation (1), which uses the streamflow forecasts and observations from the

244  previous time step *t*.

245  *Step 3*: Compute the normalized residuals $\eta_{t+1,j}$ by "de-standardizing" $v_{t+1,j}$:

246  $$\eta_{t+1,j} = \sigma_\eta^{m(t)} v_{t+1,j} + \mu_\eta^{m(t)} \tag{4}$$

247  *Step 4*: Back-transform each normalized residual $\eta_{t+1,j}$ to obtain the post-processed streamflow forecast:

248  $$Q_{t+1,j}^{PP} = Z^{-1}[Z(Q_{t+1}^F) + \eta_{t+1,j}] \tag{5}$$

249  Steps 1-4 are repeated for all ensemble members (6640 in our case).

250  Note that the above algorithm may occasionally generate negative streamflow, which is then set to zero.

251  This aspect is discussed in Section 5.6.

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

### 2.5 Transformations used in the residual error model

252

253 The observed streamflow and median streamflow forecast are transformed in Step 1 of streamflow post-

254 processing (Section 2.4.1), to account for the heteroscedasticity and skewness of the forecast residuals.

255 To achieve Aim 2 of this study, we trial three different transformations, namely the logarithmic, log-

256 sinh and Box-Cox transformations.

#### 2.5.1 Logarithmic (Log) transformation

257

258 The logarithmic (Log) transformation is

259
$$Z(Q) = \log(Q + c) \tag{6}$$

260 The offset $c$ ensures the transformed flows are defined when $Q = 0$. Here we set $c = 0.01 \times (\tilde{Q})_{ave}$

261 , where $(\tilde{Q})_{ave}$ is the average observed streamflow over the calibration period. The use of a small fixed

262 value for $c$ is common in the literature for coping with zero flow events (Wang et al., 2012).

#### 2.5.2 Log-Sinh transformation

263

264 The Log-Sinh transformation (Wang et al., 2012) is

265
$$Z(Q) = \frac{1}{b} \log \left[ \sinh(a + bQ) \right] \tag{7}$$

266 The parameters $a$ and $b$ are calibrated for each month by maximising the p-value of the Shapiro-Wilk

267 test (Shapiro and Wilk, 1965) for normality of the residuals, $\nu$. This pragmatic approach is part of the

268 existing Bureau's operational dynamic streamflow forecasting system (Lerat et al., 2015).

#### 2.5.3 Box-Cox

269

270 The Box-Cox transformation (Box and Cox, 1964) is

271
$$Z(Q; \lambda, c) = \frac{(Q + c)^{\lambda} - 1}{\lambda} \tag{8}$$

272 where $\lambda$ is a power parameter and $c = 0.01 \times (\tilde{Q})_{ave}$. Following the recommendations of McInerney et

273 al. (2017), the parameter $\lambda$ is fixed to 0.2. This avoids the need to calibrate $\lambda$, and related problems with

274 doing so.

#### 2.5.4 Rationale for selecting transformational approaches

275

276 The Log transformation is a widely used transformation that is simple to implement; McInerney et al.

277 (2017) reported that in daily scale modelling it produced the best reliability in perennial catchments

278 (from a set of eight residual error schemes, including standard least squares, weighted least squares, BC,

279    Log-Sinh and reciprocal transformation). However, the Log transformation performed poorly in

280    ephemeral catchments, where its precision was far worse than in perennial ones.

281    The Log-Sinh transformation is an alternative to the Log and BC transformations proposed by Wang et

282    al. (2012) to improve the precision at higher flows. The Log-Sinh approach has been extensively applied

283    to water forecasting problems (see for example, Del Giudice et al., 2013; Robertson et al., 2013b, Bennett

284    et al., 2016). However, McInerney et al. (2017) found that in daily scale modelling of perennial

285    catchments, when using observed rainfall, the Log-Sinh scheme did not improve on the Log

286    transformation (its parameters tend to calibrate to values for which the Log-Sinh transformation reduces

287    to the Log transformation).

288    Finally, the BC transformation with fixed $\lambda = 0.2$ is recommended by McInerney et al. (2017) as one of

289    only two schemes (from the set of eight, see above) that achieve "Pareto-optimal" (e.g., Cohon and

290    Marks, 1975) performance in terms of reliability, precision and bias, across both perennial and

291    ephemeral catchments. McInerney et al. (2017) also found that calibrating $\lambda$ did not generally improve

292    predictive performance, due to the inferred value being dominated by the fit to the low flows at the

293    expense of the high flows.

### 2.6    Summary

295    In the remainder of the paper, the term "uncorrected forecasts" refers to streamflow forecasts obtained

296    using steps in Sections 2.1-2.3, and the term "post-processed forecasts" refers to forecasts based on a

297    streamflow post-processing model, which includes the standardization and AR(1) model from Section

298    2.4, as well as a transformation (Log, Log-Sinh or BC0.2) from Section 2.5. As the streamflow residual

299    error models considered in this work differ solely in the transformation used, they will be referred to as

300    the Log, Log-Sinh and BC0.2 schemes.

## 3    Application

### 3.1    Data

303    A comprehensive set of 300 catchments representative of the diverse Australian hydro-climatic

304    conditions is used, with locations shown in Figure 2. In each catchment, data from 1980-2008 is used.

305    Observed daily rainfall data was obtained from the Australian Water Availability Project (AWAP)

306    (Jeffrey et al., 2001). Potential evaporation and observed streamflow data were obtained from the Bureau

307    of Meteorology. Rainfall forecasts from POAMA-2 were downscaled based on an analogue approach

308    (Timbal and McAvaney, 2001). These 300 sites are currently being evaluated as part of the expansion

309    of dynamic modelling for the seasonal streamflow forecasting service of the Bureau of Meteorology.

310    The figure also shows the Koppen climate zones.

### 3.2 Catchment classification

311

312 The performance of the residual error models is evaluated separately in dry versus wet catchments. In

313 this work, the classification of catchments into dry and wet is based on the aridity index (AI) according

314 to the following equation

$$\text{AI} = \frac{P}{PET} \qquad (9)$$

315

316 where P is the total rainfall volume and PET is the total potential evapotranspiration volume. The aridity

317 index has been used extensively to identify drought and wetness conditions of hydrological regimes (

318 Zhang et al., 2009; Carrillo et al., 2011; Sawicz et al., 2014).

319 Catchments with AI < 0.5 are categorised as "dry", which corresponds to hyper-arid, arid and semi-arid

320 classifications suggested by the United Nations Environment Programme (Middleton et al., 1997).

321 Conversely, catchments with AI ≥ 0.5 are classified as "wet". Overall, about 28% of catchments used in

322 this work are classified as dry.

### 3.3 Cross-validation procedure

323

324 The forecast verification is carried out using a moving-window cross-validation framework, as shown

325 in Figure 3. Suppose we are validating the streamflow forecasts in year $j$ ( $j = 1990$ in Figure 3). In this

326 case the calibration is carried out using all years except $j$, $j+1$, $j+2$, $j+3$ and $j+4$. The four-year period

327 after year $j$ are excluded to avoid the effects of memory in the hydrological model. The process is then

328 repeated for each year during 1980-2008. Once the validation has been carried out for each year, the

329 results are concatenated together to produce a single "validation" time series, for which the verification

330 metrics are calculated.

### 3.4 Verification metrics

331

332 The goal of the forecasting exercise is to maximise sharpness without sacrificing reliability (Gneiting et

333 al., 2005; Wilks, 2011; Bourdin et al., 2014). Therefore the performance of uncorrected and post-

334 processed streamflow forecasts is evaluated using reliability and sharpness metrics, as well as the

335 Continuous Ranked Probability Skill Score (CRPSS, see section 3.4.3). Note that the Bureau of

336 Meteorology uses Root Mean Squared Error (RMSE) and Root Mean Squared Error in Probability

337 (RMSEP) scores in the operational service in addition to CRPSS, however, RMSE and RMSEP results

338 have not been included in the current paper.

339 Forecast verification metrics are computed separately for each forecast month. To facilitate the

340 comparison and evaluation of streamflow forecast performance in different streamflow regimes, the high

341    and low flow months are defined using long-term average streamflow data calculated for each month –

342    "high flow" months are the 6 months with the highest average streamflow, while low flows are the 6

343    months with the lowest average streamflow. Note that although the verification metrics are computed

344    for each month separately, indices denoting the month are excluded from Equations (10), (11) and (12)

345    below to avoid cluttering the notation.

### 3.4.1 Reliability

347    The reliability of forecasts is evaluated using the Probability Integral Transform (PIT) (Dawid, 1984;

348    Laio and Tamea, 2007). To evaluate and compare reliability across 300 catchments, the p-value of the

349    Kolmogorov-Smirnov (KS) test applied to the PIT is used. In this study, forecasts with PIT plots where

350    the KS test yields a p-value ≥ 5% are classified as "reliable".

### 3.4.2 Sharpness

352    The sharpness of forecasts is evaluated using the ratio of inter-quantile ranges (IQR) of streamflow

353    forecasts and a historical reference (Tuteja et al., 2016). The following definition is used:

$$IQR_q = \frac{1}{n}\sum_{i=1}^{n} \frac{F_i(100-q) - F_i(q)}{C_i(100-q) - C_i(q)} \times 100\ \% \tag{10}$$

355    where $IQR_q$ is the $IQR$ value corresponding to percentile $q$, $F_i(q)$, and $C_i(q)$ are the $q$th percentiles of

356    forecast and historical reference for years $i = 1, 2, ..., N$, respectively.

