

1 **A strategy to overcome adverse effects of** 2 **autoregressive updating of streamflow forecasts**

3

4 M. Li¹, Q. J. Wang², J. C. Bennett² and D. E. Robertson²

5 [1] CSIRO Computational Informatics, Floreat, Western Australia, Australia

6 [2] CSIRO Land and Water, Highett, Victoria, Australia

7 Correspondence to: M. Li (Ming.Li@csiro.au)

8

9 **Abstract**

10 For streamflow forecasting, rainfall-runoff models are often augmented with updating
11 procedures that correct forecasts based on the latest available streamflow observations
12 of streamflow. A popular approach for updating forecasts is autoregressive (AR)
13 models, which exploit the “memory” in hydrological model simulation errors. AR
14 models may be applied to raw errors directly or to normalised errors. In this study, we
15 demonstrate that AR models applied in either way can sometimes cause over-
16 correction of forecasts. In using an AR model applied to raw errors, the over-
17 correction usually occurs when streamflow is rapidly receding. In applying an AR
18 model to normalised errors, the over-correction usually occurs when streamflow is
19 rapidly rising. In addition, when parameters of a hydrological model and an AR model
20 are estimated jointly, the AR model applied to normalised errors sometimes degrades
21 the stand-alone performance of the base hydrological model. This is not desirable for
22 forecasting applications, as forecasts should rely as much as possible on the base
23 hydrological model, with updating only used to correct minor errors. To overcome the
24 adverse effects of the conventional AR models, a restricted AR model applied to
25 normalised errors is introduced. We show that the new model reduces over-correction
26 and improves the performance of the base hydrological model considerably.

27

28 **1. Introduction**

29 Rainfall-runoff models are widely used to generate streamflow forecasts, which
30 provide essential information for flood warning and water resources management. For
31 streamflow forecasting, rainfall-runoff models are often augmented by updating
32 procedures that correct streamflow forecasts based on the latest available observations
33 of streamflow and their departures from model simulations. Model errors reflect
34 limitations of the hydrological models in reproducing physical processes as well as
35 inaccuracies in data used to force and evaluate the models.

36 The most popular updating approach uses autoregressive (AR) models, which exploit
37 the “memory” - more precisely the autocorrelation structure - of errors in hydrological
38 simulations (Morawietz et al., 2011). Essentially, AR updating uses a linear function
39 of the known errors at previous time steps to anticipate errors in a forecast period.
40 Forecasts are then updated according to these anticipated errors. AR updating is
41 conceptually simple and yet generally leads to significantly improved forecasts
42 (World Meteorological Organization, 1992). AR updating has been shown to provide
43 equivalent performance to more sophisticated non-linear and nonparametric updating
44 procedures (Xiong and O'Connor, 2002).

45 In rainfall-runoff modelling, model errors are generally heteroscedastic (i.e., they
46 have heterogeneous variance over time) (Xu, 2001; Kavetski et al., 2003; Pianosi and
47 Raso, 2012) and non-Gaussian (Bates and Campbell, 2001; Schaeffli et al.,
48 2007; Shrestha and Solomatine, 2008). In many applications (Seo et al., 2006; Bates
49 and Campbell, 2001; Salamon and Feyen, 2010; Morawietz et al., 2011), AR models
50 are applied to normalised errors that are considered homoscedastic and Gaussian.
51 Normalisation is often achieved through variable transformation by using, for
52 example, the Box-Cox transformation (Thyer et al., 2002; Bates and Campbell,
53 2001; Engeland et al., 2010) or, more recently, the log-sinh transformation (Wang et
54 al., 2012; Del Giudice et al., 2013). In other applications (Schoups and Vrugt,
55 2010; Schaeffli et al., 2007), AR models are applied directly to raw errors, but residual
56 errors of the AR models may be explicitly specified as heteroscedastic and non-
57 Gaussian.

