

1 Evaluating post-processing approaches for monthly and seasonal 2 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 Correspondence email: fitsum.woldemeskel@bom.gov.au

10

11 **Abstract**

12 Streamflow forecasting is prone to substantial uncertainty due to errors in meteorological forecasts,
13 hydrological model structure and parameterization, as well as in the observed rainfall and streamflow
14 data used to calibrate the models. Statistical streamflow post-processing is an important technique
15 available to improve the probabilistic properties of the forecasts. This study evaluates post-processing
16 approaches based on three transformations – logarithmic (Log), log-sinh (Log-Sinh) and Box-Cox with
17 $\lambda = 0.2$ (BC0.2) – and identifies the best performing scheme for post-processing monthly and seasonal
18 (3-months-ahead) streamflow forecasts, such as those produced by the Australian Bureau of
19 Meteorology. Using the Bureau’s operational dynamic streamflow forecasting system, we carry out
20 comprehensive analysis of the three post-processing schemes across 300 Australian catchments with a
21 wide range of hydro-climatic conditions. Forecast verification is assessed using reliability and sharpness
22 metrics, as well as the Continuous Ranked Probability Skill Score (CRPSS). Results show that the
23 uncorrected forecasts (i.e. without post-processing) are unreliable at half of the catchments. Post-
24 processing of forecasts substantially improves reliability, with more than 90% of forecasts classified as
25 reliable. In terms of sharpness, the BC0.2 scheme substantially outperforms the Log and Log-Sinh
26 schemes. Overall, the BC0.2 scheme achieves reliable and sharper-than-climatology forecasts at a larger
27 number of catchments than the Log and Log-Sinh schemes. The improvements in forecast reliability and
28 sharpness achieved using the BC0.2 post-processing scheme will help water managers and users of the
29 forecasting service to make better-informed decisions in planning and management of water resources.

30 **Keywords:** seasonal streamflow forecasts, post-processing, Box-Cox transformation

31 **Key points**

- 32 1. Uncorrected and post-processed streamflow forecasts (using three transformations, namely Log,
33 Log-Sinh and BC0.2) are evaluated over 300 diverse Australian catchments
- 34 2. Post-processing enhances streamflow forecast reliability, increasing the percentage of catchments
35 with reliable predictions from 50% to over 90%
- 36 3. The BC0.2 transformation achieves substantially better forecast sharpness than the Log-Sinh and
37 Log transformations, particularly in dry catchments

38

39 **1 Introduction**

40 Hydrological forecasts provide crucial supporting information on a range of water resource management
41 decisions, including (depending on the forecast lead-time) flood emergency response, water allocation
42 for various uses, and drought risk management (Li et al., 2016; Turner et al., 2017). The forecasts,
43 however, should be thoroughly verified and proved to be of sufficient quality to support decision-making
44 and to meaningfully benefit the economy, environment and society.

45 Sub-seasonal and seasonal streamflow forecasting systems can be broadly classified as dynamic or
46 statistical (Crochemore et al., 2016). In *dynamic* modelling systems, a hydrological model is usually
47 developed at a daily time-step and calibrated against observed streamflow using historical rainfall and
48 potential evaporation data. Rainfall forecasts from a numerical climate model are then used as an input
49 to produce daily streamflow forecasts, which are then aggregated to the time scale of interest and post-
50 processed using statistical models (e.g. Bennett et al., 2017; Schick et al., 2018). In *statistical* modelling
51 systems, a statistical model based on relevant predictors, such as antecedent rainfall and streamflow, is
52 developed and applied directly at the time scale of interest (Robertson and Wang, 2009, 2011; Lü et al.,
53 2016; Zhao et al., 2016). Hybrid systems that combine aspects of dynamic and statistical approaches
54 have also been investigated (Humphrey et al., 2016; Robertson et al., 2013a)

55 Examples of operational services based on the dynamic approach include the Australian Bureau of
56 Meteorology's dynamic modelling system (Laugesen et al., 2011; Tuteja et al., 2011; Lerat et al., 2015);
57 the Hydrological Ensemble Forecast Service (HEFS) of the US National Weather Service (NWS)
58 (Brown et al., 2014; Demargne et al., 2014); the Hydrological Outlook UK (HOUK) (Prudhomme et al.,
59 2017); and the short-term forecasting European Flood Alert System (EFAS) (Cloke et al., 2013).
60 Examples of operational services based on a statistical approach include the Bureau of Meteorology's
61 Bayesian Joint Probability (BJP) forecasting system (Senlin et al., 2017).

62 Dynamic and statistical approaches have distinct advantages and limitations. Dynamic systems can
63 potentially provide more realistic responses in unfamiliar climate situations, as it is possible to impose
64 physical constraints in such situations (Wood and Schaake, 2008). In comparison, statistical models have
65 the flexibility to include features that may lead to more reliable predictions. For example, the BJP model
66 uses climate indices (e.g. NINO3.4), which are typically not used in dynamic approaches. That said, the
67 suitability of statistical models for the analysis of non-stationary catchment and climate conditions is
68 questionable (Wood and Schaake, 2008).

69 Streamflow forecasts obtained using hydrological models are affected by uncertainties in rainfall
70 forecasts, observed rainfall and streamflow data, as well as by uncertainties in the model structure and
71 parameters. Progress has been made towards reducing biases and characterizing the sources of

72 uncertainty in streamflow forecasts. These advances include improving rainfall forecasts through post-
73 processing (Robertson et al., 2013b; Crochemore et al., 2016), accounting for input, parametric and/or
74 structural uncertainty (Kavetski et al., 2006; Kuczera et al., 2006; Renard et al., 2011; Tyralla and
75 Schumann, 2016), and using data assimilation techniques (Dechant and Moradkhani, 2011). Although
76 these steps may improve some aspects of the forecasting system, a predictive bias may nonetheless
77 remain. Such bias can only be reduced via post-processing, which, if successful, will improve forecast
78 accuracy and reliability (Madadgar et al., 2014; Lerat et al., 2015).

79 This study focuses on improving streamflow forecasting at monthly and seasonal time-scales using
80 dynamic approaches, more specifically, by evaluating several forecast post-processing approaches. Post-
81 processing of streamflow forecasts is intended to remove systemic biases in the mean, variability and
82 persistence of uncorrected forecasts, which arise due to inaccuracies in the downscaled rainfall forecasts
83 (e.g. errors in downscaling forecast rainfall from a grid with ≈ 250 km resolution to the catchment scale)
84 and in the hydrological model (e.g. due to the effects of data errors on the model calibration and due to
85 structural errors in the model itself).

86 A number of post-processing approaches have been investigated in the literature, including quantile
87 mapping (Hashino et al., 2007) and Bayesian frameworks (Pokhrel et al., 2013; Robertson et al., 2013a),
88 as well as methods based on state-space models and wavelet transformations (Bogner and Kalas, 2008).
89 Wood and Schaake (2008) used the correlation between forecast ensemble means and observations to
90 generate a conditional forecast. Compared with the traditional approach of correcting individual forecast
91 ensembles, the correlation approach improved forecast skill and reliability. In another study, Pokhrel et
92 al. (2013) implemented a Bayesian Joint Probability (BJP) method to correct biases, update predictions
93 and quantify uncertainty in monthly hydrological model predictions in 18 Australian catchments. The
94 study found that the accuracy and reliability of forecasts improved. More recently, Mendoza et al. (2017)
95 evaluated a number of seasonal streamflow forecasting approaches, including purely statistical, purely
96 dynamical, and hybrid approaches. Based on analysis of catchments contributing to five reservoirs, the
97 study concluded that incorporating catchment and climate information into post-processing improves
98 forecast skill. While the above review mainly focused on post-processing at sub-seasonal and seasonal
99 forecasts (as it is the main focus of the current study), post-processing is also commonly applied to short-
100 range forecasts (e.g. Li et al., 2016) and to long-range forecasts up to 12 months ahead (Bennett et al.,
101 2016).

102 In most streamflow post-processing approaches, a residual error model is applied to quantify forecast
103 uncertainty. Most residual error models are based on least squares techniques with weights and/or data
104 transformations (e.g. Carpenter and Georgakakos, 2001; Li et al., 2016). In order to produce post-
105 processed streamflow forecasts, a daily-scale residual error model is used in the calibration of

106 hydrological model parameters, and a monthly/seasonal-scale residual error model is used as part of
107 streamflow post-processing to quantify the forecast uncertainty. In a recent study, McInerney et al.
108 (2017) concluded that residual error models based on Box-Cox transformations with fixed parameter
109 values are particularly effective for daily scale streamflow predictions using observed rainfall, yielding
110 substantial improvements in dry catchments. This study investigates whether these findings generalize
111 to monthly and seasonal forecasts using forecast rainfall.

112 An important aspect of this work is its focus on general findings applicable over diverse hydro-
113 climatological conditions. Most of the studies in the published literature use a limited number of
114 catchments and case studies to test prospective methods. Dry catchments, characterised by intermittent
115 flows and frequent low flows, pose the greatest challenge to hydrological models (Ye et al., 1997;
116 Knoche et al., 2014). Yet the provision of good quality forecasts across a large number of catchments is
117 an essential attribute of national scale operational forecasting services, especially in large countries with
118 diverse climatic and catchment conditions, such as Australia.

119 This paper develops streamflow post-processing approaches suitable for use in an operational
120 streamflow forecasting service. We pose the following aims:

121 **Aim 1:** Evaluate the value of streamflow forecast post-processing by comparing forecasts with no post-
122 processing (hereafter called ‘uncorrected’ forecasts) against post-processed forecasts;

123 **Aim 2:** Evaluate three post-processing schemes based on residual error models with data transformations
124 recommended in recent publications, namely the Log, Box-Cox (McInerney et al., 2017) and Log-Sinh
125 (Wang et al., 2012) schemes, for monthly and seasonal streamflow post-processing;

126 **Aim 3:** Evaluate the generality of results over a diverse range of hydro-climatic conditions, in order to
127 ensure the recommendations are robust in the context of an operational streamflow forecasting service.

128 To achieve these aims, we use the operational monthly and seasonal (3-months-ahead) dynamic
129 streamflow forecasting system of the Australian Bureau of Meteorology (Lerat et al., 2015). We evaluate
130 the post-processing approaches across 300 catchments across Australia, with detailed analysis of dry and
131 wet catchments. Forecast verification is carried out using Continuous Ranked Probability Skill Score
132 (CRPSS) as well as metrics measuring reliability and sharpness, which are important aspects of a
133 probabilistic forecast (Wilks, 2011). These metrics are used by the Bureau of Meteorology to describe
134 streamflow forecast performance of the operational service.

135 The rest of the paper is organised as follows. The forecasting methodology is described in Section 2 and
136 application studies are described in Section 3. Results are presented in Section 4, followed by discussions
137 and conclusions in Sections 5 and 6 respectively.

138 **2 Seasonal streamflow forecasting methodology**

139 **2.1 Overview**

140 The streamflow forecasting system adopted in this study is based on the Bureau of Meteorology's
141 dynamic modelling system (Figure 1). Daily rainfall forecasts are input into a daily rainfall-runoff model
142 to produce "uncorrected" daily streamflow forecasts. These streamflow forecasts are then aggregated in
143 time and post-processed to produce monthly and seasonal streamflow forecasts, which are issued each
144 month. Two steps are involved: calibration and forecasting, discussed below.

145 **2.2 Uncorrected streamflow forecasts procedure**

146 **2.2.1 Rainfall-runoff model**

147 The rainfall-runoff model GR4J (Perrin et al., 2003) is used as it has been proven to provide (on average)
148 good performance across a large number of catchments ranging from semi-arid to temperate and tropical
149 humid (Perrin et al., 2003; Tuteja et al., 2011). GR4J is a lumped conceptual model with four calibration
150 parameters: maximum capacity of the production store x_1 (mm); ground water exchange coefficient x_2
151 (mm); one day ahead maximum capacity of the routing store x_3 (mm); and time base of unit hydrograph
152 x_4 (days).

153 **2.2.2 Rainfall-runoff model calibration**

154 In the calibration step, the daily rainfall-runoff model is calibrated to observed daily streamflow using
155 observed rainfall (Jeffrey et al., 2001) as forcing. The calibration of the parameters is based on the
156 weighted least squares likelihood function, similar to that outlined in Evin et al. (2014). Markov Chain
157 Monte Carlo (MCMC) analysis is used to estimate posterior parametric uncertainty (Tuteja et al., 2011).
158 Following MCMC analysis, 40 random sets of GR4J parameters are retained and used in the forecast
159 step. A cross-validation procedure is implemented to verify the forecasts, as described in Section 3.4.
160 The calibration and cross-validation is computationally intensive; therefore, we use the High
161 Performance Computing (HPC) facility at the National Computing Infrastructure (NCI) in Australia.

162 **2.2.3 Producing uncorrected streamflow forecasts**

163 Prior to the forecast period, observed rainfall is used to force the rainfall-runoff model. During the
164 forecast period, 166 replicates of daily downscaled rainfall forecasts from the Bureau of Meteorology's
165 global climate model, namely the Predictive Ocean Atmosphere Model for Australia, POAMA-2 are
166 used (see Section 3.2 for details on POAMA-2). These rainfall forecasts are input into GR4J and
167 propagated using the 40 GR4J parameter sets to obtain 6640 (166×40) daily streamflow forecasts. The
168 daily streamflow forecasts generated using GR4J are then aggregated to monthly and seasonal time
169 scales to produce ensembles of 6640 uncorrected monthly and seasonal forecasts. The computational

170 time required to generate 6640 streamflow forecast ensembles through this process is small compared
 171 with the time required to calibrate and cross-validate the hydrological model, and is easily achieved in
 172 an operational setting using HPC. Note that in this study the forecasting system does not use data
 173 assimilation technique to update the GR4J state variables. This choice is based on the limited effect of
 174 initial conditions after a number of days, which generally reduces the benefit of state-updating in the
 175 context of seasonal streamflow forecasting.

