

1 **The challenge of forecasting high streamflows 1-3 months**
2 **in advance with lagged climate indices in south-east**
3 **Australia**

4 James C. Bennett¹, Q. J. Wang¹, Prafulla Pokhrel¹ and David E. Robertson¹

5 ¹CSIRO Land and Water, Graham Road, Highett, Victoria, Australia 3190

6

7 Correspondence to:

8 James. C. Bennett

9 T: +61 3 9252 6229

10 E: james.bennett@csiro.au

11 A: CSIRO Land and Water, Graham Road, Highett, Victoria, Australia, 3190

12

13 **Abstract**

14 Skilful forecasts of high streamflows a month or more in advance are likely to be of
15 considerable benefit to emergency services and the broader community. This is particularly
16 true for mesoscale catchments ($<2000 \text{ km}^2$) with little or no seasonal snow melt, where real-
17 time warning systems are only able to give short notice of impending floods. In this study, we
18 generate forecasts of high streamflows for the coming 1-month and coming 3-month periods
19 using large-scale ocean/atmosphere climate indices and catchment wetness as predictors.
20 Forecasts are generated with a combination of Bayesian joint probability modeling and
21 Bayesian model averaging. High streamflows are defined as maximum single-day
22 streamflows and maximum 5-day streamflows that occur during each 1-month or 3-month
23 forecast period. Skill is clearly evident in the 1-month forecasts of high streamflows.
24 Surprisingly, in several catchments positive skill is also evident in forecasts of large threshold
25 events (exceedance probabilities of 25%) over the next month. Little skill is evident in
26 forecasts of high streamflows for the 3-month period. We show that including lagged climate
27 indices as predictors adds little skill to the forecasts, and thus catchment wetness is by far the
28 most important predictor. Accordingly, we recommend that forecasts may be improved by
29 using accurate estimates of catchment wetness.

30

31 **Keywords:** prediction, seasonal, Bayesian methods, high streamflows

32

33 **1 Introduction**

34 Skilful forecasts of high streamflows a month or more in advance have the potential to
35 improve the management of floods. Flood warnings in Australia are presently derived from
36 event-based forecast models that use real-time streamflow and rainfall observations to
37 forecast floods with typical lead-times from hours to a few days, depending on flood travel
38 time (Elliott et al., 2005). Real-time forecasts offer precise estimates of flood stage, but are
39 only available around the time of the flood itself. This leaves emergency services a narrow
40 window to prepare themselves and the community to mitigate flood impacts, particularly in
41 mesoscale catchments that have little or no seasonal snowmelt. In these catchments flood
42 warning systems can only give warning of floods from hours to one or two days in advance of
43 an event. Ill-preparedness for floods can have serious implications. Pfister (2002) identified
44 poor community preparedness to evacuate as the major cause of citizens' slow (and non-
45 existent) responses to a flood evacuation order issued by emergency services. Australian
46 emergency services rely heavily on volunteers for disaster response (Baxter-Tomkins and
47 Wallace, 2009), and ensuring that sufficient volunteer-labour is available during emergencies
48 is a challenge for flood-response agencies like the State Emergency Services (SES). Medium
49 range forecasts (to forecast horizons of 3 months) of high streamflows are needed to enable
50 both emergency services and the community to be better prepared for floods.

51 This study is a response to a request from the Australian Bureau of Meteorology to explore
52 the skill of real-time high streamflow forecasts at medium range forecast horizons. The
53 Bureau of Meteorology is the lead agency for flood warnings in Australia, and emergency
54 services are important users of these flood warnings. While medium range forecasts of high
55 streamflows cannot hope to be as precise as real-time flood models, forewarning of conditions
56 that could result in large or frequent flooding in the next month or more could allow
57 emergency services to better plan and prepare for the impacts of floods, for example by

58 informing volunteer emergency services personnel of heightened flood risk in the coming
59 month(s).

60 Several studies have described teleconnections between Australian runoff variability and
61 large-scale oceanic and atmospheric climate indices (hereafter, *climate indices*), particularly
62 climate indices describing the El Niño Southern Oscillation (ENSO) (Chiew et al., 1998;
63 Verdon et al., 2004; Schepen et al., 2012a). These teleconnections have been used to produce
64 forecasts of total seasonal streamflows that are skilful relative to forecasts derived from
65 streamflow climatologies (Wang et al., 2009; Piechota et al., 1998; Sharma, 2000). Flood risk
66 in south-east Australia has also been linked to ENSO (Kiem et al., 2003), but despite this no
67 attempt has yet been made to use such a teleconnection to forecast high streamflows in
68 Australia. Attempts to forecast high streamflows a month or more in advance are rarely
69 reported for other continents, and the examples that exist focus on catchments where
70 snowmelt makes a large contribution to seasonal floods (e.g. Kwon et al., 2009; Lindström
71 and Olsson, 2011). Seasonal snow-melt is rarely an important feature of Australian rivers, and
72 accordingly forecasts that rely on indicators of snow-melt have limited application in
73 Australia.

74 The aim of this study is to apply a statistical technique, the Bayesian joint probability
75 modelling approach (BJP), to the problem of forecasting high streamflows in mesoscale
76 catchments over the coming 1-month and 3-month periods. The BJP was developed to
77 forecast seasonal total volumes of streamflows (Wang et al., 2009; Wang and Robertson,
78 2011; Robertson and Wang, 2012) and is now used operationally by the Bureau of
79 Meteorology to issue forecasts for more than 70 sites across Australia (forecasts available at
80 <http://www.bom.gov.au/water/ssf/>). The BJP produces probabilistic streamflow forecasts that
81 are more accurate than climatology, and, importantly, it is able to reliably estimate uncertainty
82 in the streamflow forecasts. Knowledge of the amount of water held in storage in a catchment

83 (in the soil, as ground water, in surface stores, or as snow/ice – collectively, *catchment*
84 *wetness*) often contributes more skill to next-month/next-season forecasts of streamflow than
85 climate forecasts (Shukla and Lettenmaier, 2011; Li et al., 2009; Koster et al., 2010;
86 Mahanama et al., 2012). The BJP is able to use multiple predictors to generate forecasts,
87 meaning forecasts can be constructed from both catchment wetness and predictors of climate.
88 For example, Wang et al. (2009) used the BJP to pair the initial catchment wetness with the
89 southern oscillation index (SOI) to forecast seasonal streamflow totals.

90 A number of sets of predictors can be used to construct different forecast models, and
91 forecasts can be improved by selecting models with the best predictive power (Robertson and
92 Wang, 2012) or by weighting models according to predictive power (Wang et al., 2012a).
93 Wang et al. (2012a) showed that Bayesian model averaging (BMA) outperformed predictor
94 selection methods for merging rainfall forecast models generated with the BJP. In addition,
95 predictor selection can lead to artificially inflated estimates of cross-validation skill if the
96 predictor selection is not included in the cross-validation (DelSole and Shukla, 2009;
97 Robertson and Wang, 2013), a problem that is not present with the BMA method we use in
98 this study.

99 Our study aims to test the ability of the BJP to forecast high streamflows up to three months
100 in advance. To achieve this, we build a set of forecast models with the BJP by combining an
101 estimate of initial catchment wetness with a suite of climate indices derived from oceanic and
102 atmospheric variables. We combine the models with the BMA method described by Wang et
103 al. (2012a) to maximise predictive power.

104 We next describe the study sites and give an overview of the forecast models. This is
105 followed by descriptions of the verification measures we use to demonstrate the reliability and
106 skill of the forecasts. We present the reliability and skill of these forecasts, and discuss the

107 prospects for improving long lead forecasts of high streamflows. We conclude with a
108 summary of the paper.

109 **2 Data and methods**

110 **2.1 Study sites**

111 Forecasts are generated for six catchments in south-east Australia shown in Fig. 1.
112 Characteristics of the six catchments are summarised in Table 1 and Fig. 2. The catchments
113 are selected as they have long (>40 year) streamflow records, are free of diversions or
114 impoundments, and are minimally impacted by human activities. Streamflow data is taken
115 from the quality controlled Catchment Water Yield Estimation Tool (CWYET) dataset (Vaze
116 et al., 2011). All the catchments are of a size we describe as *mesoscale*, with drainage areas
117 between 1000 km² and 2000 km². The catchments are large enough to minimise the influence
118 of highly localised storms (e.g. localised convective storms) on the streamflow records.
119 Conversely, catchments are small enough so that flood travel times extend no more than two
120 days, making it difficult to get advance warning of floods of more than two days with a
121 forecasting model that makes use only of observed rainfalls.

122 The catchments span a range of climate and hydrological conditions. Streamflows in the two
123 north-eastern catchments, the Orara River (ORB) and the Nowendoc River (NOR), are only
124 weakly seasonal, with the highest streamflows occurring in February and March (Fig. 2). The
125 remaining catchments - Abercrombie River (ABH), Murray River (MUR), Mitta Mitta River
126 (MMH) and Tarwin River (TAW) - have more strongly seasonal streamflow regimes, with
127 high streamflows in the austral winter/spring, and low streamflows in the austral summer
128 (Fig. 2). High-elevation areas in the MUR and MMH catchments often receive snowfalls in
129 the Austral winter. However, even in these two catchments the contribution of seasonal
130 snowmelt to streamflows is relatively small.

131 2.2 Forecast model

132 2.2.1 Overview

133 Forecasts are generated on the last day of each month for two periods: the coming month (Jan,
134 Feb, ..., Dec), and the coming three months (JFM, FMA, ..., DJF). We refer to these as 1-
135 month and 3-month forecast periods.

136 Fig. 3 gives a schematic overview of how forecasts are generated. Thirteen forecast models
137 are generated with the BJP method (Fig. 3a) for each forecast period and for each predictand.
138 Forecasts from these individual models are then merged using BMA (Fig. 3b). We now
139 describe the components shown in Fig. 3 in detail.

140 2.2.2 Predictands

141 While we pursue forecasts of large streamflows in a bid to improve information available for
142 the management of floods, we employ the term *high flows* rather than *floods* in this paper.
143 This is because we seek to build monthly statistical models in catchments that often have
144 highly seasonal flow regimes. We define high flows from each month by exceedance
145 probability, and in months where mean flows are low these ‘high’ flows often do not
146 constitute what would be considered flood flows in other months.

147 We investigate two predictands to represent high streamflows:

- 148 1. The maximum 1-day streamflow (mm/d) for each forecast period (Max1D).
- 149 2. The maximum 5-day aggregated streamflow (mm/d averaged across the 5 days)
150 calculated for each forecast period (Max5D).

151 As already noted, neither Max5D nor Max1D is necessarily a large flood. For example, in the
152 catchments with strongly seasonally delineated streamflows, Max5D streamflows in summer
153 can be very low compared to Max5D winter streamflows. In low streamflow months, medians

154 of both Max1D and Max5D streamflows are sometimes not much larger than average
155 monthly streamflows (Fig. 2). For this reason, we also evaluate the performance of the
156 forecasts in terms of probabilities of events exceeding larger thresholds (see Section 2.3.3).
157 The BJP is able to generate forecasts jointly for multiple predictands. In addition to either
158 Max1D or Max5D, we also include total rainfall for the forecast period as a predictand (from
159 the Australian water availability project (AWAP) gridded rainfall dataset; Jones et al., 2009).
160 We jointly forecast rainfall and streamflow because the influence of lagged climate indices on
161 streamflow occurs mainly through rainfall (Robertson and Wang, 2012). Statistically, the
162 correlations between lagged climate indices and rainfall and between rainfall and streamflow
163 tend to be stronger, and thus easier to capture from data, than the correlation directly between
164 lagged climate indices and streamflow. By including rainfall as a co-predictand, the statistical
165 model needs to satisfy three correlations, with the two stronger correlations providing some
166 guidance on sensible values for the weaker correlation.

167 2.2.3 Predictors

168 We use lagged catchment wetness and lagged climate indices as predictors of high
169 streamflows. We approximate catchment wetness with total streamflow in the previous month
170 for both 1-month and 3-month forecast periods. Total streamflow can be a somewhat coarse
171 measure of catchment wetness, and takes no account of differences in catchment wetness
172 stores (e.g. snow *cf.* soil moisture). However, using total streamflow as an estimate of
173 catchment wetness has the virtue of simplicity, and is adequate for this exploratory study.

174 Eleven lagged climate indices are evaluated as potential predictors in this study, and these are
175 listed in Table 2. We select these climate indices as they have been linked to rainfall in south-
176 east Australia. The teleconnection between south-east Australian rainfall and ENSO has been
177 extensively described (e.g. Schepen et al., 2012a; Chiew et al., 1998; Wang et al., 2009)

178 including, as already noted, the link between flooding and ENSO (Kiem et al., 2003). We use
179 five indices to describe ENSO: NINO3, NINO3.4, NINO4, the ENSO Modoki index (EMI)
180 (Ashok et al., 2007) and the southern oscillation index (SOI) (Troup, 1965). The influence of
181 Indian Ocean sea surface temperatures has also been linked to rainfall in south-east Australia,
182 with the teleconnection being most evident in winter months (Verdon and Franks, 2005;
183 Schepen et al., 2012a; Ashok et al., 2003). We use four Indian Ocean indices as predictors:
184 the Indian Ocean west pole index (WPI), east pole index (EPI) and dipole mode index (DMI)
185 (Saji et al., 1999), as well as the Indonesia index (II) (Verdon and Franks, 2005). Finally,
186 extra-tropical sea surface temperatures and atmospheric features along Australia's east coast
187 have been linked to south-east Australian rainfall (Murphy and Timbal, 2008; Risbey et al.,
188 2009; Pook et al., 2006). We use the Tasman Sea index (TSI) (Murphy and Timbal, 2008) and
189 an index of atmospheric blocking (BI140) (Risbey et al., 2009) to represent extra-tropical
190 climatic features. The teleconnection between lagged atmospheric climate indices (e.g., the
191 Antarctic Oscillation index describing the Southern Annular Mode; Schepen et al., 2012a) and
192 Australian seasonal precipitation is often weak, as they show little persistence in comparison
193 to SST-derived indices. We note that Schepen et al. (2012a) found no evidence of a
194 relationship of lagged B140 and TSI with mean rainfall in any season. It is therefore unlikely
195 that lagged TSI or B140 will contribute skill to high streamflow forecasts, however we have
196 included them in case they have a relationship with high rainfall events. Atmospheric
197 blocking, for example, has been correlated with larger rain storms (Pook et al., 2006).