58 There is no agreement on whether it is better to apply an AR model to normalised or
59 raw errors. Recent work by Evin et al. (2013) found that an AR model applied to raw
60 errors may lead to poor performance with exaggerated uncertainty. They
61 demonstrated that such instability can be mitigated by applying an AR model to

62 standardised errors (raw errors divided by standard deviations). Here, standardisation
 63 has a similar effect to normalisation in that it homogenises the variance of the errors
 64 (but does not consider the non-Gaussian distribution of errors). Conversely, Schaefli
 65 et al. (2007) pointed out that when an AR model is jointly estimated with a
 66 hydrological model, there is a clear advantage in applying an AR model to raw errors
 67 rather than normalised (or standardised) errors. Schaefli et al. (2007) found that using
 68 raw errors leads to more reliable parameter inference and uncertainty estimation,
 69 because the mean error is close to zero and therefore the simulations are free of
 70 systematic bias. The same is not necessarily true when applying an AR model to
 71 normalised errors.

72 In this study, we evaluate AR models applied to both raw and normalised errors on
 73 four Australian catchments and three United States (US) catchments. We show that
 74 when estimated jointly with a hydrological model, the AR model applied to
 75 normalised errors sometimes degrades the stand-alone performance of the base
 76 hydrological model. We also identify that both of these conventional AR models can
 77 sometimes cause over-correction of forecasts. We introduce a restricted AR model
 78 applied to normalised errors and demonstrate its effectiveness in overcoming the
 79 adverse effects of the conventional AR models.

80 **2. Autoregressive error models**

81 **2.1 Formulations**

82 A hydrological model is a function of forcing variables (precipitation and potential
 83 evapotranspiration), initial catchment state, S_0 , and a set of hydrological model
 84 parameters, θ_H . We denote the observed streamflow and model simulated streamflow
 85 at time t by Q_t and \tilde{Q}_t , respectively. An error model is used to describe the difference
 86 between Q_t and \tilde{Q}_t . The log-sinh transformation defined by Wang et al. (2012)

$$87 \quad f(x) = b^{-1} \log\{\sinh(a + bx)\} \quad (1)$$

88 is applied to stabilise variance and normalise data.

89 In this study, we firstly examine two first-order AR error models:

90 (1) An AR error model applied to normalised errors (referred to as *AR-Norm*) defined
 91 by:

$$92 \quad Z_t = \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t, \quad (2)$$

93 where Z_t and \tilde{Z}_t are the log-sinh transformed variables of Q and \tilde{Q} ;

94 (2) An AR error model applied to raw errors (referred to as *AR-Raw*) defined by

$$95 \quad Z_t = f\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\} + \varepsilon_t. \quad (3)$$

96 For both models, ρ is the lag-1 autoregression parameter, and ε_t is an identically
 97 and independently distributed Gaussian deviate with a mean of zero and a constant
 98 standard deviation σ .

99 Both the AR-Norm and AR-Raw models represent the lag-one autocorrelation by an
 100 AR process and both employ the log-sinh transformation. However, the way the log-
 101 sinh transformation is applied differs between the two models. The AR-Norm model
 102 first applies the log-sinh transformation to the observed and model simulated
 103 streamflow, and then assumes that the error in the transformed space follows an AR(1)
 104 process. In contrast, the AR-Raw model essentially assumes that the error in the
 105 original space follows an AR(1) process and only applies the log-sinh transformation
 106 to fit the asymmetric and non-Gaussian error distribution.

107 The median of the updated streamflow forecast (referred to as *updated streamflow*)
 108 for the AR-Norm and AR-Raw models (see Appendix A for proof), denoted by \tilde{Q}_t^* ,
 109 are respectively

$$110 \quad \tilde{Q}_t^* = f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\}, \quad (4)$$

111 and

$$112 \quad \tilde{Q}_t^* = \tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1}), \quad (5)$$

113 where $f^{-1}(x)$ is the inverse of log-sinh transformation (or back-transformation). The
 114 magnitude of the error update by the AR-Raw model, $\tilde{Q}_t^* - \tilde{Q}_t$, is dependent only on
 115 the difference between Q_{t-1} and \tilde{Q}_{t-1} . In contrast, the magnitude of the error update by
 116 the AR-Norm model is dependent not only on the difference between Q_{t-1} and \tilde{Q}_{t-1} ,
 117 but also on \tilde{Q}_t . Put differently, the AR-Norm model uses errors calculated in the

118 transformed domain, and this means that the error in the original domain can be
 119 amplified (or reduced) by the back-transformation (Equation (4)). The AR-Raw model
 120 uses errors calculated in the original domain and no back-transformation is used in
 121 calculating \tilde{Q}_t^* (Equation (5)), meaning that the error in the original domain cannot be
 122 amplified (or reduced). In Appendix B, we show that the AR-Norm model gives
 123 greater error updates for larger values of \tilde{Q}_t .