176 **2.3 Streamflow post-processing procedure**

177 **2.3.1 Post-processing model**

178 The streamflow post-processing method used in this work consists of fitting a statistical model to the
 179 streamflow forecast residual errors, defined by the differences between the observed and forecast
 180 streamflow time series over a calibration period. Typically these errors are heteroscedastic, skewed and
 181 persistent. Heteroscedasticity and skew are handled using data transformations (e.g. the Box-Cox
 182 transformation), whereas persistence is represented using autoregressive models (e.g., the lag-one
 183 autoregressive model, AR(1)) (Wang et al., 2012; McInerney et al., 2017). We begin by describing the
 184 two major steps of the streamflow post-processing procedure (Sections 2.3.2 and 2.3.3), and then
 185 describe the transformations under consideration (Section 2.4).

186 **2.3.2 Post-processing model calibration**

187 The parameters of the streamflow post-processing model are calibrated as follows:

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

$$189 \quad \eta_t = Z(\widetilde{Q}_t) - Z(Q_t^F) \quad (1)$$

190 where η_t is the normalised residual, \widetilde{Q}_t is the observed streamflow, Q_t^F is the median of the uncorrected
 191 streamflow forecast ensemble, and Z is a transformation function. The transformation functions
 192 considered in this work are detailed in Section 2.4.

193 *Step 2:* Compute the standardised residuals:

$$194 \quad v_t = (\eta_t - \mu_\eta^{m(t)}) / \sigma_\eta^{m(t)} \quad (2)$$

195 where $\mu_\eta^{m(t)}$ and $\sigma_\eta^{m(t)}$ are the monthly mean and standard deviation of the residuals in the calibration
 196 period for the month $m(t)$.

197 The standardisation process in equation (2) aims to account for seasonal variations in the distribution of
 198 residuals. The quantities $\mu_\eta^{m(t)}$ and $\sigma_\eta^{m(t)}$ are calculated independently as the sample mean and standard

199 deviation of residuals for each monthly period (for a monthly forecast) or three-monthly period (for
 200 seasonal forecasts). Based on equation (2), the standardised residuals v_t are assumed to have a zero mean
 201 and unit standard deviation.

202 *Step 3:* Assume the standardised residuals are described by a first order autoregressive (AR(1)) model
 203 with Gaussian innovations:

$$204 \quad v_{t+1} = \rho v_t + y_{t+1} \quad (3)$$

205
 206 where ρ is the AR(1) coefficient and $y_{t+1} \sim N(0, \sigma_y)$ is the innovation.

207 The parameters ρ and σ_y are estimated using the method of moments (Hazelton, 2011): ρ is estimated
 208 as the sample auto-correlation of the standardized residuals \mathbf{v} , and σ_y is estimated as the sample
 209 standard deviation of the observed innovations \mathbf{y} , which in turn are calculated from the standardized
 210 residuals \mathbf{v} by re-arranging equation (3).

211 **2.3.3 Producing post-processed streamflow forecasts**

212 Once the streamflow post-processing scheme is calibrated, the post-processed streamflow forecasts for
 213 a given period are computed. For a given ensemble member j , the following steps are applied:

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

215 *Step 2:* Generate the standardized residuals $v_{t+1,j}$ using equation (3). Here $V_{t,j}$ is computed using
 216 equation (2) and $\eta_{t,j}$ is computed using equation (1), using the streamflow forecasts and observations
 217 from the previous time step t .

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

$$219 \quad \eta_{t+1,j} = \sigma_\eta^{m(t)} v_{t+1,j} + \mu_\eta^{m(t)} \quad (4)$$

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

$$221 \quad Q_{t+1,j}^{PP} = Z^{-1}[Z(Q_{t+1}^F) + \eta_{t+1,j}] \quad (5)$$

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

223 Note that the above algorithm may occasionally generate negative streamflow predictions, which we
 224 reset to zero. In addition, the algorithm can generate predictions that exceed historical maxima; such

225 predictions could in principle also be “adjusted” a posteriori, though we do not attempt such an
226 adjustment in this study. These aspects are discussed further in Section 5.6.

227 **2.4 Transformations used in the post-processing model**

228 The observed streamflow and median streamflow forecast are transformed in Step 1 of streamflow post-
229 processing (Section 2.3.2), to account for the heteroscedasticity and skewness of the forecast residuals.
230 We consider three transformations, namely the logarithmic, log-sinh and Box-Cox transformations.

231 **2.4.1 Logarithmic (Log) transformation**

232 The logarithmic (Log) transformation is

$$233 \quad Z(Q) = \log(Q + c) \quad (6)$$

234 The offset c ensures the transformed flows are defined when $Q = 0$. Here we set $c = 0.01 \times (\tilde{Q})_{ave}$
235 , where $(\tilde{Q})_{ave}$ is the average observed streamflow over the calibration period. The use of a small fixed
236 value for c is common in the literature for coping with zero flow events (Wang et al., 2012).

237 **2.4.2 Log-Sinh transformation**

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

$$239 \quad Z(Q) = \frac{1}{b} \log[\sinh(a + bQ)] \quad (7)$$

240 The parameters a and b are calibrated for each month by maximising the p-value of the Shapiro-Wilk
241 test (Shapiro and Wilk, 1965) for normality of the residuals, \mathbf{v} . This pragmatic approach is part of the
242 existing Bureau’s operational dynamic streamflow forecasting system (Lerat et al., 2015).

243 **2.4.3 Box-Cox transformation**

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

$$245 \quad Z(Q; \lambda, c) = \frac{(Q + c)^\lambda - 1}{\lambda} \quad (8)$$

246 where λ is a power parameter and $c = 0.01 \times (\tilde{Q})_{ave}$. Following the recommendations of McInerney et
247 al. (2017), the parameter λ is fixed to 0.2.

248 **2.4.4 Rationale for selecting transformational approaches**

249 The Log transformation is a simple and widely used transformation; McInerney et al. (2017) reported
250 that in daily scale modelling it produced the best reliability in perennial catchments (from a set of eight
251 residual error schemes, including standard least squares, weighted least squares, BC, Log-Sinh and

252 reciprocal transformation). However, the Log transformation performed poorly in ephemeral
253 catchments, where its precision was far worse than in perennial ones.

254 The Log-Sinh transformation is an alternative to the Log and BC transformations proposed by Wang et
255 al. (2012) to improve precision at higher flows. The Log-Sinh approach has been extensively applied to
256 water forecasting problems (see for example, Del Giudice et al., 2013; Robertson et al., 2013b, Bennett
257 et al., 2016). However, in daily scale streamflow modelling of perennial catchments using observed
258 rainfall, the Log-Sinh scheme did not improve on the Log transformation: its parameters tend to calibrate
259 to values for which the Log-Sinh transformation effectively reduces to the Log transformation
260 (McInerney et al., 2017).

261 Finally, the BC transformation with fixed $\lambda = 0.2$ is recommended by McInerney et al. (2017) as one of
262 only two schemes (from the set of eight schemes listed earlier in this section) that achieve Pareto-optimal
263 performance in terms of reliability, precision and bias, across both perennial and ephemeral catchments.
264 McInerney et al. (2017) also found that calibrating λ did not generally improve predictive performance,
265 due to the inferred value being dominated by the fit to the low flows at the expense of the high flows.

266 **2.5 Summary of key terms**

267 In the remainder of the paper, the term “uncorrected forecasts” refers to streamflow forecasts obtained
268 using steps in Section 2.2.3, and the term “post-processed forecasts” refers to forecasts based on a
269 streamflow post-processing model, which includes the standardization and AR(1) model from Section
270 2.3, as well as a transformation (Log, Log-Sinh or BC0.2) from Section 2.4. As the post-processing
271 schemes considered in this work differ solely in the transformation used, they will be referred to as the
272 Log, Log-Sinh and BC0.2 schemes.

273 **3 Application**

274 **3.1 Study catchments**

275 The empirical case study is carried out over a comprehensive set of 300 catchments with locations shown
276 in Figure 2. The figure also shows the Koppen climate zones. These catchments are selected as
277 representative of the diverse hydro-climatic conditions across Australia. The catchment areas range from
278 as small as 6 km² to as large as 232,846 km², with 90% of the catchments having areas below 6,000 km².
279 The seasonal streamflow forecasting service of the Bureau of Meteorology is currently evaluating these
280 300 catchments as part of an expansion of their dynamic modelling system.

281 **3.2 Catchment data**

282 In each catchment, data from 1980-2008 is used. Observed daily rainfall data was obtained from the
283 Australian Water Availability Project (AWAP) (Jeffrey et al., 2001). Potential evaporation and observed
284 streamflow data were obtained from the Bureau of Meteorology.

285 Catchment-scale rainfall forecasts are estimated from daily downscaled rainfall forecasts produced by
286 the Bureau of Meteorology's global climate model, namely the Predictive Ocean Atmosphere Model for
287 Australia (POAMA-2) (Hudson et al., 2013). The atmospheric component of POAMA-2 uses a spatial
288 scale of approximately 250×250 km (Charles et al., 2013). To estimate catchment-scale rainfall, a
289 statistical downscaling model based on an analogue approach (which could also be considered as rainfall
290 forecast post-processing) was applied (Timbal and McAvaney, 2001). In the analogue approach, local
291 climate information is obtained by matching analogous previous situations to the predicted climate. To
292 this end, an ensemble of 166 rainfall forecast time series (33 POAMA ensembles \times 5 replicates from
293 downscaling + 1 ensemble mean) were generated. In operation, POAMA-2 forecasts are generated every
294 week by running 33 member ensembles out to 270 days. In this study we use rainfall forecasts up to 3
295 months ahead and produce 166 rainfall forecast ensembles through the analogue downscaling procedure
296 described above.

297 **3.3 Catchment classification**

298 The performance of the post-processing schemes is evaluated separately in dry versus wet catchments.
299 In this work, the classification of catchments into dry and wet is based on the aridity index (AI) according
300 to the following equation

$$301 \quad AI = \frac{P}{PET} \quad (9)$$

302 where P is the total rainfall volume and PET is the total potential evapotranspiration volume. The aridity
303 index has been used extensively to identify and classify drought and wetness conditions of hydrological
304 regimes (Zhang et al., 2009; Carrillo et al., 2011; Sawicz et al., 2014).

305 Catchments with $AI < 0.5$ are categorised as “dry”, which corresponds to hyper-arid, arid and semi-arid
306 classifications suggested by the United Nations Environment Programme (Middleton et al., 1997).
307 Conversely, catchments with $AI \geq 0.5$ are classified as “wet”. Overall, about 28% of catchments used in
308 this work are classified as dry.

309 **3.4 Cross-validation procedure**

310 The forecast verification is carried out using a moving-window cross-validation framework, as shown
311 in Figure 3. We use 5 years data (1975-1979) to warm-up the model and apply data from 1980-2008 for

312 calibration in a cross-validation framework based on a 5-year moving window. Suppose we are
313 validating the streamflow forecasts in year j (e.g., $j=1990$ in Figure 3). In this case the calibration is
314 carried out using all years except years $j, j+1, j+2, j+3$ and $j+4$. The four-year period after year j is
315 excluded to prevent the memory of the hydrological model from affecting model performance in the
316 validation window period. The process is then repeated for each year during 1980-2008. Once the
317 validation has been carried out for each year, the results are concatenated to produce a single “validation”
318 time series, for which the performance metrics are calculated.

319 **3.5 Forecast performance (verification) metrics**

320 The performance of uncorrected and post-processed streamflow forecasts is evaluated using reliability
321 and sharpness metrics, as well as the Continuous Ranked Probability Skill Score (CRPSS, see section
322 3.5.3). Note that the Bureau of Meteorology uses Root Mean Squared Error (RMSE) and Root Mean
323 Squared Error in Probability (RMSEP) scores in the operational service in addition to CRPSS, however
324 these metrics have not been considered in this study.

325 Forecast performance (verification) metrics are computed separately for each forecast month. To
326 facilitate the comparison and evaluation of streamflow forecast performance in different streamflow
327 regimes, the high and low flow months are defined using long-term average streamflow data calculated
328 for each month. The 6 months with the highest average streamflow are classified as “high flow” months,
329 and the remaining 6 months are classified as “low flow” months. The performance metrics listed below
330 are computed for each month separately; the indices denoting the month are excluded from Equations
331 (10), (11) and (12) below to avoid cluttering the notation.

332 **3.5.1 Reliability**

333 The reliability of forecasts is evaluated using the Probability Integral Transform (PIT) (Dawid, 1984;
334 Laio and Tamea, 2007). To evaluate and compare reliability across 300 catchments, the p-value of the
335 Kolmogorov-Smirnov (KS) test applied to the PIT is used. In this study, forecasts with PIT plots where
336 the KS test yields a p-value $\geq 5\%$ are classified as “reliable”.

337 **3.5.2 Sharpness**

338 The sharpness of forecasts is evaluated using the ratio of inter-quantile ranges (IQR) of streamflow
339 forecasts and a historical reference (Tuteja et al., 2016). The following definition is used:

$$340 \quad IQR_q = \frac{1}{N} \sum_{i=1}^N \frac{F_i(100-q) - F_i(q)}{C_i(100-q) - C_i(q)} \times 100 \% \quad (10)$$

341 where IQR_q is the IQR value corresponding to percentile q , and $F_i(q)$ and $C_i(q)$ are, respectively, the
342 q^{th} percentiles of forecast and historical reference for year i .

343 An IQR_q of 100% indicates a forecast with the same sharpness as the reference, an IQR_q below 100%
 344 indicates forecasts that are sharper (tighter predictive limits) than the reference, and an IQR_q above
 345 100% indicates forecasts that are less sharp (wider predictive limits) than the reference. We report IQR_{99} ,
 346 i.e., the IQR at the 99 percentile, in order to detect forecasts with unreasonably long tails in their
 347 predictive distributions.

348 **3.5.3 CRPS skill score (CRPSS)**

349 The $CRPS$ metric quantifies the difference between a forecast distribution and observations, as follows
 350 (Hersbach, 2000),

$$351 \quad CRPS = \frac{1}{N} \times \sum_{i=1}^N \int_{-\infty}^{\infty} [F_i(y) - H_i\{y \geq y_o\}]^2 dy \quad (11)$$

352 where F_i is the cumulative distribution function (cdf) of the forecast for year i , y is the forecast variable
 353 (here streamflow) and y_o is the corresponding observed value. $H_i\{y \geq y_o\}$ is the Heaviside step function,
 354 which equals 1 when the forecast values are greater than the observed value and equals 0 otherwise.