198 We have not considered using multiple climate indices as joint predictors, which may
199 describe the effects of interactions between climate indices on high streamflows. Some
200 studies suggest that these interactions may be important in understanding concurrent
201 relationships (e.g. Kiem et al., 2003), however results from our previous work demonstrates
202 that adding a second joint predictor does not result in any improvement in forecast skill of

203 seasonal total rainfalls or streamflows when using lagged climate indices (Robertson and
204 Wang, 2012; Wang et al., 2012a).

205 Sea surface temperature climate indices are derived from the National Center for Atmospheric
206 Research (NCAR) Extended Reconstruction of Sea Surface Temperature version 3 (Smith et
207 al., 2008). B140 is derived from the National Centers for Environmental Prediction (NCEP)–
208 NCAR reanalysis data (Kalnay et al., 1996). SOI is sourced from the Australian Bureau of
209 Meteorology (BOM).

210 Mean monthly values of each climate index for the previous month are used for both 1-month
211 and 3-month forecasts; accordingly we refer to these as *lagged* climate indices. Schepen et al.
212 (2012a) showed that teleconnections between rainfall and lagged climate indices are strongest
213 at short lags, and for this study we investigate only climate indices lagged by one month to
214 establish forecast models. For example, for a 1-month forecast for June we use catchment
215 wetness and NINO3 calculated for May as predictors, while for a 3-month forecast for
216 January-February-March we use predictors calculated for December.

217 Catchment wetness is combined with each of the 11 climate indices to create 11 forecast
218 models for each predictand and for each forecast period. In addition, one forecast model is
219 developed using only catchment wetness as a predictor, and one forecast model is developed
220 based only on climatology (using no predictors). This gives a total of 13 forecast models for
221 each predictand and for each forecast period.

222 While the effect of snow on the two alpine catchments (MUR and MMH) is expected to be
223 small, we investigated the use of snow accumulation as a predictor for these two snow-
224 affected catchments. Including snow accumulation as a predictor in these two catchments
225 resulted in no increase in forecast skill and is not presented here.

226 2.2.4 Bayesian joint probability modelling

227 The BJP is used to generate the 13 individual forecast models for each predictand and each
228 forecast period (Fig. 3a), which we call *BJP forecast models*. Detailed mathematical
229 formulations of the BJP are given by Wang et al. (2009), Wang and Robertson (2011) and
230 Robertson and Wang (2012). In summary, the BJP is implemented as follows:

- 231 1. Predictands and predictors are transformed to normalise their distributions and
232 stabilise their variances. Streamflow and rainfall are transformed with a log-sinh
233 transform (Wang et al., 2012b), and climate indices are transformed with the Yeo-
234 Johnson transform (Yeo and Johnson, 2000).
- 235 2. We assume that the set of transformed predictors and predictands can be described by
236 a joint probability distribution – in this case a multivariate normal distribution.
- 237 3. The parameters of the log-sinh transform, the Yeo-Johnson transform, and the
238 multivariate normal distribution are inferred jointly. Parameter inference is performed
239 with Bayesian methods and Markov chain Monte Carlo (MCMC) sampling. Taken
240 together, the parameters of the log-sinh transform, the Yeo-Johnson transform and the
241 multivariate normal distribution define the statistical relationship between predictors
242 and predictands, and allow us to generate forecasts.

243 Mathematically, if predictors are given by vector $y(1)$ and predictands by vector $y(2)$, the
244 probabilistic forecast is given by

$$f[y(2)|y(1)] = p[y(2)|y(1); Y_{OBS}, M] = \int p[y(2)|y(1); \theta] \cdot p[\theta|Y_{OBS}, M] \cdot d\theta \quad (1)$$

245 where M is the model used, and Y_{OBS} contains the historical data of both the predictors and the
246 predictands used for model inference. θ is the vector of parameters for the log-sinh transform,
247 the Yeo-Johnson transform, and the multivariate normal distribution.

248 2.2.5 Bayesian model averaging

249 Forecasts from the thirteen BJP forecast models are merged with BMA to produce one *BJP-*
250 *BMA forecast* for each predictand and for each forecast period (Fig. 3b). The BMA method
251 we use is described in detail by Wang et al. (2012a). For a set of models M_k , $k=1, 2, \dots, K$,
252 each model is assigned a weight, w_k . The forecasts are then merged by:

$$f_{BMA}(y(2)|y(1)) = \sum_{k=1}^K w_k f_k(y(2)|y(1)) \quad (2)$$

253 We calculate w_k by maximizing the posterior distribution of the weights, which is proportional
254 to:

$$A = \prod_{k=1}^K (w_k)^{\alpha-1} \prod_{t=1}^T \sum_{k=1}^K w_k \cdot p(y_{OBS}^t(2) | y_{OBS}^t(1); Y_{OBS}^{(t)}, M_k) \quad (3)$$

255 where α is the concentration parameter, $y_{OBS}^t(1)$ and $y_{OBS}^t(2)$ are the predictors and predictands
256 for events $t=1, \dots, T$, and $Y_{OBS}^{(t)}$ is a matrix containing observed values of predictors and
257 predictands for all the events except event t . $\prod_{k=1}^K (w_k)^{\alpha-1}$ is from the symmetric Dirichlet
258 prior distribution used by Wang et al. (2012a). We use α values greater than 1 to distribute
259 weights more evenly among models, which helps to stabilise the weights when there is
260 significant sampling variability. Specifically, $\alpha=1+a/K$ with $a=1$. The remainder of the right
261 side of Eq. 3 is the cross-validation likelihood function. By using the cross-validation
262 likelihood function, we base each model weight on the predictive power of the model, rather
263 than on the fitting ability of the model. A is maximised with an iterative expectation-
264 maximization (EM) algorithm, as described by Wang et al. (2012a).

265 2.3 Forecast verification

266 Forecasts are verified using leave-one-out cross validation. Forecasts for events in year
267 $t=1, 2, \dots, n$ are generated from all available historical data except those at year t . For each

268 forecast variable y , this produces a series of forecast cumulative probability distributions
269 $y^f \sim F^f(y^f)$. Forecasts are then verified against observations y^f_{OBS} .
270 Leave-one-out cross validation ensures that a forecast model is not validated against data used
271 to build that model. We note that in this approach we use data after the forecast date to build
272 the forecast model, data which would not be available to build operational real-time forecast
273 models. The purpose of cross validation is to get an indication of model performance for
274 future events. For future events, we would use all historical events to establish the model. The
275 length of record used in model establishment in cross-validation is similar to (more precisely
276 just short of) the full record length. In this sense, cross-validation gives a good indication of
277 the skill of a true implementation for the future events.
278 Verifying the probabilistic forecasts is not straightforward, particularly when the aim is to
279 forecast rare events. Here we evaluate forecast reliability to demonstrate that the probabilistic
280 forecasts are neither too confident nor underconfident. We then assess forecast accuracy using
281 three skill scores. We now describe each of the verification measures in detail.

282 2.3.1 Forecast reliability

283 For probabilistic forecasts to be meaningful, we must first demonstrate that the forecast
284 probability distributions are reliable; that is, the uncertainty in the forecasts is reliably
285 represented, and thus the forecast distributions are neither too wide (not confident enough)
286 nor too narrow (overconfident). To achieve this, we present reliability diagrams. A reliability
287 diagram plots the observed frequency against the forecast probability and shows how well the
288 predicted probability of an event corresponds to its observed frequency (Wilks, 1995). We
289 present reliability diagrams calculated from events that are larger than the 50% exceedance
290 probability threshold of Max1D and Max5D streamflows.

291 2.3.2 Overall forecast accuracy: root mean square error in probability

292 The root mean square error in probability (RMSEP) works on the principle that if forecast and
293 observed values are of similar exceedance probabilities then the forecast should be rewarded,
294 even if the magnitudes of observed and forecast values are quite different (Wang and
295 Robertson, 2011). RMSEP is calculated as follows:

- 296 1. We represent the observed historical distribution (climatology), y , in the form of non-
297 exceedance probability, $F_{CLI}(y)$.
- 298 2. For events $t=1, 2, \dots, n$, we take the median of the forecast distribution, y_{MED}^t .
- 299 3. RMSEP is then calculated as

$$RMSEP = \left[\frac{1}{n} \sum_{t=1}^n (F_{CLI}(y_{MED}^t) - F_{CLI}(y_{OBS}^t))^2 \right]^{\frac{1}{2}} \quad (4)$$

- 300 4. We calculate $RMSEP_{REF}$ by substituting the forecast median, y_{MED}^t , in Eq. 4 with the
301 climatology median. We then calculate the RMSEP skill score:

$$SS_{RMSEP} = \frac{RMSEP_{REF} - RMSEP}{RMSEP_{REF}} \quad (5)$$

302 RMSEP (eq. 4) demonstrates the ability of the model to forecast the rank of a given event,
303 ranked in relation to historical events (i.e., the ability to forecast an event's place on a
304 cumulative distribution function generated from historical data). While this does not
305 necessarily give an indication of how well the model is able to forecast the magnitude of an
306 event, the ability to forecast an event's rank is likely to be very useful to users of the forecast,
307 who could categorise an event as, for example, 'likely to exceed the 50 percentile of high
308 flows' or similar. SS_{RMSEP} (eq. 5) measures the ability of the forecasts to outperform a naive
309 climatology forecast.

310 In addition, we calculate SS_{RMSEP} with $RMSEP_{REF}$ represented by the BJP forecast generated
311 with only catchment wetness as a predictor (i.e., no climate information is used to generate
312 $RMSEP_{REF}$). This allows us to show the relative contribution of catchment wetness and
313 climate indices to forecast skill.

314 2.3.3 Accuracy of forecasts for large threshold events

315 For a given month, we consider a subset of larger ‘high’ streamflows to assess forecast
316 performance. These larger streamflows are defined as having exceedance probabilities of 50%
317 (Q_{50}), 25% (Q_{25}) and 10% (Q_{10}) for observed Max1D and Max5D. (These streamflows
318 approximately correspond to annual exceedance probabilities (AEP) of 1:2 AEP, 1:4 AEP and
319 1:10 AEP. To keep the study as simple as possible, we have defined larger events on the basis
320 of empirical exceedance probabilities rather than fitting an extreme value distribution, so we
321 continue to refer to large streamflows in terms of exceedance probabilities.) We treat these
322 large streamflows as thresholds (we term them *large threshold events*), and measure forecast
323 skill by comparing the forecast probability of exceeding a large threshold event with the
324 corresponding observation. Q_{50} , Q_{25} , and Q_{10} thresholds are shown for 1-Month Max1D and
325 Max5D streamflows are shown in Fig. 2.

326 Use of multiple skill scores is recommended to demonstrate robustness in the results (e.g.
327 Cloke and Pappenberger, 2008). We use two measures of skill to verify forecasts at larger
328 streamflow thresholds: the Brier Score and the log-likelihood ratio.

329 **Brier Score**

330 The Brier score has been a staple for the verification of probabilistic forecasts since it was
331 proposed by Brier (1950). We use the Brier score to verify forecasts of larger streamflows in
332 order that our study can be compared to others.

333 Given forecast distributions y^t at events $t=1, 2, \dots, n$, and streamflow thresholds Q_P , with
 334 exceedance probabilities $P=50\%, 25\%, 10\%$, the forecast is presented as the probability of
 335 exceeding the streamflow threshold:

$$1 - F^t = p(y^t > Q_P) \quad (6)$$

336 We calculate the Brier score as:

$$BS = \frac{1}{n} \sum_{t=1}^n (1 - F^t - O^t) \quad (7)$$

337 where O^t takes the value of 1 if the threshold is exceeded, and 0 if it is not exceeded. We
 338 calculate BS_{REF} by substituting F^t with a forecast calculated from climatology, F^t_{REF} . We then
 339 calculate the Brier skill score:

$$SS_{BS} = \frac{BS_{REF} - BS}{BS_{REF}} \quad (8)$$

340 **Log-likelihood ratio**

341 The Brier score has been subject to criticism, particularly for producing unintuitive results for
 342 rare (and in our case, large) events when assessing very sharp forecasts (i.e., forecast
 343 probabilities of 100% or 0%) (Jewson, 2008; Benedetti, 2010). We adopt the
 344 recommendations of Benedetti (2010) and Jewson (2008), who both advocate variations on
 345 the likelihood to assess probabilistic forecasts. We term this measure the log-likelihood ratio
 346 (LLR).