357    An $IQR_q$ of 100% indicates a forecast with the same sharpness as the reference, an $IQR_q$ below 100%

358    indicates forecasts that are sharper (predictive limits that are smaller) than the reference, and an $IQR_q$

359    above 100% indicates forecasts that are less sharp (predictive limits are wider) than the reference. We

360    consider $IQR_{99}$, i.e., the $IQR$ at the 99 percentile, in order to detect forecasts with unreasonably long

361    tails in their predictive distributions.

### 3.4.3 CRPS skill score (CRPSS)

363    The $CRPS$ metric quantifies the difference between a forecast distribution and observations, as follows

364    (Hersbach, 2000):

$$CRPS = \int_{-\infty}^{\infty} \left[F_f(y) - F_o(y)\right]^2 dy \tag{11}$$

366    where $F_f$ and $F_o$ are the cumulative distribution functions (cdfs) of the streamflow forecast and

367    observation, respectively. The cdf of the observation is taken as the Heaviside step function at the

368    observed point value.

Hydrology and
Earth System
Sciences

Open Access

Discussions

369     The $CRPS$ summarises the reliability, sharpness and bias attributes of the forecast (Hersbach, 2000). A

370     "perfect" forecast – namely a point prediction that matches the actual value of the predicted quantity –

371     has $CRPS^P = 0$. In this work, we use $CRPS$ skill score, CRPSS, defined by:

372     $$CRPSS = \frac{CRPS^F - CRPS^C}{CRPS^P - CRPS^C} \times 100\%$$     (12)

373     where $CRPS^F$, $CRPS^C$ and $CRPS^P$ represent the $CRPS$ value for model forecast, climatology and

374     "perfect" forecast respectively. A higher CRPSS indicates better performance, with a value of 0

375     representing the same performance as climatology.

376     **3.4.4   Historical reference**

377     The IQR and CRPSS metrics are defined as skill scores relative to a reference forecast. In this work, we

378     use the climatology as the reference forecast, as it represents the long-term climate condition. To

379     construct these "climatological forecasts", we used the same historical reference as the operational

380     seasonal streamflow forecasting service of the Bureau of Meteorology. This reference is resampled from

381     a Gaussian probability distribution fitted to the observed streamflow data transformed using the log-sinh

382     transformation (Equation 7). This approach leads to more stable and continuous historical reference

383     estimates than sampling directly from the empirical distribution of historical streamflow, and can be

384     computed at any percentile (which facilitates comparison with forecast percentiles). Although the choice

385     of a particular reference affects the computation of skill scores, it does not affect the ranking of error

386     models when the same reference is used, which is the main aim of this paper.

387     **3.4.5   Summary Skill: Summarising forecast performance using multiple metrics**

388     When evaluating forecast performance, a focus on any single individual metric can lead to misleading

389     interpretations. For example, two forecasts might have a similar sharpness, however, if one is not

390     reliable, then it can over or underestimate risk which could lead to a sub-optimal decision by forecast

391     users (e.g. a water resources manager).

392     Given inevitable trade-offs between individual metrics (McInerney et al., 2017), it is important to

393     consider multiple metrics jointly rather than individually. Following the approach suggested by Gneiting

394     et al. (2007), we consider a forecast to have "high skill" when it is both reliable and has a better sharpness

395     than climatology. To determine the "summary skill" of the forecasts in each catchment, we evaluate the

396     total number of months (out of 12) in which forecasts are reliable (i.e., with a p-value greater than 5%)

397     and sharper than the climatology (i.e., IQR99 < 100%). Accordingly, a catchment is classified as having

398     high (low) summary skill if it has a 10-12 months (0-2 months) with reliable forecasts that are shaper

399    than climatology. Note that we do not include the CRPSS in the summary skill, because the CRPSS does

400    not provide an independent measure of forecast attribute (see Section 3.4.3 for more details).

401    A table providing the percentage of catchments with high and low summary skills is used to summarise

402    forecasts performance. In addition, to identify any geographic trends in the forecast performance, the

403    summary skills are plotted on a map. The summary skills together with individual skill score values are

404    used to evaluate the overall forecast performance.

## 4    Results

406    Results for monthly and seasonal streamflow forecasts are now presented. Section 4.1 compares the

407    uncorrected and post-processed streamflow forecast performance. Section 4.2 evaluates the performance

408    of post-processed streamflow forecasts obtained using the Log, Log-Sinh and BC0.2 schemes. The

409    CRPSS, reliability and sharpness metrics are presented in Figure 4 and Figure 5 for monthly and seasonal

410    forecasts respectively.

411    Initial inspection of results found considerable overlap in the performance metrics achieved by the error

412    models. To determine whether the differences in metrics are consistent over multiple catchments, the

413    Log and Log-Sinh schemes are compared to the BC0.2 scheme. This comparison is presented in

414    Figure 6 and Figure 7 for monthly and seasonal forecasts respectively. The BC0.2 scheme is taken as

415    the baseline because inspection of Figure 4 and Figure 5 suggests that the BC0.2 scheme has better

416    median sharpness than the Log and Log-Sinh schemes, over all the catchments and for both high and

417    low flow months individually.

418    The streamflow forecast time-series and corresponding skill for a single representative site, Dieckmans

419    Bridge, are presented in Figures 8 and 9, respectively.

420    The results are presented separately for wet and dry catchments, as well as separately for high and low

421    flow months (Sections 3.2 and 3.4). The summary skills of the monthly and seasonal forecasts are

422    presented in Figure 10 and Figure 11. The figures include a histogram of summary skills across all

423    catchments to enable comparison between the uncorrected and the post-processing approaches.

### 4.1    Comparison of uncorrected and post-processed streamflow forecasts: Individual
        metrics

426    In terms of CRPSS, largest improvement as a result of post-processing using the Log, Log-Sinh and

427    BC0.2 schemes occurs in dry catchments for both monthly (Figure 4c) and seasonal forecasts (Figure

428    5c). For example, when post-processing is used with the three transformation schemes, the median

429    CRPSS of monthly forecasts in dry catchments increases from approximately 7% (high flow months)

Hydrology and
Earth System
Sciences
Discussions

430   and -15% (low flow months) to more than 10% (Figure 4c) for both high and low flows. Visible

431   improvement is also observed in dry catchments for seasonal forecasts, however, the improvement is

432   not as pronounced as for monthly forecasts (Figure 5c).

433   In terms of reliability, the performance of uncorrected streamflow forecasts is poor, with about 50% of

434   the catchments being characterized by unreliable forecasts at both the monthly and seasonal time scales

435   (Figure 4 and Figure 5, middle row). In comparison, post-processing using the three transformation

436   approaches produces much better reliability, achieving reliable forecasts in more than 90% of the

437   catchments.

438   In terms of sharpness, the uncorrected forecasts and the BC0.2 post-processed forecasts are generally

439   sharper than forecasts generated using the other transformations (Figures 4g and 5g). The use of post-

440   processing achieves much better sharpness than uncorrected forecasts for low flow months, particularly

441   in dry catchments. For example, for low flow months in dry catchments (Figure 4i), the median IQR99

442   is greater than 200% while similar values range between 40-100% for post-processed forecasts.

443   Similarly, for seasonal forecasts, post-processing approaches improve the median sharpness from in

444   excess of 150% (uncorrected forecasts) to 50%-110% (Figure 45i).

## 445   4.2   Comparison of residual error models for post-processing: Individual metrics

446   In terms of CRPSS, Figure 4 (a, b, c) and Figure 5 (a, b, c) show considerable overlap in the boxplots

447   corresponding to all three residual error models, both in wet and dry catchments. This finding suggests

448   little difference in the performance of the residual error models, and is further confirmed by Figure 6 (a,

449   b, c) and Figure 7 (a, b, c), which show boxplots of the differences between the CRPSS of the Log and

450   Log-Sinh schemes versus the CRPSS of the BC0.2 scheme. Across all catchments, the distribution of

451   these differences is approximately symmetric with a mean close to 0. In dry catchments, the BC0.2

452   slightly outperforms the Log scheme for high flow months and the Log-Sinh scheme slightly

453   outperforms the Log scheme for low flow months. Overall, these results suggest that none of the Log,

454   Log-Sinh or BC0.2 schemes is consistently better in terms of CRPSS values.

455   In terms of reliability, post-processing using any of the three residual error models produces reliable

456   forecasts at both monthly and seasonal scales, and in the majority of the catchments (Figure 4 and Figure

457   5, middle row). The median p-value is approximately 60% for monthly forecasts compared with 45%

458   for seasonal forecasts. This indicates that better reliability is achieved at shorter lead times. Median

459   reliability is somewhat reduced when using the BC0.2 scheme compared to the Log and Log-Sinh

460   schemes in wet catchments (Figure 6e), but not so much in dry catchments (Figure 8f). Nevertheless,

461   the monthly and seasonal forecasts are reliable in 96% and 91% of the catchments, respectively. The

462 corresponding percentages for the Log scheme are 97% and 94%, and for Log-Sinh they are 95% and

463 90%.