355 The $CRPS$ summarises the reliability, sharpness and bias attributes of the forecast (Hersbach, 2000). A
 356 “perfect” forecast – namely a point prediction that matches the actual value of the predicted quantity –
 357 has $CRPS^P = 0$. In this work, we use the $CRPS$ skill score, CRPSS, defined by

$$358 \quad CRPSS = \frac{CRPS^F - CRPS^C}{CRPS^P - CRPS^C} \times 100\% \quad (12)$$

359 where $CRPS^F$, $CRPS^C$ and $CRPS^P$ represent the $CRPS$ value for model forecast, climatology and
 360 “perfect” forecast respectively. A higher CRPSS indicates better performance, with a value of 0
 361 representing the same performance as climatology.

362 **3.5.4 Historical reference**

363 The IQR and CRPSS metrics are defined as skill scores relative to a reference forecast. In this work, we
 364 use the climatology as the reference forecast, as it represents the long-term climate condition. To
 365 construct these “climatological forecasts”, we used the same historical reference as the operational
 366 seasonal streamflow forecasting service of the Bureau of Meteorology. This reference is resampled from
 367 a Gaussian probability distribution fitted to the observed streamflow transformed using the Log-Sinh
 368 transformation (Equation 7). This approach leads to more stable and continuous historical reference
 369 estimates than sampling directly from the empirical distribution of historical streamflow, and can be
 370 computed at any percentile (which facilitates comparison with forecast percentiles). Although the choice
 371 of a particular reference affects the computation of skill scores, it does not affect the ranking of post-
 372 processing models when the same reference is used, which is the main aim of this paper.

373 **3.5.5 Summary skill: Summarising forecast performance using multiple metrics**

374 When evaluating forecast performance, a focus on any single individual metric can lead to misleading
375 interpretations. For example, two forecasts might have a similar sharpness, yet if one of these forecasts
376 is unreliable it can lead to an over- or under- estimation of the risk of
377 an event of interest, which in turn can lead to a sub-optimal decision by forecast users (e.g. a water
378 resources manager).

379 Given inevitable trade-offs between individual metrics (McInerney et al., 2017), it is important to
380 consider multiple metrics jointly rather than individually. Following the approach suggested by Gneiting
381 et al. (2007), we consider a forecast to have “high skill” when it is reliable *and* sharper than climatology.
382 To determine the “summary skill” of the forecasts in each catchment, we evaluate the total number of
383 months (out of 12) in which forecasts are reliable (i.e., with a p-value greater than 5%) and sharper than
384 the climatology (i.e., $IQR_{99} < 100\%$). A catchment is classified as having high summary skill if “high
385 skill” forecasts are obtained 10-12 months per year (on average), and is classified as having low
386 summary skill otherwise. Note that CRPSS is not included in the summary skill, because it does not
387 represent an independent measure of a forecast attribute (see Section 3.5.3 for more details).

388 A table providing the percentage of catchments with high and low summary skills is used to summarise
389 forecasts performance of a given post-processing scheme. To identify any geographic trends in the
390 forecast performance, the summary skills are plotted on a map. The summary skills together with
391 individual skill score values are used to evaluate the overall forecast performance, and are presented
392 separately for wet and dry catchments, as well as separately for high and low flow months.

393 **4 Results**

394 Results for monthly and seasonal streamflow forecasts are now presented. Section 4.1 compares the
395 uncorrected and post-processed streamflow forecast performance. Section 4.2 evaluates the performance
396 of post-processed streamflow forecasts obtained using the Log, Log-Sinh and BC0.2 schemes. The
397 CRPSS, reliability and sharpness metrics are presented in Figure 4 and Figure 5 for monthly and seasonal
398 forecasts respectively.

399 Initial inspection of results found considerable overlap in the performance metrics achieved by the error
400 models. To determine whether the differences in metrics are consistent over multiple catchments, the
401 Log and Log-Sinh schemes are compared to the BC0.2 scheme. This comparison is presented in Figure
402 6 and Figure 7 for monthly and seasonal forecasts respectively. The BC0.2 scheme is taken as the
403 baseline because inspection of Figure 4 and Figure 5 suggests that the BC0.2 scheme has better median

404 sharpness than the Log and Log-Sinh schemes, over all the catchments and for both high and low flow
405 months individually.

406 The streamflow forecast time-series and corresponding skill for a single representative catchment,
407 Dieckmans Bridge, are presented in Figure 8 and Figure 9, respectively.

408 The summary skills of the monthly and seasonal forecasts are presented in Figure 10 and Figure 11. The
409 figures include a histogram of summary skills across all catchments to enable comparison between the
410 uncorrected and the post-processing approaches.

411 **4.1 Comparison of uncorrected and post-processed streamflow forecasts: Individual** 412 **metrics**

413 In terms of CRPSS, the largest improvement as a result of post-processing (using any of the
414 transformations considered here) occurs in dry catchments. This finding holds for both monthly (Figure
415 4c) and seasonal forecasts (Figure 5c). For example, when post-processing is implemented, the median
416 CRPSS of monthly forecasts in dry catchments increases from approximately 7% (high flow months)
417 and -15% (low flow months) to more than 10% (Figure 4c) for both high and low flows. Visible
418 improvement is also observed in dry catchments for seasonal forecasts, however, the improvement is
419 not as pronounced as for monthly forecasts (Figure 5c).

420 In terms of reliability, the performance of uncorrected streamflow forecasts is poor, with about 50% of
421 the catchments being characterized by unreliable forecasts at both the monthly and seasonal time scales
422 (Figure 4 and Figure 5, middle row). In comparison, post-processing using the three transformation
423 approaches produces much better reliability, achieving reliable forecasts in more than 90% of the
424 catchments.

425 In terms of sharpness, the uncorrected forecasts and the BC0.2 post-processed forecasts are generally
426 sharper than forecasts generated using the other transformations (Figure 4g and Figure 5g). The use of
427 post-processing achieves much better sharpness than uncorrected forecasts for low flow months,
428 particularly in dry catchments. For example, for low flow months in dry catchments (Figure 4i), the
429 median IQR99 is greater than 200%, while similar values range between 40-100% for post-processed
430 forecasts. Similarly, for seasonal forecasts, post-processing approaches improve the median sharpness
431 from 150% (uncorrected forecasts) to 50%-110% (Figure 5i).

432 **4.2 Comparison of post-processing schemes: Individual metrics**

433 In terms of CRPSS, Figure 4 (a, b, c) and Figure 5 (a, b, c) show considerable overlap in the boxplots
434 corresponding to all three post-processing schemes, both in wet and dry catchments. This finding
435 suggests little difference in the performance of the post-processing schemes, and is further confirmed by

436 Figure 6 (a, b, c) and Figure 7 (a, b, c), which show boxplots of the differences between the CRPSS of
437 the Log and Log-Sinh schemes versus the CRPSS of the BC0.2 scheme. Across all catchments, the
438 distribution of these differences is approximately symmetric with a mean close to 0. In dry catchments,
439 the BC0.2 slightly outperforms the Log scheme for high flow months and the Log-Sinh scheme slightly
440 outperforms the Log scheme for low flow months. Overall, these results suggest that none of the Log,
441 Log-Sinh or BC0.2 schemes is consistently better in terms of CRPSS values.

442 In terms of reliability, post-processing using any of the three post-processing schemes produces reliable
443 forecasts at both monthly and seasonal scales, and in the majority of the catchments (Figure 4 and Figure
444 5, middle row). The median p-value is approximately 60% for monthly forecasts compared with 45%
445 for seasonal forecasts. This indicates that better forecast reliability is achieved at shorter lead times.
446 Median reliability is somewhat reduced when using the BC0.2 scheme compared to the Log and Log-
447 Sinh schemes in wet catchments (Figure 6e), but not so much in dry catchments (Figure 6f).
448 Nevertheless, the monthly and seasonal forecasts are reliable in 96% and 91% of the catchments,
449 respectively. The corresponding percentages for the Log scheme are 97% and 94%, and for Log-Sinh
450 they are 95% and 90%.

451 In terms of sharpness, the BC0.2 scheme outperforms the Log and Log-Sinh schemes. This finding holds
452 in all cases (i.e., high/low flow months and wet/dry catchments), both for monthly and seasonal forecasts
453 (Figure 4 and Figure 5, bottom row). The plot of differences in the sharpness metric (Figure 6 and Figure
454 7, bottom row) highlights this improvement. In half of the catchments, during both high and low flow
455 months, the BC0.2 scheme improves the IQR99 by 30% (or more) compared to the Log and Log-Sinh
456 schemes. In dry catchments, the improvements are larger than in wet catchments. For example, in dry
457 catchments during high flow months, the BC0.2 scheme improves on the IQR99 of Log and Log-Sinh
458 by 40-60% in over a half of the catchments, and by as much as 170%-190% in a quarter of the
459 catchments.

460 To illustrate these results, a streamflow forecast time-series at Dieckmans Bridge catchment (site id:
461 145010A) is shown in Figure 8 and performance metrics calculated over six high flow months and six
462 low flow months are shown in Figure 9. This catchment is selected as it is broadly representative of
463 typical results obtained across the wide range of case study catchments. The period in Figure 8 (2003-
464 2007) is chosen because it highlights the difference in forecast interval between the uncorrected and
465 post-processing approaches. The figure indicates that in terms of reliability, the uncorrected forecast has
466 a number of observed data points outside the 99% predictive range (Figure 8a). This is an indication that
467 the forecast is unreliable. This finding can be confirmed from the corresponding p-value in Figure 9,
468 which shows that the forecast is below the reliability threshold during most of the high flow months and

469 during some low flow months. In terms of sharpness, Log and Log-Sinh schemes produce a wider 99%
470 predictive range than the BC0.2 scheme (Figure 8 and Figure 9).

471 **4.3 Comparison of summary skill between uncorrected and post-processing approaches**

472 Figure 10 and Figure 11 show the geographic distribution of the summary skill of the uncorrected and
473 post-processing approaches for monthly and seasonal forecasts respectively. Recall that the summary
474 skill represents the number of months with streamflow forecasts that are both reliable and sharper than
475 climatology. Table 1 provides a summary of the percentage of catchments with high and low summary
476 skill for the uncorrected and post-processing approaches for monthly and seasonal forecasts (see Section
477 3.5.5).

478 The findings for forecasts at monthly scale are as follows (Figure 10 and Table 1):

- 479 • Uncorrected forecasts perform worse than post-processing techniques in the sense that they have
480 low summary skill in the largest percentage of catchments (16%). The percentage of catchments
481 where high summary skill is achieved by uncorrected forecasts is 40%.
- 482 • Post-processing forecasts with the Log and Log-Sinh scheme reduces the percentage of
483 catchments with low summary skills from 16% to 2% and 7% respectively. However, the
484 percentage of catchments with high summary skill also decreases (in comparison to uncorrected
485 forecasts), from 40% to 33% for both the Log and Log-Sinh schemes.

486 Post-processing with the BC0.2 scheme provides the best performance, with the smallest percentage of
487 catchments with low summary skills (<1%) and the largest percentage of catchments with high summary
488 skills (84%). As seen in Figure 10

- 489 • Figure 10, the improvement achieved by the BC0.2 scheme (compared to the Log/Log-Sinh
490 schemes) is most pronounced in New South Wales (NSW) and in the tropical catchments in
491 Queensland (QLD) and the Northern Territory (NT). The few catchments where the BC0.2
492 scheme does not achieve a high summary skill are located in the north and north-west of
493 Australia.

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

- 495 • Log scheme has the largest percentage (19%) of catchments with low summary skill and a
496 relatively small percentage (9%) of catchments with high summary skill.
- 497 • Post-processing forecasts with the Log and Log-Sinh schemes reduces the percentages of
498 catchments with low summary skill from 19% to 18% and 17% respectively. The percentage of
499 catchments with high summary skill increases from 9% to 12% and 22% respectively.
- 500 • Post-processing with the BC0.2 scheme once again provides the best performance: it produces
501 forecasts with low summary skill in only 2% of the catchments, and achieves high summary skill

502 in 54% of the catchments. As seen in Figure 11, similar to the case of monthly forecasts, the
503 biggest improvements for seasonal forecasts occur in the NSW and Queensland regions of
504 Australia.

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

509 **4.4 Summary of empirical findings**

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

517 **5 Discussion**

518 **5.1 Benefits of forecast post-processing**

519 A comparison of uncorrected and post-processed streamflow forecasts was provided in Section 4.1.
520 Uncorrected forecasts have reasonable sharpness (except for in dry catchments), but suffer from low
521 reliability: uncorrected forecasts are unreliable at approximately 50% of the catchments. In wet
522 catchments, poor reliability is due to overconfident forecasts, which appears a common concern in
523 dynamic forecasting approaches (Wood and Schaake, 2008). In dry catchments, uncorrected forecasts
524 are both unreliable and exhibit poor sharpness. Post-processing is thus particularly important to correct
525 for these shortcomings and improve forecast skill. In this study, all post-processing models provide a
526 clear improvement in reliability and sharpness, especially in dry catchments. The value of post-
527 processing is more pronounced in dry catchments than in wet catchments (Figure 4 and Figure 5). This
528 finding can be attributed to the challenge of capturing key physical processes in dry and ephemeral
529 catchments (Ye et al., 1997), as well as the challenge of achieving accurate rainfall forecasts in arid
530 areas. In addition, the simplifications inherent in any hydrological model, including the conceptual
531 model GR4J used in this work, might also be responsible for the forecast skill being relatively lower in
532 dry catchments than in wet catchments. Whilst using a single conceptual model is attractive for practical
533 operational system, there may be gains in exploring alternative structures for ephemeral catchments (e.g.
534 Clark et al., 2008; Fenicia et al., 2011). We intend to explore such alternative model structures for

535 difficult ephemeral catchments. In such dry catchments, the hydrological model forecasts are particularly
536 poor and leave a lot of room for improvement: post-processing can hence make a big difference on the
537 quality of results.