124 We will demonstrate in Section 4 that the AR-Norm and AR-Raw models can
 125 sometimes cause over-correction of forecasts. Motivated to overcome the potential for
 126 over-correction, we introduce a modification of the AR-Norm model, called the
 127 restricted AR-Norm model (referred to as *RAR-Norm*). A condition
 128 $|\tilde{Q}_t^* - \tilde{Q}_t| \leq |Q_{t-1} - \tilde{Q}_{t-1}|$ is used to limit the correction to an amount not exceeding the
 129 raw error at the last time step. The updated streamflow is given by

$$130 \quad \tilde{Q}_t^* = \begin{cases} f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\} & \text{if } D_t \leq |Q_{t-1} - \tilde{Q}_{t-1}| \\ \tilde{Q}_t + (Q_{t-1} - \tilde{Q}_{t-1}) & \text{otherwise.} \end{cases} \quad (6)$$

131 where

$$132 \quad D_t = \left| f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\} - \tilde{Q}_t \right|. \quad (7)$$

133 The full RAR-Norm model in the transformed space is given by

$$134 \quad Z_t = \begin{cases} \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t & \text{if } D_t \leq |Q_{t-1} - \tilde{Q}_{t-1}| \\ f(\tilde{Q}_t + Q_{t-1} - \tilde{Q}_{t-1}) + \varepsilon_t & \text{otherwise.} \end{cases} \quad (8)$$

135 2.2 Estimation

136 The AR-Norm, AR-Raw and RAR-Norm models are each calibrated jointly with the
 137 hydrological model. The method of maximum likelihood is used to estimate the error
 138 model parameters θ_E and the hydrological model parameters θ_H . Using a similar
 139 derivation as given by Li et al. (2013), the likelihood functions can be written as

140 (a) for AR-Norm

$$141 \quad L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_t J_{Z_t \rightarrow Q_t} \phi(Z_t | \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}), \sigma^2), \quad (9)$$

142 (b) for AR-Raw

$$L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_t J_{Z_t \rightarrow Q_t} \phi\left(Z_t | f\left\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\right\}, \sigma^2\right),$$

(10)

(c) for RAR-Norm

$$L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_{t: D_t \leq |\tilde{Q}_{t-1} - \tilde{Q}_t|} J_{Z_t \rightarrow Q_t} \phi\left(Z_t | \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}), \sigma^2\right) \\ + \prod_{t: D_t > |\tilde{Q}_{t-1} - \tilde{Q}_t|} J_{Z_t \rightarrow Q_t} \phi\left(Z_t | f\left\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\right\}, \sigma^2\right),$$

(11)

where $J_{Z_t \rightarrow Q_t} = \{\tanh(a + bQ_t)\}^{-1}$ is the Jacobian determinant of the log-sinh transformation and $\phi(x | \mu, \sigma^2)$ is the probability density function of a Gaussian random variable x with mean μ and standard deviation σ . The probability density function is replaced by the cumulative probability function when evaluating events of zero flow occurrences (Wang and Robertson, 2011; Li et al., 2013). The Shuffled Complex Evolution (SCE) algorithm (Duan et al., 1994) is used to minimize the log likelihood.

3. Data

We use daily data from four Australian catchments and three catchments from the United States (US; Figure 1, Table 1). Australian streamflow data are taken from the Catchment Water Yield Estimation Tool (CWYET) dataset (Vaze et al., 2011). Australian rainfall and potential evaporation data are derived from the Australian Water Availability Project (AWAP) dataset (Jones et al., 2009). All data for the US catchments come from the Model Intercomparison Experiment (MOPEX) dataset (Duan et al., 2006). The selected US catchments are amongst the 12 catchments used by Evin et al. (2014) to compare joint and postprocessor approaches to estimate hydrological uncertainty, and allows us to compare results with that study (the other catchments used by Evin et al. (2014) are influenced by snowmelt, which is not considered in the hydrological model used in this study). The Abercrombie River and the Guadalupe River intermittently experience periods of very low (to zero) flow, while the other rivers flow perennially (Table 1). Such dry catchments are challenging for hydrological simulations and error modelling. All catchments have high-quality streamflow records with very few missing data.