464 In terms of sharpness, the BC0.2 scheme produces much sharper forecasts than the Log and Log-Sinh

465 schemes. This finding holds in all cases (i.e., high/low flow months and wet/dry catchments), both for

466 monthly and seasonal forecasts (Figure 4 and Figure 5, bottom row). The plot of differences in the

467 sharpness metric (Figure 6 and Figure 7, bottom row) clearly highlights this improvement. In half of the

468 catchments, during both high and low flow months, the BC0.2 scheme improves the IQR99 by 30% or

469 more compared to the Log and Log-Sinh schemes. In dry catchments, the magnitude of the

470 improvements are higher than wet catchments. For example, in dry catchments during high flow months,

471 the BC0.2 scheme improves on the IQR99 of Log and Log-Sinh by 40-60% in over a half of the

472 catchments, and by as much as ~170%-190% in a quarter of the catchments.

473 To highlight the implication of these results, a representative streamflow forecast time-series at

474 Dieckmans Bridge catchment (site id: 145010A) is shown in Figure 8 and performance metrics

475 calculated over six high and low flow months are shown in Figure 9. In terms of reliability, the

476 uncorrected forecast has a number of observed data points outside the 99% predictive range (Figure 8a).

477 This is an indication that the forecast is unreliable. This finding can also be confirmed from the

478 corresponding p-value in Figure 9, which shows that the forecast is below the reliability threshold during

479 most of the high flow months and also during some low flow months. In terms of sharpness, Log and

480 Log-Sinh schemes produce a wider 99% predictive range than BC0.2 (Figures 8 and 9).

481 **4.3  Comparison of summary skill between uncorrected and post-processing approaches**

482 Figure 10 and Figure 11 show the geographic distribution of the summary skill of the uncorrected and

483 post-processing approaches for monthly and seasonal forecasts respectively. The summary skill

484 aggregates multiple verifications metrics: it represents the number of months with streamflow forecasts

485 that are both reliable and exhibit a sharpness that is better than climatology. Table 1 provides a summary

486 of the percentage of catchments with high and low summary skill for the uncorrected and post-processing

487 approaches for monthly and seasonal forecasts. Catchments with high (low) summary skill are defined

488 as those with 10-12 months (0-2 months) with forecasts that are reliable and sharper than climatology.

489 At the monthly scale (Figure 10 and Table 1), we obtain the following key findings:

490 • Uncorrected forecasts perform worse than post-processing techniques in the sense that they have

491 low summary skill in the largest percentage of catchments (16%). The percentage of catchments

492 where high summary skill is achieved is 40%.

- Post-processing forecasts with the Log and Log Sinh scheme, reduces the percentage of catchments with low summary skills to 2% and 7% respectively. However, the percentage of catchments with high summary skill also decreases (in comparison to unprocessed forecasts), to 33% for both Log and Log-Sinh.

- Post-processing with the BC0.2 scheme provides the best performance, with the smallest percentage of catchments with low summary skills (<1%) and the largest percentage of catchments with high summary skills (84%). Figure 10 shows the improvement achieved by the BC0.2 scheme (compared to the Log/Log Sinh schemes) is most pronounced in NSW and in the tropical catchments in QLD and NT. The few catchments where the BC0.2 scheme does not achieve a high summary skill are located in the north and north-west of Australia.

The findings for seasonal forecasts (Figure 11 and Table 1) are as follows:

- Log scheme has the largest percentage (19%) of catchments with low summary skill and a relatively small percentage of catchments (9%) with high summary skill (9%).

- Post-processing forecasts with the Log and Log-Sinh schemes reduces the percentages of catchments with low summary skill to 18% and 17% respectively. The percentage of catchments with high summary skill increases to 12% and 22% respectively.

- Post-processing with the BC0.2 scheme once again provides a clear improvement: it produces forecasts with low summary skill in only 2% of the catchments, and achieves high summary skill in 54% of the catchments. Figure 11, shows that similar to monthly forecasts, the biggest improvements occur in the NSW and Queensland regions of Australia.

Overall, the summary skills of post-processing approaches are lower for seasonal forecasts than for monthly forecasts. Table 1 shows that, across all schemes, BC0.2 results in a larger percentage of catchments with low summary skill and a larger percentage of catchments with high summary skill.

## 4.4 Summary

Sections 4.1-4.3 show that post-processing produces major improvements in reliability, as well as CRPSS and sharpness, particularly in dry catchments. Although all three residual error models under consideration provide improvements in some of the performance metrics, the BC0.2 scheme consistently produces better sharpness than the Log and Log-Sinh schemes, while maintaining similar reliability and CRPSS. This finding holds for both monthly and, to a less degree, seasonal forecasts. Of the three residual error models, the BC0.2 scheme improves by the largest margin the percentage of sites and the number of months where the post-processed forecasts are reliable and sharper than climatology.

## 5   Discussion

### 5.1   Benefit of post-processing

524   A comparison of uncorrected and post-processed streamflow forecasts was provided in Section 4.1.
Uncorrected forecasts have reasonable sharpness (except for dry catchments), but suffer from low
reliability: uncorrected forecasts are unreliable at approximately 50% of the sites. In wet catchments,
poor reliability is due to overconfident forecasts, which appears a common concern in dynamic
forecasting approaches (Wood and Schaake, 2008). In dry catchments, uncorrected forecasts are both
unreliable and exhibit poor sharpness. Post-processing is thus particularly important to correct for these
shortcomings and improve forecast skill. In this study, all post-processing models provide a clear
improvement in reliability and sharpness, especially in dry catchments. The value of post-processing is
more significant in dry catchments than in wet catchments (Figure 4 and Figure 5). This finding can be
attributed to the challenge of capturing key physical processes in modelling dry and ephemeral
catchments (Ye et al., 1997) as well as the challenge of achieving accurate rainfall forecasts in arid areas.
In such cases, the hydrological model forecasts are particularly poor and leave a lot of room for
improvement: post-processing can hence make a big difference on the quality of results.

### 5.2   Interpretation of differences between residual error models

We now discuss the large differences in sharpness between the BC0.2 scheme versus the Log and Log-
Sinh schemes. The Log-Sinh residual error model was designed by Wang et al. (2012) in order to
improve the reliability and sharpness of predictions, particularly for high flows, and has worked well
when used as part of statistical modelling system for operational streamflow forecasts by the Bureau of
Meteorology. The Log-Sinh transformation corresponds to a variance stabilizing function that (for
certain parameter values) tapers off for high flows. In theory, this feature can prevent the explosive
growth of predictions for high flows that can occur with the log and Box-Cox residual error models
(especially when $\lambda < 0$).

McInerney et al. (2017) found that, when modelling perennial catchments at the daily scale, the Log-
Sinh scheme did not achieve better sharpness than the Log scheme; instead, the parameters for the Log
scheme tended to converge to values for which the tapering off of the Log-Sinh scheme occurs well
outside the range of simulated flows, and hence the Log-Sinh scheme effectively reduces to the Log
scheme. In contrast, the Box-Cox error model when using a fixed $\lambda > 0$ has a variance-stabilizing
function that gradually flattens as streamflow increases, i.e., it exhibits the "desired" tapering-off
behaviour.

Hydrology and
Earth System
Sciences
Discussions
Open Access

555 Our findings in this study confirm the insights of McInerney et al. (2017) – namely that the Log-Sinh

556 scheme produces comparable sharpness to the Log scheme – across a larger number of catchments. This

557 finding indicates that insights from modelling residual errors at the daily scale apply at least to some

558 extent to streamflow forecast post-processing at the monthly and seasonal scales. Note the minor

559 difference in the treatment of the offset parameter $c$ in equation (6): in the Log scheme used in McInerney

560 et al. (2017) this parameter is inferred, whereas in this study it is fixed a priori. This minor difference

561 does not impact on the qualitative behaviour of the error models, as described earlier in this section. The

562 BC0.2 scheme provides an opportunity to further improve forecast performance relative to what is

563 currently possible using the Log and Log-Sinh schemes when used as residual error post-processor of

564 forecasts in a dynamical modelling systems.

565 **5.3 Importance of using multiple metrics to assess forecast performance**

566 The study results show that relying on a single metric for evaluating forecast performance can lead to

567 sub-optimal conclusions. For example, if one considers the CRPSS metric alone, all post-processing

568 schemes yield comparable performance and there is no basis for favouring any single one of them.

569 However, once sharpness is taken into consideration explicitly, the BC0.2 scheme can be recommended

570 due to significantly better sharpness than the Log and Log-Sinh schemes. Similarly, comparisons based

571 solely on CRPSS might suggest reasonable performance of the uncorrected forecasts (55%-80% of

572 months have CRPSS > 0 depending on high/low flow months and monthly/seasonal forecasts), yet once

573 reliability is considered explicitly, it is found that uncorrected forecasts are unreliable at approximately

574 50% of the catchments. Note that, for example, CRPSS reflects an implicitly weighted combination of

575 reliability, sharpness and bias characteristics of the forecasts (Hersbach, 2000), whereas the reliability

576 and sharpness metrics are specifically designed to target reliability and sharpness attributes respectively.

577 These findings highlight the value of multiple independent performance metrics and diagnostics that

578 evaluate specific attributes of the forecasts, and highlight important limitations of aggregate measures

579 of performance (Clark et al., 2011).