538 **5.2 Interpretation of differences between post-processing schemes**

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

546 McInerney et al. (2017) found that, when modelling perennial catchments at the daily scale, the Log-
547 Sinh scheme did not achieve better sharpness than the Log scheme. Instead, the parameters for the Log
548 scheme tended to converge to values for which the tapering off of the Log-Sinh transformation function
549 occurs well outside the range of simulated flows, effectively reducing the Log-Sinh scheme to the Log
550 scheme. In contrast, the Box-Cox transformation function with a fixed $\lambda > 0$ gradually flattens as
551 streamflow increases, and exhibits the “desired” tapering-off behaviour within the range of simulated
552 flows. This behaviour leads to the Box-Cox scheme achieving, on average, more favourable variance-
553 stabilizing characteristics than the Log-Sinh scheme.

554 Our findings in this study confirm the insights of McInerney et al. (2017) – namely that the Log-Sinh
555 scheme produces comparable sharpness to the Log scheme – across a wider range of catchments. This
556 finding indicates that insights from modelling residual errors at the daily scale apply at least to some
557 extent to streamflow forecast post-processing at the monthly and seasonal scales. Note the minor
558 difference in the treatment of the offset parameter c in equation (6): in the Log scheme used in McInerney
559 et al. (2017) this parameter is inferred, whereas in this study it is fixed a priori. This minor difference
560 does not impact on the qualitative behaviour of the error models described earlier in this section. Overall,
561 when used for post-processing seasonal and monthly forecasts in a dynamic modelling system, the
562 BC0.2 scheme provides an opportunity to improve forecast performance further than is possible using
563 the Log and Log-Sinh schemes.

564 **5.3 Importance of using multiple metrics to assess forecast performance**

565 The goal of the forecasting exercise is to maximise sharpness without sacrificing reliability (Gneiting et
566 al., 2005; Wilks, 2011; Bourdin et al., 2014). The study results show that relying on a single metric for
567 evaluating forecast performance can lead to sub-optimal conclusions. For example, if one considers the

568 CRPSS metric alone, all post-processing schemes yield comparable performance and there is no basis
569 for favouring any single one of them. However, once sharpness is taken into consideration explicitly,
570 the BC0.2 scheme can be recommended due to substantially better sharpness than the Log and Log-Sinh
571 schemes.

572 Similarly, comparisons based solely on CRPSS might suggest reasonable performance of the
573 uncorrected forecasts: 55%-80% of months have $CRPSS > 0$ (with some variability across high/low flow
574 months and monthly/seasonal forecasts). Yet once reliability is considered explicitly, it is found that
575 uncorrected forecasts are unreliable at approximately 50% of the catchments. Note that performance
576 metrics based on the CRPSS reflect an implicitly weighted combination of reliability, sharpness and bias
577 characteristics of the forecasts (Hersbach, 2000). In contrast, the reliability and sharpness metrics are
578 specifically designed to quantify reliability and sharpness attributes individually. These findings
579 highlight the value of multiple independent performance metrics and diagnostics that evaluate specific
580 (targeted) attributes of the forecasts, and highlight important limitations of aggregate measures of
581 performance (Clark et al., 2011).

582 A number of challenges and questions remain in regards to selecting the performance verification metrics
583 for specific forecasting systems and applications. An important question is how to include user needs
584 into a forecast verification protocol. This could be accomplished by tailoring the evaluation metrics to
585 the requirements of users. Another key question is to what extent do measures of forecast skill correlate
586 to the economic and/or social value of the forecast? This challenging question was investigated by
587 Murphy and Ehrendorfer (1987) and Wandishin and Brooks (2002), who found the relationship between
588 quality and value of a forecast to be essentially nonlinear: an increase in forecast quality may not
589 necessarily lead to a proportional increase in its value. This question requires further multi-disciplinary
590 research, including human psychology, economic theory, communication and social studies (e.g. Matte
591 et al., 2017; Morss et al., 2010).

592 **5.4 Importance of performance evaluation over large numbers of catchments**

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

602 to establish the robustness of post-processing methods. Restricting the analysis to a smaller number of
603 catchments would have led to less conclusive findings.

604 **5.5 Implication of results for water resource management**

605 The empirical results clearly show that the BC0.2 post-processing scheme improves forecast sharpness
606 (precision) while maintaining forecast accuracy and reliability. As discussed below, this improvement
607 in forecast quality offers an opportunity to improve operational planning and management of water
608 resources.

609 The management of water resources, for example, deciding which water source to use for a particular
610 purpose or allocating environmental flows, requires an understanding of the current and future
611 availability of water. For water resources systems with long hydrological records, water managers have
612 devised techniques to evaluate current water availability, water demand and losses. However, one of the
613 main unknowns is the volume of future system inflows. Streamflow
614 forecasts provide crucial information to water managers and users regarding the future availability of
615 water, thus helping reduce uncertainty in decision making. This information is particularly valuable to
616 support decision during drought events. In this study, forecast performance is evaluated separately for
617 high and low flow months – providing a clearer indication of predictive ability for flows that are above
618 and below average, respectively. A detailed evaluation of forecasts for more extreme drought events is
619 challenging as these events are correspondingly rarer. Limited sample size makes it difficult to make
620 conclusive statements: e.g. if we focus on the lowest 5% of historical data with a 30 year record, we may
621 only have roughly 1.5 samples for each month/season. The uncertainty arising from limited sample size
622 requires further development of forecast verification techniques, potentially adapting some of the
623 approaches used by Hodgkins et al. (2017).

624 **5.6 Opportunities for further improvement in forecast performance**

625 There are several opportunities to further improve the seasonal streamflow forecasting system. This
626 section describes avenues related to specialised treatment of zero flows and high flow forecasts,
627 uncertainty analysis of post-processing model parameters, and the use of data assimilation (state
628 updating).

629 The post-processing approaches used in this work do not make special provision for zero flows in the
630 observed data. Robust handling of zero flows in statistical models, especially in arid and semi-arid
631 catchments, is an active research area (Wang and Robertson, 2011; Smith et al., 2015), and advances in
632 this area are certainly relevant to seasonal streamflow forecasting.

633 A similar challenge is associated with the forecasting of high flows, as the post-processing approaches
634 used in this work can produce streamflow predictions that exceed historical maxima. The IQR ratio used
635 to assess forecast sharpness will detect unreasonably long tails (i.e. extremes) in the predictive
636 distributions and hence can hence indirectly identify instances of unreasonably high flow forecasts.
637 Further research is needed to develop techniques to evaluate the realism of forecasts that exceed
638 historical maxima.

639 Another area for further investigation is the identifiability of parameters $\mu_{\eta}^{m(t)}$ and $\sigma_{\eta}^{m(t)}$ of the monthly
640 post-processing model. These parameters are estimated using monthly data (see Section 2.3.2), and
641 hence could be subject to substantial uncertainty and/or over-fitting to the calibration period. In this
642 study, 29 years of data were employed in the calibration, making these problems unlikely. Importantly,
643 the use of a cross-validation procedure (Section 3.4) is expected to detect potential overfitting. That said,
644 as many sites of potential application may lack the data length available in this work, the sensitivity of
645 forecast performance to the length of calibration period warrants further investigation.

646 Finally, the forecasting system used in this study does not employ data assimilation to update the states
647 of the GR4J hydrological model. Gibbs et al. (2018) showed that monthly streamflow forecasting
648 benefits from state updating in catchments that exhibit non-stationarity in their rainfall-runoff dynamics.
649 Note that data assimilation of ocean observations has been implemented in the climate model
650 (POAMA2) used for the rainfall forecast (Yin et al., 2011) (see Section 3.2 for additional details).

651 **6 Conclusions**

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

656 We investigated streamflow forecast post-processing schemes based on residual error models employing
657 three data transformations, namely the logarithmic (Log), log-sinh (Log-Sinh) and Box-Cox with $\lambda = 0.2$
658 (BC0.2). The Australian Bureau of Meteorology's dynamic modelling system was used as the platform
659 for the empirical analysis, which was carried out over 300 Australian catchments with diverse hydro-
660 climatic conditions.

661 The following empirical findings are obtained:

- 662 1. Uncorrected forecasts (no post-processing) perform poorly in terms of reliability, resulting in a
663 mischaracterization of forecast uncertainties;

- 664 2. All three post-processing schemes substantially improve the reliability of streamflow forecasts,
665 both in terms of the dedicated reliability metric and in terms of the summary skill given by the
666 CRPSS;
- 667 3. From the post-processing schemes considered in this work, the BC0.2 scheme is found best
668 suited for operational application. The BC0.2 scheme provides the sharpest forecasts without
669 sacrificing reliability, as measured by the reliability and CRPSS metrics. In particular, the BC0.2
670 scheme produces forecasts that are both reliable and sharper than climatology at substantially
671 more catchments than the alternative Log and Log-Sinh schemes.

672 A major practical outcome of this study is the development of a robust streamflow forecast post-
673 processing scheme that achieves forecasts that are consistently reliable and sharper than climatology.
674 This scheme is well suited for operational application, and offers the opportunity to improve decision
675 support, especially in catchments where climatology is presently used to guide operational decisions.

676 **7 Data availability**

677 The data underlying this research can be accessed from the following links: observed rainfall data
678 (<http://www.bom.gov.au/climate>), POAMA rainfall forecast (<http://poama.bom.gov.au>), and observed
679 streamflow data (<http://www.bom.gov.au/waterdata>).

680 **8 Acknowledgments**

681 Data for this study is provided by the Australian Bureau of Meteorology. This work was supported by
682 the Australian Research Council grant LP140100978 with the Australian Bureau of Meteorology and
683 South East Queensland Water. We thank the anonymous reviewers for constructive comments and
684 feedback that have helped us substantially improve the manuscript.

685

686

687 **9 References**

- 688 Bennett, J. C., Wang, Q. J., Li, M., Robertson, D. E. and Schepen, A.: Reliable long-range ensemble
689 streamflow forecasts: Combining calibrated climate forecasts with a conceptual runoff model and a
690 staged error model, *Water Resour. Res.*, 52(10), 8238–8259, doi:10.1002/2016WR019193, 2016.
- 691 Bennett, J. C., Wang, Q. J., Robertson, D. E., Schepen, A., Li, M. and Michael, K.: Assessment of an
692 ensemble seasonal streamflow forecasting system for Australia, *Hydrol. Earth Syst. Sci.*, 21(12), 6007–
693 6030, doi:10.5194/hess-21-6007-2017, 2017.
- 694 Bogner, K. and Kalas, M.: Error-correction methods and evaluation of an ensemble based hydrological
695 forecasting system for the Upper Danube catchment, *Atmos. Sci. Lett.*, 9(2), 95–102,
696 doi:10.1002/asl.180, 2008.
- 697 Bourdin, D. R., Nipen, T. N. and Stull, R. B.: Reliable probabilistic forecasts from an ensemble reservoir
698 inflow forecasting system, *Water Resour. Res.*, 50(4), 3108–3130, doi:10.1002/2014WR015462, 2014.
- 699 Box, G. E. P. and Cox, D. R.: An analysis of transformations, *J. R. Stat. Soc. Ser. B (Methodological)*,
700 211–252, doi:10.2307/2287791, 1964.
- 701 Brown, J. D., Wu, L., He, M., Regonda, S., Lee, H. and Seo, D. J.: Verification of temperature,
702 precipitation, and streamflow forecasts from the NOAA/NWS Hydrologic Ensemble Forecast Service
703 (HEFS): 1. Experimental design and forcing verification, *J. Hydrol.*, 519(PD), 2869–2889,
704 doi:10.1016/j.jhydrol.2014.05.028, 2014.
- 705 Carpenter, T. M. and Georgakakos, K. P.: Assessment of Folsom lake response to historical and potential
706 future climate scenarios: 1. Forecasting, *J. Hydrol.*, 249(1–4), 148–175,
707 doi:https://doi.org/10.1016/S0022-1694(01)00417-6, 2001.
- 708 Carrillo, G., Troch, P. A., Sivapalan, M., Wagener, T., Harman, C. and Sawicz, K.: Catchment
709 classification: hydrological analysis of catchment behavior through process-based modeling along a
710 climate gradient, *Hydrol. Earth Syst. Sci.*, 15(11), 3411–3430, doi:10.5194/hess-15-3411-2011, 2011.
- 711 Charles, A., Miles, E., Griesser, A., de Wit, R., Shelton, K., Cottrill, A., Spillman, C., Hendon, H.,
712 McIntosh, P., Nakaegawa, T., Atalifo, T., Prakash, B., Seuseu, S., Nihmei, S., Church, J., Jones, D. and
713 Kuleshov, Y.: Dynamical Seasonal Prediction of Climate Extremes in the Pacific, in 20th International
714 Congress on Modelling and Simulation (Modsim2013), pp. 2841–2847., 2013.
- 715 Clark, M. P., Slater, A. G., Rupp, D. E., Woods, R. A., Vrugt, J. A., Gupta, H. V., Wagener, T. and Hay,
716 L. E.: Framework for Understanding Structural Errors (FUSE): A modular framework to diagnose
717 differences between hydrological models, *Water Resour. Res.*, 44(12), doi:10.1029/2007WR006735,
718 2008.
- 719 Clark, M. P., Kavetski, D. and Fenicia, F.: Pursuing the method of multiple working hypotheses for

720 hydrological modeling, *Water Resour. Res.*, 47(9), n/a-n/a, doi:10.1029/2010WR009827, 2011.

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

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

725 Crochemore, L., Ramos, M. H. and Pappenberger, F.: Bias correcting precipitation forecasts to improve
726 the skill of seasonal streamflow forecasts, *Hydrol. Earth Syst. Sci.*, 20(9), 3601–3618, doi:10.5194/hess-
727 20-3601-2016, 2016.

