



Minerva Access is the Institutional Repository of The University of Melbourne

**Author/s:**

Bennett, JC;Wang, QJ;Robertson, DE;Bridgart, R;Lerat, J;Li, M;Michael, K

**Title:**

An error model for long-range ensemble forecasts of ephemeral rivers

**Date:**

2021-05-01

**Citation:**

Bennett, J. C., Wang, Q. J., Robertson, D. E., Bridgart, R., Lerat, J., Li, M. & Michael, K. (2021). An error model for long-range ensemble forecasts of ephemeral rivers. *Advances in Water Resources*, 151, <https://doi.org/10.1016/j.advwatres.2021.103891>.

**Persistent Link:**

<https://hdl.handle.net/11343/354213>



17 Key points

- 18 • No existing long-range forecasting methods produce reliable ensembles in highly  
19 ephemeral streams.
- 20 • This is a new approach to handling zero values in both observations and simulations.
- 21 • Forecasts are reliable to 12 months for even highly ephemeral streams.

22

23

24 Abstract

25 Few ensemble streamflow forecasting systems are designed to operate for ephemeral rivers.  
26 In this study, we revise our error model for generating Forecast Guided Stochastic Scenarios  
27 (FoGSS) to produce statistically reliable long-range (12-month) forecasts for ephemeral  
28 rivers. FoGSS features an error model with four stages: data transformation, bias-correction,  
29 an autoregressive error model and the statistical distribution of residuals. We revise the fourth  
30 stage of FoGSS with a parameter estimation method that uses data censoring to account for  
31 zero values in both observations and forecasts. This allows FoGSS to produce statistically  
32 reliable ensemble forecasts in even highly ephemeral streams (with >50% zero flows). We  
33 apply FoGSS to conventional ensemble hydrological prediction (ESP) forecasts for 50  
34 Australian catchments, including 26 ephemeral rivers. We show that FoGSS improves the  
35 accuracy of ESP forecasts at short lead times, while at long lead times FoGSS forecasts  
36 transition to climatology-like forecasts. FoGSS forecasts are reliable in ensemble spread at  
37 individual lead times and for volumes aggregated over lead times, even in highly ephemeral  
38 rivers. FoGSS forecasts pave the way for operational long-range forecasts in ephemeral  
39 rivers, meeting a key need for improved water management.

40 **Keywords:** ephemeral rivers; dryland streams; long-range forecasting; ESP; FoGSS;  
41 ensemble prediction.

42

43

## 44 1 Introduction

45 Ephemeral rivers pose particular problems for ensemble streamflow forecasting. They  
46 often exhibit highly non-linear responses in runoff to rainfall, making them difficult to model  
47 (e.g., Costigan et al., 2017). At the same time, statistical treatments that are necessary to  
48 generate reliable ensembles are complicated by the presence of zero values (McInerney et al.,  
49 2019; Smith et al., 2015; Smith et al., 2010). Perhaps as a consequence, few streamflow  
50 forecasting systems are designed for ephemeral rivers. This is despite the clear need for  
51 methods that can generate skillful and reliable streamflow forecasts for ephemeral rivers:  
52 ephemeral rivers drain close to half the earth’s surface (Datry et al., 2017; Tooth, 2000),  
53 providing water vital to ecosystems and humans. In Australia, the need for streamflow  
54 forecasts in ephemeral rivers is acute. Ephemeral rivers are a crucial source of water for  
55 irrigated agriculture over much of Australia, for example in the northern and western Murray-  
56 Darling Basin, a crucial region in Australia’s most productive agricultural basin. Balancing  
57 the competing needs of irrigators and ecosystems in water-stressed systems such as the  
58 Murray-Darling basin is one of Australia’s most pressing water management challenges (e.g.  
59 Grafton and Wheeler, 2018). Further, prospective agricultural development in monsoonal  
60 northern Australia (e.g. Petheram et al., 2018) will rely heavily on ephemeral streams.  
61 Streamflow forecasts - in particular long-range (12-month) forecasts that assist in water  
62 allocation decisions (e.g. Kaune et al., 2020) - are likely to be beneficial in these and other  
63 Australian drylands.

64 Long-range forecasts of streamflow are often highly uncertain, and thus require ensemble  
65 forecasting methods. To be most useful to water managers, long-range ensemble forecasts  
66 should:

- 67 1. Be in the form of hydrographs at the monthly time step. This enables water managers  
68 to i) consider forecasts of individual months and ii) accumulate forecasts for longer  
69 periods (e.g. 6-month totals) to understand forecasts at a range of temporal scales.
- 70 2. Be statistically reliable, both at individual months and for accumulated volumes, at all  
71 lead times and at all times of year – including for ephemeral rivers.
- 72 3. Be as accurate as possible when skill is available, and never less accurate than simple  
73 climatology forecasts (a property known as ‘coherence’, after Krzysztofowicz, 1999).

74 We review existing approaches to ensemble prediction in light of these properties, with  
75 particular reference to ephemeral rivers.

76 Among the longest standing methods to produce long-range streamflow forecasts are  
77 simple regression models that describe a relationship between a predictor (e.g. an estimate of  
78 moisture stored in soils) and future streamflow. These methods have been used for decades in  
79 the United States (Pagano et al., 2009). This idea has been extended to more complex  
80 Bayesian models, including the Bayesian Joint Probability (BJP) modelling approach (Wang  
81 and Robertson, 2011; Wang et al., 2009) that underpins Australia’s national seasonal  
82 streamflow forecasting service provided by the Bureau of Meteorology (BoM;  
83 <http://www.bom.gov.au/water/ssf/>). The BJP formally accounts for heteroscedasticity (the  
84 inability to assign a single variance to streamflow), parameter uncertainty and, crucially, zero  
85 values, enabling it to produce statistically reliable ensemble forecasts for ephemeral rivers.  
86 As with many statistical models, forecasts issued with the BJP are for a single lead time (total  
87 streamflow for the coming 1, 2 or 3 months), rather than hydrographs that connect  
88 streamflow forecasts at multiple lead times. While it is possible to break down 3-month  
89 streamflow totals into three 1-month forecasts with the BJP (Zhao et al., 2016), it is difficult  
90 to construct long-range hydrographs from simple predictor/predictand relationships. This is in  
91 large part because it is difficult to model the temporal properties of streamflow without some  
92 form of state variable(s) (e.g., states that mimic soil moisture stores). If temporal properties of  
93 individual ensemble members are not realistic, then aggregations of the ensemble will not be  
94 reliable - even if predictions at individual lead times are reliable (e.g. Shrestha et al., 2015).

95 Generating hydrographs is more easily performed with hydrological models. The practice  
96 of using hydrological models in seasonal streamflow prediction also has a long history:  
97 ensemble streamflow prediction (ESP) methods have been used since the 1970s (Day, 1985).  
98 ESP works by initialising a hydrological model with observed meteorological forcings (e.g.  
99 rainfall, potential evaporation) and then generating a forecast by running the model with an  
100 ensemble of these meteorological forcings taken from the historical record. An ensemble of  
101 historical forcings is simple to collate and inherently reliable (assuming a stationary climate).  
102 Forecast skill in ESP forecasts derives entirely from initial hydrological conditions (soil  
103 moisture, etc.) – i.e., forcings give no information about future meteorological conditions.  
104 Conversely, the ensemble spread in ESP forecasts derives entirely from the ensemble of  
105 forcings. ESP forecasts ignore the large uncertainties that arise in the conversion of rainfall to  
106 runoff, resulting in ensemble streamflow forecasts that are often over-confident and thus  
107 unreliable (Wood and Schaake, 2008).

108 More modern forecasting systems combine ensemble climate forecasts with hydrological  
109 models (e.g. Arnal et al., 2018; Crochemore et al., 2016). Without statistical processing,  
110 uncertainties in climate forecasts and the simulation of hydrological processes are often  
111 incorrectly specified – including in highly sophisticated ensemble climate prediction systems  
112 (Yuan et al., 2015; Zhao et al., 2017). That is, the ensembles they produce are unreliable;  
113 most usually they produce ensemble spread that is too narrow. The most common method to  
114 enforce reliable ensembles from such streamflow forecasting systems is to statistically  
115 calibrate them with methods analogous to the Model Output Statistics (MOS) approaches  
116 long used in meteorology (Hemri and Klein, 2017; Pokhrel et al., 2013; Verkade et al., 2017;  
117 Woldemeskel et al., 2018; Wood and Schaake, 2008). However, MOS approaches suffer  
118 from the same limitation described for the statistical forecasting methods above: calibration is  
119 applied separately at each lead time, completely disrupting temporal relationships of  
120 ensemble members and thus making it difficult to issue forecasts in the form of hydrographs.

121 An alternative approach is to separate forecast uncertainties into those related to climate  
122 forcings and those related to hydrological modelling. A range of methods already exists to  
123 effectively calibrate climate forecasts (e.g. Strazzo et al., 2019; Wang et al., 2019). This  
124 leaves the challenge of specifying hydrological uncertainties, which is the subject of this  
125 study. We address this challenge with error modelling. Error models have the crucial  
126 advantage over MOS methods of preserving the temporal sequences of forecast hydrographs.  
127 Most hydrological error models are developed to predict only one lead time in advance, but it  
128 is possible to propagate uncertainty through multiple lead times (Seo et al., 2006). We have  
129 previously developed the FoGSS (Forecast Guided Stochastic Scenarios) error model to  
130 generate reliable, long-range (to 12 months) ensemble forecasts in the form of hydrographs at  
131 the monthly time step (Bennett et al., 2016b; Bennett et al., 2017; Li et al., 2013). FoGSS is  
132 based on three principles:

- 133 i) Key error model parameters are adjusted by calendar month (Li et al., 2013; Liu et  
134 al., 2020)
- 135 ii) The error model is comprised of several independent stages, each addressing a  
136 particular aspect of model errors, rather than a single complex model (Bennett et  
137 al., 2016b; Bennett et al., 2017; Li et al., 2015; Li et al., 2016; Li et al., 2017).
- 138 iii) Error model parameters that are estimated for lead one predictions are used to  
139 calibrate forecasts to long lead times (Bennett et al., 2016b; Bennett et al., 2017;  
140 Li et al., 2017)

141 We have shown that FoGSS is able to generate reliable streamflow forecasts to lead times of  
142 12 months across a wide variety of perennial catchments (Bennett et al., 2016b; Bennett et  
143 al., 2017). FoGSS reduces forecast error with an autoregressive model, particularly at shorter  
144 lead times. Further, the forecasts it produces are ‘coherent’: that is, they are always at least as  
145 skillful as climatology forecasts, even at long lead times. However, the success of FoGSS in  
146 ephemeral catchments was mixed, and depended on the degree of ephemerality. It worked as  
147 expected in moderately ephemeral catchments (defined as never having any calendar month  
148 with >50% zeros), but performed poorly in highly ephemeral catchments (where some  
149 months had >50% zeros) (Bennett et al., 2017). FoGSS was unable to generate reliable  
150 ensembles in highly ephemeral rivers, sometimes resulting in negative forecast skill with  
151 respect to climatology.

152 Accounting for zeros in hydrological error models is both technically challenging and  
153 essential. Smith et al. (2010) showed that for parameter estimation methods that use a  
154 likelihood, it is not possible to correctly optimise error model parameters for ephemeral rivers  
155 without explicitly accounting for zeros. They developed a likelihood that accounted for zero  
156 values by treating residuals as a mixed discrete-continuous distribution, conditioned on  
157 observations. They later complemented this work by treating residuals as autocorrelated  
158 (Ammann et al., 2019; Smith et al., 2015). Other error models, including FoGSS, use data  
159 censoring rather than a mixed discrete-continuous distribution to treat the presence of zeros in  
160 observations (Li et al., 2016; Li et al., 2017; McInerney et al., 2019). Error models that  
161 explicitly handle the presence of zeros in observations generally perform well for modelling  
162 errors in cases where i) catchments are not more than moderately ephemeral and ii)  
163 hydrological simulations do not equal zero. However, because these error models treat  
164 residuals as symmetrically distributed around the hydrological model simulation (usually  
165 after transformation to normalise data), they cannot generate >50% zeros (Wang et al., 2020).  
166 That is, error models that handle zeros only in observations are structurally incapable of  
167 reliably simulating errors in highly ephemeral catchments, as we describe in detail in Section  
168 2.2. For this reason, FoGSS and similar models are unsuitable for generating reliable  
169 predictions in highly ephemeral catchments.

170 In a recent paper, we addressed the problems of error models in highly ephemeral streams  
171 (Wang et al., 2020). We used data censoring to establish a likelihood that applies censoring to  
172 hydrological simulations as well as observations. This likelihood was applied to a simple  
173 Gaussian error model that assumed uncorrelated residuals and produced reliable simulations

174 even in highly ephemeral rivers. However, the aims of Wang et al. (2020) were to improve  
175 hydrological simulations and to correctly specify uncertainty in highly ephemeral streams,  
176 not to produce forecasts at multiple lead times.

177 In the present study, we adapt the likelihood proposed by Wang et al. (2020) to work with the  
178 FoGSS error model to produce long-range streamflow forecasts in ephemeral rivers. We test  
179 the revised FoGSS error model on ESP forecasts generated with an experimental  
180 hydrological-model based streamflow forecasting system developed by the Bureau of  
181 Meteorology (Woldemeskel et al., 2018). To demonstrate the FoGSS error model's more  
182 general applicability, we test the forecasts on perennial, moderately ephemeral and highly  
183 ephemeral rivers.

184 The paper is structured as follows. We give a brief overview of the selected terms used in this  
185 paper in Section 1.1. The FoGSS error model is described in Section 2, dividing the existing  
186 components (Section 2.1) from the changes made for ephemeral rivers (Section 2.2). Section  
187 3 describes the ESP forecasting system to which we apply FoGSS, and the process fo  
188 generating FoGSS forecasts is described in Section 4. Catchments and data are described in  
189 Section 5. Forecast verification methods, including our cross-validation scheme, are  
190 described in Section 6, and the results of our experiments are presented in Section 7. We  
191 discuss our findings in Section 8 and summarise and conclude the study in Section 9.

## 192 1.1. Terminology

193 Several terms used in this paper are used in specific ways or are not widely known, and we  
194 define these here to assist the reader.

195 **FoGSS versions:** We designate the previous FoGSS error model (Bennett et al., 2017) as 'o-  
196 censored', to indicate that it only applies censoring to observations, and the new FoGSS  
197 presented in this paper as 'os-censored', to indicate that it applies censoring to both  
198 observations and simulations. This follows the terminology used by Wang et al. (2020).

199 **Lead times:** FoGSS produces forecasts monthly forecasts to a horizon of 12 months into the  
200 future. We refer to each forecast lead time as 'lead 1', 'lead 2', ..., 'lead 12'. A lead 1 forecast  
201 is a forecast for the next month: e.g., if a forecast is issued on 1 Feb 2000, the lead 1 forecast  
202 will predict streamflow for the period 2 Feb 2000-1 Mar 2000.

203 **Hydrological simulations/forecasts:** We designate hydrological 'simulations' as those  
204 generated by forcing the hydrological model with observed rainfall and potential evaporation.

205 ‘Forecasts’ are generated by forcing an initialised hydrological model with forecast  
206 rainfall/potential evaporation.

## 207 2 FoGSS error model

### 208 2.1 Previous FoGSS error model: o-censored

209 We briefly review the structure of o-censored FoGSS error model, which we have developed  
210 previously through a series of papers (Bennett et al., 2016b; Bennett et al., 2017; Li et al.,  
211 2013; Li et al., 2015). These papers may be consulted for further detail. Note that FoGSS is  
212 designed to produce forecasts at the monthly time step. Forecasts are issued at the beginning  
213 of each calendar month (i.e., 12 forecasts per year) to a lead time of 12 months.

#### 214 2.1.1 Stage 1 – Data transformation

215 In previous studies using FoGSS (Bennett et al., 2016b; Bennett et al., 2017), Stage 1 has  
216 encompassed the estimation of hydrological model parameters. However, it is possible to  
217 apply the FoGSS error model to a hydrological model that is already calibrated. In this study  
218 we wish to apply FoGSS to an existing forecasting setup (Section 3), and thus Stage 1  
219 concerns only the transformation.

220 We use the log-sinh transformation (Wang et al., 2012) to allow us to treat our errors as  
221 normally distributed. It is given by:

$$222 \quad z(t) = Tf(q(t)) = \frac{1}{b} \log(\sinh(a + cq(t)b)) \quad (1)$$

223 where  $q(t)$  is the observed or simulated monthly streamflow at time  $t$ , and  $a$  and  $b$  are  
224 parameters.  $c = \frac{5}{\max(\mathbf{q}_o)}$  is a standardisation constant where  $\max(\mathbf{q}_o)$  is the maximum of the  
225 time series of available streamflow observations. The standardisation constant allows  $a$  and  $b$   
226 to take comparable values across all catchments, simplifying the application of Bayesian  
227 priors for parameter estimation (Appendix A). Note that when Eq 1 is reversed (to back-  
228 transform  $z$  to  $q$ ), any values of  $z(t) < Tf(0)$  (the transformed value of zero) are first forced  
229 to  $z(t) = Tf(0)$  before back-transformation, to ensure streamflow values cannot be negative.

230 Transformation parameters are estimated only from observations using a Maximum A  
231 Posteriori (MAP) estimation (Section 2.2.2), and these parameters are applied to both  
232 observed and simulated streamflow.

### 233 2.1.2 Stage 2 – Bias-correction

234 A bias-correction is applied to transformed simulations at each calendar month:

$$235 \quad z_2(t) = d(i)z_1(t) + \mu(i) \quad (2)$$

236 where  $i = 1, 2, \dots, 12$  is the calendar month corresponding to  $t$ ,  $z_1(t)$  is the log-sinh transformed  
237 simulation (Stage 1), and  $d(i)$  and  $\mu(i)$  are parameters that vary by month. We limit  $d$  to  
238  $0 \leq d < 2$ . Values less than zero imply negative correlations (i.e., very poor hydrological  
239 model performance), meaning that it is more sensible to ignore simulations ( $d=0$ ). The upper  
240 limit is arbitrary, and avoids overly large corrections. The bias-correction parameters are  
241 estimated using a method we term ‘Least Squares after Transformation’ (LST), detailed in  
242 Section 2.2.2.