580 A number of challenges and questions remain in regards to selecting the verification metrics for specific

581 forecasting systems and applications. An important question is how to include user needs into a forecast

582 verification protocol. This could be accomplished by tailoring the evaluation metrics to the requirements

583 of users. Another key question is to what extent do measures of forecast skill correlate to the economic

584 and/or social value of the forecast? This question was investigated by Murphy and Ehrendorfer (1987)

585 and Wandishin and Brooks (2002), who found the relationship between quality and value of a forecast

586 to be essentially nonlinear: an increase in forecast quality may not necessarily lead to a proportional

587 increase in its value.

### 5.4    Importance of performance evaluation over large numbers of catchments

When designing an operational forecast service for locations with streamflow regimes as diverse and variable as in Australia (Taschetto and England, 2009), it is essential to thoroughly evaluate multiple modelling methods over multiple locations to ensure the findings are sufficiently robust and general. This was the major reason for considering the large set of 300 catchments in our study. This setup also yields valuable insights into spatial patterns in forecast performance. For example, the Log and Log-Sinh schemes perform relatively well in catchments in South-Eastern Australia, and relatively worse in catchments in Northern and North-Eastern Australia (Figure 10 and Figure 11). In contrast, the BC0.2 scheme performs well across the majority of the catchments in all regions included in the evaluation. The evaluation over a large number of catchments in different hydro-climatic regions is clearly beneficial to establish the robustness of post-processing methods. Restricting the analysis to a smaller number of catchments would have led to less conclusive findings.

### 5.5    Implication of results for water resource management

The management of water resources, for example, deciding which water source to use for a particular purpose or allocating environmental flows, requires an understanding of the current and future availability of water. For water resources systems with long hydrological records, water managers have devised techniques to evaluate current water availability, water demand and losses. However, one of the main unknowns is the volume of future system inflows. Streamflow forecasts thus provide crucial information to water managers and users regarding the future availability of water, thus helping reduce uncertainty in decision making. The ability of the BC0.2 post-processing scheme to improve forecast sharpness (precision) while maintaining forecast accuracy and reliability can hence lead to improved operational planning and management of water resources.

### 5.6    Treatment of zero flows

The post-processing approach using the three residual error models described above does not make special provision for zero flows in the calibration approach. Robust handling of zero flows in statistical models is an active research area (Wang and Robertson, 2011; Smith et al., 2015), and advances in this area are certainly relevant to seasonal streamflow forecasting.

## 6    Conclusions

This study focused on developing robust streamflow forecast post-processing schemes for an operational forecasting service at the monthly and seasonal time scales. For such forecasts to be useful to water managers and decision-makers, they should be reliable and exhibit sharpness that is better than climatology.

620 We investigated streamflow forecast postprocessor schemes employing residual error models based on

621 three data transformations, namely the logarithmic (Log), log-sinh (Log-Sinh) and Box-Cox

622 transformation with $\lambda = 0.2$ (BC0.2). The Australian Bureau of Meteorology's dynamic modelling

623 system was used as the platform for the empirical analysis, which was carried out over 300 Australian

624 catchments with diverse hydro-climatic conditions.

625 The outcomes of this study are:

626     1. Uncorrected forecasts (no post-processing) perform poorly in terms of reliability, which is an

627        indication that forecast uncertainties are misrepresented. All three post-processing schemes

628        substantially improve the reliability of streamflow forecasts, both in terms of the dedicated

629        reliability metric and in terms of the summary skill given by the CRPSS;

630     2. From the post-processing schemes considered in this work, the BC0.2 scheme is found best

631        suited for operational application. The BC0.2 scheme provides the sharpest forecasts without

632        sacrificing reliability, as measured by the reliability and CRPSS metrics. In particular, the BC0.2

633        scheme produces forecasts that are both reliable and sharper than climatology at substantially

634        more sites than the alternative Log and Log-Sinh schemes.

635 In conclusion, this study developed a robust streamflow forecast post-processing scheme that achieves

636 reliable and consistently sharper-than-climatology streamflow forecasts. This scheme is well suited for

637 operational application, and offers the opportunity to improve decision support, especially at sites where

638 climatology is presently used to guide operational decisions.

639 **7   Data Availability**

640 The data underlying this research can be accessed from the following links: Observed rainfall data

641 (http://www.bom.gov.au/climate/); POAMA rainfall forecast (http://poama.bom.gov.au/); and observed

642 streamflow data (http://www.bom.gov.au/waterdata/).

643 **8   Acknowledgments**

647

648

Hydrology and
Earth System
Sciences

Open Access

EGU

Discussions

## 9 References

649

650 Bennett, J. C., Wang, Q. J., Li, M., Robertson, D. E. and Schepen, A.: Reliable long-range ensemble
651 streamflow forecasts: Combining calibrated climate forecasts with a conceptual runoff model and a
652 staged error model, Water Resour. Res., 52(10), 8238–8259, doi:10.1002/2016WR019193, 2016.

653 Bogner, K. and Kalas, M.: Error-correction methods and evaluation of an ensemble based hydrological
654 forecasting system for the Upper Danube catchment, Atmos. Sci. Lett., 9(2), 95–102,
655 doi:10.1002/asl.180, 2008.

656 Bourdin, D. R., Nipen, T. N. and Stull, R. B.: Reliable probabilistic forecasts from an ensemble reservoir
657 inflow forecasting system, Water Resour. Res., 50(4), 3108–3130, doi:10.1002/2014WR015462, 2014.

658 Box, G. E. P. and Cox, D. R.: An analysis of transformations, J. R. Stat. Soc. Ser. B (Methodological,
659 211–252, doi:10.2307/2287791, 1964.

660 Brown, J. D., Wu, L., He, M., Regonda, S., Lee, H. and Seo, D. J.: Verification of temperature,
661 precipitation, and streamflow forecasts from the NOAA/NWS Hydrologic Ensemble Forecast Service
662 (HEFS): 1. Experimental design and forcing verification, J. Hydrol., 519(PD), 2869–2889,
663 doi:10.1016/j.jhydrol.2014.05.028, 2014.

664 Carpenter, T. M. and Georgakakos, K. P.: Assessment of Folsom lake response to historical and potential
665 future climate scenarios: 1. Forecasting, J. Hydrol., 249(1–4), 148–175,
666 doi:https://doi.org/10.1016/S0022-1694(01)00417-6, 2001.

667 Carrillo, G., Troch, P. A., Sivapalan, M., Wagener, T., Harman, C. and Sawicz, K.: Catchment
668 classification: hydrological analysis of catchment behavior through process-based modeling along a
669 climate gradient, Hydrol. Earth Syst. Sci., 15(11), 3411–3430, doi:10.5194/hess-15-3411-2011, 2011.

670 Charles, A., Miles, E., Griesser, A., de Wit, R., Shelton, K., Cottrill, A., Spillman, C., Hendon, H.,
671 McIntosh, P., Nakaegawa, T., Atalifo, T., Prakash, B., Seuseu, S., Nihmei, S., Church, J., Jones, D. and
672 Kuleshov, Y.: Dynamical Seasonal Prediction of Climate Extremes in the Pacific, in 20th International
673 Congress on Modelling and Simulation (Modsim2013), pp. 2841–2847., 2013.

674 Clark, M. P., Kavetski, D. and Fenicia, F.: Pursuing the method of multiple working hypotheses for
675 hydrological modeling, Water Resour. Res., 47(9), n/a-n/a, doi:10.1029/2010WR009827, 2011.

676 Cloke, H., Pappenberger, F., Thielen, J. and Thiemig, V.: Operational European Flood Forecasting, in
677 Environmental Modelling, pp. 415–434, John Wiley & Sons, Ltd., 2013.

678 Cohon, J. L. and Marks, D. H.: A review and evaluation of multiobjective programing techniques, Water
679 Resour. Res., 11(2), 208–220, doi:10.1029/WR011i002p00208, 1975.

680 Crochemore, L., Ramos, M. H. and Pappenberger, F.: Bias correcting precipitation forecasts to improve
681 the skill of seasonal streamflow forecasts, Hydrol. Earth Syst. Sci., 20(9), 3601–3618, doi:10.5194/hess-

Hydrology and
Earth System
Sciences

Open Access

EGU

Discussions

682    20-3601-2016, 2016.

683    Dawid,  a P.: Present Position and Potential Developments: Some Personal Views: Statistical theory: the

684    prequential approach (with discussion), J. R. Stat. Soc. Ser. A, 147(2), 278–292, doi:10.2307/2981683,

685    1984.

686    Dechant, C. M. and Moradkhani, H.: Improving the characterization of initial condition for ensemble

687    streamflow prediction using data assimilation, Hydrol. Earth Syst. Sci., 15(11), 3399–3410,

688    doi:10.5194/hess-15-3399-2011, 2011.

689    Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D. J., Hartman, R., Herr, H.

690    D., Fresch, M., Schaake, J. and Zhu, Y.: The science of NOAA's operational hydrologic ensemble

691    forecast service, Bull. Am. Meteorol. Soc., 95(1), 79–98, doi:10.1175/BAMS-D-12-00081.1, 2014.

692    Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint versus

693    postprocessor approaches for hydrological uncertainty estimation accounting for error autocorrelation

694    and heteroscedasticity, Water Resour. Res., 50(3), 2350–2375, doi:10.1002/2013WR014185, 2014.