728 Dawid, a P.: Present Position and Potential Developments: Some Personal Views: Statistical theory: the
729 prequential approach (with discussion), *J. R. Stat. Soc. Ser. A*, 147(2), 278–292, doi:10.2307/2981683,
730 1984.

731 Dechant, C. M. and Moradkhani, H.: Improving the characterization of initial condition for ensemble
732 streamflow prediction using data assimilation, *Hydrol. Earth Syst. Sci.*, 15(11), 3399–3410,
733 doi:10.5194/hess-15-3399-2011, 2011.

734 Demargne, J., Wu, L., Regonda, S. K., Brown, J. D., Lee, H., He, M., Seo, D. J., Hartman, R., Herr, H.
735 D., Fresch, M., Schaake, J. and Zhu, Y.: The science of NOAA’s operational hydrologic ensemble
736 forecast service, *Bull. Am. Meteorol. Soc.*, 95(1), 79–98, doi:10.1175/BAMS-D-12-00081.1, 2014.

737 Evin, G., Thyer, M., Kavetski, D., McInerney, D. and Kuczera, G.: Comparison of joint versus
738 postprocessor approaches for hydrological uncertainty estimation accounting for error autocorrelation
739 and heteroscedasticity, *Water Resour. Res.*, 50(3), 2350–2375, doi:10.1002/2013WR014185, 2014.

740 Fenicia, F., Kavetski, D. and Savenije, H. H. G.: Elements of a flexible approach for conceptual
741 hydrological modeling: 1. Motivation and theoretical development, *Water Resour. Res.*, 47(11), 1–13,
742 doi:10.1029/2010WR010174, 2011.

743 Gibbs, M. S., McInerney, D., Humphrey, G., Thyer, M. A., Maier, H. R., Dandy, G. C. and Kavetski,
744 D.: State updating and calibration period selection to improve dynamic monthly streamflow forecasts
745 for an environmental flow management application, *Hydrol. Earth Syst. Sci.*, 22(1), 871–887,
746 doi:10.5194/hess-22-871-2018, 2018.

747 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P. and Rieckermann, J.: Improving
748 uncertainty estimation in urban hydrological modeling by statistically describing bias, *Hydrol. Earth*
749 *Syst. Sci.*, 17(10), 4209–4225, doi:10.5194/hess-17-4209-2013, 2013.

750 Gneiting, T., Raftery, A. E., Westveld, A. H. and Goldman, T.: Calibrated Probabilistic Forecasting
751 Using Ensemble Model Output Statistics and Minimum CRPS Estimation, *Mon. Weather Rev.*, 133(5),
752 1098–1118, doi:10.1175/MWR2904.1, 2005.

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

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

755 Hashino, T., Bradley, a. a. and Schwartz, S. S.: Evaluation of bias-correction methods for ensemble
756 streamflow volume forecasts, Hydrol. Earth Syst. Sci., 11, 939–950, doi:10.5194/hess-11-939-2007,
757 2007.

758 Hazelton, M. L.: Methods of Moments Estimation BT - International Encyclopedia of Statistical
759 Science, edited by M. Lovric, pp. 816–817, Springer Berlin Heidelberg, Berlin, Heidelberg., 2011.

760 Hersbach, H.: Decomposition of the Continuous Ranked Probability Score for Ensemble Prediction
761 Systems, Weather Forecast., 15(5), 559–570, doi:10.1175/1520-
762 0434(2000)015<0559:DOTCRP>2.0.CO;2, 2000.

763 Hudson, D., Marshall, A. G., Yin, Y., Alves, O. and Hendon, H. H.: Improving Intraseasonal Prediction
764 with a New Ensemble Generation Strategy, Mon. Weather Rev., 141(12), 4429–4449,
765 doi:10.1175/MWR-D-13-00059.1, 2013.

766 Humphrey, G. B., Gibbs, M. S., Dandy, G. C. and Maier, H. R.: A hybrid approach to monthly
767 streamflow forecasting: Integrating hydrological model outputs into a Bayesian artificial neural network,
768 J. Hydrol., 540, 623–640, doi:10.1016/j.jhydrol.2016.06.026, 2016.

769 Jeffrey, S. J., Carter, J. O., Moodie, K. B. and Beswick, A. R.: Using spatial interpolation to construct a
770 comprehensive archive of Australian climate data, Environ. Model. Softw., 16(4), 309–330,
771 doi:10.1016/S1364-8152(01)00008-1, 2001.

772 Kavetski, D., Kuczera, G. and Franks, S. W.: Bayesian analysis of input uncertainty in hydrological
773 modeling: 1. Theory, Water Resour. Res., 42(3), n/a-n/a, doi:10.1029/2005WR004368, 2006.

774 Knoche, M., Fischer, C., Pohl, E., Krause, P. and Merz, R.: Combined uncertainty of hydrological model
775 complexity and satellite-based forcing data evaluated in two data-scarce semi-arid catchments in
776 Ethiopia, J. Hydrol., 519, 2049–2066, doi:https://doi.org/10.1016/j.jhydrol.2014.10.003, 2014.

777 Kuczera, G., Kavetski, D., Franks, S. and Thyer, M.: Towards a Bayesian total error analysis of
778 conceptual rainfall-runoff models: Characterising model error using storm-dependent parameters, J.
779 Hydrol., 331(1–2), 161–177, doi:10.1016/j.jhydrol.2006.05.010, 2006.

780 Laio, F. and Tamea, S.: Verification tools for probabilistic forecasts of continuous hydrological
781 variables, Hydrol. Earth Syst. Sci., 11(4), 1267–1277, doi:10.5194/hess-11-1267-2007, 2007.

782 Laugesen, R., Tuteja, N. K., Shin, D., Chia, T. and Khan, U.: Seasonal Streamflow Forecasting with a
783 workflow-based dynamic hydrologic modelling approach, in MODSIM 2011 - 19th International
784 Congress on Modelling and Simulation - Sustaining Our Future: Understanding and Living with
785 Uncertainty, pp. 2352–2358. [online] Available from: [http://www.scopus.com/inward/record.url?eid=2-](http://www.scopus.com/inward/record.url?eid=2-s2.0-84858823270&partnerID=tZOtx3y1)
786 [s2.0-84858823270&partnerID=tZOtx3y1](http://www.scopus.com/inward/record.url?eid=2-s2.0-84858823270&partnerID=tZOtx3y1), 2011.

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

788 G., Thyer, M. and Kavetski, D.: Dynamic streamflow forecasts within an uncertainty framework for 100
789 catchments in Australia, in In: 36th Hydrology and Water Resources Symposium: The art and science
790 of water, pp. 1396–1403, Barton, ACT: Engineers Australia., 2015.

791 Li, M., Wang, Q. J., Bennett, J. C. and Robertson, D. E.: Error reduction and representation in stages
792 (ERRIS) in hydrological modelling for ensemble streamflow forecasting, *Hydrol. Earth Syst. Sci.*, 20(9),
793 3561–3579, doi:10.5194/hess-20-3561-2016, 2016.

794 Lü, H., Crow, W. T., Zhu, Y., Ouyang, F. and Su, J.: Improving streamflow prediction using remotely-
795 sensed soil moisture and snow depth, *Remote Sens.*, 8(6), doi:10.3390/rs8060503, 2016.

796 Madadgar, S., Moradkhani, H. and Garen, D.: Towards improved post-processing of hydrologic forecast
797 ensembles, *Hydrol. Process.*, 28(1), 104–122, doi:10.1002/hyp.9562, 2014.

798 Matte, S., Boucher, M. A., Boucher, V. and Fortier Filion, T. C.: Moving beyond the cost-loss ratio:
799 Economic assessment of streamflow forecasts for a risk-Averse decision maker, *Hydrol. Earth Syst. Sci.*,
800 21(6), 2967–2986, doi:10.5194/hess-21-2967-2017, 2017.

801 McInerney, D., Thyer, M., Kavetski, D., Lerat, J. and Kuczera, G.: Improving probabilistic prediction
802 of daily streamflow by identifying Pareto optimal approaches for modeling heteroscedastic residual
803 errors, *Water Resour. Res.*, 53(3), 2199–2239, doi:10.1002/2016WR019168, 2017.

804 Mendoza, P. A., Wood, A. W., Clark, E., Rothwell, E., Clark, M. P., Nijssen, B., Brekke, L. D. and
805 Arnold, J. R.: An intercomparison of approaches for improving predictability in operational seasonal
806 streamflow forecasting, *Hydrol. Earth Syst. Sci. Discuss.*, 2017, 1–37, doi:10.5194/hess-2017-60, 2017.

807 Middleton, N., Programme, U. N. E. and Thomas, D. S. G.: *World Atlas of Desertification*, Arnold.,
808 1997.

809 Morss, R. E., Lazo, J. K. and Demuth, J. L.: Examining the use of weather forecasts in decision scenarios:
810 Results from a us survey with implications for uncertainty communication, *Meteorol. Appl.*, 17(2), 149–
811 162, doi:10.1002/met.196, 2010.

812 Murphy, A. H. and Ehrendorfer, M.: On the relationship between the accuracy and value of forecasts in
813 the cost–loss ratio situation, *Weather Forecast.*, 2(3), 243–251, doi:10.1175/1520-
814 0434(1987)002<0243:OTRBTA>2.0.CO;2, 1987.

815 Perrin, C., Michel, C. and Andréassian, V.: Improvement of a parsimonious model for streamflow
816 simulation, *J. Hydrol.*, 279(1–4), 275–289, doi:10.1016/S0022-1694(03)00225-7, 2003.

817 Pokhrel, P., Robertson, D. E. and Wang, Q. J.: A Bayesian joint probability post-processor for reducing
818 errors and quantifying uncertainty in monthly streamflow predictions, *Hydrol. Earth Syst. Sci.*, 17(2),
819 795–804, doi:10.5194/hess-17-795-2013, 2013.

820 Prudhomme, C., Hannaford, J., Harrigan, S., Boorman, D., Knight, J., Bell, V., Jackson, C., Svensson,
821 C., Parry, S., Bachiller-Jareno, N., Davies, H., Davis, R., Mackay, J., McKenzie, A., Rudd, A., Smith,

822 K., Bloomfield, J., Ward, R. and Jenkins, A.: Hydrological Outlook UK: an operational streamflow and
823 groundwater level forecasting system at monthly to seasonal time scales, *Hydrol. Sci. J.*, 62(16), 2753–
824 2768, doi:10.1080/02626667.2017.1395032, 2017.

825 Renard, B., Kavetski, D., Leblois, E., Thyer, M., Kuczera, G. and Franks, S. W.: Toward a reliable
826 decomposition of predictive uncertainty in hydrological modeling: Characterizing rainfall errors using
827 conditional simulation, *Water Resour. Res.*, 47(11), n/a-n/a, doi:10.1029/2011WR010643, 2011.

828 Robertson, D. E. Wang, Q. J.: Selecting predictors for seasonal streamflow predictions using a Bayesian
829 joint probability (BJP) modelling approach, 18th World IMACS/MODSIM Congr. Cairns, Aust. 13-
830 17 July 2009, (July), 376–382, 2009.

831 Robertson, D. E. and Wang, Q. J.: A Bayesian Approach to Predictor Selection for Seasonal Streamflow
832 Forecasting, *J. Hydrometeorol.*, 13(1), 155–171, doi:10.1175/JHM-D-10-05009.1, 2011.

833 Robertson, D. E., Pokhrel, P. and Wang, Q. J.: Improving statistical forecasts of seasonal streamflows
834 using hydrological model output, *Hydrol. Earth Syst. Sci.*, 17(2), 579–593, doi:10.5194/hess-17-579-
835 2013, 2013a.

836 Robertson, D. E., Shrestha, D. L. and Wang, Q. J.: Post-processing rainfall forecasts from numerical
837 weather prediction models for short-term streamflow forecasting, *Hydrol. Earth Syst. Sci.*, 17(9), 3587–
838 3603, doi:10.5194/hess-17-3587-2013, 2013b.

839 Sawicz, K. A., Kelleher, C., Wagener, T., Troch, P., Sivapalan, M. and Carrillo, G.: Characterizing
840 hydrologic change through catchment classification, *Hydrol. Earth Syst. Sci.*, 18(1), 273–285,
841 doi:10.5194/hess-18-273-2014, 2014.

842 Schick, S., Rössler, O. and Weingartner, R.: Monthly streamflow forecasting at varying spatial scales in
843 the Rhine basin, *Hydrol. Earth Syst. Sci.*, 22(2), 929–942, doi:10.5194/hess-22-929-2018, 2018.

844 Senlin, Z., Feikema, P., Shin, D., Tuteja, N. K., MacDonald, A., Sunter, P., Kent, D., Le, B., Pipunic,
845 R., Wilson, T., Pickett-Heaps, C. and Lerat, J.: Operational efficiency measures of the national seasonal
846 streamflow forecast service in Australia, edited by G. Syme, D. H. MacDonald, B. Fulton, and J.
847 Piantadosi, the Modelling and Simulation Society of Australia and New Zealand Inc, Hobart, Australia.,
848 2017.

849 Seo, D.-J., Herr, H. D. and Schaake, J. C.: A statistical post-processor for accounting of hydrologic
850 uncertainty in short-range ensemble streamflow prediction, *Hydrol. Earth Syst. Sci. Discuss.*, 3(4),
851 1987–2035, doi:10.5194/hessd-3-1987-2006, 2006.

852 Shapiro, S. S. and Wilk, M. B.: An Analysis of Variance Test for Normality (Complete Samples),
853 *Biometrika*, 52(3–4), 591–611, doi:10.2307/1267427, 1965.

854 Smith, T., Marshall, L. and Sharma, A.: Modeling residual hydrologic errors with Bayesian inference,
855 *J. Hydrol.*, 528(SEPTEMBER 2015), 29–37, doi:10.1016/j.jhydrol.2015.05.051, 2015.

856 Tang, Q. and Lettenmaier, D. P.: Use of satellite snow-cover data for streamflow prediction in the
857 Feather River Basin, California, *Int. J. Remote Sens.*, 31(14), 3745–3762,
858 doi:10.1080/01431161.2010.483493, 2010.