347 The LLR is based on the likelihood ratio described by Jewson (2008). For all exceedance
 348 forecasts $1-F^t$, let all the cases of t where $1-F^t$ exceeds a streamflow threshold Q be given by
 349 the set A , and all cases of t where the streamflow threshold is not exceeded be given by B . The
 350 log-likelihood for a forecast is calculated by:

$$LL = \log_e \left(\prod_A (1 - F^t) \prod_B F^t \right) \quad (9)$$

351 The log-likelihood of the reference forecast, LL_{REF} , is calculated by substituting F^t_{REF} (again,
352 based on climatology) for F^t in Eq. 9. The LLR is then calculated by:

$$LLR = LL - LL_{REF} \quad (10)$$

353 The LLR differs from skill scores like RMSEP or the Brier score in that it does not show
354 proportional improvement over a reference forecast on a normalised scale (often $-\infty\%$ -
355 100%), making direct comparisons to other skill scores difficult. However, the LLR is
356 essentially identical to the natural logarithm of the pseudo Bayes factor ($\log_e(\text{PsBF})$)
357 presented by Robertson and Wang (2012) and Schepen et al. (2012a). Robertson and Wang
358 (2012) showed that values of the $\log_e(\text{PsBF})$ up to 2 are indistinguishable from statistical
359 noise, while there is a 95% chance that the relationship between a forecast model and
360 observations is true if the $\log_e(\text{PsBF})$ is greater than 4. We adopt the qualitative categories for
361 the LLR presented by Schepen et al. (2012a) for our study: little evidence of skill where
362 $LLR < 2$; positive evidence of skill where $2 < LLR < 4$; strong evidence of skill where $4 < LLR < 6$;
363 very strong evidence of skill where $LLR > 6$.

364 **3 Results**

365 **3.1 Suitability of BJP for modelling high streamflows**

366 The log-sinh transform used to normalise streamflows has been shown to be well-suited to
367 hydrological data in general (Wang et al., 2012b; Del Giudice et al., 2013), but its ability to
368 adequately describe high streamflows needs to be established. In Fig. 4 we show the log-sinh
369 transformed normal distributions fitted to observed Max1D values for two example months,
370 February and September (other months give very similar results). These two months represent
371 low and high streamflow regimes: February is a month of low mean streamflows in MMH,
372 MUR, ABH and TAW, and a month of high mean streamflows in ORB and NOR, while
373 September is a month of high mean streamflows in MMH, MUR, ABH and TAW and a

374 month of low mean streamflows in ORB and NOR. In general, the assumed log-sinh
375 transformed normal distributions appear to adequately represent the marginal distribution of
376 observations. Almost all observations fall within the confidence bounds of the fitted
377 distributions, including large Max1D events. The log-sinh transformed normal distributions
378 represent observed events well even in catchments with highly variable streamflows, such as
379 ORB and ABH. In summary, the log-sinh transform is flexible enough to normalise the events
380 we are attempting to forecast.

381 **3.2 Forecast reliability**

382 In general, forecast uncertainty is reliably represented by the forecasts after cross-validation.
383 Fig. 5 shows reliability diagrams for the NOR and MUR catchments for Max1D 1-Month
384 forecasts (the other catchments, not shown, produce similar results). In these diagrams,
385 forecast probabilities are divided into five bins (see inserts). The [0.05, 0.95] uncertainty
386 interval of the observed relative frequency is calculated through bootstrap resampling of the
387 forecasts and observed streamflows. For the majority of forecast probability ranges, the
388 uncertainty interval of the observed relative frequency intersects the theoretical 1:1 line,
389 indicating that the forecasts of high streamflows are reliable. Similar results are obtained for
390 the other catchments for all predictands and forecast periods (not shown). These results
391 support the findings of Wang et al. (2009) and Wang and Robertson (2011), who showed the
392 BJP produces reliable forecasts of seasonal streamflows.

393 **3.3 Overall forecast skill**

394 Fig. 6 shows BJP-BMA cross-validated hindcasts of Max1D for an example 20-year period
395 for all catchments. Visual inspection of the hindcasts shows that the credible prediction
396 intervals largely encompass the range of observations. In catchments with strongly seasonal

397 streamflows (e.g. MUR, MMH), the mean of the ensemble forecast often gives realistic
398 predictions of Max1D streamflows during seasons of high streamflows. Accuracy of forecasts
399 in more variable catchments (e.g. NOR, ABH) is much more difficult to ascertain from these
400 time series, and we now turn to formal measures of skill to assess these.

401 RMSEP skill scores are positive for Max5D forecasts for the 1-month forecast period for most
402 months and catchments (Fig. 7b). Skill in Max5D 1-month forecasts is particularly strong in
403 the winter-spring months (June-November). Skill in Max1D 1-month forecasts is generally
404 lower than for Max5D 1-month forecasts (Fig. 7a, 7b). Max1D streamflows are inherently
405 more variable than Max5D streamflows, as Max5D streamflows are smoothed by the greater
406 number of data included in their calculation. This makes forecasting Max1D streamflows
407 more challenging. Nonetheless, RMSEP skill scores for Max1D 1-month forecasts are
408 positive for most catchments and seasons (Fig. 7a). Max1D 1-month forecast skill is strongest
409 in the winter-spring months. For the 3-month forecast period, RMSEP scores are generally
410 lower for both Max1D and Max5D forecasts, although positive skill scores occur in winter-
411 spring for the MUR, MMH, and ABH catchments, and the NOR catchment shows skill
412 intermittently through the year (Fig. 7c, 7d).

413 The reason for the reduced performance of the 3-month forecasts becomes evident when we
414 review the contribution of climate indices to forecast skill. Fig. 8 shows RMSEP skill scores
415 calculated relative to BJP forecasts generated using only streamflow as a predictor. The plot
416 shows the skill gained by the inclusion of climate indices for Max1D 1-month forecasts. Fig.
417 8 shows that almost no skill is gained in any month or catchment by including climate indices,
418 meaning the forecasts depend heavily on catchment wetness for skill. Results are similar for
419 Max5D (not shown). This finding is also supported by Robertson and Wang (2013), who
420 found that climate indices made only weak contributions to the skill of forecasts of seasonal
421 streamflow totals in the MMH and MUR catchments. The contribution of catchment wetness

422 to forecast skill declines over longer forecast periods (Mahanama et al., 2012; Shukla and
423 Lettenmaier, 2011; Li et al., 2009). Thus forecasts for longer periods are less accurate than for
424 shorter forecast periods. This effect is also evident in individual catchments. The TAW
425 catchment, for example, has the lowest autocorrelation of monthly streamflows of the six
426 catchments (not shown), and forecasts for this catchment show poor skill in relation to
427 streamflow climatology.

428 Nonetheless, 3-month forecasts can be skilful in certain catchments at times of the year when
429 the influence of catchment wetness on high streamflows is strong. The influence of catchment
430 wetness on streamflows is generally strongest on the receding limb of the annual hydrograph
431 (Robertson and Wang, 2013). For the ORB and NOR catchments the annual hydrograph
432 recedes in March-May, while in the ABH, MMH and MUR catchments the annual
433 hydrograph recedes in August-November. This results in positive RMSEP skill scores for 3-
434 month forecasts of these catchments during these months (Fig. 7c, 7d).

435 Overall, RMSEP generally shows positive skill scores for 1-month forecasts for both Max1D
436 and Max5D streamflows, while 3-month forecasts are substantially less skilful. However, the
437 positive RMSEP skill scores may be the result of good agreement of forecasts with lower
438 'high' streamflows, and not reflect forecasts at larger streamflows. We now turn to forecast
439 skill at higher streamflows to determine the size of streamflows for which forecasts are
440 skilful.

441 **3.4 Forecast skill for large threshold events**

442 In general, forecast skill declines as streamflows get larger (Figs. 9-12). Brier scores show
443 more instances of positive skill than LLR scores, particularly for streamflows larger than Q_{10} .
444 Because the Brier score has known problems with infrequent events (Benedetti, 2010), we
445 focus on the LLR score to discuss forecast skill at larger streamflows.

446 Substantial skill is evident in forecasts where observed Max1D streamflows are larger than
447 Q_{50} for 1-month forecasts, in both the Brier score (Fig. 9) and the LLR (Fig. 10). LLR scores
448 are higher for Max5D streamflows than for Max1D streamflows, and the highest LLR scores
449 generally occur in July-November. Skill is not related to seasonal changes in high or low
450 Max1D/Max5D streamflows. The ARB, MUR, MMH and catchments show high skill during
451 months of high streamflow (winter-spring, Fig. 10, Fig. 2) while the ORB and NOR
452 catchments only exhibit skill during months of low streamflow (Jul-Nov, Fig. 10, Fig. 2). As
453 with the RMSEP scores, the TAW catchment shows the lowest skill. Four of the six
454 catchments show positive LLR scores in 6 or more months of the year for 1-month forecasts
455 of Max5D streamflows above Q_{25} (Fig. 10). For Max1D streamflows greater than Q_{25} , three
456 catchments show positive LLR scores in six or more months of the year (Fig. 10). Little skill
457 is evident in any catchment or season for either Max1D or Max5D streamflows above Q_{10} .
458 Skill for 3-month forecasts of larger streamflows is generally low (Figs. 11, 12). Except for
459 one catchment (MUR), catchments show little forecast skill in the majority of months for any
460 of the streamflow thresholds tested for either Max1D or Max5D streamflows. We find
461 positive skill scores for 3-month forecasts in the MUR catchment of Max5D streamflows
462 above Q_{50} and Q_{25} for six or more months, and also for Max1D streamflows above Q_{50} (Fig.
463 12). Indeed, forecasts for MUR performed best in most measures and skill scores. It is not
464 clear why this should be so. MUR receives reliable rainfall in the winter and spring, resulting
465 in relatively low variability and strong autocorrelation in monthly streamflows. However
466 these characteristics also apply to the nearby MMH catchment, for which forecasts perform
467 no better than for ABH, ORB or NOR in a number of measures (e.g. Fig. 10).
468 Overall, forecast skill is positive to very strong for 1-month exceedance forecasts of
469 streamflows exceeding Q_{50} for a majority of months in all but the TAW catchment. Skill is not
470 related to seasonal cycles of high and low streamflows. Positive skill scores are also found in

471 several catchments for 1-month exceedance forecasts of streamflows exceeding Q_{25} . The
472 remaining large streamflow forecasts tested here showed little skill in most catchments.

473 **4 Discussion**

474 RMSEP skill scores reported here show the 1-month forecasts to be superior to climatology in
475 forecasting high streamflows. Further, the skill in forecasts is not limited to the lowest of the
476 'high' streamflows - forecasts of the probability of exceeding Q_{50} Max1D streamflows one
477 month in advance show robust skill in a number of catchments. We note, however, that the
478 Q_{50} Max1D streamflows are still not necessarily very large streamflows. Skill in forecasting
479 large threshold events in two catchments, ORB and NOR, is restricted to months where 'high'
480 streamflows are small, and in which damaging floods are unlikely to occur. Conversely, skill
481 in the MUR, ABH and MMH catchments is evident during periods of high streamflow.
482 Accordingly, forecast skill in these catchments may be valuable to the Bureau of Meteorology
483 when they are seeking to answer more general questions about the risks of high streamflows
484 in a coming month. We note that the usefulness of the forecast is likely to vary with
485 catchment in any case, both because forecast skill varies between catchments and because the
486 prospect of flood damage varies greatly between catchments (i.e., in one catchment a common
487 high streamflow event may damage property or have other deleterious impacts, in another
488 catchment large floods may be of little consequence).

489 The 1-month forecasts rely heavily on catchment wetness for skill. This supports the many
490 studies that have demonstrated the preeminent contribution of catchment wetness to the skill
491 of seasonal streamflow forecasts for catchments (or seasons) where seasonal snow-melt does
492 not occur (e.g. Mahanama et al., 2012; Shukla and Lettenmaier, 2011; Li et al., 2009; Koster
493 et al., 2010; Robertson and Wang, 2013). Accordingly, improving estimates of catchment
494 wetness is likely to be a simple way of improving forecasts. Accumulated streamflow for a

495 month can be a poor measure of catchment wetness. For example, a high value of total
496 streamflow may be caused by a single intense rainfall event that causes infiltration-excess
497 overland flow, resulting in a large streamflow but little infiltration. In this example the
498 catchment wetness is overestimated by total streamflow. Catchment wetness can be modelled
499 more effectively for forecasting with so-called ‘dynamical’ approaches (Rosenberg et al.,
500 2011; Robertson et al., 2013a) that use soil-moisture accounting models (e.g. conceptual
501 rainfall-runoff models forced by observed rainfall and evaporation) to improve estimates of
502 catchment wetness and thereby improve forecasts.

503 The ability of the BJP-BMA models to forecast high streamflows a month or more in advance
504 is limited by knowledge of climate during the forecast period. This problem is not likely to be
505 easily surmountable. The high variability of larger rainfall events makes their prediction
506 inherently difficult. In addition, climate indices that have the potential to forecast particular
507 types of rain-bearing weather patterns may have little persistence from month to month. This
508 is particularly so for climate indices calculated from atmospheric variables, which tend to be
509 less persistent than oceanic variables. For example, we have used the atmospheric blocking
510 index (B140, see Table 2) to attempt to account for atmospheric blocking and associated
511 cutoff lows in our forecasts. Cutoff lows associated with atmospheric blocking bring a
512 substantial proportion of rainfall to south-east Australia (Pook et al., 2006), and may
513 counteract the drying associated with very strong El Niño years (Brown et al., 2009).
514 However, we find that B140 adds little skill to forecasts of high streamflows, supporting
515 Schepen et al. (2012a) who showed that lagged B140 had no significant statistical relationship
516 to mean rainfall anywhere in Australia. Similarly, this would very likely apply to other
517 atmospheric indices, e.g. those used to describe the Southern Annular Mode or the
518 Subtropical Ridge of high pressure (position or intensity).