695    Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P. and Rieckermann, J.: Improving

696    uncertainty estimation in urban hydrological modeling by statistically describing bias, Hydrol. Earth

697    Syst. Sci., 17(10), 4209–4225, doi:10.5194/hess-17-4209-2013, 2013.

698    Gneiting, T., Raftery, A. E., Westveld, A. H. and Goldman, T.: Calibrated Probabilistic Forecasting

699    Using Ensemble Model Output Statistics and Minimum CRPS Estimation, Mon. Weather Rev., 133(5),

700    1098–1118, doi:10.1175/MWR2904.1, 2005.

701    Gneiting, T., Balabdaoui, F. and Raftery, A. E.: Probabilistic forecasts, calibration and sharpness, J. R.

702    Stat. Soc. Ser. B Stat. Methodol., 69(2), 243–268, doi:10.1111/j.1467-9868.2007.00587.x, 2007.

703    Hashino, T., Bradley,  a. a. and Schwartz, S. S.: Evaluation of bias-correction methods for ensemble

704    streamflow volume forecasts, Hydrol. Earth Syst. Sci., 11, 939–950, doi:10.5194/hess-11-939-2007,

705    2007.

706    Hersbach, H.: Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction

707    Systems,        Weather        Forecast.,        15(5),        559–570,        doi:10.1175/1520-

708    0434(2000)015<0559:DOTCRP>2.0.CO;2, 2000.

709    Hudson, D., Marshall, A. G., Yin, Y., Alves, O. and Hendon, H. H.: Improving Intraseasonal Prediction

710    with a New Ensemble Generation Strategy, Mon. Weather Rev., 141(12), 4429–4449,

711    doi:10.1175/MWR-D-13-00059.1, 2013.

712    Humphrey, G. B., Gibbs, M. S., Dandy, G. C. and Maier, H. R.: A hybrid approach to monthly

713    streamflow forecasting: Integrating hydrological model outputs into a Bayesian artificial neural network,

714    J. Hydrol., 540, 623–640, doi:10.1016/j.jhydrol.2016.06.026, 2016.

715    Jeffrey, S. J., Carter, J. O., Moodie, K. B. and Beswick, A. R.: Using spatial interpolation to construct a

716    comprehensive archive of Australian climate data, Environ. Model. Softw., 16(4), 309–330,

717    doi:10.1016/S1364-8152(01)00008-1, 2001.

718    Kavetski, D., Kuczera, G. and Franks, S. W.: Bayesian analysis of input uncertainty in hydrological

719    modeling: 1. Theory, Water Resour. Res., 42(3), n/a-n/a, doi:10.1029/2005WR004368, 2006.

720    Knoche, M., Fischer, C., Pohl, E., Krause, P. and Merz, R.: Combined uncertainty of hydrological model

721    complexity and satellite-based forcing data evaluated in two data-scarce semi-arid catchments in

722    Ethiopia, J. Hydrol., 519, 2049–2066, doi:https://doi.org/10.1016/j.jhydrol.2014.10.003, 2014.

723    Kuczera, G., Kavetski, D., Franks, S. and Thyer, M.: Towards a Bayesian total error analysis of

724    conceptual rainfall-runoff models: Characterising model error using storm-dependent parameters, J.

725    Hydrol., 331(1–2), 161–177, doi:10.1016/j.jhydrol.2006.05.010, 2006.

726    Laio, F. and Tamea, S.: Verification tools for probabilistic forecasts of continuous hydrological

727    variables, Hydrol. Earth Syst. Sci., 11(4), 1267–1277, doi:10.5194/hess-11-1267-2007, 2007.

728    Laugesen, R., Tuteja, N. K., Shin, D., Chia, T. and Khan, U.: Seasonal Streamflow Forecasting with a

729    workflow-based dynamic hydrologic modelling approach, in MODSIM 2011 - 19th International

730    Congress on Modelling and Simulation - Sustaining Our Future: Understanding and Living with

731    Uncertainty, pp. 2352–2358. [online] Available from: http://www.scopus.com/inward/record.url?eid=2-

732    s2.0-84858823270&partnerID=tZOtx3y1, 2011.

733    Lerat, J., Pickett-Heaps, C., Shin, D., Zhou, S., Feikema, P., Khan, U., Laugesen, R., Tuteja, N., Kuczera,

734    G., Thyer, M. and Kavetski, D.: Dynamic streamflow forecasts within an uncertainty framework for 100

735    catchments in Australia, in In: 36th Hydrology and Water Resources Symposium: The art and science

736    of water, pp. 1396–1403, Barton, ACT: Engineers Australia., 2015.

737    Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: Error reduction and representation in stages

738    (ERRIS) in hydrological modelling for ensemble streamflow forecasting, Hydrol. Earth Syst. Sci., 20(9),

739    3561–3579, doi:10.5194/hess-20-3561-2016, 2016.

740    Lü, H., Crow, W. T., Zhu, Y., Ouyang, F. and Su, J.: Improving streamflow prediction using remotely-

741    sensed soil moisture and snow depth, Remote Sens., 8(6), doi:10.3390/rs8060503, 2016.

742    Madadgar, S., Moradkhani, H. and Garen, D.: Towards improved post-processing of hydrologic forecast

743    ensembles, Hydrol. Process., 28(1), 104–122, doi:10.1002/hyp.9562, 2014.

744    McInerney, D., Thyer, M., Kavetski, D., Lerat, J. and Kuczera, G.: Improving probabilistic prediction

745    of daily streamflow by identifying Pareto optimal approaches for modeling heteroscedastic residual

746    errors, Water Resour. Res., 53(3), 2199–2239, doi:10.1002/2016WR019168, 2017.

747    Mendoza, P. A., Wood, A. W., Clark, E., Rothwell, E., Clark, M. P., Nijssen, B., Brekke, L. D. and

748    Arnold, J. R.: An intercomparison of approaches for improving predictability in operational seasonal

749    streamflow forecasting, Hydrol. Earth Syst. Sci. Discuss., 2017, 1–37, doi:10.5194/hess-2017-60, 2017.

750    Middleton, N., Programme, U. N. E. and Thomas, D. S. G.: World Atlas of Desertification, Arnold.

751    [online] Available from: https://books.google.com.au/books?id=aNqtQgAACAAJ, 1997.

752    Murphy, A. H. and Ehrendorfer, M.: On the relationship between the accuracy and value of forecasts in

753    the cost–loss ratio situation, Weather Forecast., 2(3), 243–251, doi:10.1175/1520-

754    0434(1987)002<0243:OTRBTA>2.0.CO;2, 1987.

755    Perrin, C., Michel, C. and Andréassian, V.: Improvement of a parsimonious model for streamflow

756    simulation, J. Hydrol., 279(1–4), 275–289, doi:10.1016/S0022-1694(03)00225-7, 2003.

757    Pokhrel, P., Robertson, D. E. and Wang, Q. J.: A Bayesian joint probability post-processor for reducing

758    errors and quantifying uncertainty in monthly streamflow predictions, Hydrol. Earth Syst. Sci., 17(2),

759    795–804, doi:10.5194/hess-17-795-2013, 2013.

760    Prudhomme, C., Hannaford, J., Harrigan, S., Boorman, D., Knight, J., Bell, V., Jackson, C., Svensson,

761    C., Parry, S., Bachiller-Jareno, N., Davies, H., Davis, R., Mackay, J., McKenzie, A., Rudd, A., Smith,

762    K., Bloomfield, J., Ward, R. and Jenkins, A.: Hydrological Outlook UK: an operational streamflow and

763    groundwater level forecasting system at monthly to seasonal time scales, Hydrol. Sci. J., 62(16), 2753–

764    2768, doi:10.1080/02626667.2017.1395032, 2017.

765    Renard, B., Kavetski, D., Leblois, E., Thyer, M., Kuczera, G. and Franks, S. W.: Toward a reliable

766    decomposition of predictive uncertainty in hydrological modeling: Characterizing rainfall errors using

767    conditional simulation, Water Resour. Res., 47(11), n/a-n/a, doi:10.1029/2011WR010643, 2011.

768    Robertson, D. E. and Wang, Q. J.: Selecting predictors for seasonal streamflow predictions using a

769    Bayesian joint probability ( BJP ) modelling approach, 18th World IMACS/MODSIM Congr. Cairns,

770    Aust. 13-17 July 2009, (July), 376–382, 2009.

771    Robertson, D. E. and Wang, Q. J.: A Bayesian Approach to Predictor Selection for Seasonal Streamflow

772    Forecasting, J. Hydrometeorol., 13(1), 155–171, doi:10.1175/JHM-D-10-05009.1, 2011.

773    Robertson, D. E., Pokhrel, P. and Wang, Q. J.: Improving statistical forecasts of seasonal streamflows

774    using hydrological model output, Hydrol. Earth Syst. Sci., 17(2), 579–593, doi:10.5194/hess-17-579-

775    2013, 2013a.

776    Robertson, D. E., Shrestha, D. L. and Wang, Q. J.: Post-processing rainfall forecasts from numerical

777    weather prediction models for short-term streamflow forecasting, Hydrol. Earth Syst. Sci., 17(9), 3587–

778    3603, doi:10.5194/hess-17-3587-2013, 2013b.