859 Taschetto, A. S. and England, M. H.: An analysis of late twentieth century trends in Australian rainfall,
860 *Int. J. Climatol.*, 29(6), 791–807, doi:10.1002/joc.1736, 2009.

861 Timbal, B. and McAvaney, B. J.: An Analogue based method to downscale surface air temperature:
862 Application for Australia, *Clim. Dyn.*, 17, 947–963, doi:10.1007/s003820100156, 2001.

863 Turner, S. W. D., Bennett, J., Robertson, D. and Galelli, S.: Value of seasonal streamflow forecasts in
864 emergency response reservoir management, *Hydrol. Earth Syst. Sci. Discuss.*, 2017, 1–26,
865 doi:10.5194/hess-2016-691, 2017.

866 Tuteja, N. K., Shin, D., Laugesen, R., Khan, U., Shao, Q., Wang, E., Li, M., Zheng, H., Kuczera, G.,
867 Kavetski, D., Evin, G., Thyer, M., MacDonald, A., Chia, T. and Le, B.: Experimental evaluation of the
868 dynamic seasonal streamflow forecasting approach, Melbourne., 2011.

869 Tuteja, N. K., Zhou, S., Lerat, J., Wang, Q. J., Shin, D. and Robertson, D. E.: Overview of
870 Communication Strategies for Uncertainty in Hydrological Forecasting in Australia, in *Handbook of*
871 *Hydrometeorological Ensemble Forecasting*, edited by Q. Duan, F. Pappenberger, J. Thielen, A. Wood,
872 H. L. Cloke, and J. C. Schaake, pp. 1–19, Springer Berlin Heidelberg, Berlin, Heidelberg., 2016.

873 Tyralla, C. and Schumann, A. H.: Incorporating structural uncertainty of hydrological models in
874 likelihood functions via an ensemble range approach, *Hydrol. Sci. J.*, 02626667.2016.1164314,
875 doi:10.1080/02626667.2016.1164314, 2016.

876 Wandishin, M. S. and Brooks, H. E.: On the relationship between Clayton’s skill score and expected
877 value for forecasts of binary events, *Meteorol. Appl.*, 9(4), 455–459, doi:10.1017/S1350482702004085,
878 2002.

879 Wang, Q. J. and Robertson, D. E.: Multisite probabilistic forecasting of seasonal flows for streams with
880 zero value occurrences, *Water Resour. Res.*, 47(2), doi:10.1029/2010WR009333, 2011.

881 Wang, Q. J., Robertson, D. E. and Chiew, F. H. S.: A Bayesian joint probability modeling approach for
882 seasonal forecasting of streamflows at multiple sites, *Water Resour. Res.*, 45(5),
883 doi:10.1029/2008WR007355, 2009.

884 Wang, Q. J., Shrestha, D. L., Robertson, D. E. and Pokhrel, P.: A log-sinh transformation for data
885 normalization and variance stabilization, *Water Resour. Res.*, 48(5), doi:10.1029/2011WR010973,
886 2012.

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

888 Wood, A. W. and Schaake, J. C.: Correcting Errors in Streamflow Forecast Ensemble Mean and Spread,
889 *J. Hydrometeorol.*, 9(1), 132–148, doi:10.1175/2007JHM862.1, 2008.

890 Ye, W., Bates, B. C., Viney, N. R., Sivapalan, M. and Jakeman, A. J.: Performance of conceptual
891 rainfall-runoff models in low-yielding ephemeral catchments, *Water Resour. Res.*, 33(1), 153–166,
892 doi:10.1029/96WR02840, 1997.

893 Yin, Y., Alves, O., Oke, P. R., Yin, Y., Alves, O. and Oke, P. R.: An ensemble ocean data assimilation
894 system for seasonal prediction, *Mon. Weather Rev.*, 139(3), 786–808, doi:10.1175/2010MWR3419.1,
895 2011.

896 Zhang, Q., Xu, C.-Y. and Zhang, Z.: Observed changes of drought/wetness episodes in the Pearl River
897 basin, China, using the standardized precipitation index and aridity index, *Theor. Appl. Climatol.*, 98(1),
898 89–99, doi:10.1007/s00704-008-0095-4, 2009.

899 Zhao, T., Schepen, A. and Wang, Q. J.: Ensemble forecasting of sub-seasonal to seasonal streamflow by
900 a Bayesian joint probability modelling approach, *J. Hydrol.*, 541, 839–849,
901 doi:<https://doi.org/10.1016/j.jhydrol.2016.07.040>, 2016.

902

903

904

905

906

907

908 **Tables**

909

910 Table 1. Performance of post-processing schemes, expressed as the percentage of catchments with high
911 and low summary skill. Results shown for monthly and seasonal forecasts. A catchment with “high
912 summary skill” is defined as a catchment where “high skill” forecasts are achieved in 10-12 months out
913 of the year; “high skill” forecasts are defined as forecasts that are reliable and sharper than climatology.

	Post-processing scheme			
	Uncorrected forecasts	Log	Log-Sinh	BC0.2
<i>Monthly Forecasts</i>				
High Summary Skill	40%	33%	33%	84%
Low Summary Skill	16%	2%	7%	<1%
<i>Seasonal Forecasts</i>				
High Summary Skill	46%	9%	20%	54%
Low Summary Skill	14%	19%	17%	2%

914

915

916

917

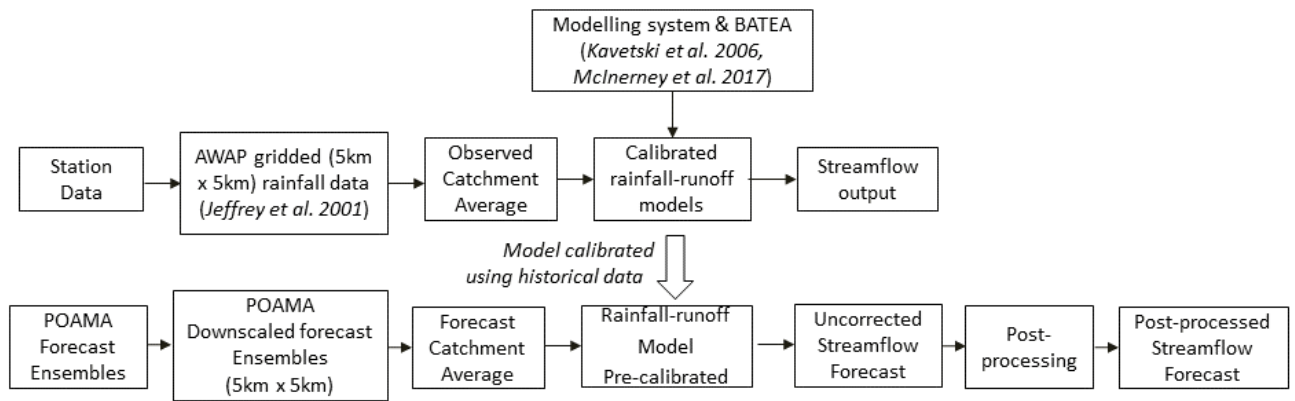
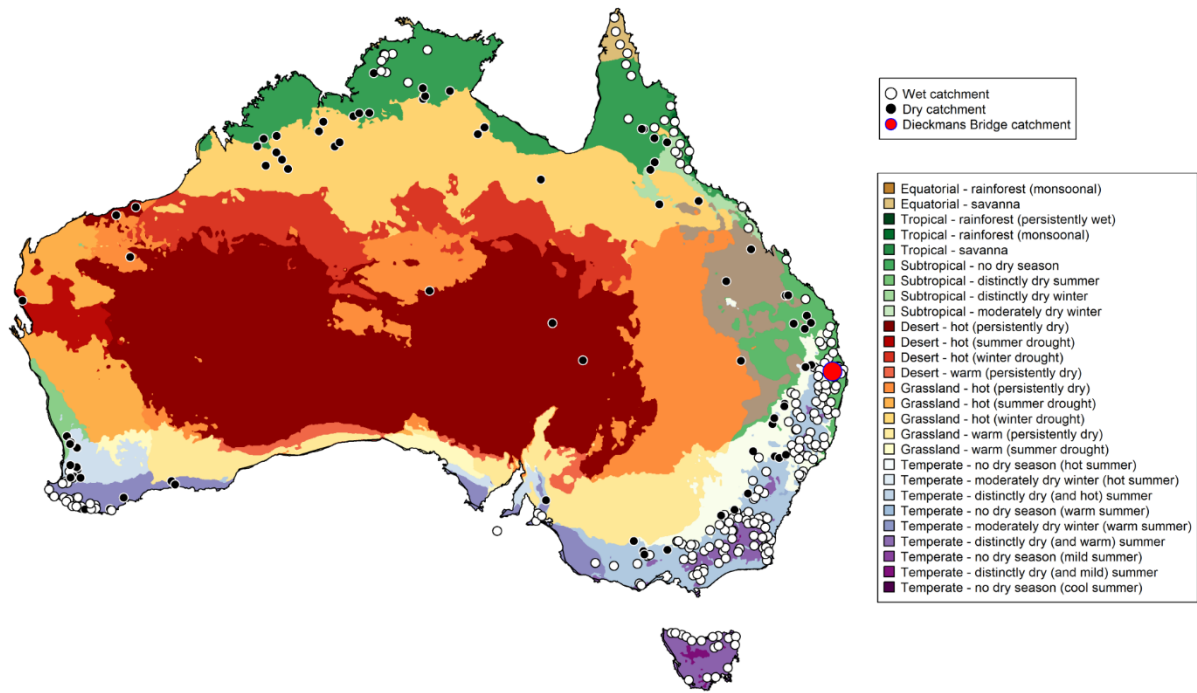


Figure 1: Schematic of the dynamic streamflow forecasting system used in this study. A similar approach is used by the Australian Bureau of Meteorology for its monthly and seasonal streamflow forecasting service.

919

920

921



922

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

927

928

929

930

931

932

933

934

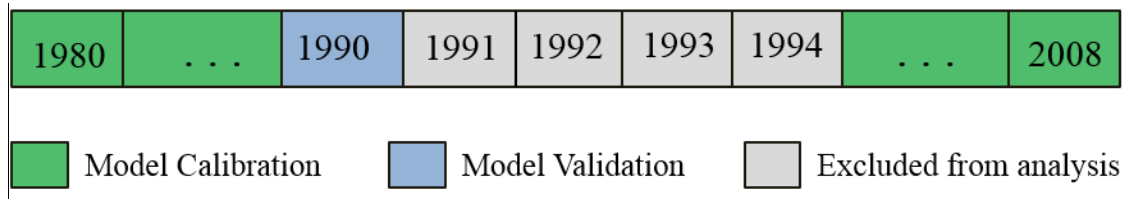
935

936

937

938

939



940

941 Figure 3: Schematic of the cross-validation framework used for forecast verification, applied with the 5-
942 year validation period window beginning in year 1990 (after Tuteja et al., 2016).