519 As we noted in the introduction, several studies have shown positive relationships between
520 climate indices and streamflow/rainfall in south-east Australia. However, our work shows that
521 the benefit of using lagged climate indices to forecast high streamflows in south-east
522 Australia is negligible. This can be explained in four ways:

523 1. Many studies examine teleconnection between concurrent climate indices and
524 streamflow/rainfall (e.g. Verdon and Franks, 2005; Ashok et al., 2003; Pook et al., 2006).

525 The teleconnection between lagged climate indices and rainfall may be weaker than for
526 concurrent indices as implied by the often weak relationships between lagged climate
527 indices and Australian rainfall found by Schepen et al. (2012a).

528 2. Even if a significant teleconnection exists between a lagged climate index and high
529 streamflows, this information may still not contribute skill to forecasts of high
530 streamflows when we include catchment wetness as a predictor because:

531 a. even if the teleconnection between high rainfalls and lagged climate indices is
532 strong, the influence of catchment wetness on high streamflows is so much more
533 powerful that the predictive information provided by lagged climate indices is
534 rendered negligible;

535 b. the catchment wetness predictor implicitly contains information about the current
536 state of the climate (e.g., a very wet October), and any information provided by
537 lagged indices may be subsumed by the climate information implicit in catchment
538 wetness.

539 3. Even in areas where lagged climate indices show a significant teleconnection to seasonal
540 rainfalls (Schepen et al., 2012a), the high variability of large rainfalls associated with high
541 streamflows means that any positive relationships that have been shown to exist between
542 lagged climate indices and seasonal rainfall totals may not apply to high rainfall events.

543 4. Some studies (e.g. Kiem et al., 2003) use an index describing the Interdecadal Pacific
544 Oscillation (IPO) to relate rainfall/streamflow to climate indices. If we limit our
545 assessment of forecasts only to periods where IPO was in the negative phase, it is possible
546 that ENSO SST indices may add more skill to the forecasts (as suggested by Kiem et al.,
547 2003). However, we sought to assess forecast skill in the context of generating forecasts in
548 real-time. Describing the IPO is not particularly useful for real-time forecasting because it
549 is only possible to define an IPO phase with certainty in retrospect (although informed
550 speculation about the present IPO phase is possible; see, e.g., Cai and van Rensch, 2012).
551 That is, it is often not possible to know with certainty which IPO phase we are in at the
552 present time, so it cannot be used to inform real-time forecasts.

553 Using conceptual rainfall runoff models forced by rainfall forecasts from dynamical climate
554 models to forecast high streamflows at long lead times is an attractive alternative to the
555 statistical models we have presented here. Statistical models require large volumes of data to
556 characterise relationships between predictors and predictands, and this is particularly
557 important when forecasting rare events. If dynamical climate and hydrological processes can
558 be accurately simulated, fewer data may be required to generate skilful forecasts. Further,
559 dynamical climate models should, in theory, be able to account for complex interactions
560 between different climate drivers, which may influence rainfall. At present dynamical climate
561 models do not necessarily exhibit more skill than statistical forecasts of seasonal precipitation
562 (e.g. Schepen et al., 2012b). Future improvements in dynamical climate models used for
563 forecasting weeks to months advance (e.g. Marshall et al., 2011) may ultimately improve
564 forecasts of high rainfalls. In addition, we note that the skill of statistical forecasts may
565 complement that of dynamical rainfall forecasts (e.g. the statistical rainfall forecasts may
566 exhibit skill in different seasons or locations to dynamical forecasts; Schepen et al., 2012b),
567 and that merging forecasts of high rainfalls from dynamical and statistical models may

568 improve overall skill. Using climate indices derived from SST forecasts from coupled ocean-
569 atmosphere dynamical climate models shows promise in improving forecasts of monthly
570 rainfall totals at lead-times of more than six months (Hawthorne et al., 2013), and avoids the
571 use of lagged climate indices for forecasting.

572 Our forecast method could be adapted to catchments in different regions by including
573 predictors that are relevant to a given region. In colder regions, seasonal snow melt has been
574 shown to be a very important predictor of seasonal streamflows (e.g. Mahanama et al., 2012),
575 and indicators of future snowmelt (e.g. temperature) could be included as predictors in this
576 model. In addition, climate indices that are important to a given region may also be included,
577 although their utility for forecasting high streamflows may be negligible, as we have shown
578 here.

579 The high streamflow forecasts we have developed here may be bolstered in future by the
580 inclusion of Numerical Weather Prediction (NWP) models in hydrological forecasting. The
581 Australian Bureau of Meteorology does not presently use NWP forecasts to quantify flood
582 forecasts, although they are used qualitatively to inform flood warnings (Elliott et al., 2005).
583 Very high resolution NWP forecasts have been shown to improve flood forecasts (Roberts et
584 al., 2008). At present, however, NWP forecasts are skilful only for a few days (typically <6
585 days); and even skilful NWP forecasts are often not accurate enough for use in hydrological
586 forecasting systems, even in catchments substantially larger than those tested here (Cloke and
587 Pappenberger, 2009; Shrestha et al., 2013; Cuo et al., 2011). As NWP models and post-
588 processing of NWP forecasts improve (e.g. Robertson et al., 2013b), NWP forecasts may
589 complement the simpler forecasts we have generated in this study.

590 **5 Summary and conclusions**

591 We have explored the ability of existing statistical forecasting methods to produce forecasts
592 for high streamflows for the coming month and the coming three months. Forecast models are
593 built from a combination of climate predictors and catchment wetness. Models are
594 constructed with a Bayesian joint probability method, and the models are then weighted based
595 on their predictive power using Bayesian model averaging.

596 Skill is clearly evident in forecasts of high streamflows for the coming 1-month period.
597 Forecasts of larger events, including maximum 1-day streamflows of exceedance probabilities
598 as low as 25%, are also skilful in comparison to long-term climatologies. Our 1-month high
599 streamflow forecasts have the potential to complement existing real-time flood warnings
600 currently used in Australia, to give emergency services and the community more warning of
601 impending high streamflows.

602 Almost all forecast skill derives from the catchment wetness predictor. If the forecasts are to
603 be extended to additional catchments, they are likely to be poor in catchments that have little
604 month-to-month memory in streamflows. Forecasts in skilful catchments may be improved
605 somewhat by using more refined estimates of catchment wetness.

606 We find substantially lower skill in forecasts of high streamflows for the coming 3-month
607 period. The influence of catchment wetness on streamflows diminishes over longer periods,
608 and climate predictors add little skill to the forecasts. Future improvements in forecasts of
609 extreme rainfalls from dynamical climate models may be able to improve longer range
610 forecasts of high streamflows.

611 **Acknowledgements**

612 This research has been supported by the Water Information Research and Development
613 Alliance between the Australian Bureau of Meteorology and CSIRO Water for a Healthy

614 Country Flagship. Thanks to Yong Song (CSIRO Land and Water), Christopher J. White
615 (Bureau of Meteorology) and Senlin Zhou (Bureau of Meteorology) for their comments on
616 earlier drafts. Thanks to Ben Livneh and two anonymous reviewers for comments that have
617 improved this manuscript.

618 **References**

- 619 Ashok, K., Guan, Z., and Yamagata, T.: Influence of the Indian Ocean dipole on the
620 Australian winter rainfall, *Geophysical Research Letters* 30, 1821, 10.1029/2003GL017926,
621 2003.
- 622 Ashok, K., Nakamura, H., and Yamagata, T.: Impacts of ENSO and Indian Ocean dipole
623 events on the southern hemisphere storm-track activity during austral winter, *Journal of*
624 *Climate*, 20, 3147-3163, 10.1175/jcli4155.1, 2007.
- 625 Baxter-Tomkins, T., and Wallace, M.: Recruitment and retention of volunteers in emergency
626 services, *Australian Journal on Volunteering*, 14, 1–11, 2009.
- 627 Benedetti, R.: Scoring rules for forecast verification, *Monthly Weather Review*, 138, 203-211,
628 10.1175/2009MWR2945.1, 2010.
- 629 Brier, G. W.: Verification of forecasts expressed in terms of probability, *Monthly Weather*
630 *Review*, 78, 1-3, 10.1126/science.27.693.594, 1950.
- 631 Brown, J. N., McIntosh, P. C., Pook, M. J., and Risbey, J. S.: An investigation of the links
632 between ENSO flavors and rainfall processes in southeastern Australia, *Monthly Weather*
633 *Review*, 137, 3786-3795, 10.1175/2009MWR3066.1, 2009.
- 634 Cai, W., and van Rensch, P.: The 2011 southeast Queensland extreme summer rainfall: a
635 confirmation of a negative Pacific Decadal Oscillation phase?, *Geophysical Research Letters*,
636 39, L08702, 10.1029/2011GL050820, 2012.
- 637 Chiew, F. H. S., Zhou, S. L., and McMahan, T. A.: Use of seasonal streamflow forecasts in
638 water resources management, *Journal of Hydrology*, 270, 135–144, 10.1016/S0022-
639 1694(02)00292-5, 1998.
- 640 Cloke, H. L., and Pappenberger, F.: Evaluating forecasts of extreme events for hydrological
641 applications: an approach for screening unfamiliar performance measures, *Meteorological*
642 *Applications*, 15, 10.1002/met.58, 2008.
- 643 Cloke, H. L., and Pappenberger, F.: Ensemble flood forecasting: a review, *Journal of*
644 *Hydrology*, 375, 613–626, 10.1016/j.jhydrol.2009.06.005, 2009.
- 645 Cuo, L., Pagano, T. C., and Wang, Q. J.: A review of quantitative precipitation forecasts and
646 their use in short-to medium range streamflow forecasting, *Journal of Hydrometeorology*, 12,
647 713–728, 10.1175/2011JHM1347.1, 2011.
- 648 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P., and Rieckermann, J.:
649 Improving uncertainty estimation in urban hydrological modeling by statistically describing
650 bias, *Hydrology and Earth System Sciences*, 17, 4209-4225, 10.5194/hess-17-4209-2013,
651 2013.
- 652 DelSole, T., and Shukla, J.: Artificial skill due to predictor screening, *Journal of Climate*, 22,
653 331–345, 10.1175/2008JCLI2414.1, 2009.
- 654 Elliott, J., Catchlove, R., Sooriyakumaran, S., and Thompson, R.: Recent advances in the
655 development of flood forecasting and warning services in Australia, *International conference*
656 *on innovation, advances and implementation of flood forecasting technology*, Tromsø,
657 Norway, 2005, 1-10, 2005.

658 Hawthorne, S., Wang, Q. J., Schepen, A., and Robertson, D. E.: Effective use of GCM
659 outputs for forecasting monthly rainfalls to long lead times, *Water Resources Research*, 49,
660 5427–5436, 10.1002/wrcr.20453, 2013.

661 Jewson, S.: The problem with the Brier score, arXiv: physics/0401046v1 [physics.ao-ph],
662 available at <http://arxiv.org/abs/physics/0401046v1> (last accessed June 2013), 2008.

663 Jones, D. A., Wang, W., and Fawcett, R.: High-quality spatial climate data-sets for Australia,
664 *Australian Meteorological and Oceanographic Journal*, 58, 233-248, 2009.

665 Kalnay, E., M. Kanamitsua, R. Kistlera, W. Collinsa, D. Deavena, L. Gandina, M. Iredella, S.
666 Sahaa, G. Whitea, J. Woollena, Y. Zhua, A. Leetmaa, R. Reynolds, M. Chelliah, W.
667 Ebisuzakib, W. Higgins, J. Janowiak, K.C. Mob, C. Ropelewskib, J. Wang, R. Jenne, and
668 Joseph, D.: The NCEP/NCAR 40-year reanalysis project, *Bull. Amer. Meteor. Soc.*, 77, 34,
669 1996.

670 Kiem, A. S., Franks, S. W., and Kuczera, G.: Multi-decadal variability of flood risk,
671 *Geophysical Research Letters*, 30, 1035, 10.1029/2002GL015992, 2003.

672 Koster, R. D., P., S. P., Mahanama, Livneh, B., Lettenmaier, D. P., and Reichle, R. H.: Skill
673 in streamflow forecasts derived from large-scale estimates of soil moisture and snow, *Nature*
674 *Geoscience*, 3, 613-616, 10.1038/NGEO944, 2010.

675 Kwon, H.-H., Brown, C., Xu, K., and Lall, U.: Seasonal and annual maximum streamflow
676 forecasting using climate information: application to the Three Gorges Dam in the Yangtze
677 River basin, China, *Hydrological Sciences Journal*, 54, 582-595, 10.1623/hysj.54.3.582, 2009.

678 Li, H., Luo, L., Wood, E. F., and Schaake, J.: The role of initial conditions and forcing
679 uncertainties in seasonal hydrologic forecasting, *Journal of Geophysical Research:*
680 *Atmospheres*, 114, D04114, 10.1029/2008jd010969, 2009.