779    Sawicz, K. A., Kelleher, C., Wagener, T., Troch, P., Sivapalan, M. and Carrillo, G.: Characterizing

780    hydrologic change through catchment classification, Hydrol. Earth Syst. Sci., 18(1), 273–285,

781    doi:10.5194/hess-18-273-2014, 2014.

782    Senlin, Z., Feikema, P., Shin, D., Tuteja, N. K., MacDonald, A., Sunter, P., Kent, D., Le, B., Pipunic,

783    R., Wilson, T., Pickett-Heaps, C. and Lerat, J.: Operational efficiency measures of the national seasonal

784    streamflow forecast service in Australia, edited by G. Syme, D. H. MacDonald, B. Fulton, and J.

785    Piantadosi, the Modelling and Simulation Society of Australia and New Zealand Inc, Hobart, Australia.,

786    2017.

787    Seo, D.-J., Herr, H. D. and Schaake, J. C.: A statistical post-processor for accounting of hydrologic

788    uncertainty in short-range ensemble streamflow prediction, Hydrol. Earth Syst. Sci. Discuss., 3(4),

789    1987–2035, doi:10.5194/hessd-3-1987-2006, 2006.

790    Shapiro, S. S. and Wilk, M. B.: An Analysis of Variance Test for Normailty (Complete Samples),

791    Biometrika, 52(3–4), 591–611, doi:10.2307/1267427, 1965.

792    Smith, T., Marshall, L. and Sharma, A.: Modeling residual hydrologic errors with Bayesian inference,

793    J. Hydrol., 528(SEPTEMBER 2015), 29–37, doi:10.1016/j.jhydrol.2015.05.051, 2015.

794    Tang, Q. and Lettenmaier, D. P.: Use of satellite snow-cover data for streamflow prediction in the

795    Feather    River    Basin,    California,    Int.    J.    Remote    Sens.,    31(14),    3745–3762,

796    doi:10.1080/01431161.2010.483493, 2010.

797    Taschetto, A. S. and England, M. H.: An analysis of late twentieth century trends in Australian rainfall,

798    Int. J. Climatol., 29(6), 791–807, doi:10.1002/joc.1736, 2009.

799    Timbal, B. and McAvaney, B. J.: An Analogue based method to downscale surface air temperature:

800    Application for Australia, Clim. Dyn., 17, 947–963, doi:10.1007/s003820100156, 2001.

801    Turner, S. W. D., Bennett, J., Robertson, D. and Galelli, S.: Value of seasonal streamflow forecasts in

802    emergency response reservoir management, Hydrol. Earth Syst. Sci. Discuss., 2017, 1–26,

803    doi:10.5194/hess-2016-691, 2017.

804    Tuteja, N. K., Shin, D., Laugesen, R., Khan, U., Shao, Q., Wang, E., Li, M., Zheng, H., Kuczera, G.,

805    Kavetski, D., Evin, G., Thyer, M., MacDonald, A., Chia, T. and Le, B.: Experimental evaluation of the

806    dynamic seasonal streamflow forecasting approach, Melbourne., 2011.

807    Tuteja, N. K., Zhou, S., Lerat, J., Wang, Q. J., Shin, D. and Robertson, D. E.: Overview of

808    Communication Strategies for Uncertainty in Hydrological Forecasting in Australia, in Handbook of

809    Hydrometeorological Ensemble Forecasting, edited by Q. Duan, F. Pappenberger, J. Thielen, A. Wood,

810    H. L. Cloke, and J. C. Schaake, pp. 1–19, Springer Berlin Heidelberg, Berlin, Heidelberg., 2016.

811    Tyralla, C. and Schumann, A. H.: Incorporating structural uncertainty of hydrological models in

812    likelihood functions via an ensemble range approach, Hydrol. Sci. J., 02626667.2016.1164314,

813    doi:10.1080/02626667.2016.1164314, 2016.

814    Wandishin, M. S. and Brooks, H. E.: On the relationship between Clayton's skill score and expected

815    value for forecasts of binary events, Meteorol. Appl., 9(4), 455–459, doi:10.1017/S1350482702004085,

816    2002.

817    Wang, Q. J. and Robertson, D. E.: Multisite probabilistic forecasting of seasonal flows for streams with

818    zero value occurrences, Water Resour. Res., 47(2), doi:10.1029/2010WR009333, 2011.

819    Wang, Q. J., Robertson, D. E. and Chiew, F. H. S.: A Bayesian joint probability modeling approach for

820    seasonal     forecasting     of     streamflows     at     multiple     sites,     Water     Resour.     Res.,     45(5),

821    doi:10.1029/2008WR007355, 2009.

822    Wang, Q. J., Shrestha, D. L., Robertson, D. E. and Pokhrel, P.: A log-sinh transformation for data

823    normalization and variance stabilization, Water Resour. Res., 48(5), doi:10.1029/2011WR010973,

824    2012.

825    Wilks, D. S.: Statistical methods in the atmospheric sciences., 2011.

826    Wood, A. W. and Schaake, J. C.: Correcting Errors in Streamflow Forecast Ensemble Mean and Spread,

827    J. Hydrometeorol., 9(1), 132–148, doi:10.1175/2007JHM862.1, 2008.

828    Ye, W., Bates, B. C., Viney, N. R., Sivapalan, M. and Jakeman, A. J.: Performance of conceptual

829    rainfall-runoff models in low-yielding ephemeral catchments, Water Resour. Res., 33(1), 153–166,

830    doi:10.1029/96WR02840, 1997.

831    Zhang, Q., Xu, C.-Y. and Zhang, Z.: Observed changes of drought/wetness episodes in the Pearl River

832    basin, China, using the standardized precipitation index and aridity index, Theor. Appl. Climatol., 98(1),

833    89–99, doi:10.1007/s00704-008-0095-4, 2009.

834    Zhao, T., Schepen, A. and Wang, Q. J.: Ensemble forecasting of sub-seasonal to seasonal streamflow by

835    a     Bayesian     joint     probability     modelling     approach,     J.     Hydrol.,     541,     839–849,

836    doi:https://doi.org/10.1016/j.jhydrol.2016.07.040, 2016.

837

838

839

840

841

842

843

844

845

846

847

848

Table 1. Percentage of catchments with high and low summary skill for the different residual error schemes for both monthly and seasonal forecasts. High (low) summary skill is defined as the percentage of catchments with 10-12 months (0-2 months) reliable forecasts that are sharper than climatology.

| Residual Error Scheme | Uncorrected forecasts | Log | Log-Sinh | BC0.2 |
|---|---|---|---|---|
| *Monthly Forecasts* | | | | |
| High Summary Skill | 40% | 33% | 33% | 84% |
| Low Summary Skill | 16% | 2% | 7% | <1% |
| *Seasonal Forecasts* | | | | |
| High Summary Skill | 46% | 9% | 20% | 54% |
| Low Summary Skill | 14% | 19% | 17% | 2% |

852

853

854

855

856    **Figures**



Figure 1: Schematic of the dynamic streamflow forecasting system used in this study. A similar approach is used by the Australian Bureau of Meteorology for its monthly and seasonal streamflow forecasting service.

857

858

859



860

Figure 2: Locations of the 300 catchments used in this study. The catchments are classified as dry or wet based on the aridity index. The Koppen climate classification for Australia are shown. The Dieckmans Bridge catchment (site id: 145010A), used as a representative site in Figure 8, is indicated by the red circle.

865

866

867

868

869

870

871

872

873

874

875

876

877



Figure 3: Schematic of the cross-validation framework used for forecast verification as an example for model validation year 1990 (after Tuteja et al., 2016).