943

944

945

946

947

948

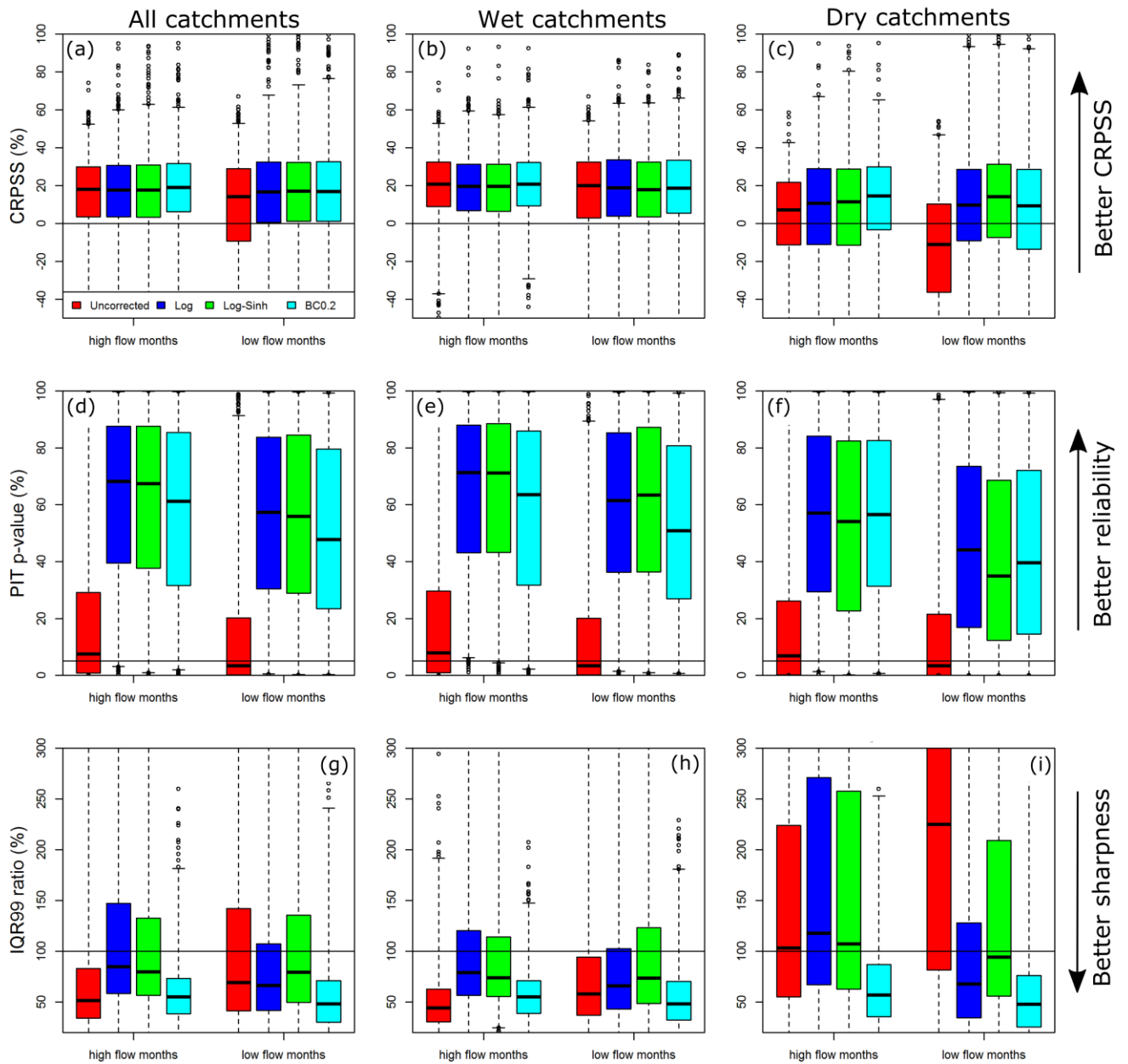
949

950

951

952

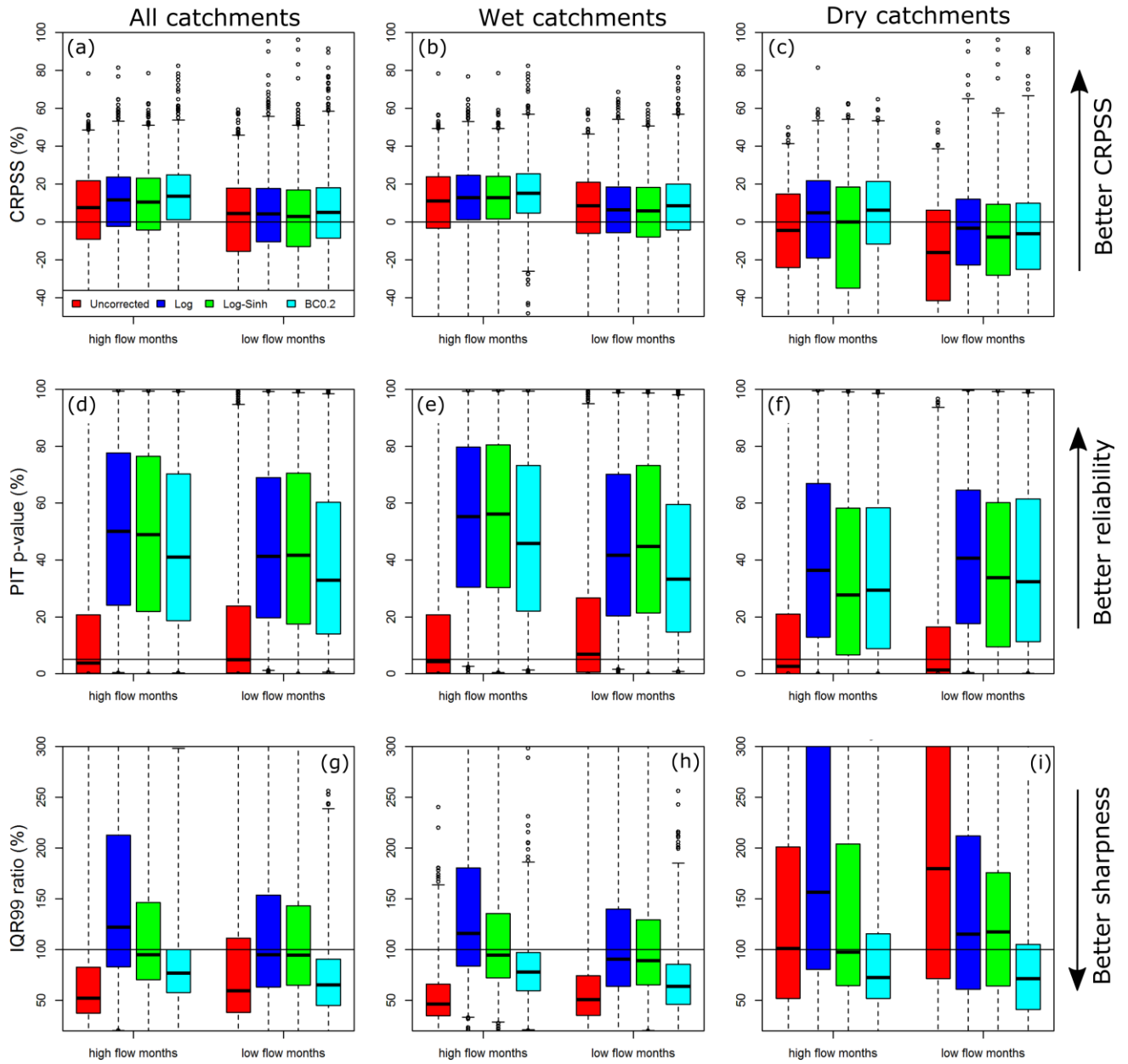
953



954

955 Figure 4: Performance of monthly forecasts in terms of CRPSS, reliability (PIT p-value) and sharpness
 956 (IQR99 ratio).
 957