243 The bias-correction performs two important tasks. First, it allows us to treat our residuals as  
244 normally distributed with a mean of zero (see Stage 4). Second, as  $d(i)$  approaches zero  
245 (which can occur when hydrological simulations are very poor), the bias-correction tends to  
246  $z_2(t) \approx \mu(i)$ ; that is, it returns a constant, akin to climatology. This is similar in concept to the  
247 statistical calibration of forecasts with MOS methods. As noted in the introduction, MOS  
248 methods apply separate regressions to each lead time, breaking the temporal sequence of  
249 ensemble members. For our application (12-month forecasts issued 12 time per year) using  
250 MOS would mean applying a regression at each forecast issue time and each lead time (  
251  $12 \times 12 = 144$  regressions). By contrast, in FoGSS we estimate the bias-correction for each  
252 calendar month at lead 1, and apply the bias-correction to that calendar month at all lead  
253 times (at total of 12 regressions). This allows us to preserve the temporal sequence from the  
254 hydrological model in the forecast (Section **Error! Reference source not found.**). We have  
255 shown in previous papers that this method acts like a meteorological calibration at all lead  
256 times, returning climatology forecasts when predictions are not skilful (Bennett et al., 2016b;  
257 Bennett et al., 2017).

### 258 2.1.3 Stage 3 – Autoregressive updating

259 In Stage 3 we apply a first-order autoregressive (AR1) that follows the general form

$$260 \quad z_3(t) = z_2(t) + \rho(i)(z_2(t-1) - z_2(t-1)) \quad (3)$$

261 where  $z_2$  is the transformed bias-corrected simulation from Eq 2,  $z_o(t-1) = Tf(q_o(t-1))$  is the  
 262 log-sinh transformed value of observed flow, and  $\rho(i)$  is a parameter that varies by calendar  
 263 month, which can take values between 0 and 1. The transformation (Eq 1) can amplify values  
 264 of  $z_3(t)$  to be unrealistically large after back transformation and to avoid this problem, we  
 265 restrict  $z_3(t)$  (Li et al., 2015):

$$266 \quad z_3(t) = \begin{cases} Tf(\min(q_3(t), q_2(t) + (q_o(t-1) - q_2(t-1)))) & \text{when } q_o(t-1) \geq q_2(t-1) \\ Tf(\max(q_3(t), q_2(t) + (q_o(t-1) - q_2(t-1)))) & \text{when } q_o(t-1) < q_2(t-1) \end{cases} \quad (4)$$

267 where  $q_2(t) = Tf^{-1}(z_2(t))$  and  $q_3(t) = Tf^{-1}(z_3(t))$  are the back-transformed values of  $z_2(t)$  and  
 268  $z_3(t)$ . In other words, we limit the size of the update to the error in the original domain. The  
 269 AR1 model reduces errors and plays a pivotal role in propagating uncertainty through  
 270 multiple lead times, as described in Section **Error! Reference source not found.** As with  
 271 the bias-correction parameters, the AR1 parameters are estimated using a method we term  
 272 ‘Least Squares after Transformation’ (LST), detailed in Section 2.2.2.

## 273 2.2 Revising FoGSS for ephemeral rivers: os-censored

### 274 2.2.1 Stage 4 - residual modelling

275 Previous (o-censored) versions of FoGSS described only 3 stages, with the statistical  
 276 modelling of residuals included in Stage 3. For os-censoring, we separate out the statistical  
 277 model of residuals into Stage 4 to highlight the improvements to FoGSS for ephemeral  
 278 catchments.

279 As with previous incarnations of FoGSS, we assume residuals follow a normal distribution:

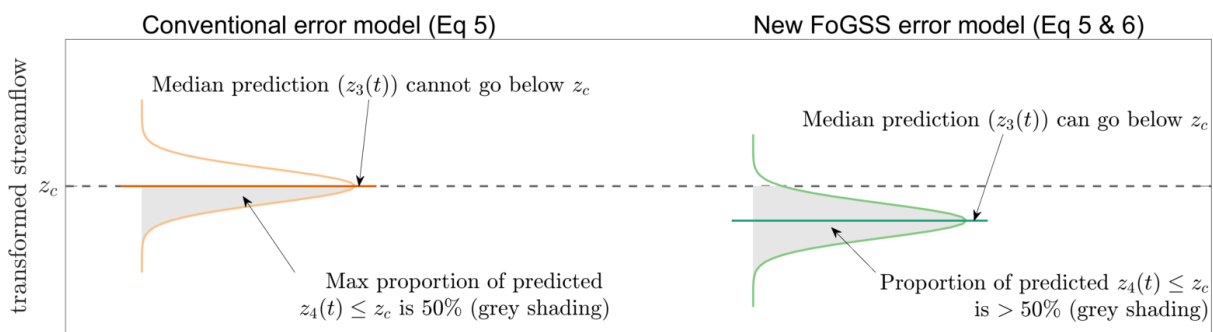
$$280 \quad \begin{aligned} z_4(t) &= z_3(t) + \varepsilon(t) \\ \varepsilon(t) &\sim N(0, \sigma^2(i)) \end{aligned} \quad (5)$$

281 where  $z_3$  is the transformed, bias-corrected and updated simulation from Eq 4, and  $\sigma^2(i)$  is  
 282 the standard deviation of residuals, varied by calendar month. To generate predictive  
 283 uncertainty, we revise FoGSS to follow Wang et al. (2020)’s use of data censoring of both  
 284 observations and simulations (os-censoring). Denote a censoring threshold of  $q_c = 0$ , and let  
 285  $z_c = Tf(q_c)$  be the transformed censoring threshold. If  $z_3(t) \leq z_c$ , a new value is assigned to  
 286  $z_3(t)$  by

287 
$$z_3(t) = \Phi^{-1}\left(\Phi\left(z_c \mid m_3(i), s_3(i)^2\right) \times r_0(t)\right)$$
 (6)  
 288 
$$r_0(t) \sim U(0,1)$$

289 where  $\Phi$  is the normal cumulative probability,  $m_3(i)$  and  $s_3(i)$  are the mean and standard  
 290 deviation of  $z_3(t)$  for the calendar month  $i$  corresponding to  $t$ , and  $U(0,1)$  is a uniform  
 291 distribution (see Section 2.2.2 for the estimation of these parameters). We then add noise to  
 292  $z_3(t)$  to give  $z_4(t)$  by sampling from Eq (5). Note that when we generate a forecast, Eq (5) and  
 293 Eq (6) are applied at each lead time through a process we call *stochastic updating*. We  
 294 describe the process of generating a forecast, including stochastic updating, in Section **Error!**  
**Reference source not found.**

295 Eq (6) is a crucial advance over the o-censored FoGSS as it allows the error model to produce  
 296 >50% zero values, which we show with a schematic in Fig 1. In conventional error models,  
 297  $z_3(t)$  can only take values  $z_3(t) \geq z_c$ . If we assume a symmetrical error distribution (e.g., Eq 5)  
 298 at most 50% of  $z_4(t)$  can be less than or equal to  $z_c$ . (Remember that any values of  $z_4(t) < z_c$   
 299 are forced to  $z_4(t) = z_c$  before back-transformation.) In the os-censored FoGSS error model,  
 300 Eq (6) re-assigns values of  $z_3(t) = z_c$  to  $z_3(t) \leq z_c$ . When noise is added according to Eq (5),  
 301 the predictive distribution can now produce >50% zeros after back-transformation. This  
 302 enables reliable predictive uncertainty to be generated for highly ephemeral streams.



303  
 304 Figure 1 Schematic of FoGSS error model compared to conventional error models, adapted  
 305 from Wang et al. (2020). Left hand example shows a conventional error model (Eq 5), where  
 306 the median transformed simulation  $z_3(t) \geq z_c$ , and thus at most 50% of  $z_4(t)$  values can be less  
 307 than or equal to  $z_c$ . Right hand shows the FoGSS error model, where Eq. 6 re-assigns a value  
 308 of  $z_3(t) \leq z_c$ , allowing >50% of  $z_4(t)$  to be less than or equal to  $z_c$ .

309 Wang et al. (2020) noted that because of the way many hydrological models are structured,  
 310 their simulations can never reach exactly zero. The GR4J model we use in this study can

311 suffer from this issue. In such cases  $z_3(t)$  is always greater than  $Tf(0)$ , meaning Eq (6) is  
312 never enacted if  $q_c = 0$ . This makes reliable predictive uncertainty impossible to generate for  
313 streams with >50% zeros. Accordingly, Wang et al. (2020) suggested using a censoring  
314 threshold slightly above zero to ensure Eq (6) is enacted. In FoGSS, however, this is  
315 unnecessary. The bias-correction (Stage 2) and AR1 model (Stage 3) can correct simulations  
316 to zero, ensuring Eq (6) is enacted frequently with a censoring threshold of  $q_c = 0$ .

317 Stage 4 parameters are estimated by Maximum Likelihood

### 318 2.2.2 Revised parameter inference for os-censored FoGSS

319 FoGSS error model parameters are estimated at lead one, as with a conventional hydrological  
320 model calibration. That is, one set of parameters is applied to all lead times when forecasts  
321 are generated. Error model parameters are inferred independently at each stage. In all cases, a  
322 numerical search method is used to find optimal parameters (Duan et al., 1993). We provide a  
323 basic overview of the parameter estimation procedure here. A more technical description,  
324 including equations and priors, is given in Appendix A.

#### 325 **Stage 1**

326 As noted in Section **Error! Reference source not found.**, in previous versions of FoGSS we  
327 have included the estimation of hydrological model parameters with the error model at Stage  
328 1. In this study, we apply FoGSS to an existing forecasting setup, where the hydrological  
329 model is already calibrated (Section 3.1). Thus in Stage 1 only the transformation parameters  
330 are estimated.

331 The Stage 1 transformation parameters are estimated by *Maximum A Posteriori* (MAP)  
332 estimation, detailed in Appendix A. Note that Stage 1 applies a single set of transformation  
333 parameters for all data, rather than varying transformations by month. This means that all  
334 forecasts and observations take a common range of values in the transform domain, greatly  
335 simplifying the AR1 modelling in Stage 3.

#### 336 **Stage 2 and Stage 3**

337 It is possible to use maximum likelihood estimation (MLE) for stages 2 and 3, using a  
338 likelihood similar to that described for Stage 4 (below, and Appendix A). However, the  
339 likelihood is quite computationally intensive. Further, we do not require an estimate of  
340 residual uncertainty at stages 2 or 3, as that is the role of Stage 4. We can therefore use a

341 simple short-cut to reduce computation at Stage 2 and Stage 3. Rather than the likelihood, we  
 342 use least squares of residuals on transformed data (abbreviated here to LST) to achieve  
 343 similar results (Chatterjee and McLeish, 1986; Li et al., 2020). As we normalize our data at  
 344 Stage 1, the assumption of normal residuals that underlies least squares regression will be  
 345 satisfied. LST reduces computation considerably. For an observed time series of length  
 346  $t = 1, 2, \dots, T$ , the LST objective minimizes

$$347 \quad L = \sum_{t=1}^T \left( \max(z_o(t), z_c) - \max(z_s(t), z_c) \right)^2 \quad (7)$$

348 where  $z_s(t)$  is the transformed simulation at stage  $S=2$  or  $S=3$ .

#### 349 **Stage 4**

350 Stage 4 parameters are estimated by MLE, using the likelihood proposed by Wang et al.  
 351 (2020). The os-censored likelihood handles four cases, depending on whether the simulation  
 352 and/or observation are equal to zero:

- 353 • case = 1: observed and simulated flow are both greater than zero, i.e.  $q_o(t) > 0$  and  $q_s(t) > 0$

$$354 \quad p(q_4(t) = q_o(t) | q_3(t)) \propto N(z_o(t) | z_3(t), \sigma^2(i)) \quad (8)$$

- 355 • case = 2: observed flow is zero and simulated flow is greater than zero, i.e.  
 356  $q_o(t) = 0$  and  $q_s(t) > 0$

$$357 \quad p(q_4(t) = q_o(t) = 0 | q_3(t)) \propto \Phi(z_c | z_3(t), \sigma^2(i)) \quad (9)$$

- 358 • case = 3: observed flow is greater than zero and simulated flow equals zero, i.e.  
 359  $q_o(t) > 0$  and  $q_s(t) = 0$

$$360 \quad p(q_4(t) = q_o(t) | q_3(t) = 0) \propto \frac{N(z_c | m_3(i), s_3^2(i) + \sigma^2(i)) \Phi\left(z_c \left| \frac{s_3^2(i) z_3(t) + \sigma^2(i) m_3(i)}{s_3^2(i) + \sigma^2(i)}, \frac{\sigma^2(i) s_3^2(i)}{s_3^2(i) + \sigma^2(i)} \right.\right)}{\Phi(z_c | m_3(i), s_3^2(i))} \quad (10)$$

- 362 • case = 4: observed and simulated flow are both zero, i.e.  $q_o(t) = 0$  and  $q_s(t) = 0$

$$363 \quad p(q_4(t) = q_o(t) = 0 | q_3(t) = 0) \propto \frac{\int_{-\infty}^{z_c} \Phi(z_c | z_3(t), \sigma^2) N(z_3(t) | m_3(i), s_3^2(i)) dz_3(t)}{\Phi(z_3(t) | m_3(i), s_3^2(i))} \quad (11)$$

365 where  $q_s(t)$  and  $z_3(t) = Tf(q_3(t))$  are hydrological model simulations at Stage 3 before and after  
 366 transformation, and  $q_o(t)$  and  $z_o(t) = Tf(q_o(t))$  are observations before and after transformation.  
 367  $z_c = Tf(q_c)$  is the transformed censoring threshold,  $m_3(i)$  and  $s_3(i)$  are the mean and standard

368 deviation of  $\mathbf{z}_3$  and  $\Phi$  is the normal cumulative probability. Note that each of  
 369  $m_3(i)$ ,  $s_3^2(i)$ , and  $\sigma^2(i)$  take different values for different calendar months (see Eq 5 and Eq  
 370 6).  $m_3(i)$  and  $s_3^2(i)$  are estimated from  $q_3(t)$  simulations generated by applying the first 3  
 371 stages of FoGSS to the hydrological model simulation (see Appendix A).

372 We denote the residual distribution parameters estimated at Stage 4 by

373  $\theta_4 = \{\sigma^2(i): i=1,2,\dots,12\}$ . The likelihood is given by:

374 
$$L(\theta_4) \propto \prod_{t: \text{case}=1} p(q_4(t)=q_o(t)|q_3(t)) \prod_{t: \text{case}=2} p(q_4(t)=q_o(t)=0|q_3(t)) \prod_{t: \text{case}=3} p(q_4(t)=q_o(t)|q_3(t)=0) \prod_{t: \text{case}=4} p(q_4(t)=q_o(t)=0|q_3(t)=0)$$

375 (12)

376 where the cases are described by equations (8)-(11). Eq (10) and Eq (11) make this likelihood  
 377 unique in the literature: they treat simulated values as censored data. As described in Section  
 378 2.2.1, the treatment of simulations as censored data at  $q_3(t)=0$  enables Eq (6) to be enacted,  
 379 making it possible to produce reliable predictive uncertainty in highly ephemeral streams  
 380 (streams that have >50% zero flow in some months).

381 We have now described the structure of FoGSS and parameter estimation procedure for os-  
 382 censoring. Before we describe the forecast generation procedure for FoGSS (Section 4) we  
 383 describe the ESP forecasts that are the key input to the FoGSS error model (Section 3).

### 384 3 ESP forecasts

385 ESP forecasts are generated with the BoM's experimental 'dynamical' seasonal streamflow  
 386 forecasting system (Feikema et al., 2018; Woldemeskel et al., 2018). This system is run at a  
 387 daily time step, and in its usual configuration outputs are aggregated to 1-month and 3-month  
 388 totals before MOS post-processing is applied. For our study, we aggregate daily forecasts to a  
 389 monthly time step to produce streamflow forecasts as 12-month time series.

390 The dynamical forecast system usually uses calibrated climate forecasts as forcing, to lead  
 391 times of ~90 days. While it is possible to generate calibrated daily forcings from climate  
 392 prediction systems (Schepen et al., 2017), these methods have not been tested to the long lead  
 393 times (365 days) we require in our study. ESP forcings are inherently reliable and simple to  
 394 generate, and as this study focusses on the development of a hydrological error model, ESP  
 395 forcings suffice. We discuss the prospects for using calibrated climate forecasts in  
 396 combination with the FoGSS error model in Section 8.

### 397 3.1 Hydrological model

398 Hydrological modelling in the BoM's dynamical forecasting system is carried out with the  
399 daily GR4J model (Perrin et al., 2003), a simple four-parameter conceptual hydrological  
400 model that has performed strongly in Australian catchments in model intercomparison studies  
401 (Bennett et al., 2016a; Coron et al., 2012). It is forced by rainfall and potential evaporation  
402 (PE).

403 For each gauge, GR4J is calibrated by minimising the sum of squared errors computed on  
404 Box-Cox transformed flows (McInerney et al., 2017), following the BoM's existing practice.  
405 Calibration is carried out under cross-validation, as described in Section 6.1. To enable  
406 FoGSS parameters to be estimated, we run GR4J in simulation mode – i.e., forced with  
407 observed rainfall and PE - and aggregate daily GR4J simulations to monthly data to generate  
408  $q_1$ . When running these simulations, we warm up GR4J for a minimum of 5 years (Section  
409 3.3).

### 410 3.2 Rainfall and potential evaporation sampling

411 Rainfall and PE forcings in ESP forecasts are taken from historical observations. We sample  
412 ESP forcings with a leave-4-years-out cross-validation procedure (i.e., including the target  
413 year), which we illustrate by an example. To generate a forecast for January-Dec in 1980, we  
414 sample daily January-Dec rainfall and PE sequences from 1985, 1986, ..., 2008. Each  
415 sequence is 1 year long. We construct sequences from 29 years of data (1980-2008). We  
416 choose 1980-2008 because this is the standard period used by the BoM to assess seasonal  
417 streamflow forecasting products. The choice of leaving out 4 years is less stringent than the  
418 cross-validation used for the hydrological component of the system (Section 5.1). Rainfall  
419 has far less interannual memory than streamflow, meaning that a less stringent cross-  
420 validation is acceptable. (The specific choice of leaving out four years was taken to generate  
421 a 25-member ensemble for both rainfall and PE. We ultimately generate 1000-member  
422 ensembles with FoGSS, and this process was simplified for an ESP ensemble of 25 members,  
423 because 1000 is evenly divisible by 25.)