681 Lindström, G., and Olsson, J.: A systematic review of sensitivities in the Swedish flood-
682 forecasting system, *Atmospheric Research*, 100, 275-284, 10.1016/j.atmosres.2010.09.013,
683 2011.

684 Mahanama, S., Livneh, B., Koster, R., Lettenmaier, D., and Reichle, R.: Soil moisture, snow,
685 and seasonal streamflow forecasts in the United States, *Journal of Hydrometeorology*, 13,
686 189-203, 10.1175/jhm-d-11-046.1, 2012.

687 Marshall, A. G., Hudson, D., Wheeler, M. C., Hendon, H. H., and Alves, O.: Assessing the
688 simulation and prediction of rainfall associated with the MJO in the POAMA seasonal
689 forecast system, *Climate Dynamics*, 37, 2129–2141, DOI 10.1007/s00382-010-0948-2, 2011.

690 Murphy, B. F., and Timbal, B.: A review of recent climate variability and climate change in
691 southeastern Australia, *International Journal of Climatology*, 28, 859-879, 10.1002/joc.1627,
692 2008.

693 Pfister, N.: The case of an evacuation from Grafton, *The Australian Journal of Emergency*
694 *Management*, 17, 19-29, 2002.

695 Piechota, T. C., Chiew, F. H. S., Dracup, J. A., and McMahon, T. A.: Seasonal streamflow
696 forecasting in eastern Australia and the El Niño–Southern Oscillation, *Water Resources*
697 *Research*, 34, 3035–3044, 10.1029/98WR02406, 1998.

698 Pook, M. J., McIntosh, P. C., and Meyers, G. A.: The synoptic decomposition of cool-season
699 rainfall in the southeastern Australian cropping region, *Journal of Applied Meteorology and*
700 *Climatology*, 45, 1156-1170, 10.1175/JAM2394.1, 2006.

701 Risbey, J. S., Pook, M. J., McIntosh, P. C., Wheeler, M. C., and Hendon, H. H.: On the
702 remote drivers of rainfall variability in Australia, *Monthly Weather Review*, 137, 3233–3253,
703 10.1175/2009MWR2861.1, 2009.

704 Roberts, N. M., Cole, S. J., Forbes, R. M., Mooreb, R. J., and Boswell, D.: Use of high-
705 resolution NWP rainfall and river flow forecasts for advance warning of the Carlisle flood,
706 north-west England, *Meteorological Applications*, 16, 23–34, 10.1002/met.94, 2008.

707 Robertson, D. E., and Wang, Q. J.: A Bayesian approach to predictor selection for seasonal
708 streamflow forecasting, *Journal of Hydrometeorology*, 13, 155–171, 10.1175/JHM-D-10-
709 05009.1, 2012.

710 Robertson, D. E., Pokhrel, P., and Wang, Q. J.: Improving statistical forecasts of seasonal
711 streamflows using hydrological model output, *Hydrology and Earth System Sciences*, 17,
712 579–593, 10.5194/hess-17-579-2013, 2013a.

713 Robertson, D. E., Shrestha, D. L., and Wang, Q. J.: Post-processing rainfall forecasts from
714 numerical weather prediction models for short-term streamflow forecasting, *Hydrology and*
715 *Earth System Sciences*, 17, 3587–3603, doi:10.5194/hess-17-3587-2013, 2013b.

716 Robertson, D. E., and Wang, Q. J.: Seasonal Forecasts of Unregulated Inflows into the
717 Murray River, Australia, *Water Resour Manage*, 27, 2747–2769, 10.1007/s11269-013-0313-4,
718 2013.

719 Rosenberg, E. A., Wood, A. W., and Steinemann, A. C.: Statistical applications of physically
720 based hydrologic models to seasonal streamflow forecasts, *Water Resources Research*, 47,
721 W00H14, 10.1029/2010WR010101, 2011.

722 Saji, N. H., Goswami, B. N., Vinayachandran, P. N., and Yamagata, T.: A dipole mode in the
723 tropical Indian Ocean, *Nature*, 401, 360–363, 1999.

724 Schepen, A., Wang, Q. J., and Robertson, D.: Evidence for using lagged climate indices to
725 forecast Australian seasonal rainfall, *Journal of Climate*, 25, 1230–1246, 10.1175/JCLI-D-11-
726 00156.1, 2012a.

727 Schepen, A., Wang, Q. J., and Robertson, D. E.: Combining the strengths of statistical and
728 dynamical modeling approaches for forecasting Australian seasonal rainfall, *Journal of*
729 *Geophysical Research*, 117, D20107, 10.1029/2012JD018011, 2012b.

730 Sharma, A.: Seasonal to interannual rainfall probabilistic forecasts for improved water supply
731 management: part 3 — a nonparametric probabilistic forecast model, *Journal of Hydrology*,
732 239, 249–258, 10.1016/S0022-1694(00)00348-6, 2000.

733 Shrestha, D. L., Robertson, D. E., Wang, Q. J., Pagano, T. C., and Hapuarachchi, H. A. P.:
734 Evaluation of numerical weather prediction model precipitation forecasts for short-term
735 streamflow forecasting purpose, *Hydrology and Earth System Sciences*, 17, 1913–1931,
736 10.5194/hess-17-1913-2013, 2013.

737 Shukla, S., and Lettenmaier, D. P.: Seasonal hydrologic prediction in the United States:
738 understanding the role of initial hydrologic conditions and seasonal climate forecast skill,
739 *Hydrology and Earth System Sciences*, 15, 3529–3538, 10.5194/hess-15-3529-2011, 2011.

740 Smith, T. M., Reynolds, R. W., T.C.Peterson, and Lawrimore, J.: Improvements to NOAA’s
741 historical merged land–ocean surface temperature analysis (1880–2006), *Journal of Climate*,
742 21, 2283–2296, 2008.

743 Troup, A. J.: The southern oscillation, *Quarterly Journal of the Royal Meteorological Society*,
744 91, 490–506, 10.1002/qj.49709139009, 1965.

745 Vaze, J., Perraud, J., Teng, J., Chiew, F., Wang, B., and Yang, Z.: Catchment Water Yield
746 Estimation Tools (CWYET), 34th World Congress of the International Association for
747 Hydro-environment Research and Engineering and 33rd Hydrology and Water Resources
748 Symposium and the 10th Conference on Hydraulics in Water Engineering, Brisbane, 2011,
749 1554-1561, 2011.

750 Verdon, D. C., Wyatt, A. M., Kiem, A. S., and Franks, S. W.: Multidecadal variability of
751 rainfall and streamflow: Eastern Australia, *Water Resources Research*, 40, W10201,
752 10.1029/2004WR003234., 2004.

753 Verdon, D. C., and Franks, S. W.: Indian Ocean sea surface temperature variability and winter
754 rainfall: Eastern Australia, *Water Resources Research*, 41, W09413,
755 10.1029/2004WR003845, 2005.

756 Wang, Q. J., Robertson, D. E., and Chiew, F. H. S.: A Bayesian joint probability modeling
757 approach for seasonal forecasting of streamflows at multiple sites, *Water Resources Research*,
758 45, W05407, 10.1029/2008WR007355, 2009.

759 Wang, Q. J., and Robertson, D. E.: Multisite probabilistic forecasting of seasonal flows for
760 streams with zero value occurrences, *Water Resources Research*, 47, W02546,
761 10.1029/2010WR009333, 2011.

762 Wang, Q. J., Schepen, A., and Robertson, D. E.: Merging seasonal rainfall forecasts from
763 multiple statistical models through Bayesian model averaging, *Journal of Climate*, 25, 5524-
764 5537, 10.1175/JCLI-D-11-00386.1, 2012a.

765 Wang, Q. J., Shrestha, D. L., Robertson, D. E., and Pokhrel, P.: A log-sinh transformation for
766 data normalization and variance stabilization, *Water Resources Research*, 48, W05514,
767 10.1029/2011WR010973., 2012b.

768 Wilks, D. S.: *Statistical Methods in the Atmospheric Sciences*, Elsevier, New York, 648 pp.,
769 1995.

770 Yeo, I. K., and Johnson, R. A.: A new family of power transformations to improve normality
771 or symmetry, *Biometrika*, 87, 954–959, 10.1093/biomet/87.4.954, 2000.

772

773

1 Table 1 Characteristics of catchments used in this study.

Name	Short name	Streamflow record used	Fraction of record missing	Area (km ²)	Annual rainfall (mm)	Annual runoff (mm)	Runoff coefficient
Orara River at Bawden Bridge	ORB	1956-2006	4.2%	1823	1396	407	0.29
Nowendoc River at Rocks Crossing	NOR	1950-2006	3.9%	1898	1155	258	0.22
Abercrombie River at Hadley No. 2	ABH	1960-2005	0.5%	1626	842	117	0.14
Murray River at Biggara	MUR	1950-2005	2.5%	1254	1178	446	0.38
Mitta Mitta River at Hinnomunjie	MMH	1950-2006	2.6%	1528	1343	297	0.22
Tarwin River at Meeniyah	TAW	1955-2006	3.1%	1066	1084	233	0.21

2

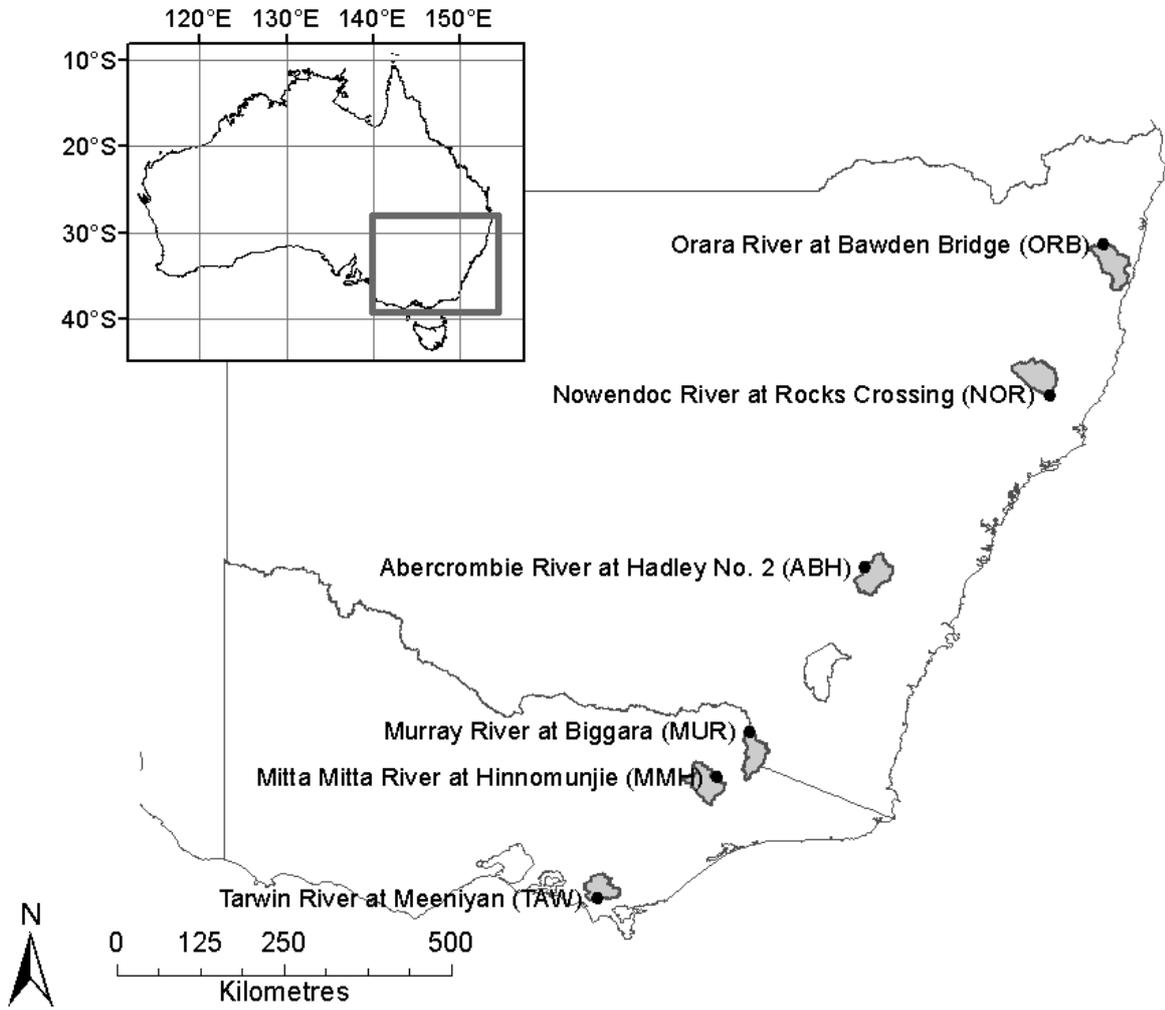
1 Table 2 List of oceanic and atmospheric climate indices used as predictors.