881

882

883

884

885

886

887

888

889

890

891

Figure 4: Performance of monthly forecasts in terms of CRPSS, reliability (PIT p-value) and sharpness (IQR99 ratio).
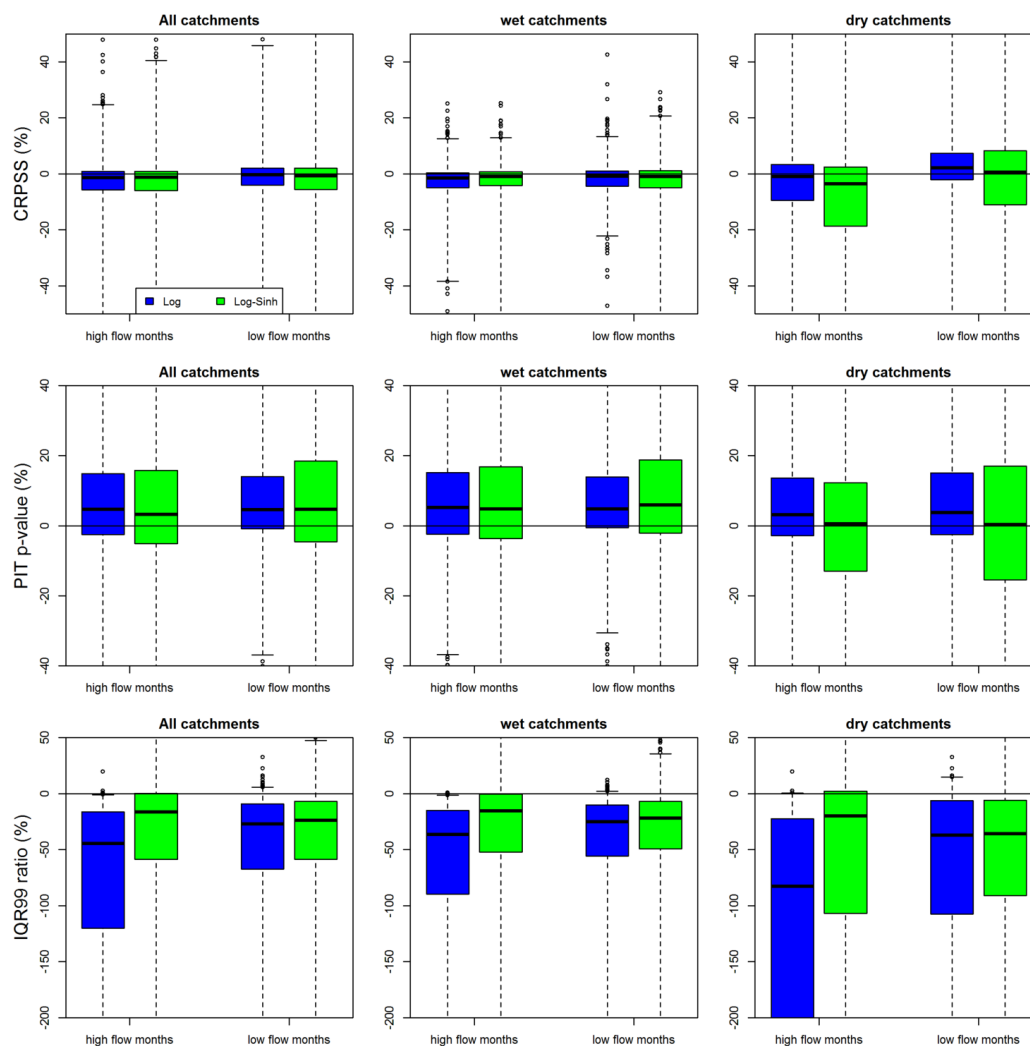
897



898 Figure 5: Performance of seasonal forecasts in terms of CRPSS, reliability (PIT p-value) and sharpness
899 (IQR99 ratio).

900
901
902

903
904

905   Figure 6: Distributions of differences in the monthly forecast performance metrics of the Log and Log-
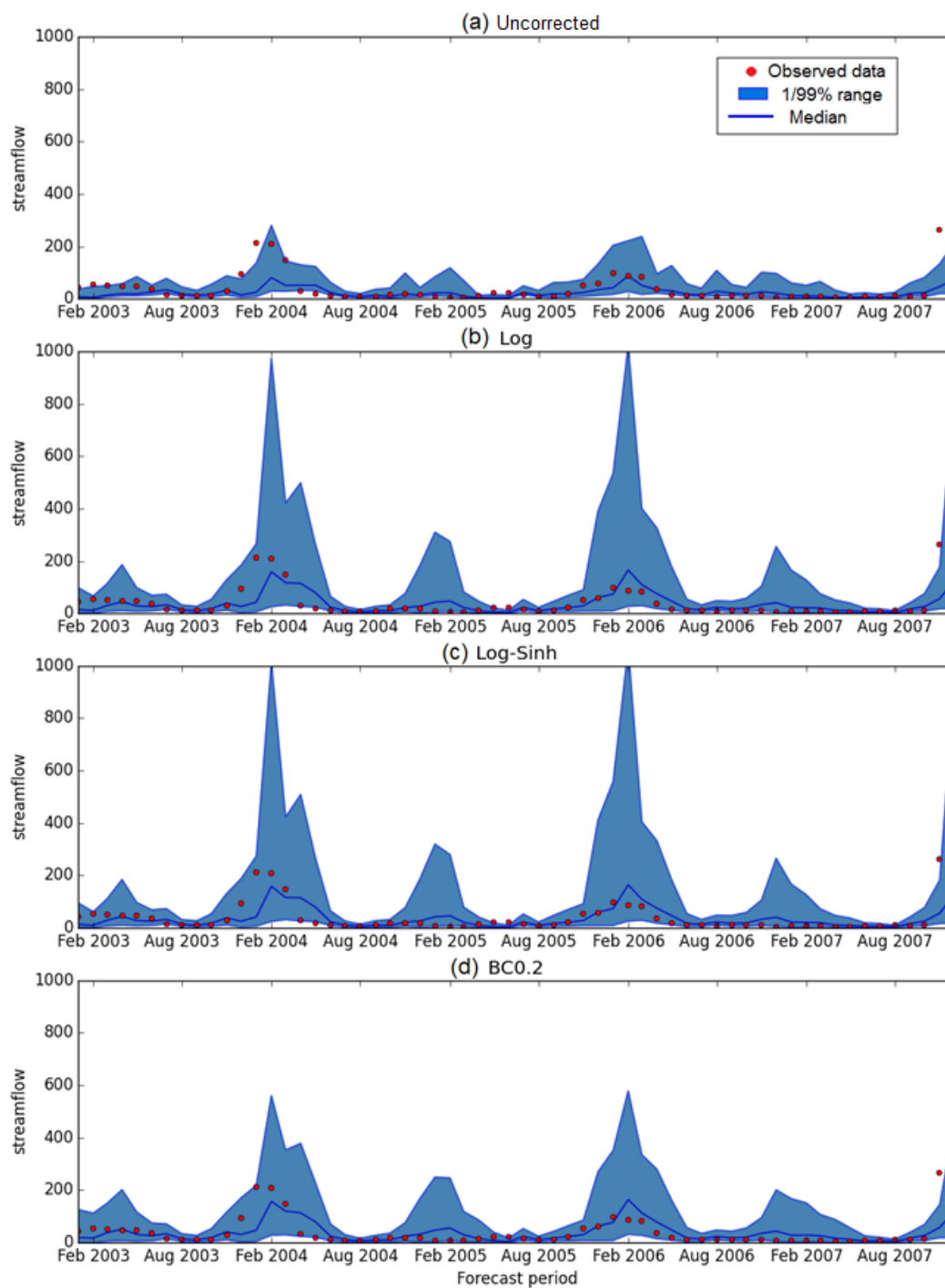906   Sinh schemes compared to the BC0.2 scheme.

907

908
909
910

911

912



Figure 7: Distributions of differences in the seasonal forecast performance metrics of the Log and Log-Sinh schemes compared to the BC0.2 scheme.