424 PE and rainfall ensemble members are paired: i.e., if a rainfall ensemble member was  
425 sampled from 1985, the matching PE member is also from 1985. Rainfall and PE are often  
426 anticorrelated, so pairing rainfall and PE ensemble members is necessary to ensure realistic  
427 outputs from the hydrological model.

### 428 3.3 Generating ESP forecasts

429 We issue a forecast at the beginning of every calendar month. To generate a forecast, we  
 430 initialise GR4J states by running it with observed forcings from 1/1/1975 onwards (a  
 431 minimum of 5 years warmup for our reforecast period of 1980-2008). At the forecast issue  
 432 time, the hydrological model is then forced with a single member of the rainfall/PE ensemble  
 433 to produce a 1-year forecast of streamflow at the daily time step. This process is repeated for  
 434 all 25 ensemble members to produce  $\mathbf{q}_{F1}$ . We then aggregate the daily forecasts to the  
 435 monthly time step before applying the FoGSS error model.

### 436 4 Generating an os-censored FoGSS forecast

437 For a given forecast issue time,  $t$ , we assume we have available an initialised hydrological  
 438 model simulation  $q_i(t)$  at the forecast issue date, together with an observation  $q_o(t)$ . We also  
 439 have available an uncorrected hydrological forecast  $\mathbf{q}_{F1} = \{q_{F1}(t+1), \dots, q_{F1}(t+\tau), \dots, q_{F1}(t+12)\}$   
 440 for lead times  $\tau = 1, 2, \dots, 12$  (Section 3). We apply FoGSS to each ensemble member  
 441 separately, where  $\mathbf{q}_{F1}$  takes hydrological states from the hydrological simulation  $q_i(t)$  and is  
 442 then forced with a single member from the ensemble climate forecast. We apply the os-  
 443 censored FoGSS error model to  $\mathbf{q}_{F1}$  as follows:

- 444 1.  $\mathbf{q}_{F1}$  is transformed (Stage 1 - Eq 1) and bias-corrected (Stage 2 - Eq 2) to produce  
 445  $\mathbf{z}_{F2} = \{z_{F2}(t+1), z_{F2}(t+2), \dots, z_{F2}(t+12)\}$ . In the same way,  $q_i(t)$  is transformed and bias-  
 446 corrected to produce  $z_2(t)$ . Then  $q_o(t)$  is transformed to produce  $z_o(t)$ .
- 447 2. For  $\tau = 1$ , we apply Eq (3) (Stage 3) as

$$448 \quad z_{F3}(t+\tau) = z_{F2}(t+\tau) + \rho(i)(z_o(t) - z_2(t)) \quad (13)$$

449 where  $i$  is the calendar month at time  $t+\tau$ . We also apply the restriction (Eq 4) but  
 450 we omit this equation for brevity.

- 451 3. If  $z_{F3}(t+\tau) \leq z_c$  we assign a new value to  $z_{F3}(t+\tau)$  by applying Eq (6):

$$452 \quad \begin{aligned} z_{F3}(t+\tau) &= \Phi^{-1}\left(\Phi\left(0 \mid m_3(i), s_3(i)^2\right) \times r_o(t+\tau)\right) \\ r_o(t+\tau) &\sim U(0,1) \end{aligned} \quad (14)$$

- 453 4. We then draw a single realisation of noise according to Eq (5) (Stage 4):

454 
$$\begin{aligned} z_{F_4}(t+\tau) &= z_{F_3}(t+\tau) + \varepsilon(t+\tau) \\ \varepsilon(t+\tau) &\sim N(0, \sigma^2(i)) \end{aligned} \quad (15)$$

455 5. For lead times  $\tau = 2, 3, \dots, 12$  we no longer have an observation at the previous time  
456 step. We resolve this by substituting  $z_{F_4}$  from Eq (15) for  $z_o$  in Eq (13) to give:

457 
$$z_{F_3}(t+\tau) = z_{F_2}(t+\tau) + \rho(i)(z_{F_4}(t+\tau-1) - z_{F_2}(t+\tau-1)) \quad (16)$$

458 where  $i$  is the calendar month at time  $t+\tau$ . (We again apply the restriction given by  
459 Eq 4, again omitted for brevity.) We term the substitution of  $z_{F_4}$  for  $z_o$  *stochastic*  
460 *updating*, which we discuss in more detail below.

461 6. Apply Eq (14) and Eq (15) to  $z_{F_3}(t+\tau)$  from Step 5 (i.e., steps 3 and 4 are repeated for  
462 values of  $\tau > 1$ ).

463 7. Repeat steps 5 and 6 for each  $\tau$  to produce the streamflow forecast ensemble member

464 
$$\mathbf{z}_{F_4} = \{z_{F_4}(t+1), \dots, z_{F_4}(t+\tau), \dots, z_{F_4}(t+12)\}.$$

465 8.  $\mathbf{z}_{F_4}$  is converted to the original domain by reversing Eq (1):

466 
$$\mathbf{q}_{F_4} = Tf^{-1}(\max(\mathbf{z}_{F_4}, z_c)) \quad (17)$$

467 We repeat steps 1-8 for each uncorrected ensemble member  $\mathbf{q}_{F_1}$  that has been generated with  
468 in the ESP ensemble. Note that it is straightforward to generate large ensembles by repeating  
469 steps 1-8 as many times as necessary for each  $\mathbf{q}_{F_1}$  - i.e., we can have many streamflow  
470 forecast ensemble members ( $\mathbf{q}_{F_4}$ ) for each climate forecast ensemble member ( $\mathbf{q}_{F_1}$ ). As noted  
471 in Section 3, the ESP forecasts have 25 ensemble members. Steps 1-8 are repeated 40 times  
472 for each ensemble member to produce a 1000-member streamflow forecast ensemble.

473 Stochastic updating (Step 5) efficiently propagates uncertainty through the forecast (Bennett  
474 et al., 2016b; Li et al., 2017). Estimating forecast uncertainty correctly at lead 1 is  
475 straightforward, because FoGSS parameters are estimated at lead 1. However, uncertainty  
476 should grow appropriately through lead times. The interaction of the random noise (Eq 15)  
477 and the AR model (Eq 16) in stochastic updating allows uncertainty to grow through the  
478 forecast. Uncertainty tends to grow rapidly at short lead times, but this growth slows and  
479 ultimately stops as the influence of the AR model recedes at longer lead times. This is what  
480 we are seeking in our uncertainty range: narrow uncertainty at short lead times when  
481 forecasts are skillful, widening uncertainty as skill recedes with lead time, and an

482 approximately static uncertainty range (approaching a climatological distribution) when skill  
483 is absent. Stochastic updating is able to achieve this whilst retaining the temporal features of  
484 the hydrograph generated by GR4J.

## 485 5 Data

486 We assess ESP and FoGSS forecasts at 50 sites from a range of climates around Australia  
487 (Fig 2). We classify these sites into three categories of ephemerality based on observed flows  
488 at the monthly time step:

- 489 1) perennial (no month has >5% zero flows) – 24 sites
- 490 2) moderately ephemeral (some months have >5% zero flow, with no months >50%  
491 zero flow) – 20 sites
- 492 3) highly ephemeral (some months have >50% zero flow) – 6 sites

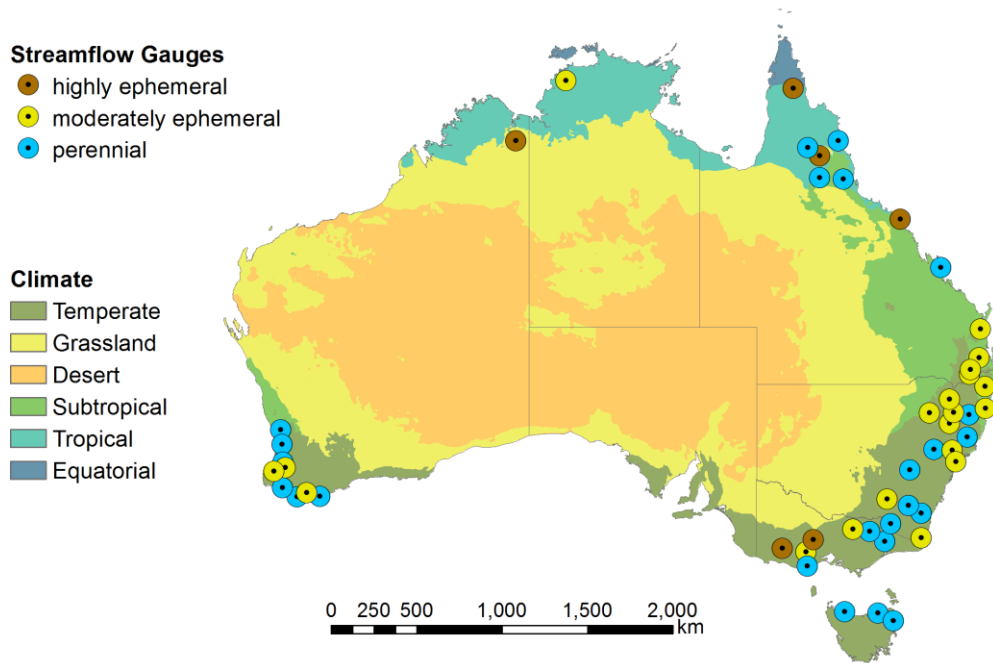
493 From these 50 sites, we select the 3 sites (one from each category) shown in Fig 3 to illustrate  
494 model performance in more detail. Note that the degree of ephemerality – in the sense of zero  
495 flows – is not replicated in the GR4J simulations. While it is technically possible for GR4J to  
496 produce zero flows for particular combinations of states/parameters, in practice this rarely  
497 occurs – i.e., like a lot of hydrological models, it tends not to produce zeros. For the 26  
498 ephemeral catchments tested here, simulations were never zero.

499 Daily flow data are extracted from Water Data Online (<http://www.bom.gov.au/waterdata/>).

500 Rainfall and potential evaporation (PE) are taken from the gridded AWAP data set  
501 (Australian Water Availability Project; Jones et al., 2009; <http://www.csiro.au/awap/>).

502 AWAP produces daily estimates of rainfall interpolated from gauges to a ~5-km grid. PE  
503 estimates from AWAP are at a monthly time step; we disaggregate these to daily estimates by  
504 simple linear interpolation. Catchment estimates of rainfall and PE are calculated by areal  
505 averaging of AWAP grid cells that intersect with each catchment.

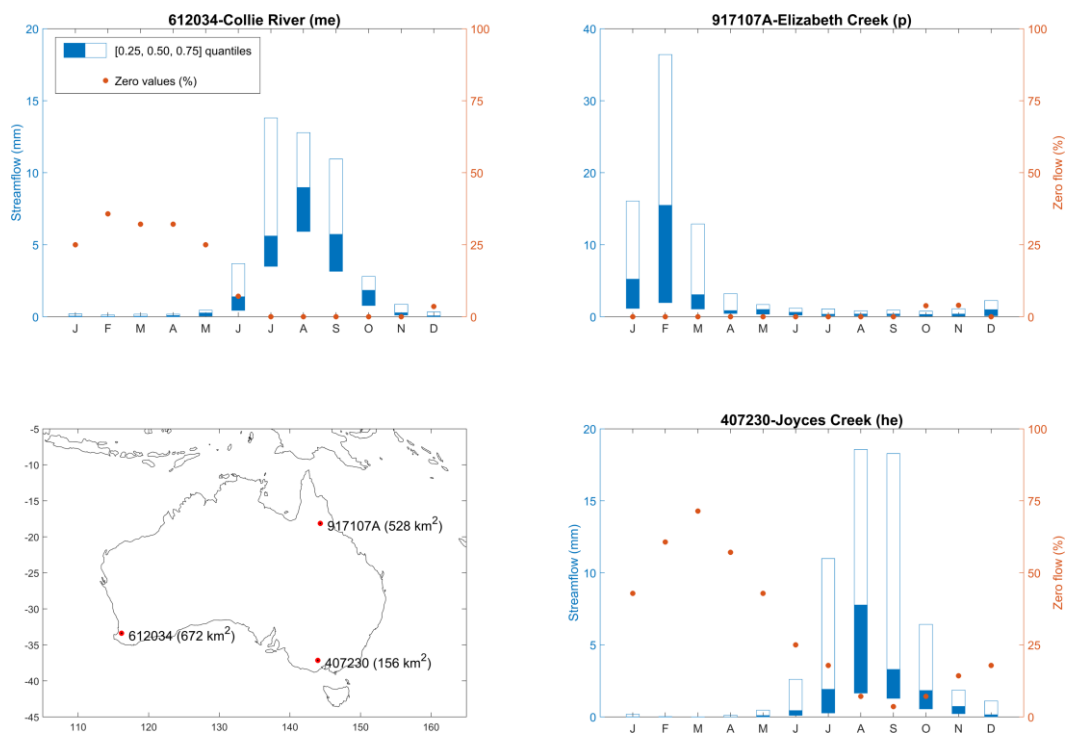
506 We use data for the period 01/01/1980-31/12/2012, which covers the 1980-2008 period used  
507 by the Bureau of Meteorology to assess all its seasonal streamflow forecasting products. We  
508 use the 01/01/1980-31/12/2012 period to estimate parameters (Section 2.2.2) and generate  
509 climatology reference forecasts (Section 6.2) under buffered cross-validation (Section 6.1) to  
510 ensure a consistent 4-year buffer for all years.



511

512 Figure 2 Location map of 50 sites used in this study, classified by degree of ephemerality (see  
 513 text for details). Climate zones are from Köppen climate classifications (Stern et al., 2000)

514



515

516 Figure 3 Exemplar catchments for moderately ephemeral (me), perennial (p) and highly  
 517 ephemeral (he) flow regimes

## 518 6 Forecast verification

### 519 6.1 Cross-validation scheme

520 All forecast verification is carried out under a buffered leave-one-year-out cross-validation  
521 scheme, with a buffer of 4 years. The cross-validation scheme is best described with an  
522 example. To attain GR4J and FoGSS parameters for 2000, we omit data from 2000 and the  
523 succeeding 4 years (2000-2004). The buffer is necessary to avoid informing the parameter  
524 estimation with rainfall information from the target year (2000), which can influence  
525 observed streamflow in subsequent years through catchment memory.

### 526 6.2 Verification scores

527 We focus on two probabilistic verification metrics: the probability integral transform (PIT) to  
528 measure ensemble reliability and the Continuous Ranked Probability Score (CRPS) to  
529 measure forecast accuracy (Gneiting and Katzfuss, 2014). Reliability measures the  
530 appropriateness of the ensemble spread: it should not be too wide (underconfident) nor too  
531 narrow (overconfident). CRPS measures both accuracy and reliability, but in the context of  
532 this study it is not strongly sensitive to reliability. This is because the benefits of the new (os-  
533 censored) FoGSS method are mainly that it produces more reliable forecasts at very low  
534 flow. Because CRPS is averaged over a range of forecasts, it tends to be more sensitive to  
535 errors in larger events than to the reliability of low flow months.

536 A PIT value is calculated for each forecast by

$$537 \quad p(t) = \begin{cases} F(t, q_o(t)) & q_o(t) > 0 \\ U(0,1) \times F(0) & q_o(t) = 0 \end{cases} \quad (18)$$

538 where  $F(t, [ ])$  is the cumulative distribution function (CDF) of the forecast ensemble at time  
539  $t$ . A set of forecasts at  $t = 1, 2, \dots, T$  is reliable if PIT values  $\mathbf{p} = \{p(1), p(2), \dots, p(T)\}$  are  
540 uniformly distributed. We check this by plotting PIT values against a standard uniform  
541 variate (we refer to these plots as ‘PIT plots’ for brevity). When PIT values follow the  
542 diagonal in PIT plots, the forecasting system is perfectly reliable. The treatment of PIT values  
543 at  $q_o = 0$  in Eq (18) is necessary to allow PIT values to follow a uniform distribution in the  
544 presence of zero values (Wang and Robertson, 2011). We refer to PIT values calculated when  
545  $q_o = 0$  as ‘pseudo-PIT’ values.

546 Uniformity of PIT values can be summarised with the  $\alpha$ -index (Renard et al., 2010). The  $\alpha$ -  
 547 index describes the tendency of PIT values to deviate from the diagonal in PIT plots:

$$548 \quad \alpha = 1 - \frac{2}{N} \sum_{t=1}^T |p(t) - p_U(t)| \quad (19)$$

549  $p_U(t)$  is the theoretical value corresponding to  $p(t)$ .  $\alpha$  ranges between 0 (unreliable) and 1  
 550 (perfectly reliable). Because it reduces PIT diagrams to a single value,  $\alpha$  allows easy  
 551 comparison between catchments and lead times.

552 CRPS is given by

$$553 \quad C_M = \frac{1}{T} \sum_{t=1}^T \int_{-\infty}^{\infty} (F(t, x) - H(q_o(t) \leq x))^2 dx \quad (20)$$

554 where  $F(t, [ ])$  is the CDF of the forecast ensemble at time  $t$ , and  $H$  is the Heaviside step  
 555 function. In this study we wish to compare the performance of ESP forecasts with  $m = 25$   
 556 ensemble members to FoGSS forecasts with  $M = 1000$  ensemble members. When the CDF of  
 557 the forecast is estimated empirically from an ensemble, CRPS calculations are sensitive to  
 558 ensemble size, with greater errors occurring in smaller ensembles (e.g. Zamo and Naveau,  
 559 2018). Ferro et al. (2008) derived an unbiased estimator for CRPS for cases where ensemble  
 560 members are exchangeable:

$$561 \quad C_m = \frac{M(m+1)}{m(M+1)} C_M \quad (21)$$

562 where  $C_M$  is calculated on ensembles of  $M = 1000$  with Eq (20), and  $C_m$  is an estimate for the  
 563 smaller ensemble size  $m = 25$ .