Index	Description
Southern Oscillation Index (SOI)	Troup (1965)
NINO3	Mean SST anomaly over 150–90°W and 5°N–5°S
NINO3.4	Mean SST anomaly over 170–120°W and 5°N–5°S
NINO4	Mean SST anomaly over 150–160°E and 5°N–5°S
ENSO Modoki Index (EMI)	Ashok et al. (2003)
Indian Ocean Dipole Mode Index (DMI)	Saji et al. (1999)
Indian Ocean West Pole Index (WPI)	Saji et al. (1999)
Indian Ocean East Pole Index (EPI)	Saji et al. (1999)
Indonesia Index (II)	Verdon and Franks (2005)
Tasman Sea Index (TSI)	Murphy and Timbal (2008)
140°E Blocking Index (B140)	Risbey et al. (2009)

2

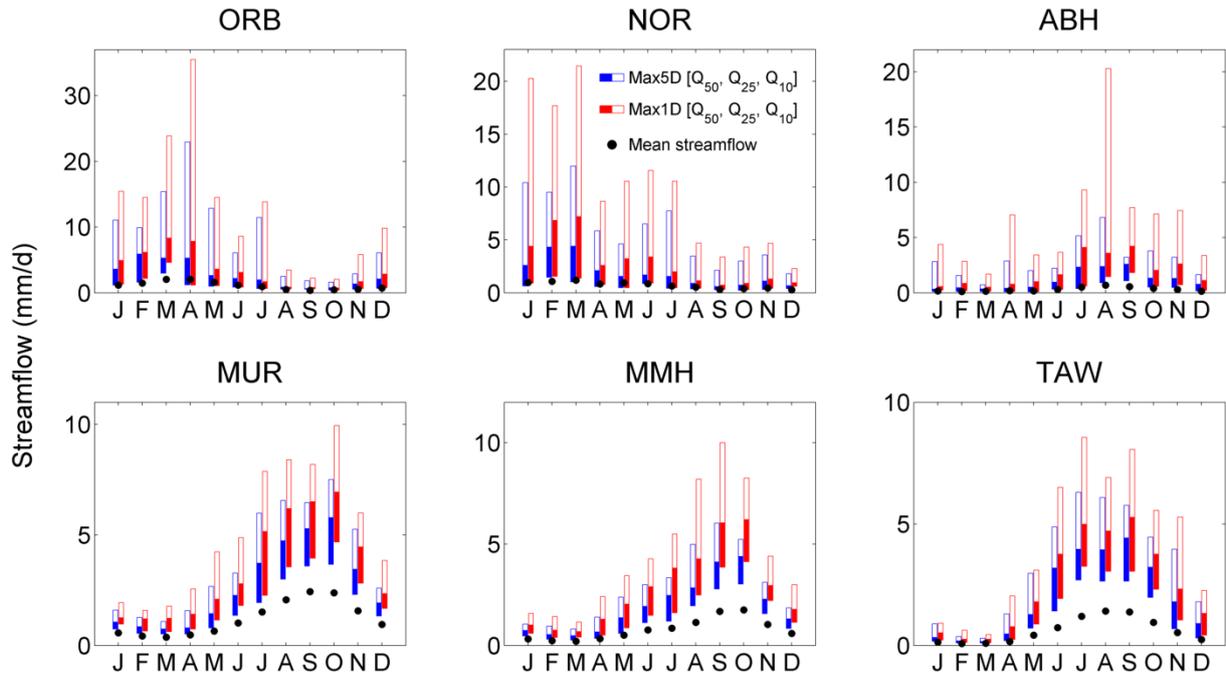
3

1 Fig. 1 Catchments (shaded) and streamflow gauge sites (black dots) used in this study.



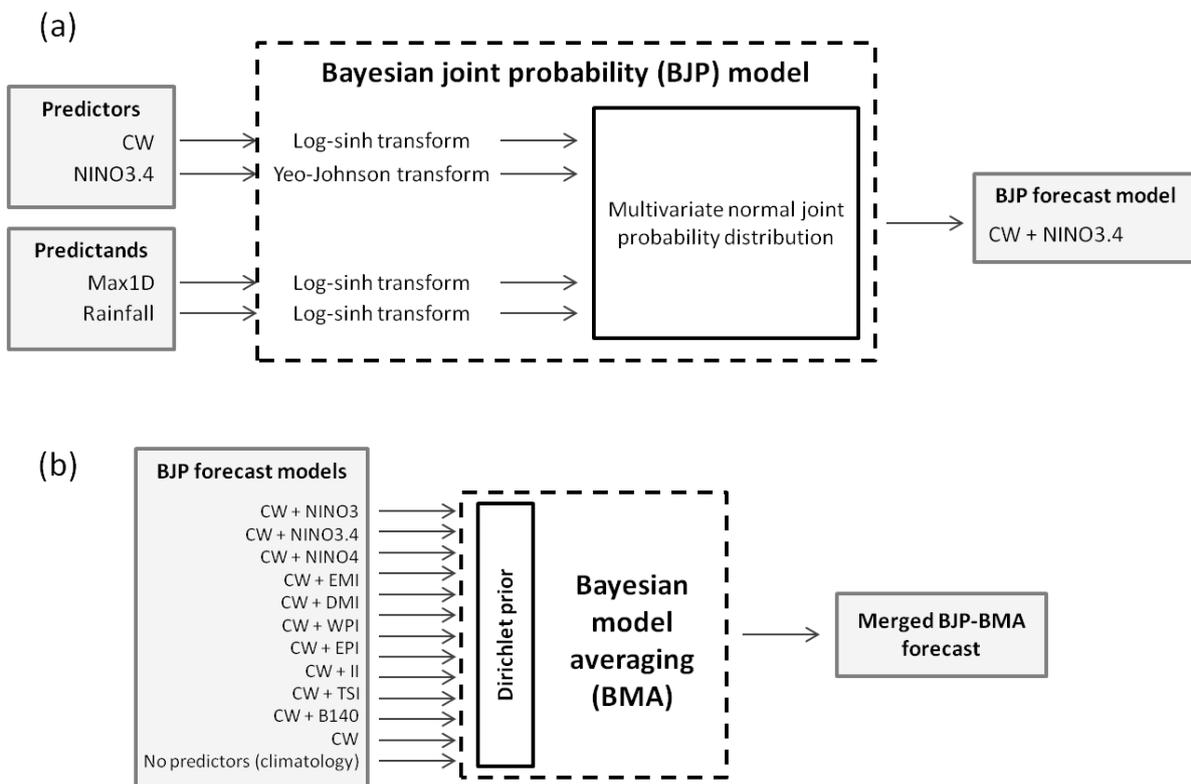
2
3

- 1 Fig. 2 Catchment streamflow characteristics. Black dots show average monthly streamflows.
- 2 Boxes show maximum five-day streamflow (Max5D - blue) and maximum 1-day streamflow
- 3 (Max1D - red) occurring during each month for exceedance probabilities of 50% (Q_{50} , bottom
- 4 edge) to 10% (Q_{10} , top edge), with box centreline showing Max5D/Max1D streamflows of
- 5 exceedance probability of 25% (Q_{25}).



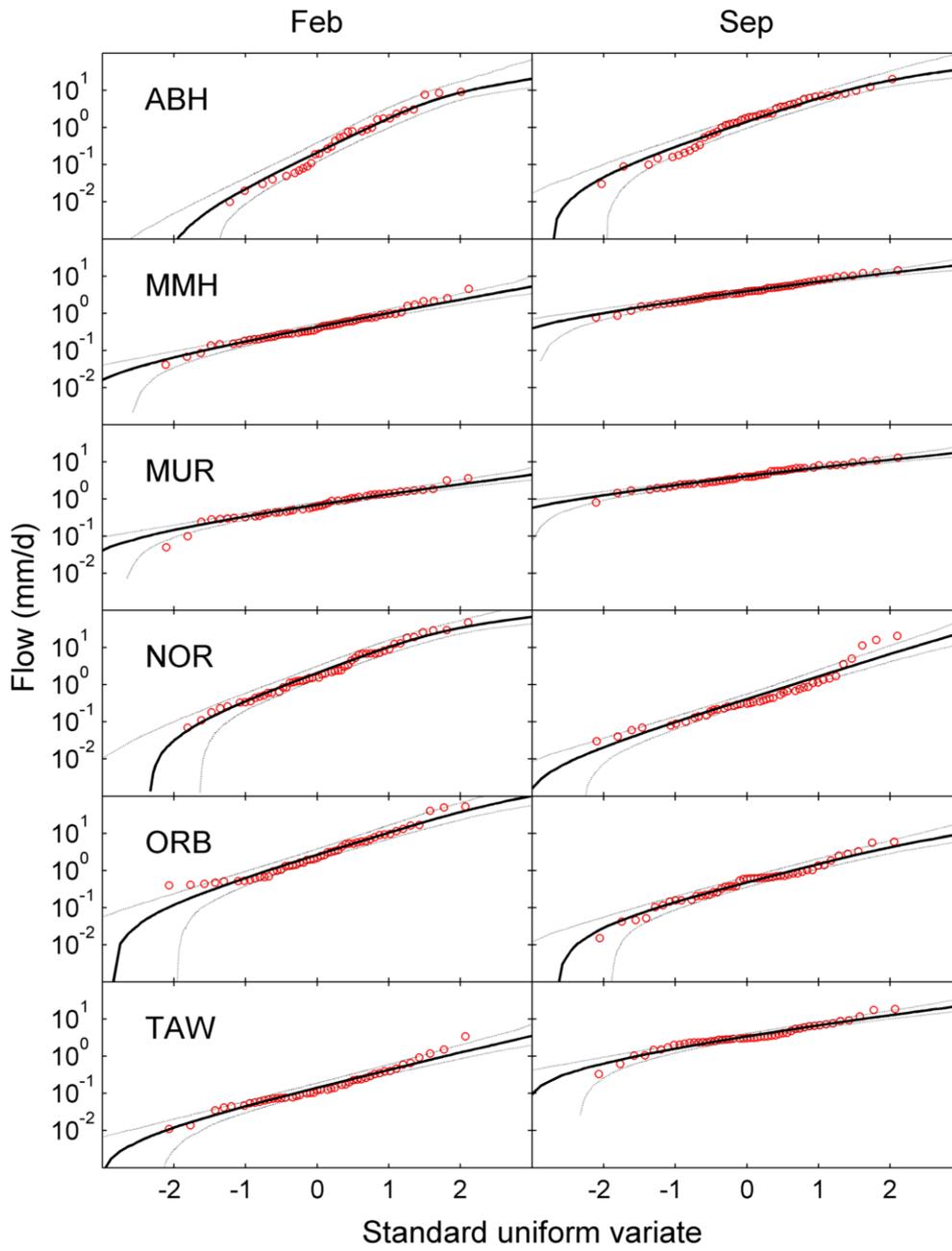
6
7

1 Fig. 3 Schematic of forecast model. (a) Example of individual forecast model generated with
 2 the Bayesian joint probability method. In this example, catchment wetness (CW) and
 3 NINO3.4 predictors are used to predict Max1D streamflows. Rainfall is included as a joint
 4 predictand to elicit more information from the climate indices. Parameters for the transforms
 5 and joint probability distribution are inferred jointly. This process is repeated for thirteen
 6 different predictor-sets. (b) The forecasts from thirteen BJP models are weighted based on
 7 cross-validated predictive performance with Bayesian model averaging (BMA) to produce a
 8 merged BJP-BMA forecast. The use of a symmetric Dirichlet prior encourages even weights
 9 in instances of high sampling uncertainty. See text for details.



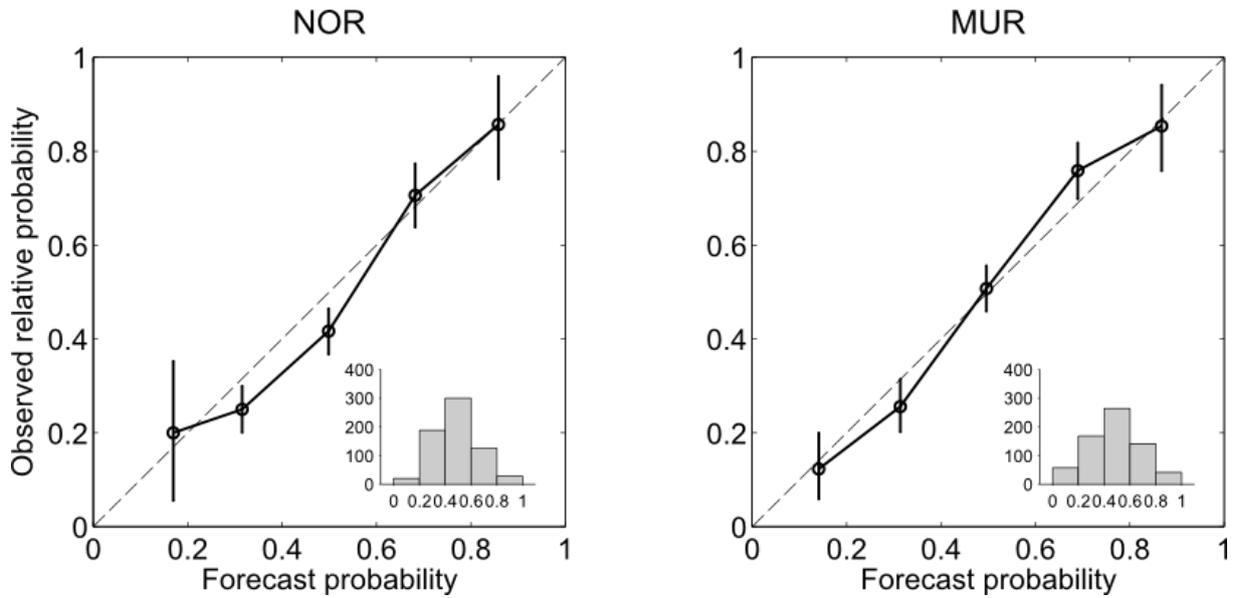
10
 11

- 1 Fig. 4 Fit of log-sinh transformed normal distributions to Max1D values for two months.
- 2 Red circles show actual values, black solid line shows fitted log-sinh transform, dashed lines show
- 3 [0.1, 0.9] confidence intervals.