958



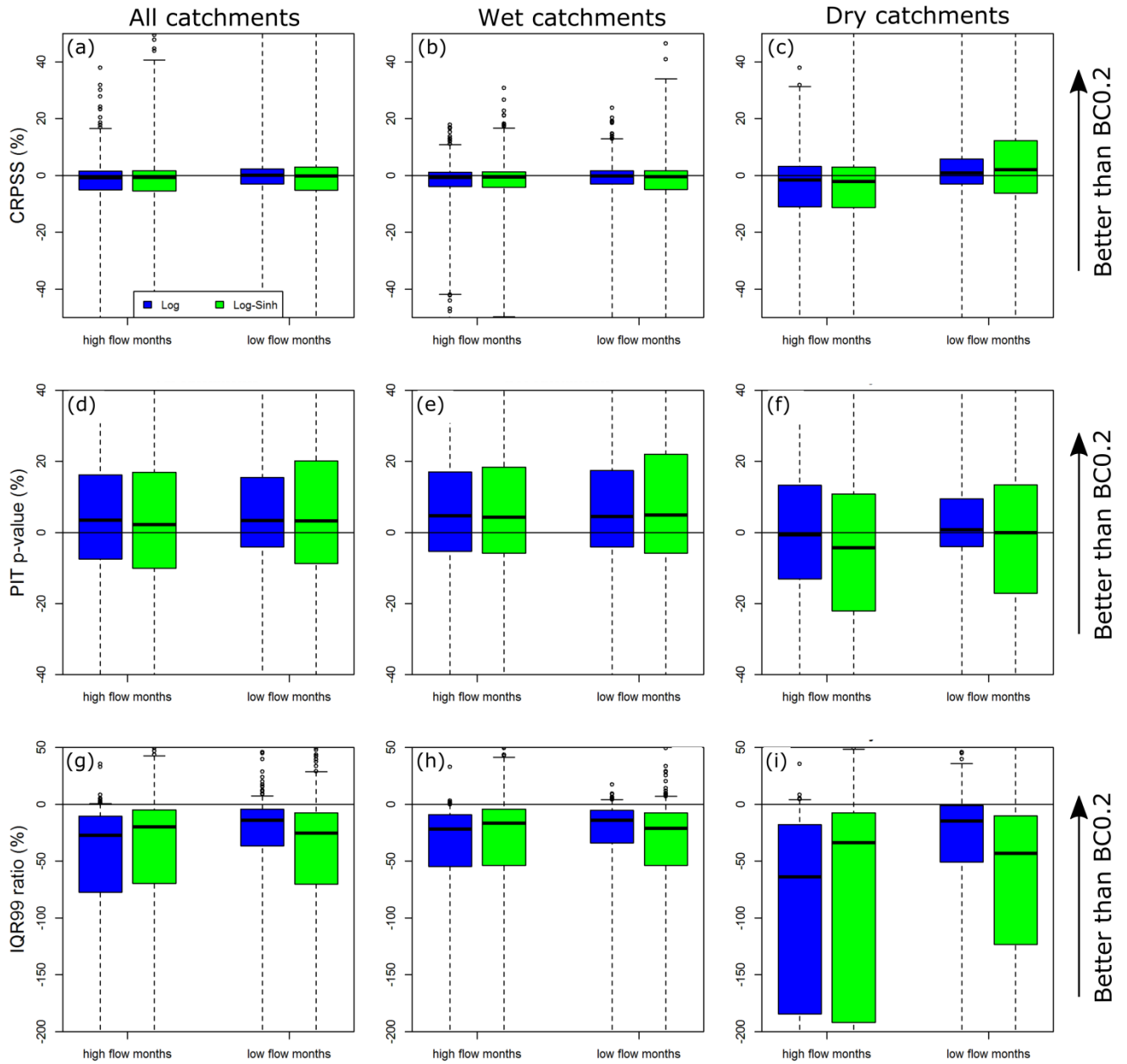
959

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

962

963

964



965
966

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

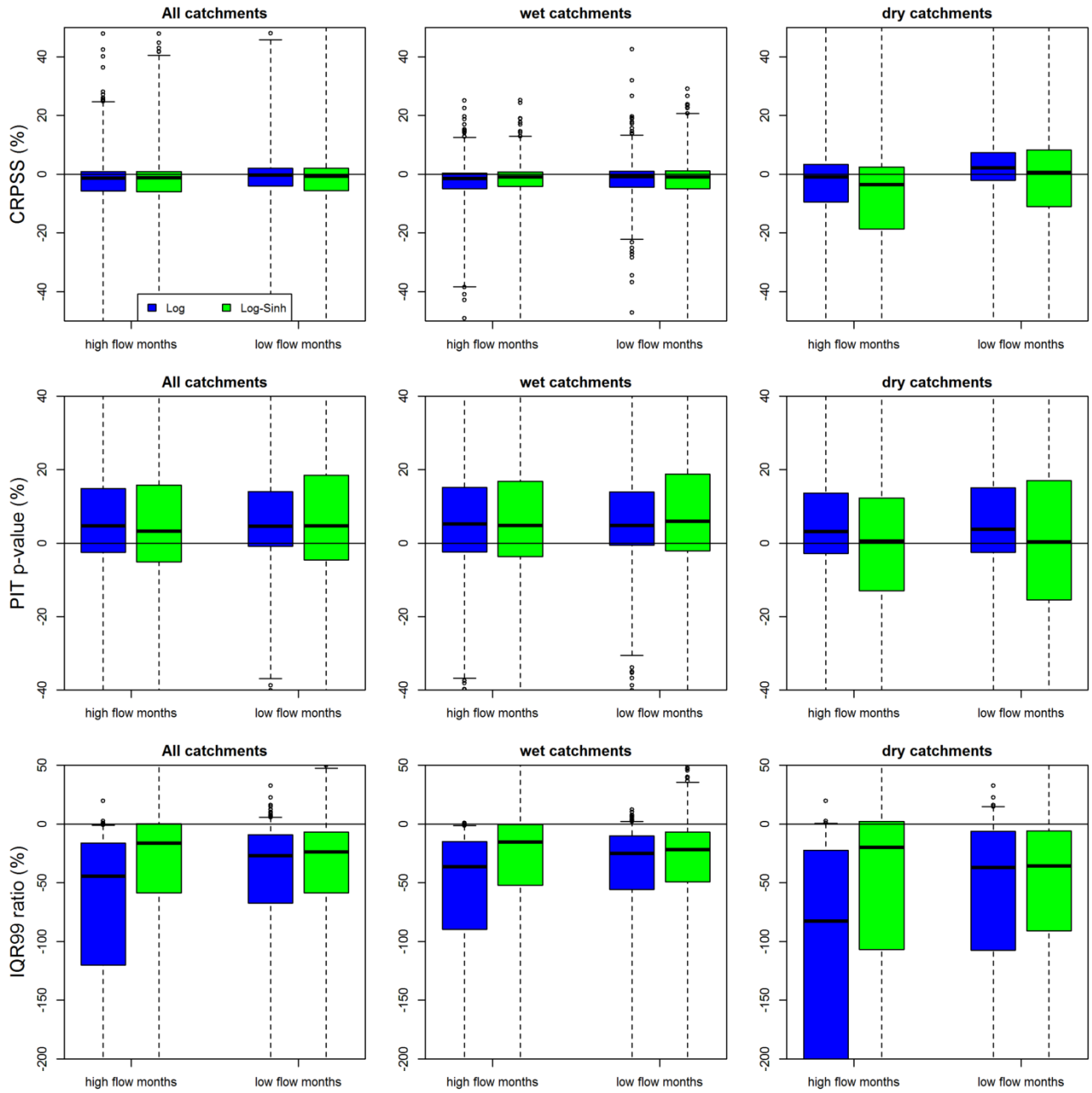
969

970

971

972

973

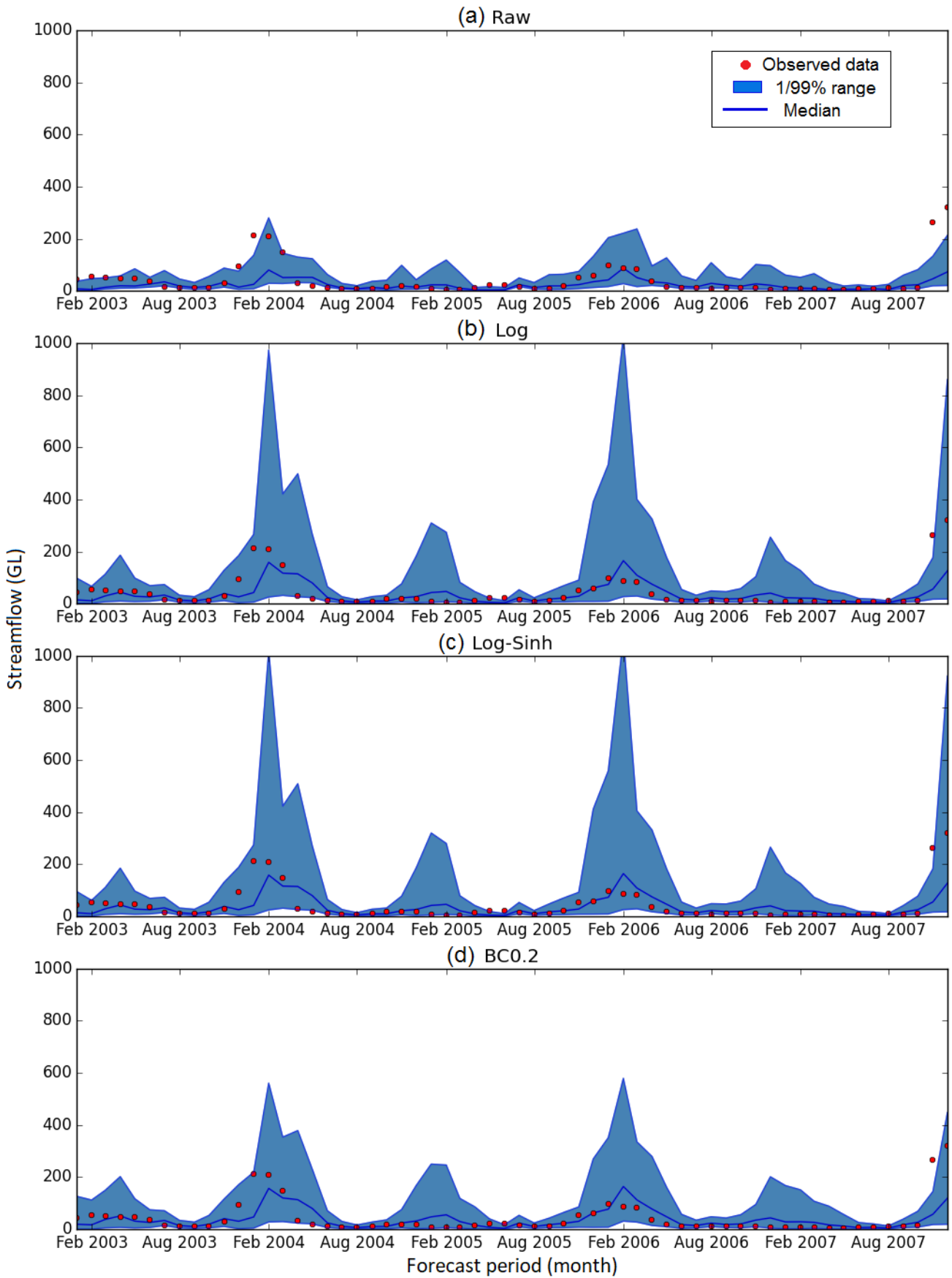


974

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

977

978



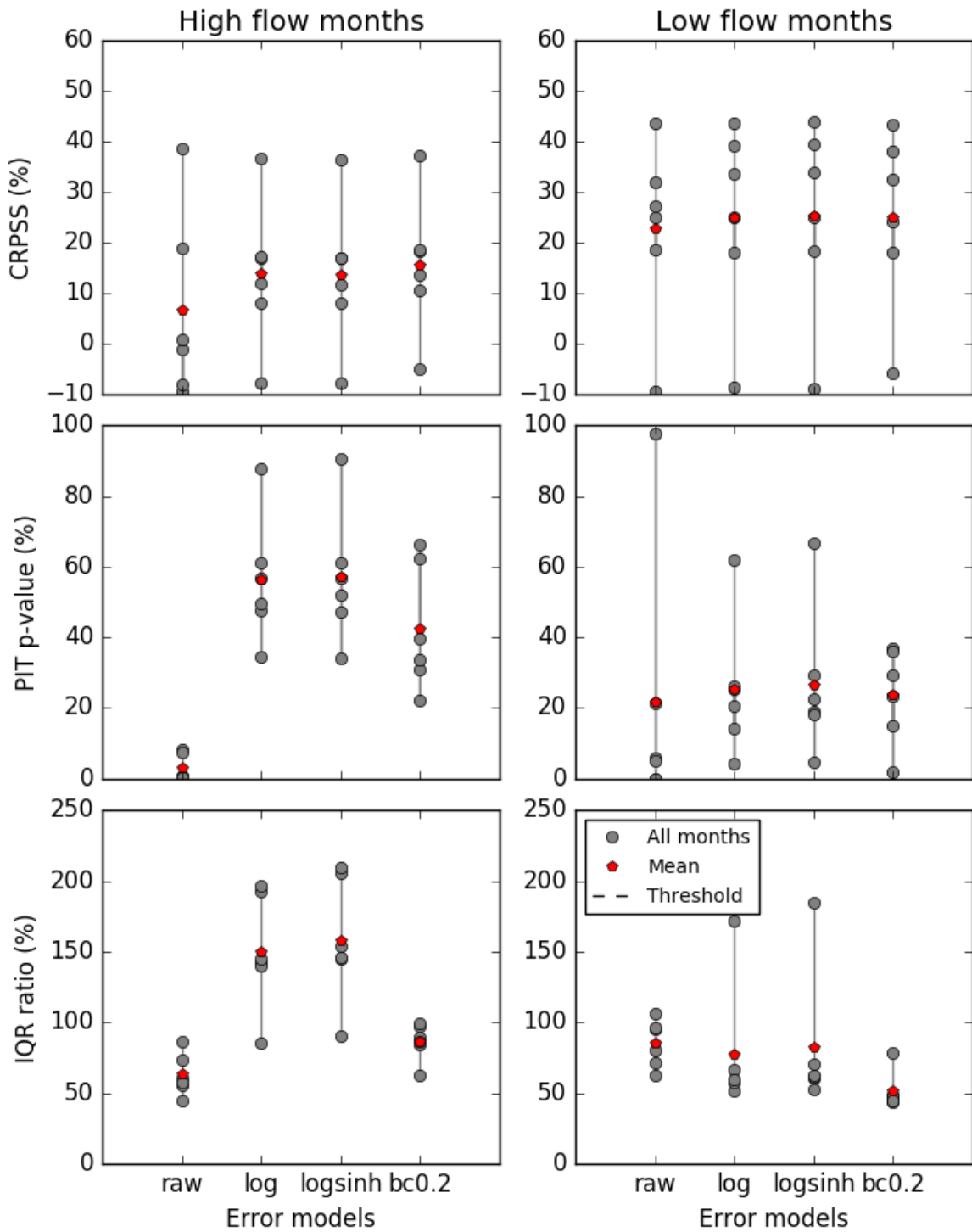
980

981

982

983

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



985

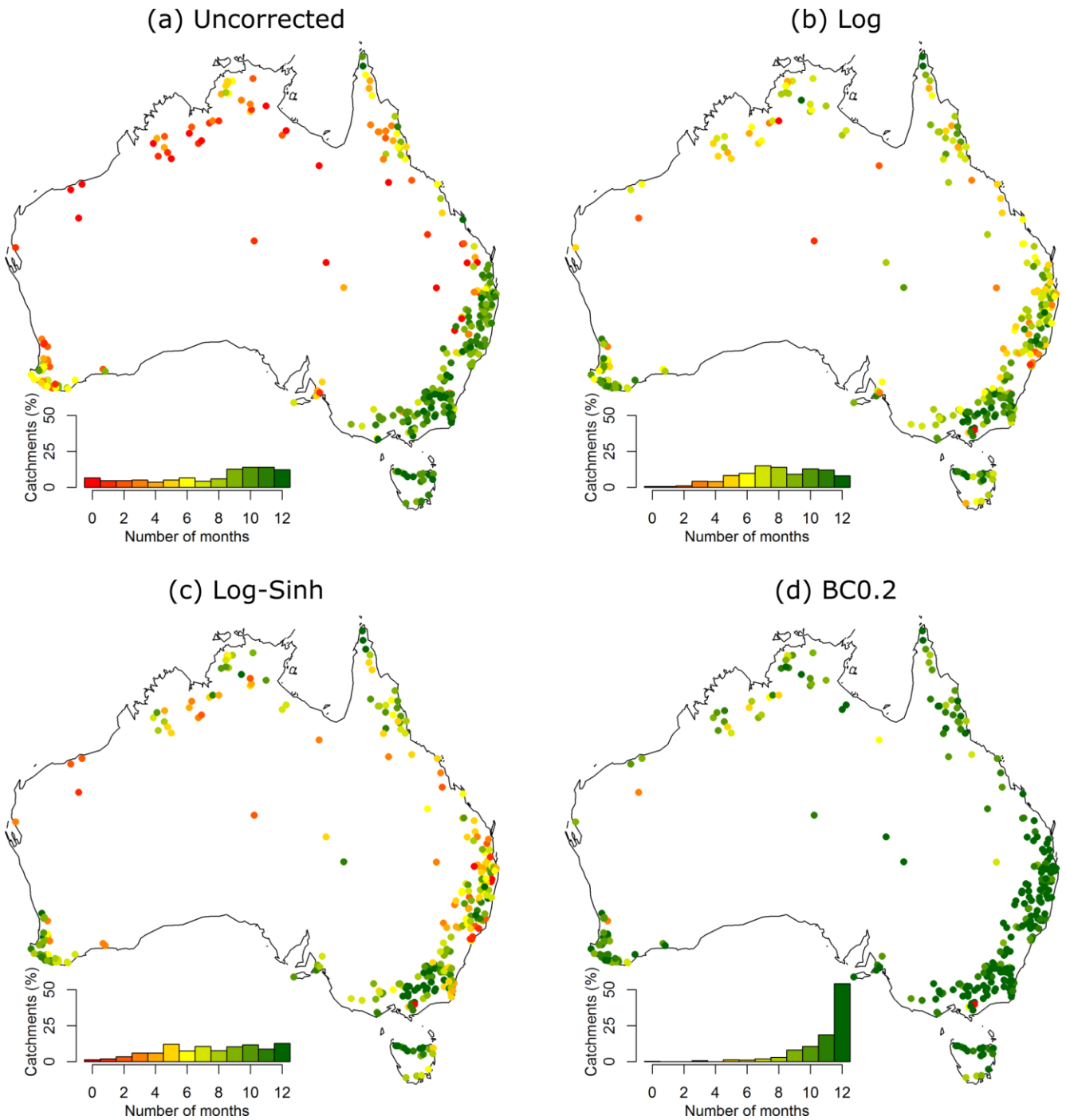
986 Figure 9: Seasonal streamflow forecast skill scores at Dieckmans Bridge catchment, computed from the
 987 time series shown in Figure 8 for six high flow months and six low flow months.

988

989

990

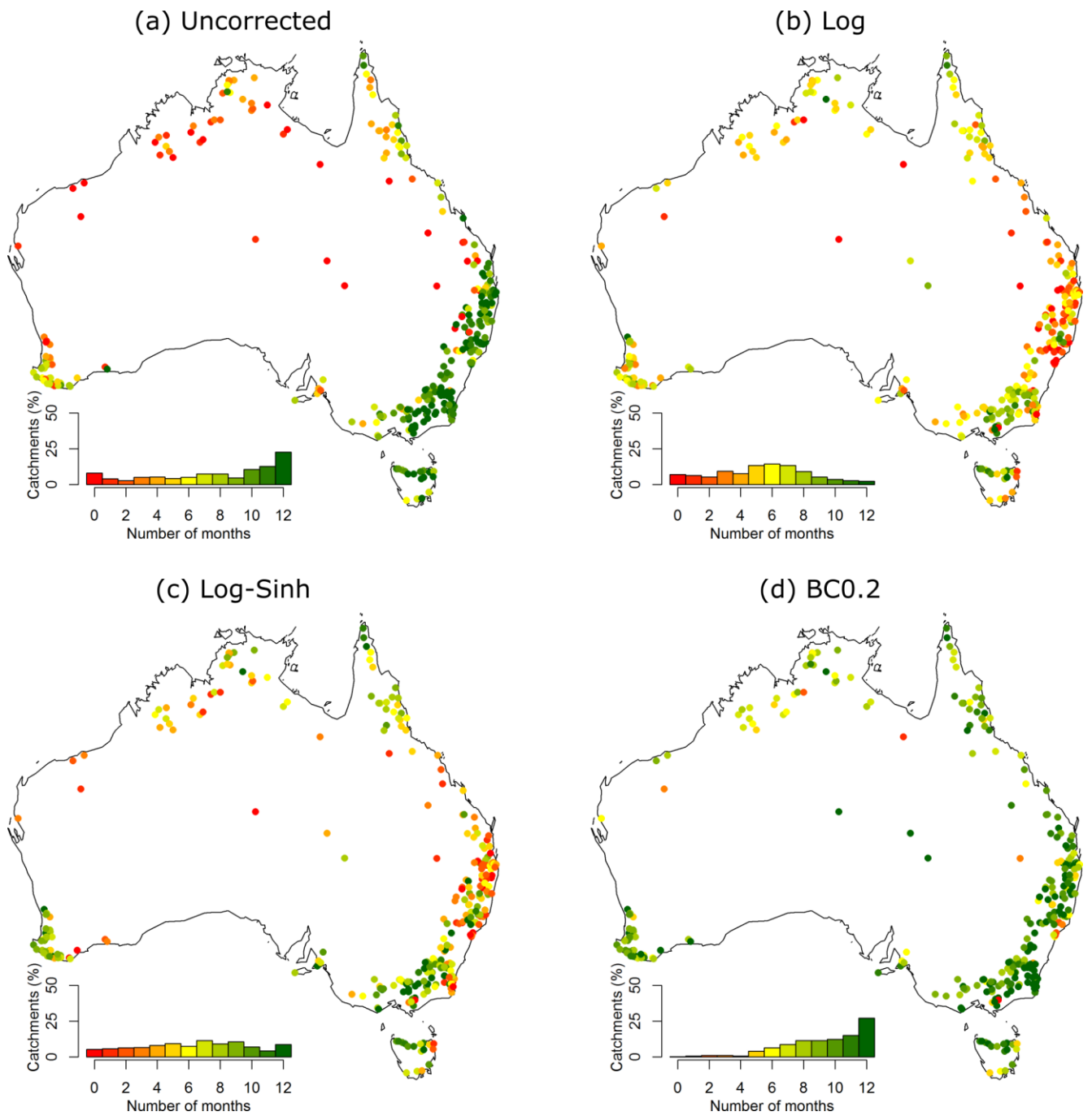
991
992



993
994

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

1000



1001
1002

1003 Figure 11: Summary skill of seasonal forecasts obtained using the Log, Log-Sinh and BC0.2 schemes
1004 across 300 Australian catchments. See

1005 Figure 10 caption for details.