564 CRPS is commonly presented as a skill score, where forecast accuracy in a forecasting  
 565 system is calculated against a reference forecast:

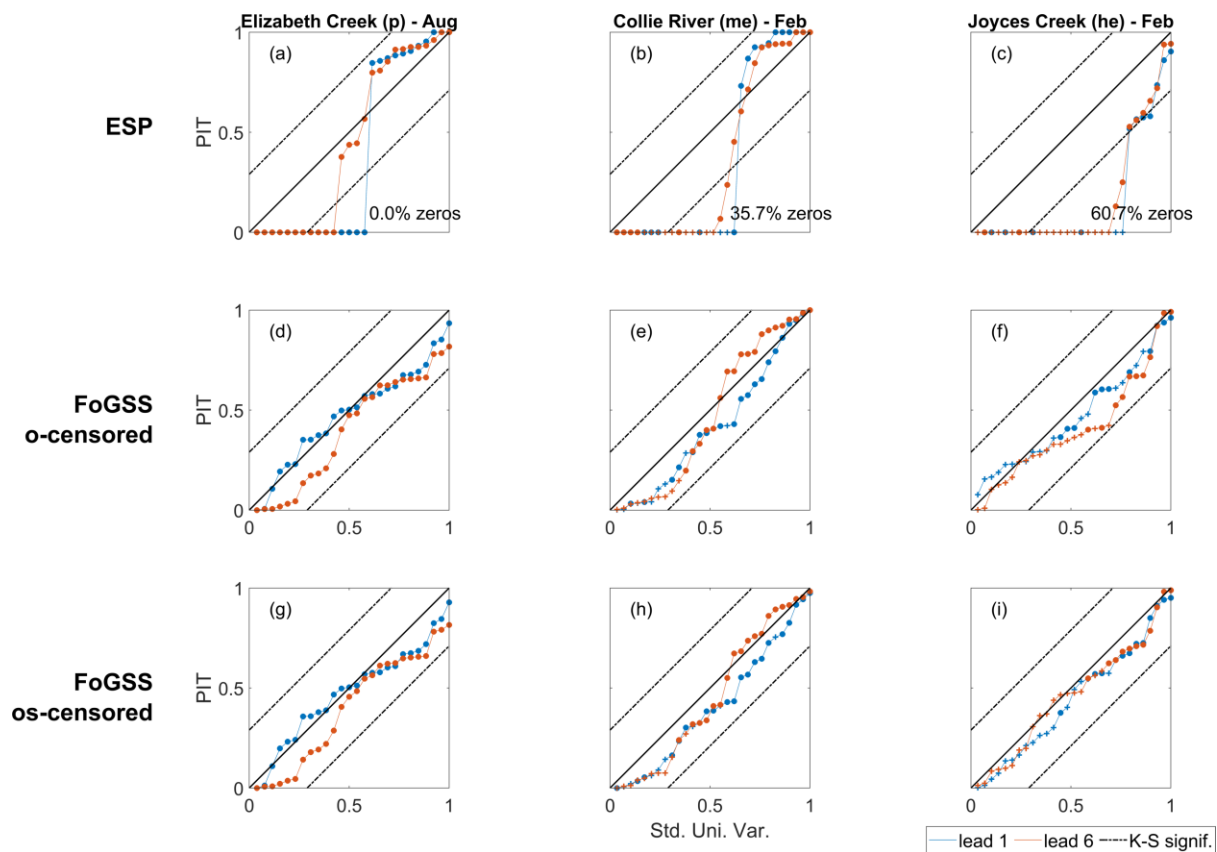
$$566 \quad CRPSS = 1 - \frac{C_{m,Fct}}{C_{m,Ref}} \quad (22)$$

567 CRPSS ranges from  $-\infty$  (worst performance) to  $\infty$  (best performance), with values near zero  
 568 indicating the forecast performs similarly to the reference. We use climatology forecasts as  
 569 our reference. Climatology forecasts are generated by randomly drawing 1000 realisations  
 570 from a log-sinh transformed normal distribution fitted to each calendar month. We use the

571 BJP to fit the log-sinh transformation, following the data and cross-validation scheme  
 572 described in sections 5 and 6.1. The BJP generates reliable distributions even where many  
 573 zeros are present (Wang and Robertson, 2011). We also assess the skill of volume forecasts  
 574 accumulated over multiple months. To calculate the climatology reference forecasts in these  
 575 cases, we first accumulate volumes from the observed record, and then generate a  
 576 climatology as described above.

577 To assess the significance of skill, we bootstrap equations (20)-(22) with 500 repeats. When  
 578 97.5% of bootstrapped CRPSS values are positive (negative), we consider the forecasts to be  
 579 significantly skilful (negatively skilful).

580



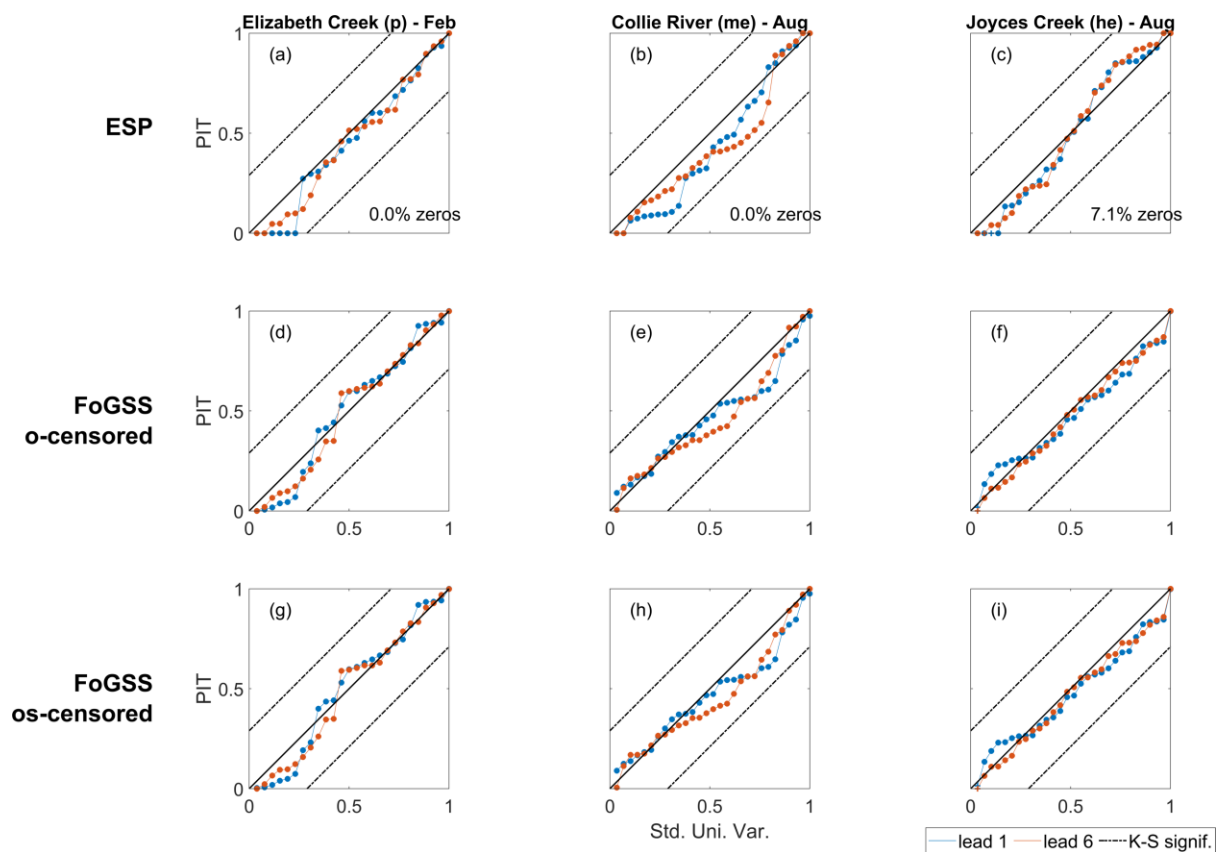
581

582 Figure 4 PIT uniform probability plots for forecasts for a dry month in example catchments  
 583 (columns). Insets in (a)-(c) show proportion of zero flows in that month. Points show PIT  
 584 values, crosses show pseudo-PIT values. Dashed lines give 95% Kolmogorov-Smirnoff  
 585 confidence intervals. Top row shows ESP forecasts, middle row shows forecasts generated  
 586 with the previous FoGSS method (o-censored) and bottom row shows new FoGSS method  
 587 (os-censored).

588 7 Results

589 7.1 Reliability

590 Reliability of ESP and FoGSS forecasts for a dry month is shown in Fig 4 for our three  
 591 example catchments (all 50 catchments are presented in Fig S1). Perfectly reliable forecasts  
 592 follow the 1-1 line in the PIT diagrams. As foreshadowed in the introduction, ESP forecasts  
 593 are often overconfident (the PIT diagram is s-shaped) because they do not consider  
 594 uncertainties in the conversion of rainfall to runoff (e.g., uncertainties in model states,  
 595 structure and/or streamflow observations). Overconfidence is particularly prevalent for the  
 596 dry months when rainfall is low and streamflow results primarily from catchment stores  
 597 draining rainfall accrued in previous (often wetter) months. This means forcing uncertainty  
 598

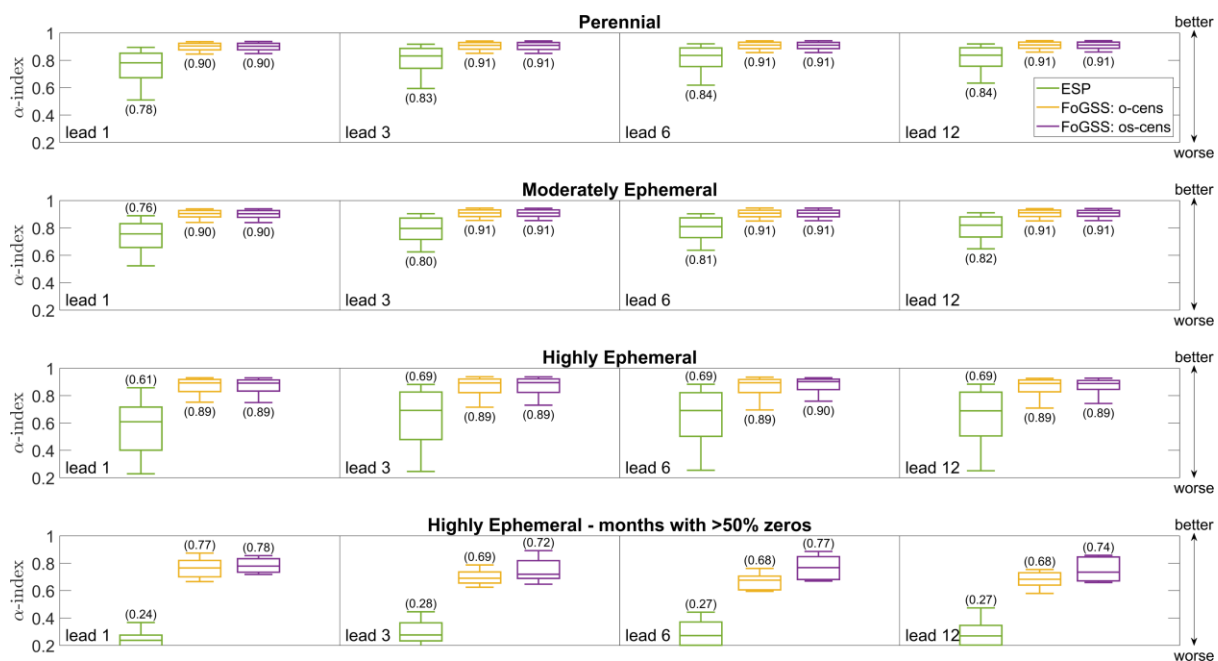


599

600 Figure 5 PIT uniform probability plots for forecasts for a wet month in example catchments  
 601 (columns). Insets in (a)-(c) show proportion of zero flows in that month. Points show PIT  
 602 values, crosses show pseudo-PIT values. Dashed lines give 95% Kolmogorov-Smirnoff  
 603 confidence intervals. Top row shows ESP forecasts, middle row shows forecasts generated  
 604 with the previous FoGSS method (o-censored) and bottom row shows new FoGSS method  
 605 (os-censored).

606 has little bearing on overall uncertainty. At shorter lead times, forecast uncertainty in dry  
 607 months is dominated by uncertainty in initial hydrological conditions. ESP reliability is  
 608 considerably better for wetter months (Fig 5) than for drier months. In wetter months rainfall  
 609 forcings are often a dominant source of uncertainty, and thus ESP ensembles do a reasonably  
 610 good job of representing total uncertainty.

611 Both previous and new FoGSS methods (o-censoring and os-censoring, respectively) correct  
 612 overconfidence in the ESP ensembles in many cases, producing reliable forecasts at short and  
 613 long lead times in most instances. Improvements are most striking in dry months (Fig 4),  
 614 where poor reliability in ESP is markedly improved by FoGSS. Even in wetter months,  
 615 however, we can see small improvements in reliability after FoGSS is applied. At lead 1 in  
 616 the perennial Elizabeth Creek catchment (Fig 5a), six observations are below the ESP  
 617 forecast ensemble (these points fall on the bottom of the vertical axis). While these are  
 618 technically within our confidence intervals, to have six forecasts (of 29) that completely miss  
 619 an observation is evidence of a poorly performing forecast. FoGSS completely corrects this  
 620 problem (Fig 5d).

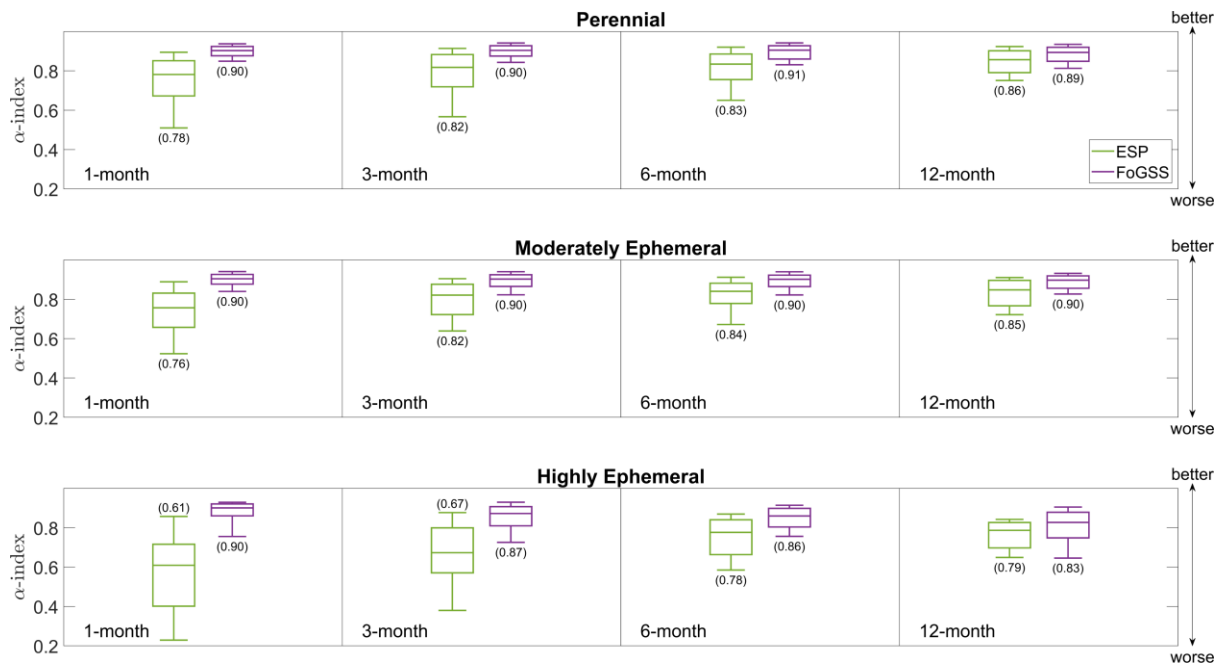


621  
 622 Figure 6 Alpha index with lead time summarised in box plots for all month in perennial (top),  
 623 moderately ephemeral (upper middle) and highly ephemeral (lower middle) catchments.  
 624 Bottom panel shows performance only for months with >50% zeros. Columns show lead  
 625 times. Boxes give interquartile range and median, whiskers show 90% range. Numbers in  
 626 brackets give median value for the corresponding box.

627 The benefits of os-censoring over o-censoring are apparent at longer lead times in dry months  
628 in the highly ephemeral Joyce catchment (Fig 4f,i). FoGSS with os-censoring produces  
629 reliable ensemble forecasts in the highly ephemeral Joyce River catchment in a month with  
630 >60% zeros (Fig 4f), regardless of lead time. For o-censoring, at lead 1 the forecasts are  
631 slightly underconfident (ensembles are too wide), signified by the transposed s-shape of the  
632 PIT diagram. At lead 6 the forecasts are increasingly positively biased; a result of  
633 accumulated instances of underconfident forecasts at multiple lead times. We note that o-  
634 censoring still performs reasonably well in this catchment, but this is less true of other  
635 catchments, as discussed in the next paragraph.

636 The improved performance of os-censoring over o-censoring becomes more evident in highly  
637 ephemeral months when we summarise reliability at a range of lead times (Fig 6). Both o-  
638 censored and os-censored versions of FoGSS markedly outperform ESP at short-lead times,  
639 especially in highly ephemeral rivers. ESP forecasts tend to become more reliable (closer to  
640 the ideal value of 1) with lead time because uncertainties from the forcings explain a large  
641 proportion of total uncertainty at long lead times. The benefits of os-censoring over o-  
642 censoring are marked in months with >50% zeros at longer lead times (Fig 6, bottom panel).  
643 Os-censoring allows FoGSS to maintain reliability in even highly ephemeral months to long  
644 lead times, whereas the reliability of o-censored FoGSS forecasts drops with lead time, as  
645 instances of underconfidence/positive bias accumulate. Consistent with Wang et al. (2020),  
646 we do not expect (and do not see) strong differences in skill or the reliability of accumulated  
647 volumes between o-censored and os-censored FoGSS forecasts; we therefore concentrate on  
648 os-censored forecasts in the remaining figures.