4
5

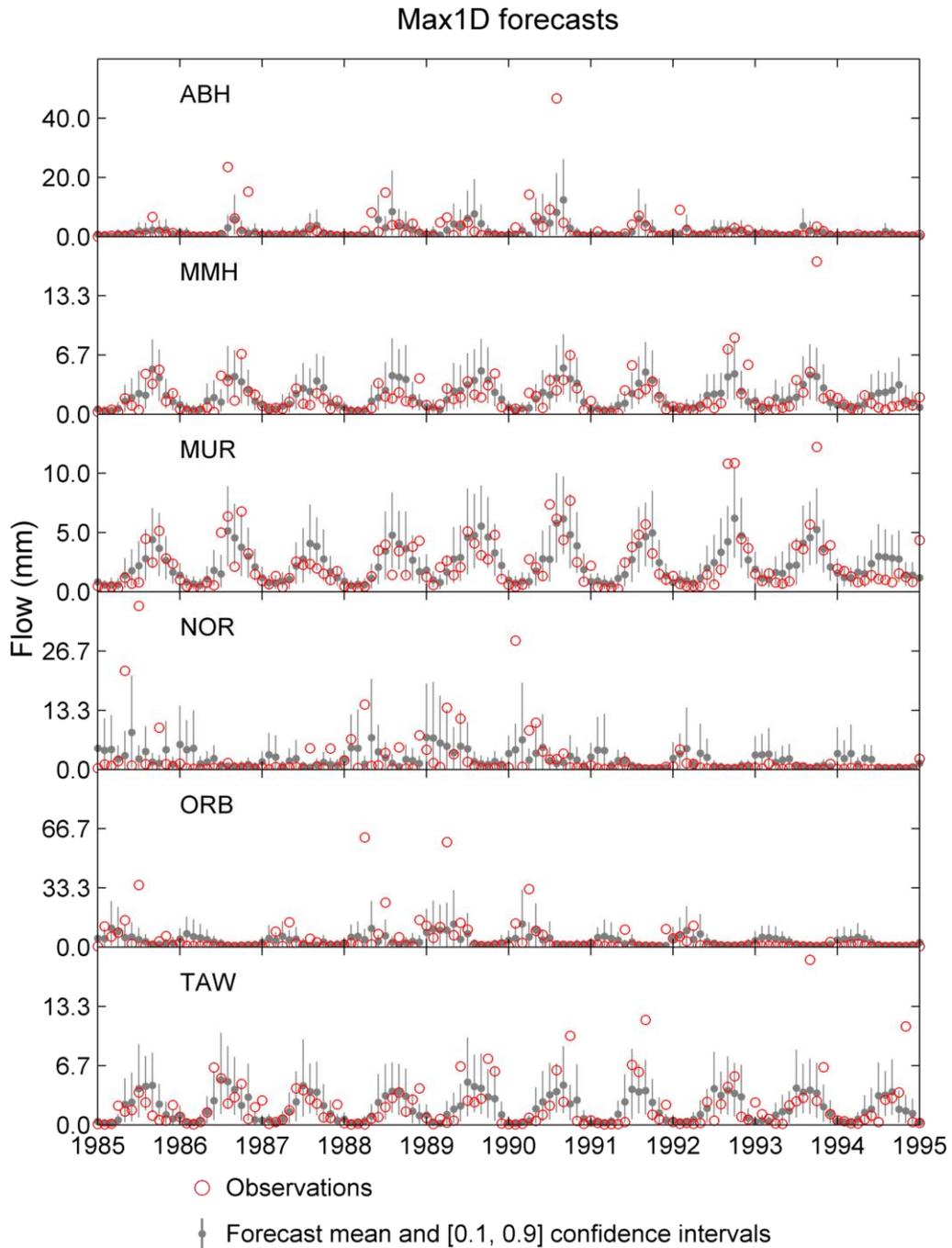
1 Fig. 5 Forecast reliability diagrams at two catchments for Max1D streamflows of exceedance
 2 probability $\leq 50\%$. (Forecasts are divided into five bins. 1:1 dashed lines, perfectly reliable
 3 forecast; circles, observed relative frequency; vertical lines, [0.05, 0.95] uncertainty interval
 4 of observed relative frequency; inserts, number of events in the different forecast probability
 5 bins.)



6

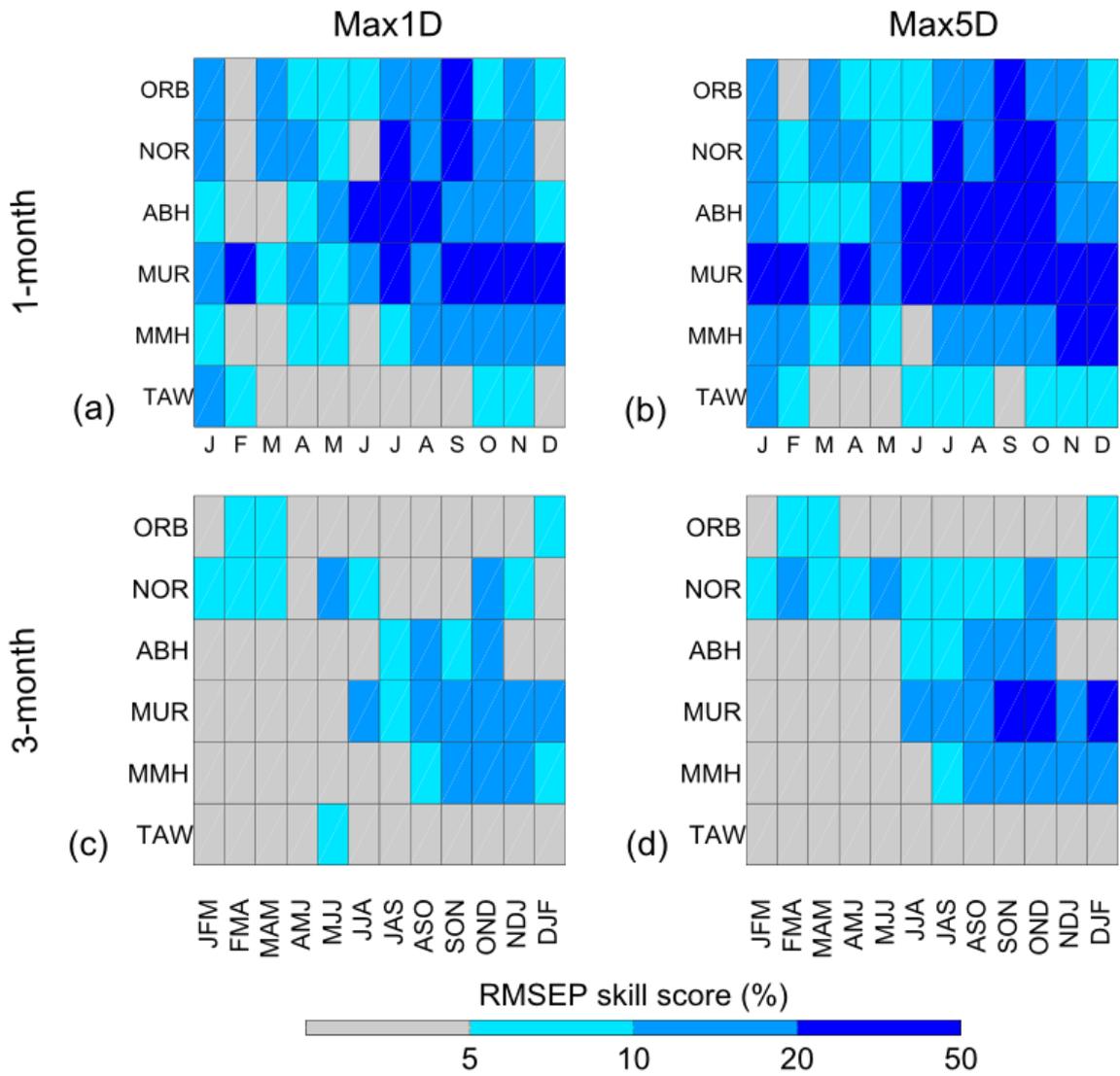
7

- 1 Fig. 6 Example forecast time series of cross-validated BJP-BMA for Max1D. Red circles
- 2 show observed Max1D values, black points and lines show mean forecast and [0.1, 0.9]
- 3 credible prediction intervals.



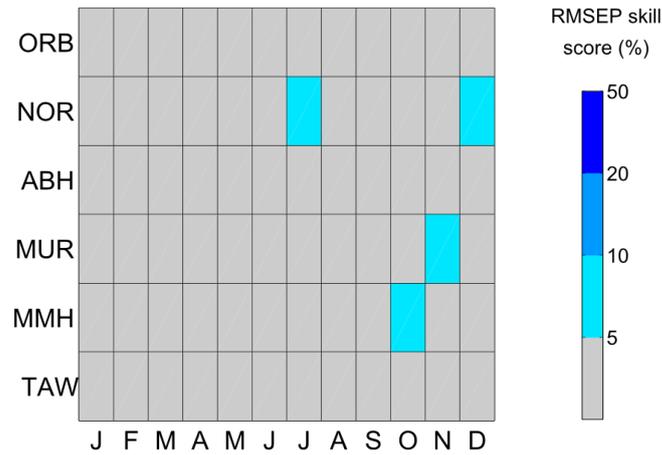
- 4
- 5
- 6

1 Fig. 7 RMSEP skill scores. Catchments are ordered by their location, from northernmost (top)
 2 to southernmost (bottom). (a) Max1D streamflows for 1-month forecasts, (b) Max5D
 3 streamflows for 1-month forecasts, (c) Max1D streamflows for 3-month forecasts, and (d)
 4 Max5D streamflows at 3-month forecasts. Scores show proportional improvement of
 5 forecasts over climatology forecasts.



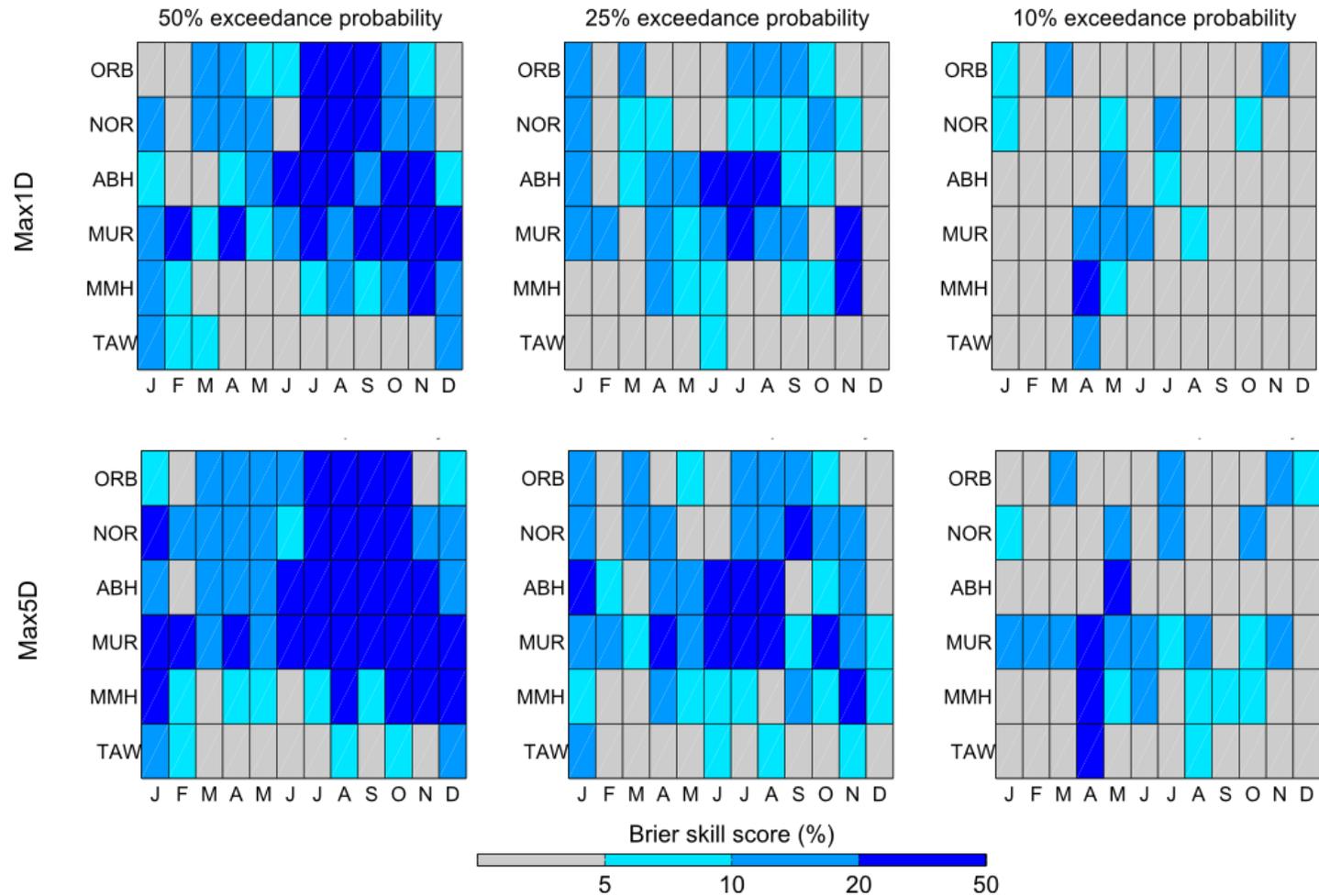
6

- 1 Fig. 8 Skill added by climate indices to forecasts. Plot shows RMSEP skill scores for Max1D
- 2 1-month forecasts calculated with respect to BJP forecasts generated with only catchment
- 3 wetness as a predictor. Scores show proportional improvement of BJP-BMA forecasts over
- 4 BJP forecasts generated with only catchment wetness as a predictor.