913
914

915

916

36

917
918
919   Figure 8: Seasonal streamflow forecast time series (blue line) and observations (red dots) at Dieckmans
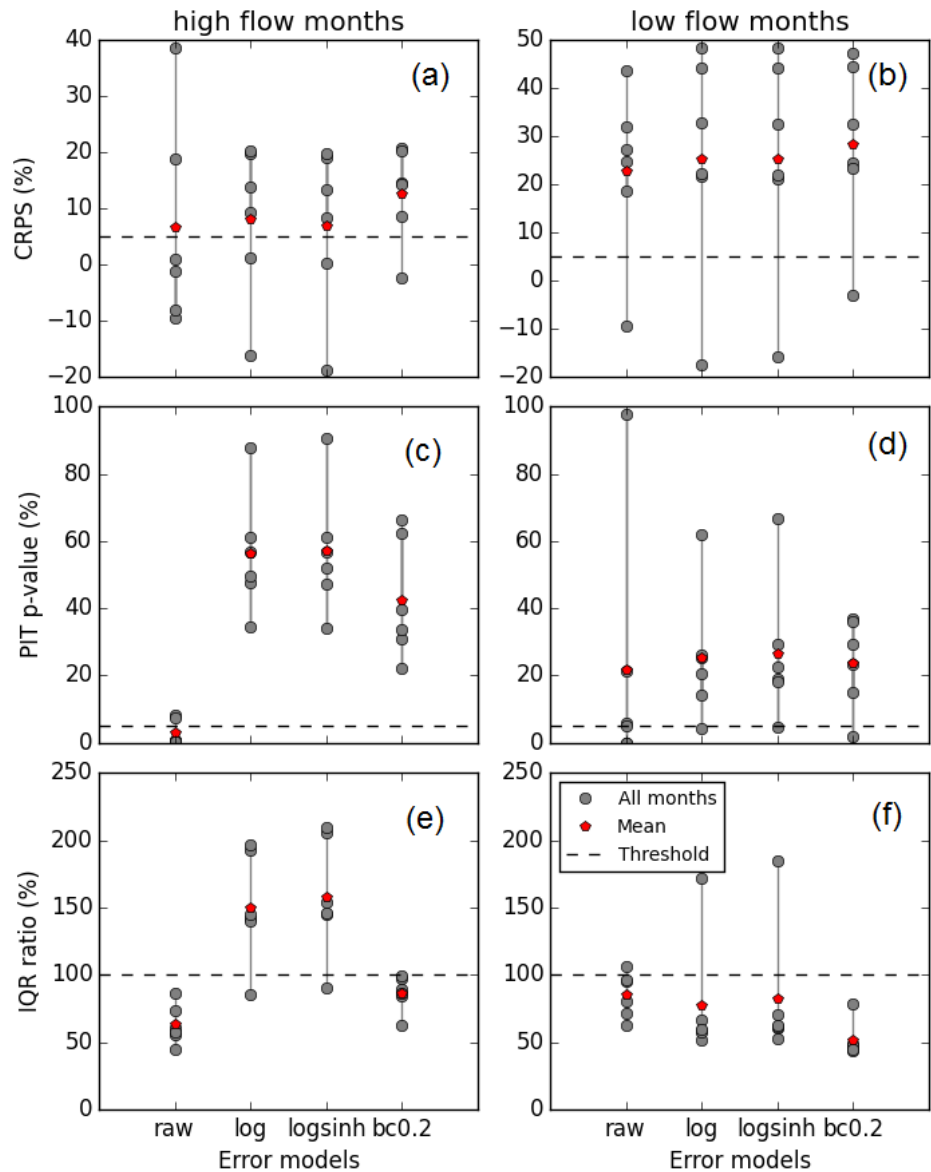920   Bridge catchment (site id: 145010A). The shaded area shows the 99% prediction limits.
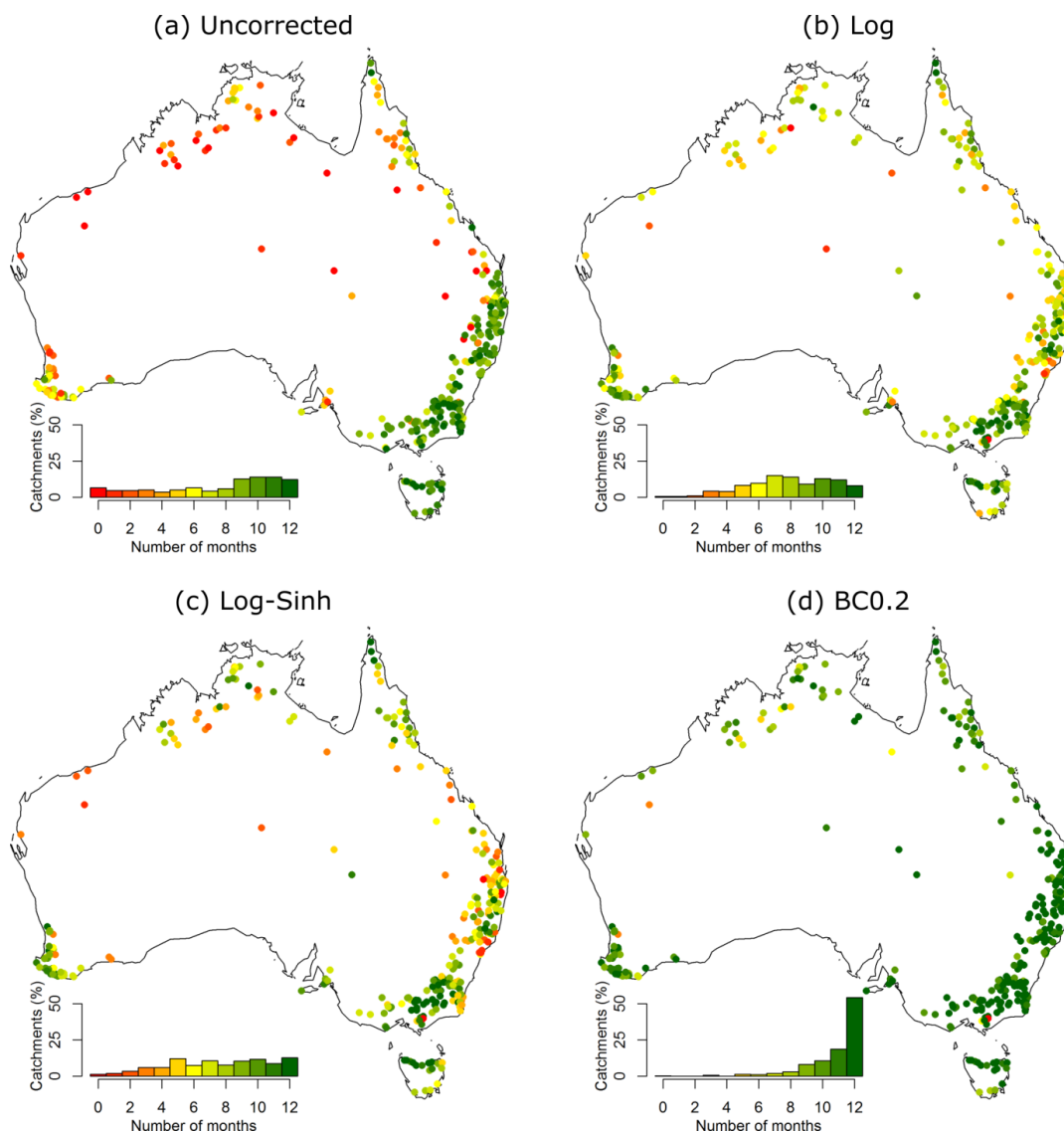
921

922  Figure 9: Seasonal streamflow forecast skill-score at the Dieckmans Bridge catchment corresponding to
923  the time series shown in Figure 8 for six high flow months and six low flow months. Note that skill-
924  score values of 5%, 5% and 100% are indicated for CRPSS, p-value and IQR ratio respectively, using
925  dashed lines.

926
927
928
929
930
931

933



Figure 10: Summary skill of monthly forecasts obtained using the Log, Log-Sinh and BC0.2 schemes across 300 Australian catchments. The performance of uncorrected forecasts is also shown. The summary skill is defined as the number of months where the forecasts are reliable and sharper than climatology. The inset histogram shows the percentage of catchments in each performance category and also serves as the color legend.
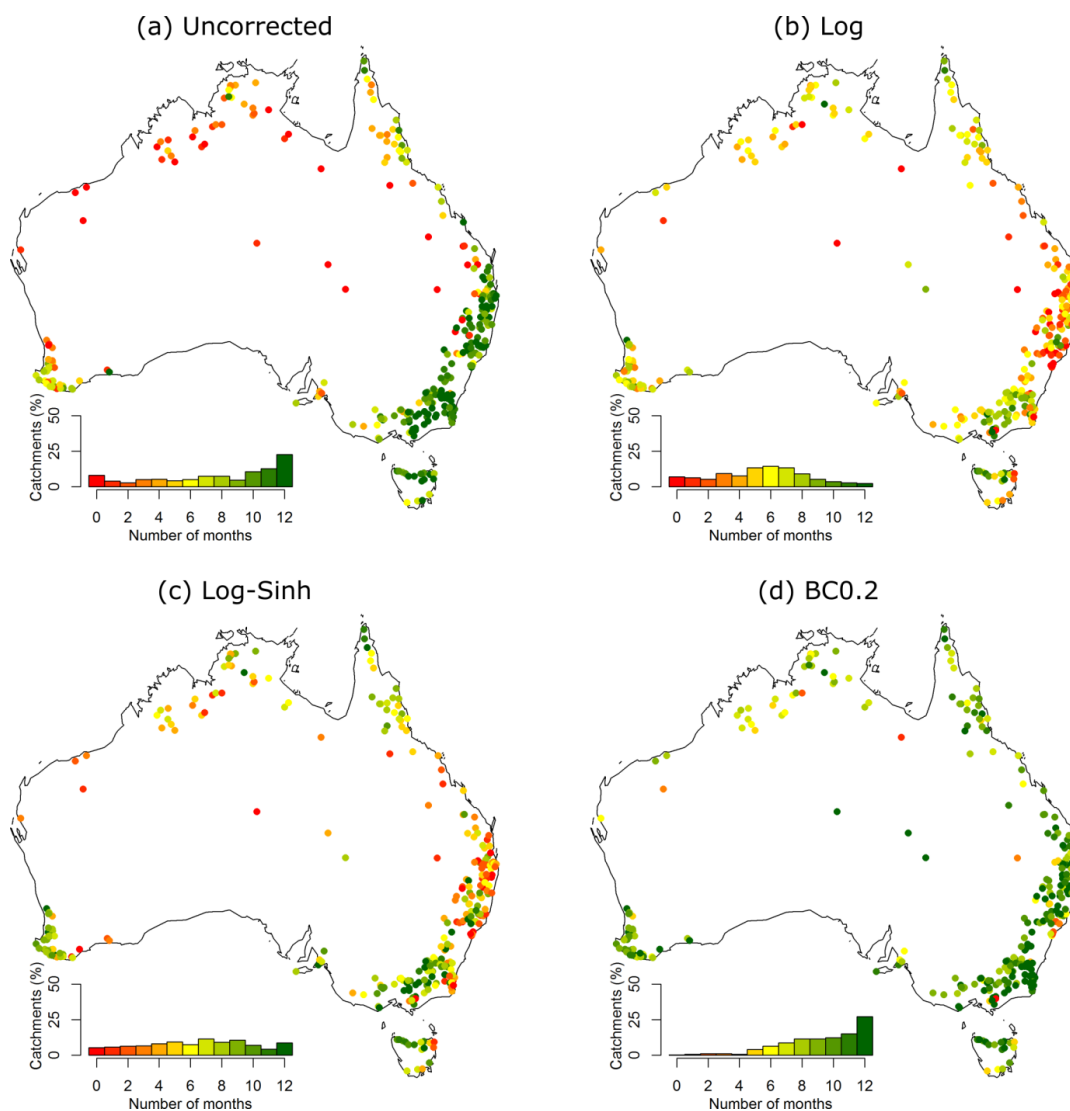
940

Figure 11: Summary skill of seasonal forecasts obtained using the Log, Log-Sinh and BC0.2 schemes across 300 Australian catchments. See Figure 10 for details.