649 As stated in the introduction, a key advantage of forecasts in the form of time series is that  
650 they can be summed to produce total volume forecasts. Reliability of total volume forecasts is  
651 not ensured by reliability at individual lead times; it can only be guaranteed if hydrographs in  
652 the ensemble have realistic temporal properties (Demargne et al., 2014). Fig 7 shows that  
653 ESP forecasts can achieve reasonably reliable 12-month aggregations, but reliability falls  
654 away for shorter aggregation periods, irrespective of ephemerality. FoGSS forecasts of  
655 aggregated values are strongly reliable, and reliability is consistent across different  
656 aggregation periods (and across different catchments). Reliability of aggregated volumes in  
657 highly ephemeral catchments is slightly less than that of perennial and moderately ephemeral  
658 rivers, but is still high ( $\alpha > 0.8$  in a large majority of cases, where  $\alpha = 1$  is perfect  
659 reliability).



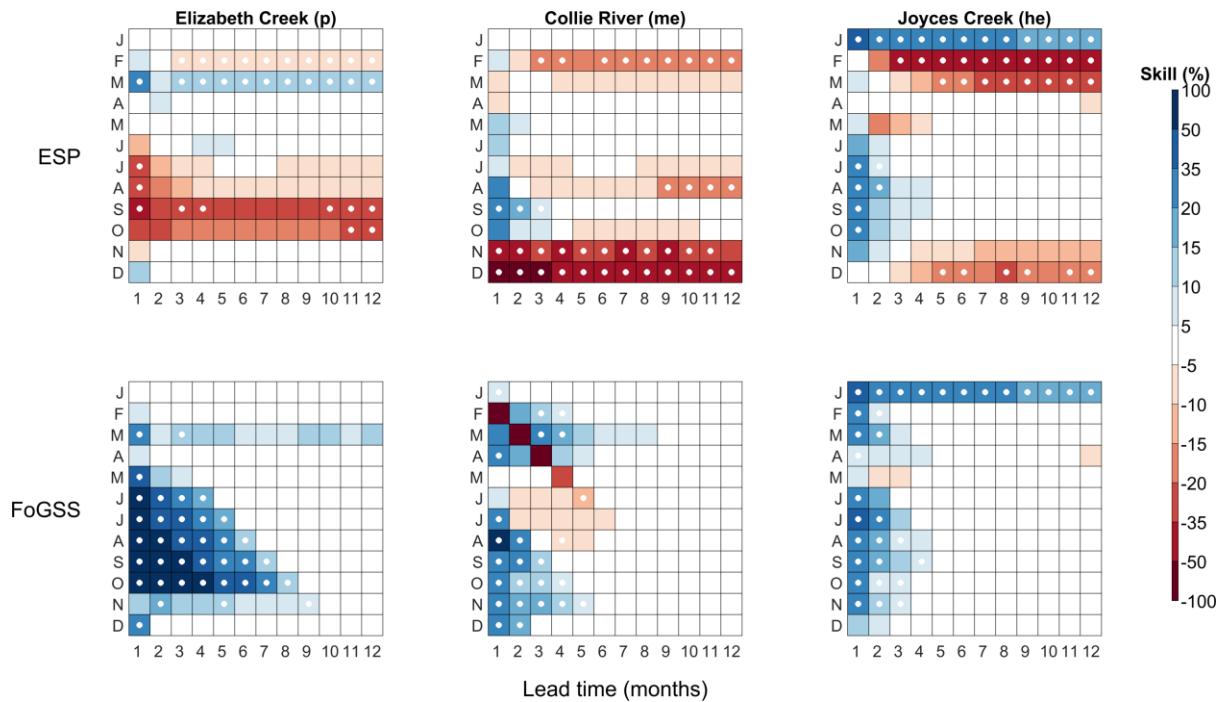
660

661 Figure 7 Alpha index of lead 1 forecasts of aggregated volumes. Alpha indices are  
 662 summarised in box plots for forecasts in perennial (top), moderately ephemeral (middle) and  
 663 highly ephemeral (bottom) catchments. Columns show aggregated volumes over different  
 664 periods. Colours show ESP (green) and os-censored FoGSS (purple) forecasts. Boxes give  
 665 interquartile range and median, whiskers show 90% range. Numbers in brackets give median  
 666 value for the corresponding box.

667

## 668 7.2 Accuracy and skill

669 While ESP forecasts can be skilful to multiple lead times, they are beset by statistically  
 670 significant negative skill in many months and lead times, irrespective of catchment  
 671 ephemerality (Fig 8). Significant negative skills are present in Jul-Oct and Feb in Elizabeth  
 672 Creek; Aug, Nov, Dec and Feb in the Collie River; and Dec, Feb and Mar in Joyces Creek.  
 673 FoGSS forecasts, by contrast, virtually always produce positively skillful forecasts at short  
 674 lead times and neutrally skilful forecasts (i.e., similarly skilful to climatology) at longer lead  
 675 times – the one exception in June at lead 5 in the Collie catchment, where os-censored  
 676 FoGSS produces significant negative skill but the ESP forecasts do not. This is consistent  
 677 with the findings of our previous work in perennial and moderately ephemeral catchments  
 678 (Bennett et al., 2016b; Bennett et al., 2017). It confirms that stages 2 and 3 of FoGSS  
 679 substantially improve biases and reduce errors at short lead times, respectively.

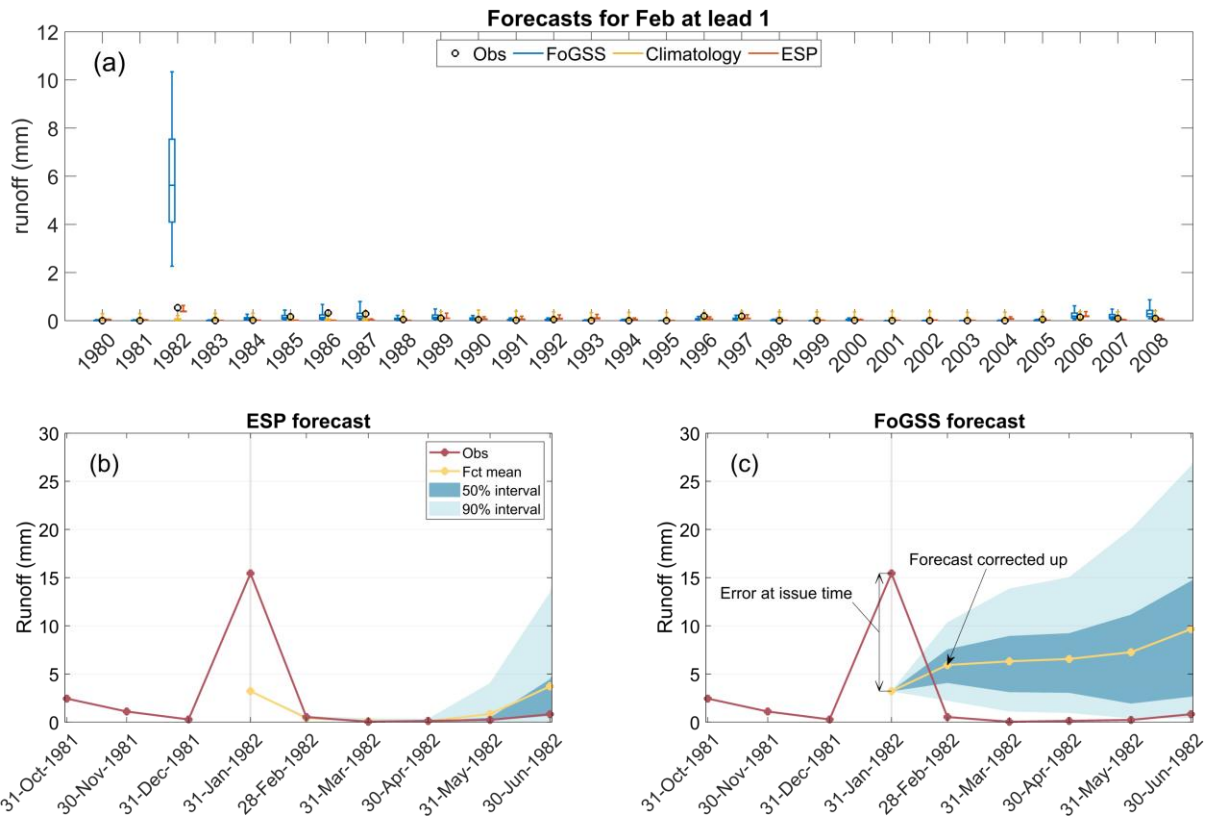


680

681 Figure 8 Continuous ranked probability skill scores of target months for three example  
 682 catchments. Ephemerality is denoted by p-perennial (left column), me-moderately ephemeral  
 683 (middle column) and he-highly ephemeral (right column). White points indicating statistical  
 684 significance at the 5% level under bootstrapping (see text for details). Top row shows ESP;  
 685 bottom row shows os-censored FoGSS. See text for discussion of prominent features.

686

687 Fig 8 also shows that os-censored FoGSS forecasts for the highly ephemeral Joyces Creek  
 688 catchment are now positively or neutrally skilful in all months. This is also true of other  
 689 highly ephemeral catchments: FoGSS improves skill of ESP forecasts for >85% of highly  
 690 ephemeral months used in this study (Fig S4). Indeed FoGSS almost always improves ESP  
 691 forecasts, often substantially, with few exceptions (Fig 8, Figs S3 & S4). For example,  
 692 negative skills in the perennial Elizabeth Creek in Aug-Oct are completely removed (Fig 8),  
 693 and earlier lead times for these months are now significantly skilful. Similar improvements  
 694 are evident for Nov-Dec in the moderately ephemeral Collie River and for Feb-Mar in the  
 695 highly ephemeral Joyce Creek.



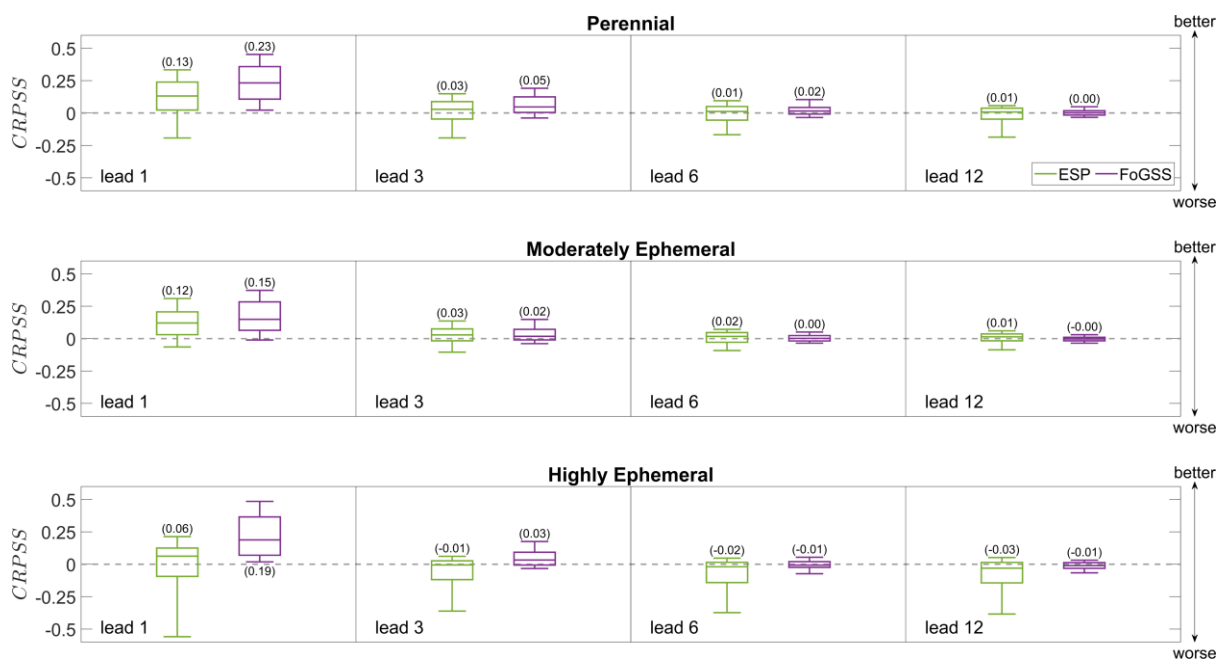
696

697 Figure 9 Cause of negative skill scores for Lead 1 forecasts for February in the Collie River.  
 698 (a) Shows that negative forecast skill is caused by a single poor forecast in 1982; (b) shows  
 699 ESP forecast issued for Feb 1982 and (c) shows os-censored FoGSS forecast issued for Feb  
 700 1982, and how autoregressive updating causes this poor forecast (see text for details).

701 Some positive skills are caused by poor performance of the reference forecast. A notable  
 702 example is the significant positive skill of both ESP and FoGSS forecasts for January for  
 703 Joyces Creek, even to very long lead times. In this case, the BJP produces a very wide  
 704 climatology ensemble for January. This is caused by one very large event in the streamflow  
 705 record (in what is usually a dry month), which affects the fit of the log-sinh transformation by  
 706 the BJP. FoGSS is less affected, because in FoGSS we fit the transformation to data from all  
 707 months. The BJP's poor performance in this case is highly unusual: more often, fitting  
 708 climatology distributions by month leads to more accurate forecasts. Further, fitting the  
 709 transformation to all months in FoGSS sometimes leads to a poorly fitting transformation at  
 710 individual months. In exceptional cases this can result in negative skill in FoGSS forecasts  
 711 (e.g. forecasts for Nov in the 922001A gauge, Fig S4).

712 There is one notable exception to FoGSS improving ESP forecasts in the Collie River:  
 713 strongly negative skill has been introduced by FoGSS at lead 1 in Feb and follows a diagonal

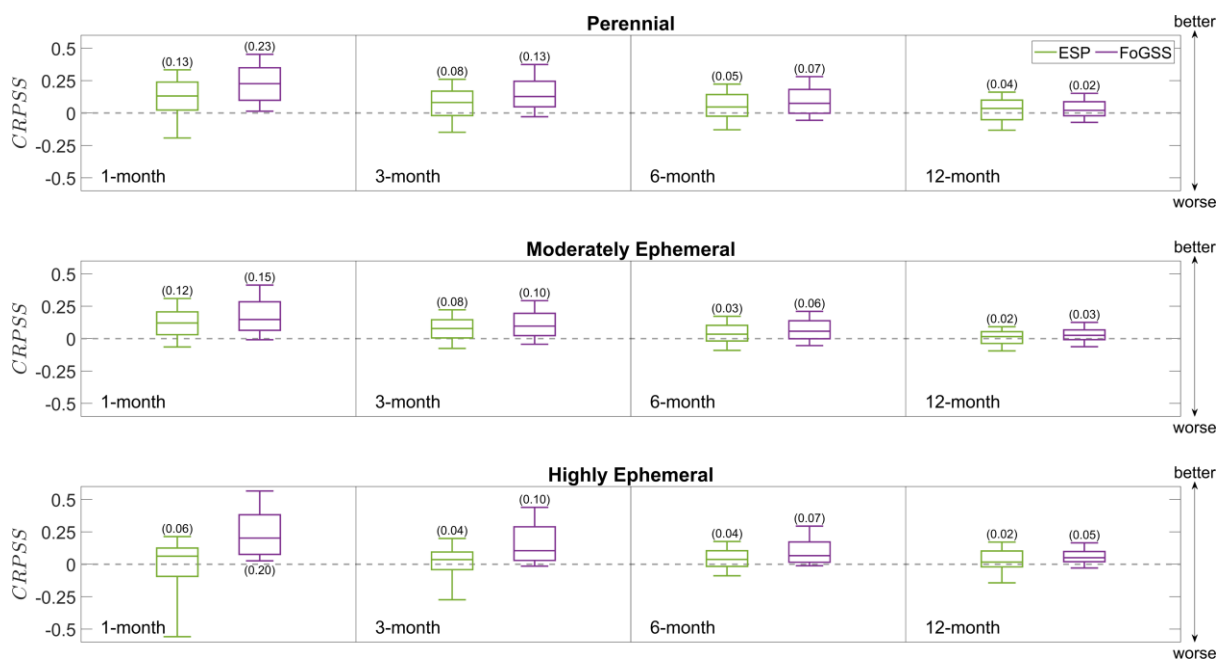
714 to Lead 4 in May. These negative skills are not statistically significant, but demonstrating  
 715 their cause is instructive. Fig 9 shows how negative forecast skill is introduced by FoGSS for  
 716 the Collie River. Negative skill for the lead 1 Feb forecasts is caused by a single large  
 717 overestimation of streamflow in 1982 (Fig 9a), and this also causes negative skill at lead 2  
 718 Mar forecasts through to lead 4 May forecasts. The large overestimation in Feb 1982 is due to  
 719 the underestimation of a very large and unseasonal streamflow event in Jan 1982 driven by  
 720 Tropical Cyclone Bruno (<http://www.bom.gov.au/cyclone/history/bruno.shtml>). Fig 9b shows  
 721 that this has virtually no impact on the ESP forecast for lead 1 Feb (or the subsequent lead  
 722 times). The AR1 model in FoGSS, however, corrects forecasts for Feb 1982 up dramatically  
 723 (Fig 9c), following the mismatch in simulations and observations in Jan 1982. Under strict  
 724 cross-validation, poor performance of a calibrated hydrological model can be difficult to  
 725 avoid for isolated extreme events. In addition, the AR1 model in FoGSS exacerbates the poor  
 726 performance of the hydrological model in this instance. However, while the AR1 model  
 727 exacerbates poor model performance in this instance (as well as a few other isolated instances  
 728 in Fig S4), the AR1 model overwhelmingly improves forecast accuracy in the vast majority  
 729 of cases (e.g. in Elizabeth Creek for May-Oct).