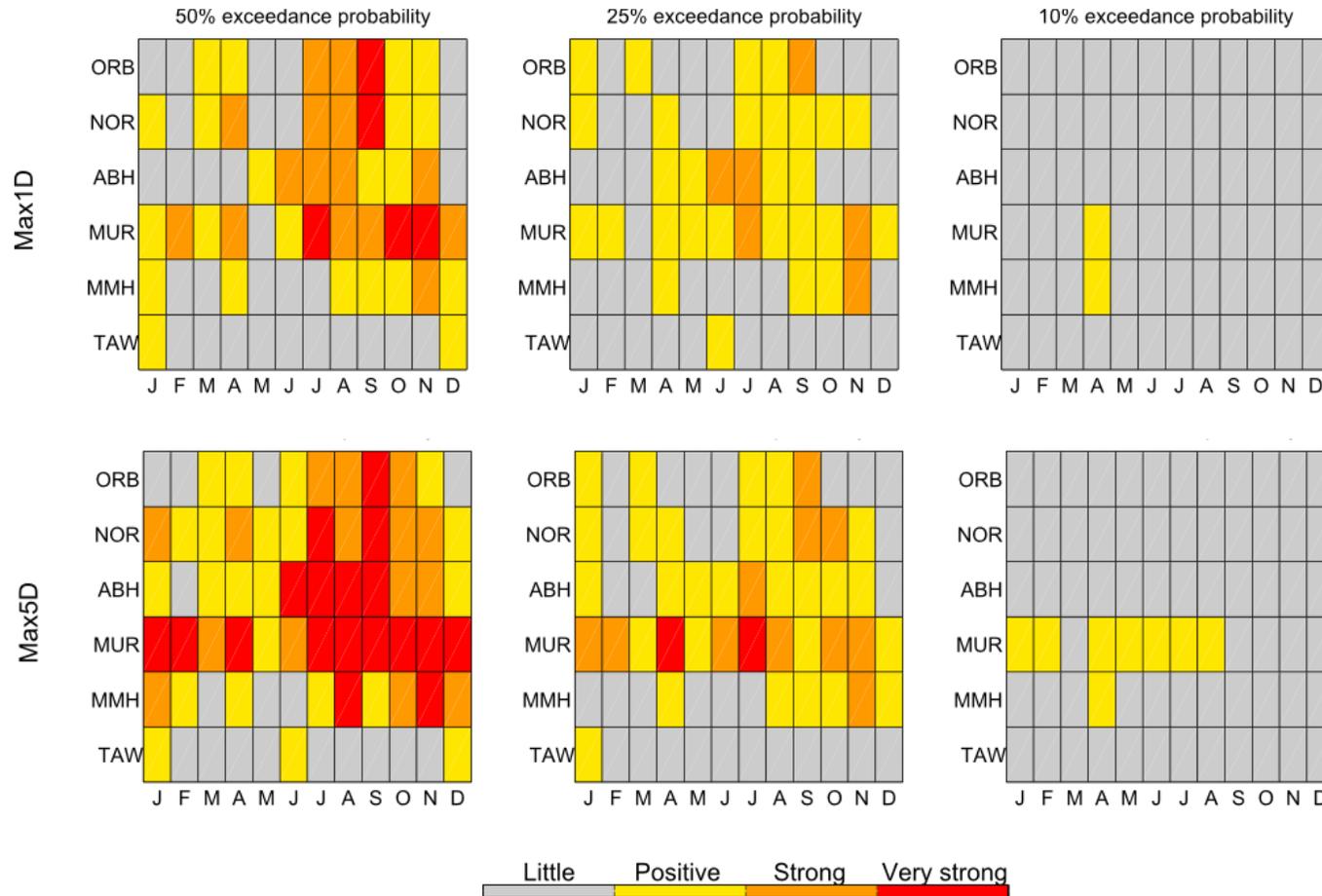
5

- 1 Fig. 9 Brier skill scores calculated at three streamflow thresholds for 1-month forecasts. Scores show proportional improvement of BJP-BMA
- 2 forecasts over climatology forecasts.



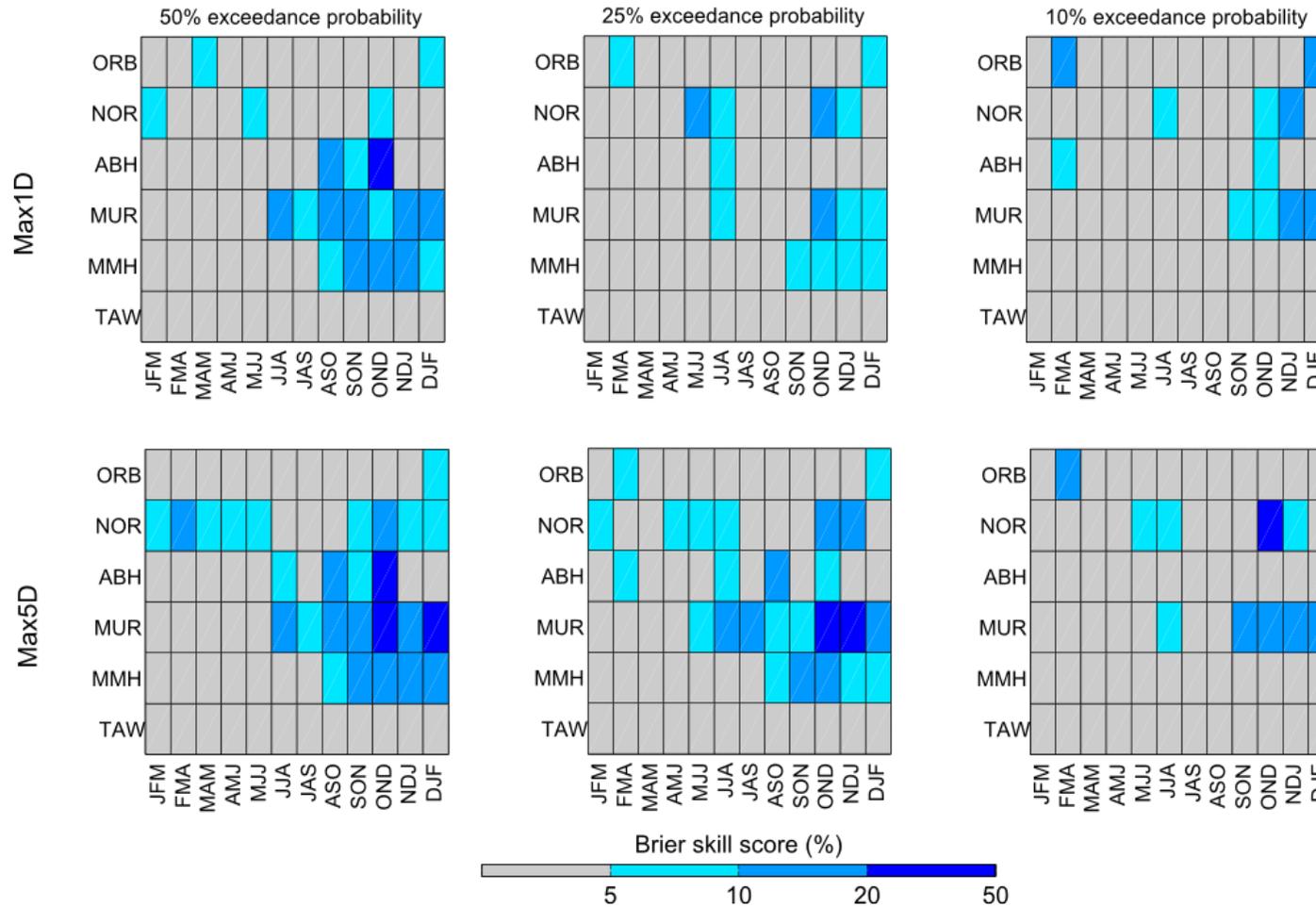
3

- 1 Fig. 10 Evidence of skill from the log-likelihood ratio (LLR) at three streamflow thresholds for 1-month forecasts. Scores show evidence of
- 2 skill of BJP-BMA forecasts over climatology forecasts. Categories are taken from Schepen et al. (2012a): little evidence of skill where
- 3 $LLR < 2$; positive evidence where $2 < LLR < 4$; strong evidence where $4 < LLR < 6$; very strong evidence where $LLR > 6$.



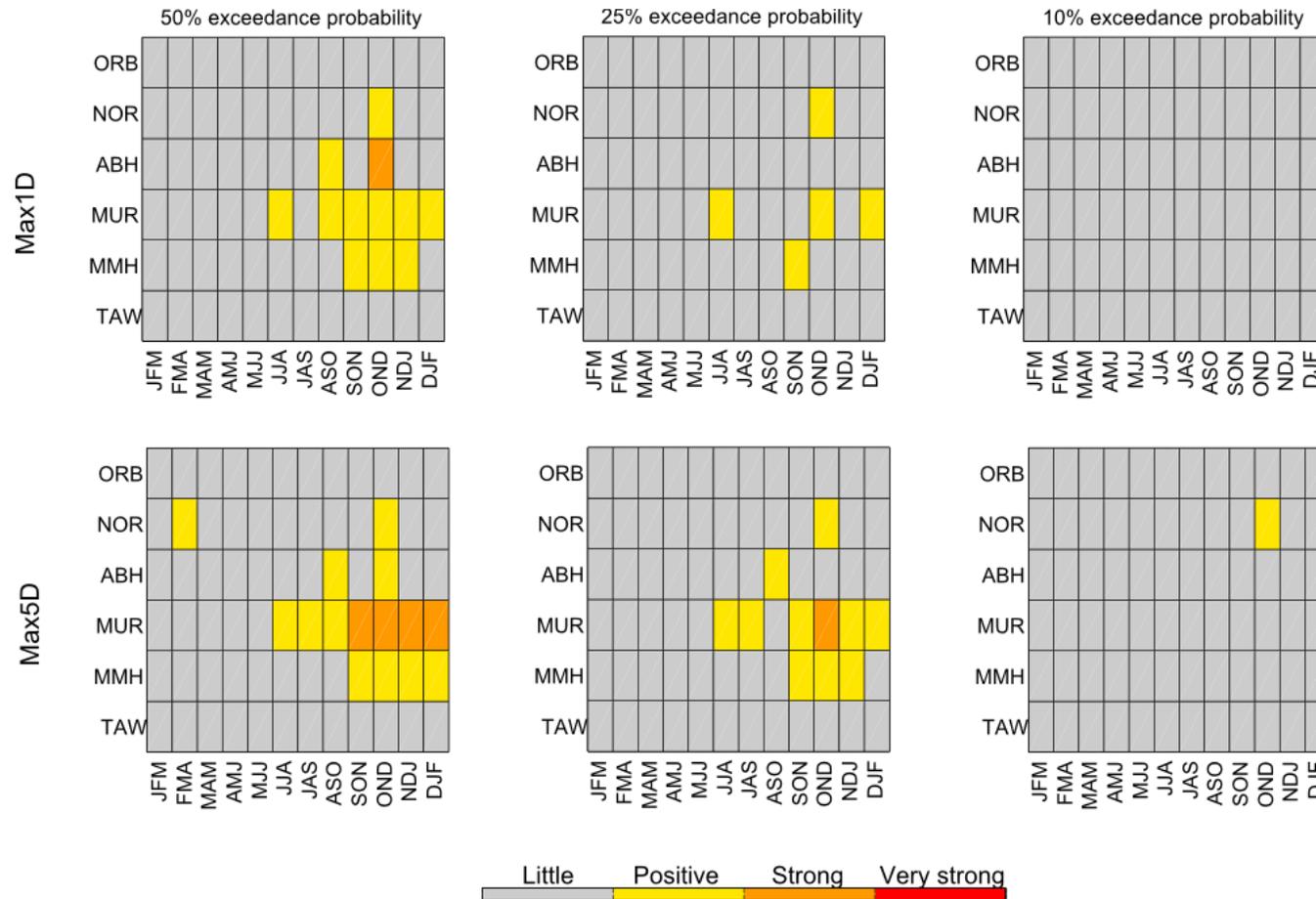
4

- Fig. 11 Brier skill scores calculated at three streamflow thresholds for 3-month forecasts. Scores show proportional improvement of BJP-
- BMA forecasts over climatology forecasts.



3

- 1 Fig. 12 Evidence of skill from the log-likelihood ratio at three streamflow thresholds for 3-month forecasts. Scores show evidence of skill of
- 2 BJP-BMA forecasts over climatology forecasts. Categories are taken from Schepen et al. (2012a): little evidence of skill where $LLR < 2$;
- 3 positive evidence where $2 < LLR < 4$; strong evidence where $4 < LLR < 6$; very strong evidence where $LLR > 6$.



4