730  
 731 Figure 10 Continuous Ranked Probability Skill scores with lead time summarised in box  
 732 plots for all months in perennial (top), moderately ephemeral (middle) and highly ephemeral  
 733 (bottom) catchments. Columns show lead times. Colours show ESP (green) and os-censored  
 734 FoGSS (purple) forecasts. Boxes give interquartile range and median, whiskers show 90%  
 735 range. Numbers in brackets give median value of each box.

736 FoGSS clearly improves the accuracy of ESP forecasts across the 50 catchments tested in this  
 737 study (Fig 10). ESP forecasts for highly ephemeral rivers are noticeably less accurate than  
 738 forecasts for perennial or moderately ephemeral rivers at all lead times. Conversely, FoGSS  
 739 forecasts for highly ephemeral catchments tend to be similarly accurate as those for perennial  
 740 and moderately ephemeral rivers. Crucially, at long lead times FoGSS forecasts tend to be  
 741 neutrally skillful (i.e., similarly skillful to climatology), whereas ESP forecasts can be  
 742 negatively skillful, in particular for highly ephemeral catchments.

743 Removing negative skills at long lead times has considerable benefits for the accuracy of  
 744 aggregated volume forecasts (Fig 11). Skill at short lead times can carry through the entire  
 745 volume forecast, resulting in FoGSS forecasts that can be skillful even for 12-month  
 746 aggregations, including in highly ephemeral catchments. FoGSS forecasts for 6-month  
 747 aggregations are very often skillful, irrespective of ephemerality. As with individual lead  
 748 times, FoGSS forecasts of aggregated volumes are almost always more skillful than ESP  
 749 forecasts of aggregated volumes.



750

751 Figure 11 Continuous Ranked Probability Skill scores of lead 1 forecasts of aggregated  
 752 volumes summarised for all months for perennial (top), moderately ephemeral (middle) and  
 753 highly ephemeral (bottom) catchments. Columns show aggregation periods. Colours show  
 754 ESP (green) and os-censored FoGSS (purple) forecasts. Boxes give interquartile range and  
 755 median, whiskers show 90% range. Numbers in brackets give median value for the  
 756 corresponding box.

## 757 8 Discussion

758 With the addition of the data censoring treatment described in this study, FoGSS now  
759 produces 12-month streamflow forecasts at the monthly time step that are reliable in highly  
760 ephemeral rivers. Forecasts are in the form of an ensemble of time series, where each  
761 ensemble member can be aggregated to produce ensemble forecasts of aggregated volumes.  
762 Forecasts are reliable at individual lead times and for aggregated volumes. As with previous  
763 versions of FoGSS, forecasts are almost always more skilful than ESP forecasts. Crucially, as  
764 skill declines with lead time, forecasts become neutrally skilful. This allows aggregated  
765 volume forecasts to be skillful, even to very long (e.g. 6-month) aggregation periods. These  
766 properties remove many of the barriers to the use of long-range ensemble streamflow  
767 forecasts in ephemeral rivers. We have shown in other work that FoGSS forecasts can  
768 improve the management of dams (Turner et al., 2017) and the allocation of water (Kaune et  
769 al., 2020), and these benefits can now be realised for highly ephemeral rivers.

770 This study reaffirms that it is possible to use FoGSS error model parameters estimated at lead  
771 1 – analogous to a hydrological model calibration – to calibrate forecasts to multiple lead  
772 times. Crucially, this allows forecasts to be represented as a hydrograph, rather than a set of  
773 discrete probability distributions at each lead time. This is possible through the use of  
774 stochastic updating. Stochastic updating assumes that for a given month, the autocorrelation  
775 properties and error distributions at lead 1 hold at all lead times. This assumption is not  
776 guaranteed to hold, and accordingly FoGSS cannot enforce reliability at each month/lead  
777 time in the way that statistical methods can – although we have shown that in practice it  
778 almost always does. It is possible, then, that statistical forecasts that account for zero values  
779 (like the BJP; Wang and Robertson, 2011) will produce more reliable forecasts at discrete  
780 lead times (noting again that there is no clear method to generate hydrographs with statistical  
781 methods). In future work, we plan to compare climate-forced os-censored FoGSS forecasts  
782 with statistical forecasts generated with the BJP, where we adapt BJP forecasting methods to  
783 produce monthly forecasts to 12 months.

784 We note that there are aspects of the FoGSS error model that may still be improved. For  
785 example, McInerney et al. (2019) found that they could achieve slightly sharper predictions  
786 using a fixed-parameter Box-Cox transformation instead of a log-sinh transformation in  
787 ephemeral rivers. Conversely, the log-sinh transformation performs strongly in perennial  
788 catchments in comparison to other transformations (McInerney et al., 2017; Wang et al.,

789 2012). It is therefore possible that a different transformation may improve the properties of  
790 the FoGSS ensemble in ephemeral rivers.

791 In this study we have used uninformative forcings to drive our streamflow forecasts. FoGSS  
792 completely separates uncertainties in hydrological modelling from uncertainties in forcings,  
793 so it is relatively straightforward to include informative forcings from seasonal climate  
794 forecasts. However, to ensure reliable streamflow forecasts, forcings must also be reliable.  
795 Only then can uncertainties from forcings and hydrological models sum to correctly represent  
796 total forecast uncertainty. In addition, forcings should be coherent (i.e., at least as skilful as  
797 climatology) to ensure that streamflow forecast skill does not become negative at long lead  
798 times. ESP forcings are inherently reliable and coherent, but this is not necessarily true of  
799 climate forecasts (e.g., Peng et al., 2014; Schepen et al., 2016; Strazzo et al., 2019; Wang et  
800 al., 2019). Thus climate forecasts should only be used with FoGSS after they have been  
801 formally statistically calibrated, for which a number of methods are available (Manzanas et  
802 al., 2019; Sansom et al., 2016; Schepen and Wang, 2014; Siegert and Stephenson, 2019).  
803 Simple bias-corrections do not ensure reliability or coherence, and are not suitable for  
804 calibrating forecasts (Zhao et al., 2017).

805 In our previous applications of FoGSS, we have used calibrated climate forecasts as forcings,  
806 and forecasts were reliable in perennial and moderately ephemeral catchments (Bennett et al.,  
807 2016b; Bennett et al., 2017). However, these studies differed from the present one in that they  
808 used a hydrological model run at the monthly time step, not at the daily time step. Calibrating  
809 climate forecasts at a monthly time step is much more straightforward than at the daily time  
810 step, as monthly data are far less noisy. While it is technically possible to generate calibrated  
811 daily climate forecasts (Schepen et al., 2017), these methods have never been tested to the  
812 very long lead times (i.e. 365 days) used in our study. In addition, propagating uncertainty  
813 with stochastic updating – as used in FoGSS – is more difficult over many lead times  
814 (Bennett et al., 2021; Li et al., 2020) and is likely to require further development to produce  
815 forecasts for 365 lead times compared to the 12 lead times in this study. Testing and further  
816 development of daily climate forecast calibration methods and error modelling to very long  
817 lead times is thus a clear target for future research.

818 FoGSS is not yet in operational use, but offers an attractive means for a future operational  
819 long-range streamflow forecasting system. As we have shown, censoring works across all  
820 catchment types; in perennial catchments censoring is not enacted, while in ephemeral

821 catchments it plays an important role. This means that a single method can be applied to all  
822 catchment types – a clear advantage for operations. Further, we have shown that FoGSS is  
823 effective as a ‘bolt-on’ corrective to ESP forecasts – in this case for a hydrological model  
824 calibrated using an independent procedure. As ESP forecasts are widely used in operational  
825 forecasting, FoGSS can be easily applied to improve these forecasts. Correcting ESP  
826 forecasts is also attractive as ESP can be readily adapted to accept calibrated climate forecast  
827 inputs. In other words, because FoGSS is premised on the idea of separating the uncertainties  
828 in climate forecasts from uncertainties in hydrological modelling, this makes any FoGSS  
829 operational system highly modular: climate forecasts can be updated/improved independently  
830 of hydrological models, and vice versa. We note also that FoGSS is highly computationally  
831 efficient, often a key consideration for operationalisation: on a standard desktop computer, it  
832 takes <2 minutes to estimate parameters for the 29-year estimation period used in this study,  
833 while it takes at most a few seconds to generate a 12-month forecast with 1000 ensemble  
834 members.

## 835 9 Summary and Conclusions

836 This study further develops the FoGSS (forecast guided stochastic scenarios) method to  
837 generate long-range (12-month) streamflow forecasts for use in highly ephemeral rivers.  
838 Forecasts are in the form of an ensemble of time series at the monthly time step. Each  
839 ensemble member can be summed to produce an ensemble of aggregated volume forecasts.  
840 We combine the basic staged structure of the FoGSS error model with a data censoring  
841 method that is applied to both observations and simulations. The data censoring method treats  
842 both modelled and simulated flow as censored data when they are equal to zero. This allows  
843 FoGSS to generate reliable ensemble forecasts in highly ephemeral rivers, which can cease to  
844 flow >50% of the time in some months.

845 We test the new version of FoGSS on 50 catchments, of which 26 are ephemeral. FoGSS  
846 forecasts are generally highly skillful at short lead times compared to climatology, and very  
847 often more skillful than conventional ESP (Ensemble Streamflow Prediction) forecasts. At  
848 longer lead times FoGSS forecasts are ‘coherent’ – that is, never less skillful than  
849 climatology – which cannot be guaranteed by ESP forecasts. Forecasts are reliable  
850 irrespective if they are issued for highly ephemeral (>50% zeros) or perennial (<5% zeros)  
851 months. Forecasts perform similarly in highly ephemeral, moderately ephemeral and  
852 perennial catchments in both skill and reliability. These improvements pave the way for

853 operational long-range forecasts in ephemeral rivers, meeting a key need for improved water  
854 management.

## 855 Acknowledgements

856 This work was conducted on the traditional lands of the Boonwurrung and Wurundjeri  
857 peoples of the Kulin Nation. We acknowledge their continuing custodianship of these lands  
858 and the rivers that flow through them, and pay our respects to their elders, past and present.  
859 We also acknowledge the traditional custodians of the catchments and rivers used in this  
860 study.

861 This research was supported by the Water Information Research And Development Alliance  
862 (WIRADA) between the Bureau of Meteorology and CSIRO Land & Water, and ARC  
863 linkage project LP170100922. Thanks to Elisabeth Vogel and Richard Laugesen (both  
864 Bureau of Meteorology) for helpful comments on the manuscript.

865 All data and model runs used in this paper are available from CSIRO's data access portal at  
866 [http://dx.doi.org/\[DOI NO STILL TO COME\]](http://dx.doi.org/[DOI NO STILL TO COME]). Matlab and C++ code used to generate  
867 simulations and forecasts, and to verify forecasts, are available on request; license conditions  
868 apply.

869

870 References

- 871 Ammann, L., Fenicia, F., Reichert, P., 2019. A likelihood framework for deterministic  
872 hydrological models and the importance of non-stationary autocorrelation. *Hydrol.  
873 Earth Syst. Sci.*, 23(4): 2147-2172. DOI:10.5194/hess-23-2147-2019
- 874 Arnal, L. et al., 2018. Skilful seasonal forecasts of streamflow over Europe? *Hydrol. Earth  
875 Syst. Sci.*, 22(4): 2057-2072. DOI:10.5194/hess-22-2057-2018
- 876 Bennett, J.C., Robertson, D.E., Wang, Q.J., Li, M., Perraud, J.-M., 2021. Propagating reliable  
877 estimates of hydrological forecast uncertainty to many lead times. Submitted to  
878 *Journal of Hydrology*.
- 879 Bennett, J.C., Robertson, D.E., Ward, P.G.D., Hapuarachchi, H.A.P., Wang, Q.J., 2016a.  
880 Calibrating hourly rainfall-runoff models with daily forcings for streamflow  
881 forecasting applications in meso-scale catchments. *Environmental Modelling &  
882 Software*, 76: 20-36. DOI:10.1016/j.envsoft.2015.11.006
- 883 Bennett, J.C., Wang, Q.J., Li, M., Robertson, D.E., Schepen, A., 2016b. Reliable long-range  
884 ensemble streamflow forecasts: Combining calibrated climate forecasts with a  
885 conceptual runoff model and a staged error model. *Water Resources Research*, 52:  
886 8238–8259. DOI:10.1002/2016wr019193
- 887 Bennett, J.C. et al., 2017. Assessment of an ensemble seasonal streamflow forecasting system  
888 for Australia. *Hydrol. Earth Syst. Sci.*, 21(12): 6007-6030. DOI:10.5194/hess-21-  
889 6007-2017
- 890 Chatterjee, S., McLeish, D.L., 1986. Fitting linear regression models to censored data by least  
891 squares and maximum likelihood methods. *Communications in Statistics - Theory and  
892 Methods*, 15(11): 3227-3243. DOI:10.1080/03610928608829305
- 893 Coron, L. et al., 2012. Crash testing hydrological models in contrasted climate conditions: An  
894 experiment on 216 Australian catchments. *Water Resources Research*, 48(5):  
895 W05552. DOI:10.1029/2011wr011721
- 896 Costigan, K.H. et al., 2017. Chapter 2.2 - Flow Regimes in Intermittent Rivers and  
897 Ephemeral Streams. In: Datry, T., Bonada, N., Boulton, A. (Eds.), *Intermittent Rivers  
898 and Ephemeral Streams*. Academic Press, pp. 51-78. DOI:10.1016/B978-0-12-  
899 803835-2.00003-6
- 900 Crochemore, L., Ramos, M.H., Pappenberger, F., 2016. Bias correcting precipitation  
901 forecasts to improve the skill of seasonal streamflow forecasts. *Hydrology and Earth  
902 System Sciences*, 20: 3601-3618. DOI:10.5194/hess-2016-78
- 903 Datry, T., Bonada, N., Boulton, A.J., 2017. Chapter 1 - General Introduction. In: Datry, T.,  
904 Bonada, N., Boulton, A. (Eds.), *Intermittent Rivers and Ephemeral Streams*.  
905 Academic Press, pp. 1-20. DOI:10.1016/B978-0-12-803835-2.00001-2
- 906 Day, G.N., 1985. Extended streamflow forecasting using NWSRFS. *Journal of Water  
907 Resources Planning and Management*, 111(2): 157–170. DOI:10.1061/(ASCE)0733-  
908 9496(1985)111:2(157)
- 909 Demargne, J. et al., 2014. The Science of NOAA's Operational Hydrologic Ensemble  
910 Forecast Service. *Bulletin of the American Meteorological Society*, 95: 79–98.  
911 DOI:10.1175/bams-d-12-00081.1
- 912 Duan, Q.Y., Gupta, V.K., Sorooshian, S., 1993. Shuffled complex evolution approach for  
913 effective and efficient global minimization. *Journal of Optimization Theory and  
914 Applications*, 76(3).
- 915 Feikema, P.M. et al., 2018. Service and Research on Seasonal Streamflow Forecasting in  
916 Australia. In: Jung, H.-S., Wang, B. (Eds.), *Bridging Science and Policy Implication  
917 for Managing Climate Extremes*. World Scientific Series on Asia-Pacific Weather and  
918 Climate. World Scientific, pp. 157-175. DOI:doi:10.1142/9789813235663\_0010

919 10.1142/9789813235663\_0010

920 Ferro, C.A.T., Richardson, D.S., Weigel, A.P., 2008. On the effect of ensemble size on the  
921 discrete and continuous ranked probability scores. *Meteorological Applications*,  
922 15(1): 19-24. DOI:10.1002/met.45

923 Gneiting, T., Katzfuss, M., 2014. Probabilistic forecasting. *Annual Review of Statistics and  
924 Its Application*, 1(1): 125-151. DOI:10.1146/annurev-statistics-062713-085831

925 Grafton, R.Q., Wheeler, S.A., 2018. Economics of Water Recovery in the Murray-Darling  
926 Basin, Australia. *Annual Review of Resource Economics*, 10(1): 487-510.  
927 DOI:10.1146/annurev-resource-100517-023039

928 Hemri, S., Klein, B., 2017. Analog-Based Postprocessing of Navigation-Related  
929 Hydrological Ensemble Forecasts. *Water Resources Research*, 53(11): 9059-9077.  
930 DOI:10.1002/2017wr020684

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

933 Kaune, A., Chowdhury, F., Werner, M., Bennett, J., 2020. The benefit of using an ensemble  
934 of seasonal streamflow forecasts in water allocation decisions. *Hydrol. Earth Syst.  
935 Sci. Discuss.*, 2020: 1-32. DOI:10.5194/hess-2020-60

936 Krzysztofowicz, R., 1999. Bayesian theory of probabilistic forecasting via deterministic  
937 hydrologic model. *Water Resources Research*, 35(9): 2739-2750.  
938 DOI:doi:10.1029/1999WR900099

939 Li, M., Robertson, D.E., Wang, Q.J., Bennett, J.C., Perraud, J.-M., 2020. Reliable hourly  
940 streamflow forecasting with emphasis on ephemeral rivers. *Journal of Hydrology*,  
941 125739. DOI:10.1016/j.jhydrol.2020.125739

942 Li, M., Wang, Q.J., Bennett, J., 2013. Accounting for seasonal dependence in hydrological  
943 model errors and prediction uncertainty. *Water Resources Research*, 49: 5913–5929.  
944 DOI:10.1002/wrcr.20445

945 Li, M., Wang, Q.J., Bennett, J.C., Robertson, D.E., 2015. A strategy to overcome adverse  
946 effects of autoregressive updating of streamflow forecasts. *Hydrology and Earth  
947 System Sciences*, 19(1): 1-15. DOI:10.5194/hess-19-1-2015

948 Li, M., Wang, Q.J., Bennett, J.C., Robertson, D.E., 2016. Error reduction and representation  
949 in stages (ERRIS) in hydrological modelling for ensemble streamflow forecasting.  
950 *Hydrology and Earth System Sciences*, 20(9): 3561-3579. DOI:10.5194/hess-20-  
951 3561-2016

952 Li, M., Wang, Q.J., Robertson, D.E., Bennett, J.C., 2017. Improved error modelling for  
953 streamflow forecasting at hourly time steps by splitting hydrographs into rising and  
954 falling limbs. *Journal of Hydrology*, 555: 586-599.  
955 DOI:10.1016/j.jhydrol.2017.10.057

956 Liu, L., Wang, Q.J., Xu, Y.-P., 2020. Temporally varied error modelling for improving  
957 simulations and quantifying uncertainty. *Journal of Hydrology*, 586: 124914.  
958 DOI:10.1016/j.jhydrol.2020.124914

959 Manzanas, R. et al., 2019. Bias adjustment and ensemble recalibration methods for seasonal  
960 forecasting: a comprehensive intercomparison using the C3S dataset. *Climate  
961 Dynamics*, 53(3): 1287-1305. DOI:10.1007/s00382-019-04640-4

962 McInerney, D., Kavetski, D., Thyer, M., Lerat, J., Kuczera, G., 2019. Benefits of explicit  
963 treatment of zero flows in probabilistic hydrological modelling of ephemeral  
964 catchments. *Water Resources Research*, 0(ja). DOI:10.1029/2018WR024148

965 McInerney, D., Thyer, M., Kavetski, D., Lerat, J., Kuczera, G., 2017. Improving probabilistic  
966 prediction of daily streamflow by identifying Pareto optimal approaches for modeling  
967 heteroscedastic residual errors. *Water Resources Research*, 53(3): 2199-2239.  
968 DOI:10.1002/2016wr019168

- 969 Pagano, T.C., Garen, D.C., Perkins, T.R., Pasteris, P.A., 2009. Daily Updating of Operational  
970 Statistical Seasonal Water Supply Forecasts for the western U.S.1. JAWRA Journal of  
971 the American Water Resources Association, 45(3): 767-778. DOI:10.1111/j.1752-  
972 1688.2009.00321.x
- 973 Peng, Z. et al., 2014. Statistical calibration and bridging of ECMWF System4 outputs for  
974 forecasting seasonal precipitation over China. Journal of Geophysical Research  
975 (Atmospheres), 119: 7116–7135. DOI:10.1002/2013JD021162.
- 976 Perrin, C., Michel, C., Andréassian, V., 2003. Improvement of a parsimonious model for  
977 streamflow simulation. Journal of Hydrology, 279: 275-289. DOI:10.1016/S0022-  
978 1694(03)00225-7
- 979 Petheram, C., Gallant, J., Stone, P., Wilson, P., Read, A., 2018. Rapid assessment of potential  
980 for development of large dams and irrigation across continental areas: application to  
981 northern Australia. The Rangeland Journal, 40(4): 431-449. DOI:10.1071/RJ18012
- 982 Pokhrel, P., Robertson, D.E., Wang, Q.J., 2013. A Bayesian joint probability post-processor  
983 for reducing errors and quantifying uncertainty in monthly streamflow predictions.  
984 Hydrology and Earth System Sciences, 17(2): 795-804. DOI:10.5194/hess-17-795-  
985 2013
- 986 Renard, B., Kavetski, D., Kuczera, G., Thyer, M., Franks, S.W., 2010. Understanding  
987 predictive uncertainty in hydrologic modeling: The challenge of identifying input and  
988 structural errors. Water Resources Research, 46(5): W05521.  
989 DOI:10.1029/2009wr008328
- 990 Sansom, P.G., Ferro, C.A.T., Stephenson, D.B., Goddard, L., Mason, S.J., 2016. Best  
991 Practices for Postprocessing Ensemble Climate Forecasts. Part I: Selecting  
992 Appropriate Recalibration Methods. Journal of Climate, 29(20): 7247-7264.  
993 DOI:10.1175/jcli-d-15-0868.1
- 994 Schepen, A., Wang, Q.J., 2014. 12 month out forecasts of catchment rainfall by post-  
995 processing ECMWF System 4 and POAMA M2.4 forecasts, CSIRO Land and Water  
996 for a Healthy Country Flagship, Highett.
- 997 Schepen, A., Wang, Q.J., Everingham, Y., 2016. Calibration, bridging and merging to  
998 improve GCM seasonal temperature forecasts in Australia. Monthly Weather Review,  
999 144(6): 2421-2441. DOI:10.1175/MWR-D-15-0384.1
- 1000 Schepen, A., Zhao, T., Wang, Q.J., Robertson, D.E., 2017. A new method for post-processing  
1001 daily sub-seasonal to seasonal rainfall forecasts from GCMs and evaluation for 12  
1002 Australian catchments. Submitted to Hydrology and Earth System Sciences.
- 1003 Seo, D.J., Herr, H.D., Schaake, J.C., 2006. A statistical post-processor for accounting of  
1004 hydrologic uncertainty in short-range ensemble streamflow prediction. Hydrology and  
1005 Earth System Sciences Discussions, 3: 1987-2035. DOI:10.5194/hessd-3-1987-2006
- 1006 Shrestha, D.L., Robertson, D.E., Bennett, J.C., Wang, Q.J., 2015. Improving precipitation  
1007 forecasts by generating ensembles through postprocessing. Monthly Weather Review,  
1008 143: 3642-3663. DOI:10.1175/mwr-d-14-00329.1
- 1009 Siegert, S., Stephenson, D.B., 2019. Chapter 15 - Forecast Recalibration and Multimodel  
1010 Combination. In: Robertson, A.W., Vitart, F. (Eds.), Sub-Seasonal to Seasonal  
1011 Prediction. Elsevier, pp. 321-336. DOI:[https://doi.org/10.1016/B978-0-12-811714-  
1012 9.00015-2](https://doi.org/10.1016/B978-0-12-811714-9.00015-2)
- 1013 Smith, T., Marshall, L., Sharma, A., 2015. Modeling residual hydrologic errors with Bayesian  
1014 inference. Journal of Hydrology, 528: 29-37.  
1015 DOI:<https://doi.org/10.1016/j.jhydrol.2015.05.051>
- 1016 Smith, T., Sharma, A., Marshall, L., Mehrotra, R., Sisson, S., 2010. Development of a formal  
1017 likelihood function for improved Bayesian inference of ephemeral catchments. Water  
1018 Resources Research, 46(12): W12551. DOI:10.1029/2010wr009514

- 1019 Stern, H., Hoedt, G., Ernst, J., 2000. Objective classification of Australian climates.  
 1020 Australian Meteorological Magazine, 49: 87-96.
- 1021 Strazzo, S. et al., 2019. Application of a Hybrid Statistical–Dynamical System to Seasonal  
 1022 Prediction of North American Temperature and Precipitation. Monthly Weather  
 1023 Review, 147(2): 607-625. DOI:10.1175/mwr-d-18-0156.1
- 1024 Tooth, S., 2000. Process, form and change in dryland rivers: a review of recent research.  
 1025 Earth-Science Reviews, 51(1): 67-107. DOI:10.1016/S0012-8252(00)00014-3
- 1026 Turner, S.W.D., Bennett, J.C., Robertson, D.E., Galelli, S., 2017. Complex relationship  
 1027 between seasonal streamflow forecast skill and value in reservoir operations. Hydrol.  
 1028 Earth Syst. Sci., 21(9): 4841-4859. DOI:10.5194/hess-21-4841-2017
- 1029 Verkade, J.S., Brown, J.D., Davids, F., Reggiani, P., Weerts, A.H., 2017. Estimating  
 1030 predictive hydrological uncertainty by dressing deterministic and ensemble forecasts;  
 1031 a comparison, with application to Meuse and Rhine. Journal of Hydrology,  
 1032 555(Supplement C): 257-277. DOI:10.1016/j.jhydrol.2017.10.024
- 1033 Wang, Q., Shrestha, D.L., Robertson, D., Pokhrel, P., 2012. A log-sinh transformation for  
 1034 data normalization and variance stabilization. Water Resources Research, 48(5).  
 1035 DOI:<http://dx.doi.org/10.1029/2011WR010973>
- 1036 Wang, Q.J., Bennett, J.C., Robertson, D.E., Li, M., 2020. A data censoring approach for  
 1037 predictive error modelling of flow in ephemeral rivers. Water Resources Research,  
 1038 56: e2019WR026128. DOI:10.1029/2019WR026128
- 1039 Wang, Q.J., Robertson, D.E., 2011. Multisite probabilistic forecasting of seasonal flows for  
 1040 streams with zero value occurrences. Water Resources Research, 47: W02546.  
 1041 DOI:10.1029/2010wr009333
- 1042 Wang, Q.J., Robertson, D.E., Chiew, F.H.S., 2009. A Bayesian joint probability modeling  
 1043 approach for seasonal forecasting of streamflows at multiple sites. Water Resources  
 1044 Research, 45: W05407. DOI:10.1029/2008WR007355
- 1045 Wang, Q.J. et al., 2019. An evaluation of ECMWF SEAS5 seasonal climate forecasts for  
 1046 Australia using a new forecast calibration algorithm. Environmental Modelling &  
 1047 Software, 122: 104550. DOI:<https://doi.org/10.1016/j.envsoft.2019.104550>
- 1048 Woldemeskel, F. et al., 2018. Evaluating post-processing approaches for monthly and  
 1049 seasonal streamflow forecasts. Hydrol. Earth Syst. Sci., 22(12): 6257-6278.  
 1050 DOI:10.5194/hess-22-6257-2018
- 1051 Wood, A.W., Schaake, J.C., 2008. Correcting Errors in Streamflow Forecast Ensemble Mean  
 1052 and Spread. J. Hydrometeorol., 9(1): 132-148. DOI:10.1175/2007JHM862.1
- 1053 Yuan, X., Wood, E.F., Ma, Z., 2015. A review on climate-model-based seasonal hydrologic  
 1054 forecasting: physical understanding and system development. Wiley Interdisciplinary  
 1055 Reviews: Water, 2(5): 523-536. DOI:10.1002/wat2.1088
- 1056 Zamo, M., Naveau, P., 2018. Estimation of the Continuous Ranked Probability Score with  
 1057 Limited Information and Applications to Ensemble Weather Forecasts. Mathematical  
 1058 Geosciences, 50(2): 209-234. DOI:10.1007/s11004-017-9709-7
- 1059 Zhao, T. et al., 2017. How suitable is quantile mapping for post-processing GCM  
 1060 precipitation forecasts? Journal of Climate, 30(9): 3185-3196. DOI:10.1175/jcli-d-16-  
 1061 0652.1
- 1062 Zhao, T., Schepen, A., Wang, Q.J., 2016. Ensemble forecasting of sub-seasonal to seasonal  
 1063 streamflow by a Bayesian joint probability modelling approach. Journal of  
 1064 Hydrology, 541, Part B: 839-849.  
 1065 DOI:<http://dx.doi.org/10.1016/j.jhydrol.2016.07.040>

1066

1067 **Appendix A: Parameter estimation for each stage**

1068 **Stage 1**

1069 Maximum a posteriori (MAP) estimation is used to estimate the transformation parameters in  
 1070 Eq (1). We assume that transformed observed streamflow follows a normal distribution

$$1071 \quad z_o \sim N(m, s^2) \quad (A1)$$

1072 In Eq (1) and Eq (A1) a total of four parameters -  $a$ ,  $b$ ,  $m$  and  $s$  - are estimated. The four  
 1073 parameters are reparameterised to  $\log(a)$ ,  $\log(b)$ ,  $m/s$  and  $\log(s)$ , respectively, to make it  
 1074 more numerically efficient to find a MAP solution. Our parameter set is then

1075  $\theta_1 = \{\log(a), \log(b), m/s, \log(s)\}$ . The likelihood is given by:

$$1076 \quad L(\theta_1) \propto p(\theta_1) \prod_t J_{z_o(t) \rightarrow c q_o(t)} N(z_o(t) | m, s^2) \quad (A2)$$

1077 where  $J_{z_o(t) \rightarrow c q_o(t)}$  is the Jacobian

$$1078 \quad J_{z_o(t) \rightarrow c q_o(t)} = \coth(a + b c q_o(t)) \quad (A3)$$

1079 with a standardization constant  $c = 5 / \max(\mathbf{q}_o)$ , and  $p(\theta_1)$  is a prior

$$1080 \quad p(\theta_1) \propto p(\log(a)) p(m/s) p(\log(s)) p(\log(b)), \quad (A4)$$

1081 where

$$1082 \quad p(\log(a)) \propto 1, \quad \log(a) \leq 0, \quad (A5)$$

$$1083 \quad p(m/s) \propto 1, \quad (A6)$$

$$1084 \quad p(\log(s)) \propto 1 \quad (A7)$$

1085 and

$$1086 \quad p(\log(b)) \sim N(0, 1^2). \quad (A8)$$

1087 The priors in equations (A5)-(A7) are uninformative. The prior on  $\log(b)$  (Eq A8) encourages  
 1088 values closer to zero, which in our experience is a good starting point. If the datum at time  $t$   
 1089 is a censored value (i.e.,  $q_o(t) = 0$ ), the term  $J_{z_o(t) \rightarrow c q_o(t)} N(z_o(t) | m, s^2)$  in Eq (A2) is replaced with

1090 the normal cumulative probability  $\Phi(z_c | m, s^2)$ , where  $z_c = Tf(0)$  is the log-sinh transformed  
1091 value of zero.

1092 **Stage 2**

1093 Parameters in Eq (2),  $\theta_2 = \{d(i), \mu(i): i=1,2,\dots,12\}$ , are estimated by minimizing the LST  
1094 objective (Eq 7).

1095 **Stage 3**

1096 The  $\theta_3 = \{\rho(i): i=1,2,\dots,12\}$  parameter in Eq (3) is estimated by minimizing the LST  
1097 objective (Eq 7).

1098 **Stage 4**

1099 We first estimate  $m_3(i)$  and  $s_3^2(i)$ . To do this, we generate a streamflow simulation  $\mathbf{Q}_3$  by  
1100 applying stages 1-3 to the hydrological model simulation. We extract the simulated  
1101 streamflow for each calendar month from  $\mathbf{Q}_3$ , and then use the MAP procedure from Stage 1  
1102 to estimate  $m_3(i)$ ,  $s_3^2(i)$ , keeping  $a$  and  $b$  fixed from the Stage 1 estimation.

1103 We can then use maximum likelihood estimation to infer  $\theta_4 = \{\sigma^2(i): i=1,2,\dots,12\}$  with the  
1104 likelihood specified in Eq (12). Note that no closed form is available for  $case=4$  (Eq 11), and  
1105 thus it requires a Monte Carlo integration. For details of this integration, see Wang et al.  
1106 (2020).

1107

1108 **Supplementary Materials**1109 **Table S1: Gauge site information**

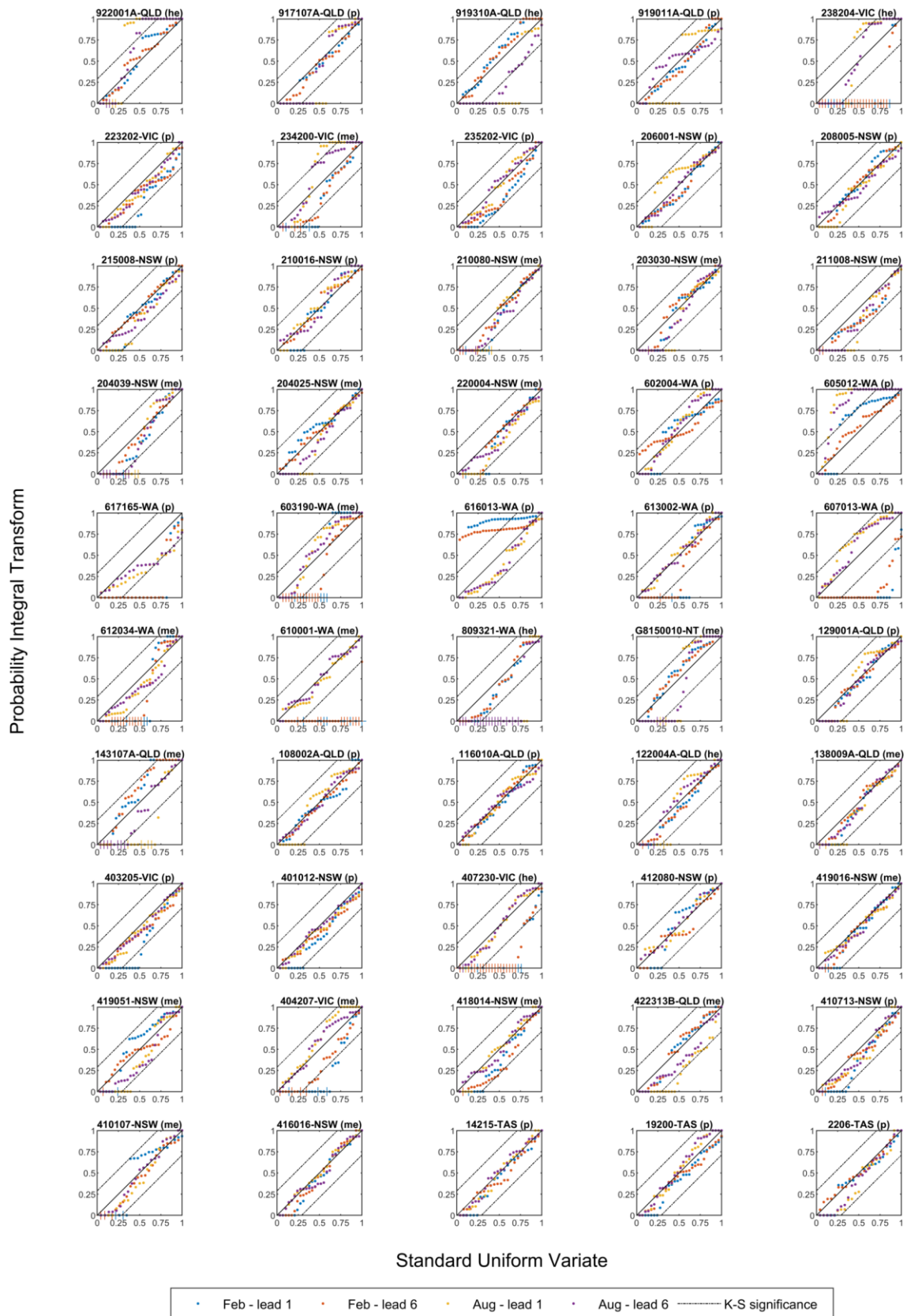
Gauge Name	Gauge Id	State	Area (km <sup>2</sup> )	Lon.	Lat.	Missing (%)	Ephemerality	No. months % zero >5	No. months % zero >50
Cockburn River at Mulla Crossing	419016	NSW	897	151.13	-31.06	0.9	Moderate	4	0
Flyers Creek at Beneree	412080	NSW	89	149.04	-33.51	20.5	Perennial	0	0
Goulburn River at Kerrabee	210016	NSW	4981	150.32	-32.42	14	Perennial	0	0
Gwydir River at Yarrowyck	418014	NSW	835	151.36	-30.47	0.9	Moderate	2	0
Jigadee Creek at Avondale	211008	NSW	64	151.47	-33.07	9.2	Moderate	2	0
Macintyre River at Inverell	416016	NSW	751	151.13	-29.79	0	Moderate	1	0
Maryland River downstream of Wylie Creek	204039	NSW	377	152.2	-28.47	4.8	Moderate	10	0
Maules Creek at Avoca East	419051	NSW	663	150.08	-30.5	0.9	Moderate	3	0
Mountain Creek at Mountain Creek	410107	NSW	185	147.85	-35.03	20.2	Moderate	6	0
Murray River at Biggara	401012	NSW	1165	148.05	-36.32	1.5	Perennial	0	0
Myrtle Creek at Rappville	203030	NSW	392	153	-29.11	0.3	Moderate	5	0
Nowendoc River at Rocks Crossing	208005	NSW	1870	152.08	-31.78	1.8	Perennial	0	0
Orara River at Karangi	204025	NSW	136	153.03	-30.25	4.2	Moderate	2	0
Paddy's River at Riverlea	410713	NSW	224	148.97	-35.38	2.1	Perennial	0	0
Shoalhaven River at Kadoona	215008	NSW	282	149.64	-35.79	8.3	Perennial	0	0
Styx River at Jeogla	206001	NSW	163	152.16	-30.59	3	Perennial	0	0
Towamba River at Towamba	220004	NSW	766	149.66	-37.07	1.2	Moderate	2	0
West Brook upstream of Glendon Brook	210080	NSW	72	151.28	-32.47	6.3	Moderate	6	0
Finniss River at Batchelor Dam Site	G8150010	NT	363	130.95	-13.02	0	Moderate	4	0

Archer River at Telegraph Crossing	922001A	QLD	2928	142.92	-13.42	14.3	High	5	1
Blencoe Creek at Blencoe Falls	116010A	QLD	224	145.54	-18.2	1.8	Perennial	0	0
Bremer River at Walloon	143107A	QLD	398	152.69	-27.6	10.7	Moderate	9	0
Daintree River at Bairds	108002A	QLD	907	145.28	-16.18	5.4	Perennial	0	0
Elizabeth Creek at Mount Surprise	917107A	QLD	528	144.31	-18.13	9.8	Perennial	0	0
Emu Creek at Emu Vale	422313B	QLD	153	152.23	-28.23	0.6	Moderate	1	0
Gregory River at Lower Gregory	122004A	QLD	46	148.55	-20.3	8.3	High	5	0
Mitchell River at Gamboola	919011A	QLD	20315	143.68	-16.54	8.9	Perennial	0	0
Tinana Creek at Tagigan Road	138009A	QLD	102	152.78	-26.08	6.3	Moderate	3	0
Walsh River at Rookwood	919310A	QLD	4927	144.29	-16.98	18.5	High	4	1
Waterpark Creek at Byfield	129001A	QLD	212	150.67	-22.84	5.4	Perennial	0	0
Brid River upstream of tidal limit	19200	TAS	138	147.37	-41.02	2.4	Perennial	0	0
Flowerdale River at Moorleah	14215	TAS	156	145.61	-40.97	10.7	Perennial	0	0
Scamander River upstream of Scamander water supply intake	2206	TAS	261	148.18	-41.45	20.8	Perennial	0	0
Gellibrand at Upper Gellibrand	235202	VIC	53	143.66	-38.56	0.3	Perennial	0	0
Holland Creek at Kelfeera	404207	VIC	448	146.06	-36.61	0	Moderate	4	0
Joyces Creek at Strathlea	407230	VIC	156	143.96	-37.16	0	High	11	3
Ovens River at Bright	403205	VIC	495	146.95	-36.73	0	Perennial	0	0
Tambo River at Swifts Creek	223202	VIC	943	147.72	-37.26	0	Perennial	0	0
Wannon River at Dunkeld	238204	VIC	385	142.34	-37.63	0	High	8	4
Woody Yaloak River Pitfield	234200	VIC	317	143.59	-37.81	0	Moderate	3	0
Collie River at South Branch	612034	WA	672	116.16	-33.39	0	Moderate	6	0

Dunham River at Dunham Gorge	809321	WA	1637	128.3	-16.19	0	High	8	2
Frankland River at Mount Frankland	605012	WA	4467	116.79	-34.91	0	Perennial	0	0
Harvey River at Dingo Road	613002	WA	148	116.04	-33.09	0	Perennial	0	0
Helena River at Ngangaguringurin	616013	WA	316	116	-32.15	0.6	Perennial	0	0
Kalgan River at Stevens Farm	602004	WA	2433	118	-34.89	0	Perennial	0	0
Lefroy Brook River at Rainbow Trail	607013	WA	249	116.02	-34.43	0	Perennial	0	0
Lennard Brook at Molecap Hill	617165	WA	58	115.92	-31.38	25.9	Perennial	0	0
Margaret River at Willmots Farm	610001	WA	435	115.58	-33.56	0	Moderate	3	0
Yale Flat Creek at Woonanup	603190	WA	53	117.29	-34.7	0	Moderate	6	0

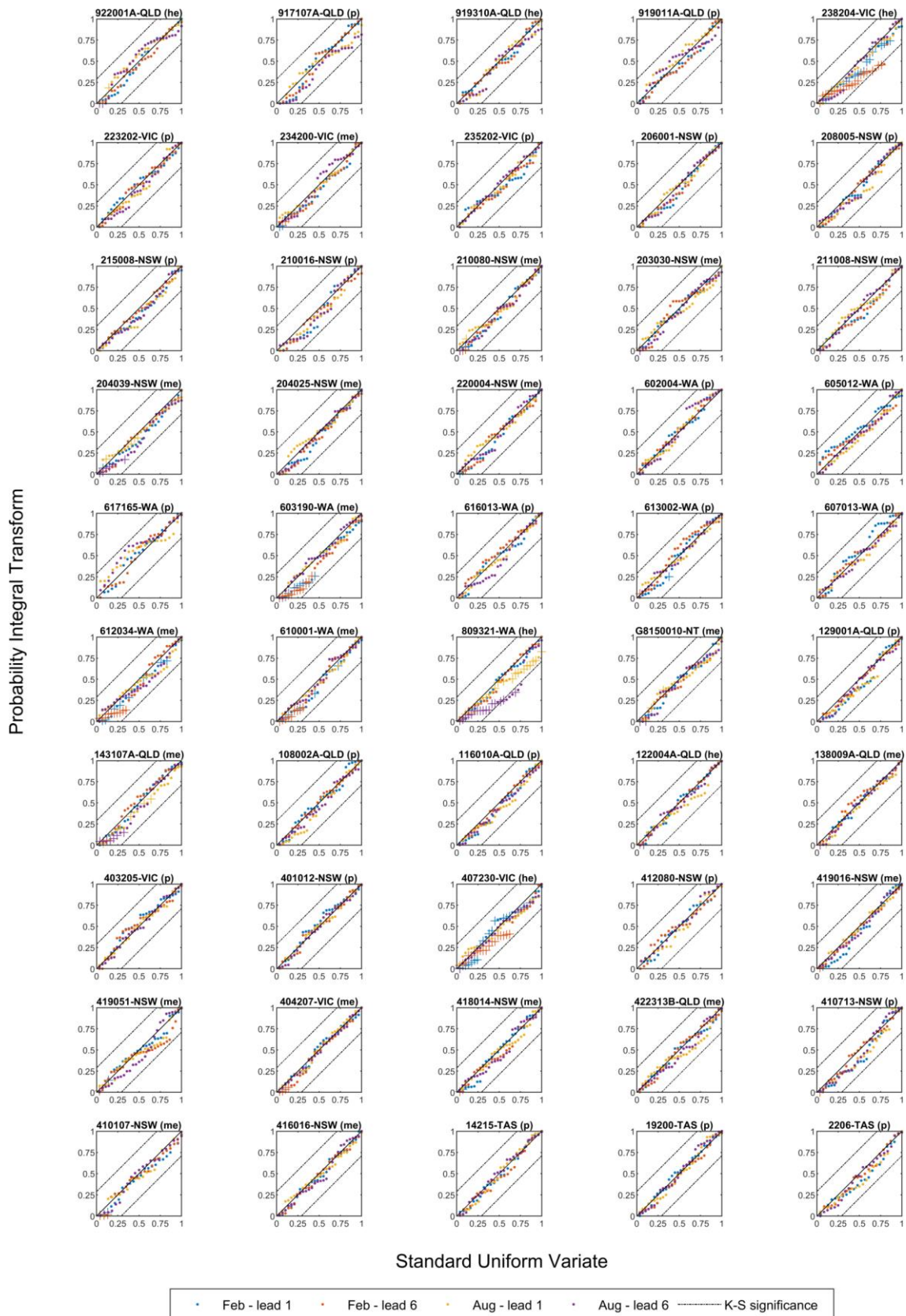
1110

1111



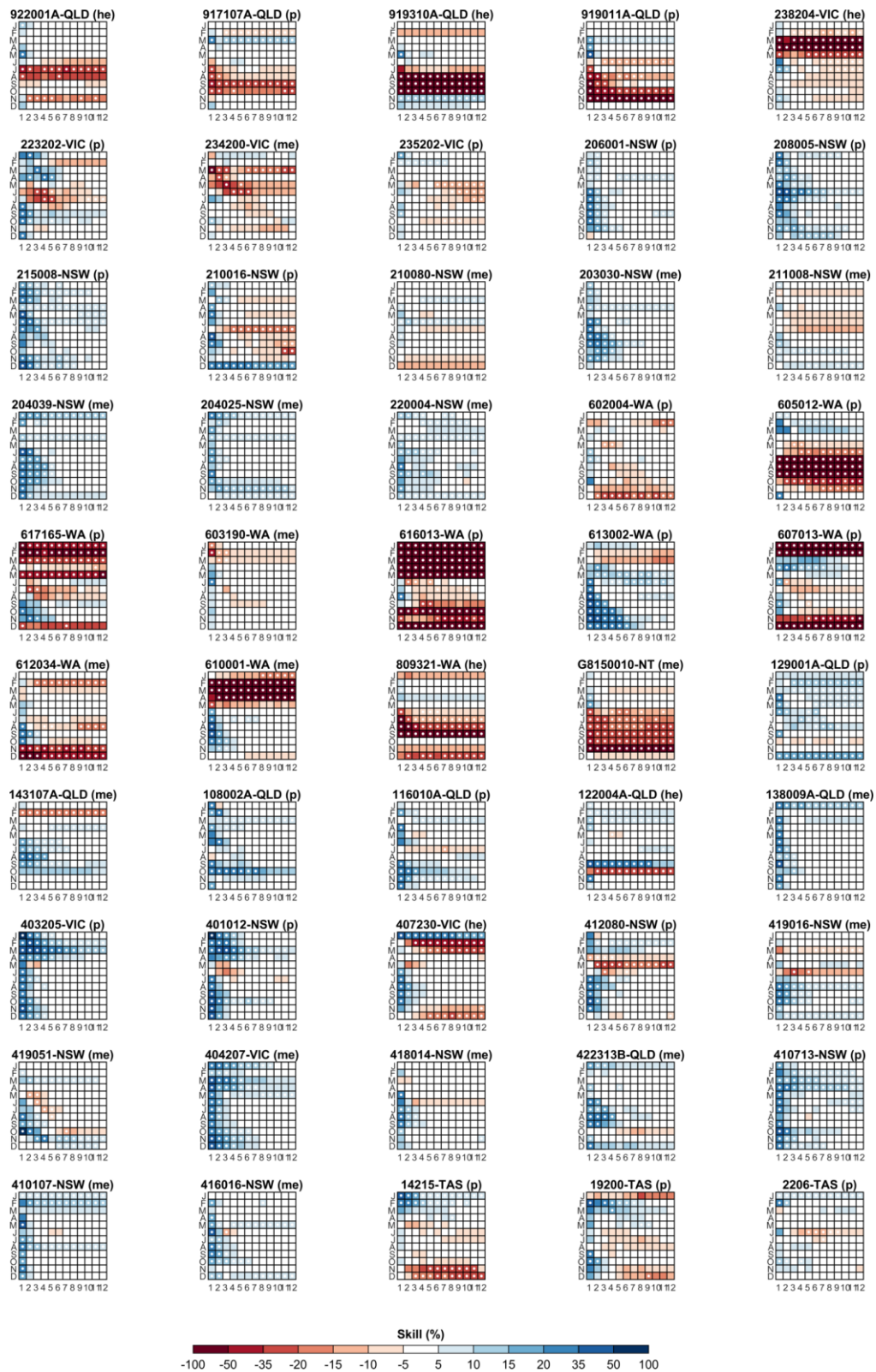
1112

1113 Figure S1 PIT-uniform probability plots for ESP forecasts for selected months and lead times  
 1114 catchments. Degree of ephemerality is denoted by (p)-perennial; (me)-moderately ephemeral;  
 1115 (he)-highly ephemeral.



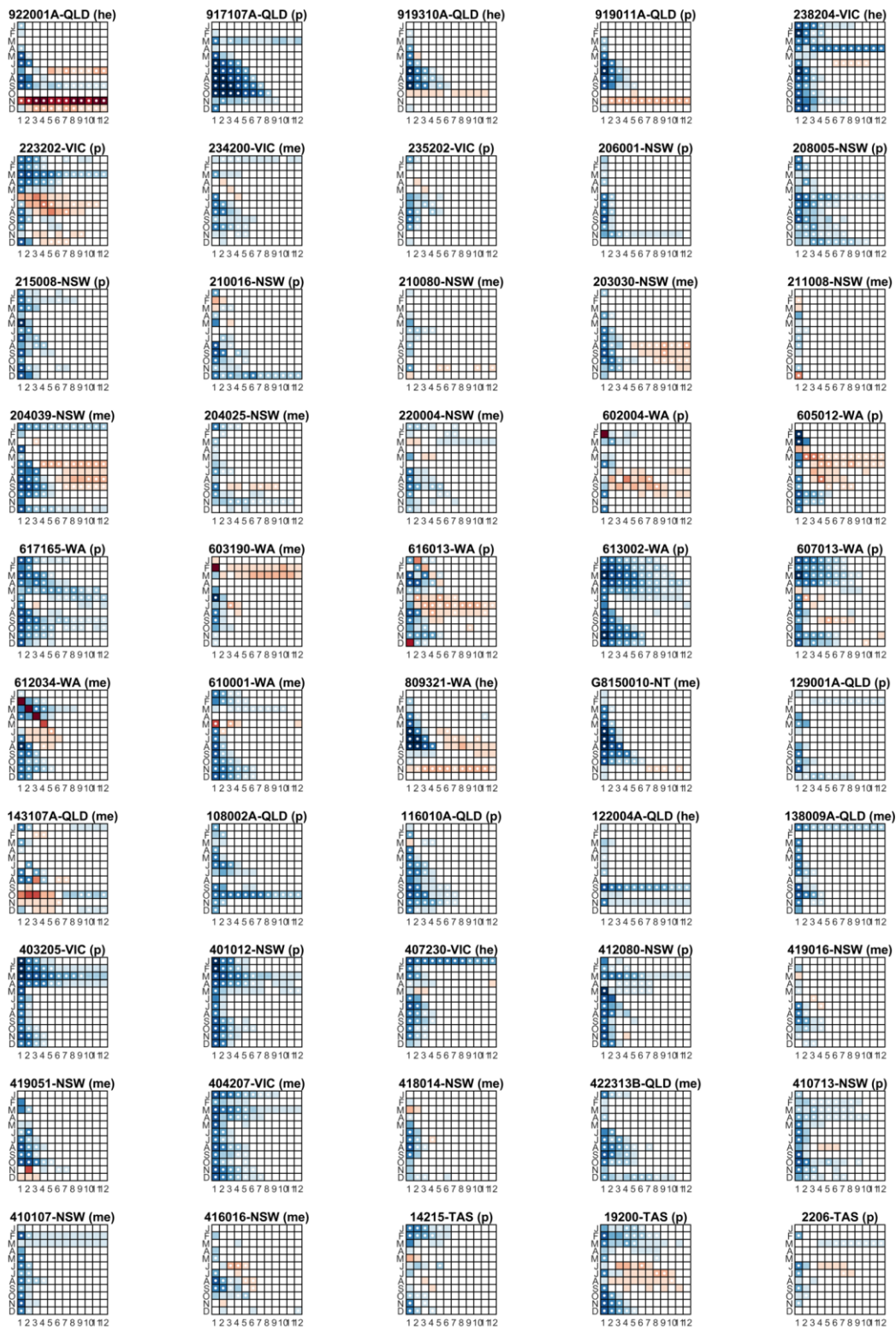
1116

1117 Figure S2 PIT-uniform probability plots for FoGSS forecasts for selected months and lead  
 1118 times catchments. Degree of ephemerality is denoted by (p)-perennial; (me)-moderately  
 1119 ephemeral; (he)-highly ephemeral.



1120

1121 Figure S3 Continuous ranked probability skill scores for ESP forecasts of target months for  
 1122 all catchments. White points indicating statistical significance at the 5% level under  
 1123 bootstrapping (see text for details). Degree of ephemerality is denoted by (p)-perennial; (me)-  
 1124 moderately ephemeral; (he)-highly ephemeral.



1125

1126 Figure S4 Continuous ranked probability skill scores for FoGSS forecasts of target months  
 1127 for all catchments. Degree of ephemerality is denoted by (p)-perennial; (me)-moderately  
 1128 ephemeral; (he)-highly ephemeral.