171 We forecast daily streamflow with the GR4J rainfall-runoff model (Perrin et al.,
172 2003). We apply updating procedures to correct these forecasts. All results presented
173 in this paper are based on cross-validation to ensure the results can be generalised to
174 independent data. We use different cross-validation schemes for the Australian and
175 US catchments, because of the shorter streamflow records available for the Australian
176 catchments:

- 177 i. For the Australian catchments we use data from 1992 to 2005 (14 years) for
178 these catchments. We then generate 14-fold cross-validated streamflow
179 forecasts. The data from 1990-1991 are only used to warm up the GR4J model.
180 For a given year, we leave out the data from that year and the following year
181 when estimating the parameters of GR4J and error models. For example, if we
182 wish to forecast streamflows at any point in 1999, we leave out data from 1999
183 and 2000 when we estimate parameters. The removal of data from the
184 following year (2000) is designed to minimise the impact of hydrological
185 memory on model parameter estimation. We then generate streamflow
186 forecasts in that year (1999) with model parameters estimated from the
187 remaining data.
- 188 ii. For the US catchments we follow the split-sampling validation scheme
189 suggested by Evin et al. (2014) to make our results comparable to that study:
190 (1) an 8-year calibration (09/09/1973- 26/11/1981) (i.e. 3000 days) with an 8-
191 year warm-up period and (2) a 17-year validation (27/11/1981-01/05/1998)
192 (i.e. 6000 days) with an 8-year warm-up period.

193 To demonstrate the problems of over-correction of errors in updating and poor stand-
194 alone performance of the base hydrological model, we consider only streamflow
195 forecasts for one time step ahead. We will consider longer lead times in future work.
196 Forecasts are generated using observed rainfall (i.e., a 'perfect' rainfall forecast) as
197 input. In streamflow forecasting, forecasts may be generated from rainfall information
198 that comes from a different source (e.g., a numerical weather prediction model). Our
199 study is aimed at streamflow forecasting applications, so we preserve the distinction
200 between observed and forecast forcings by referring to streamflows modelled with
201 observed rainfall as *simulations* and those modelled with forecast rainfall as *forecasts*.
202 In this study the forecast rainfall is observed rainfall, so the terms *forecast* and
203 *simulation* are interchangeable.

204 4. Results

205 4.1 Over-correction of forecasts as the hydrograph rises

206 The first adverse effect of the conventional AR models is over-correction of errors in
 207 updating as streamflows are rising. By over-correction, we mean that the AR model
 208 updates the hydrological model simulations too much. Over-correction is difficult to
 209 define precisely, however we will demonstrate the concept with two examples in the
 210 Mitta Mitta catchment: the first example illustrates over-correction by the AR-Norm
 211 model, and the second example illustrates over-correction by the AR-Raw model.

212 To illustrate the problem of over-correction caused by the AR-Norm model, Figure 2
 213 presents a 1-week time series for the Mitta Mitta catchment, showing streamflow
 214 forecasts with GR4J before error updating (referred to as streamflows forecast with
 215 the *base hydrological model*) and after error updating. Figure 2 shows that the base
 216 hydrological models consistently under-estimate the streamflow from 23/09/2000 to
 217 25/09/2000, and the corresponding updating procedures successfully identify the need
 218 to compensate for this under-estimation. For the AR-Norm model, however, the
 219 correction for 26/09/2000 is unreasonably large. Because the forecast streamflow on
 220 26/09/2000 is much higher than that of the previous day, the correction is greatly
 221 amplified by the back-transformation, leading to the over-correction. In contrast, the
 222 AR-Raw model works better in this situation because the magnitude of the error
 223 update never exceeds the simulation error on the previous day regardless of whether
 224 the forecast streamflow is high or low. The RAR-Norm model behaves similarly to
 225 the AR-Raw model for correcting the peak on 26/09/2000 and avoids the over-
 226 correction made by the AR-Norm model.

227 Figure 3 shows instances of possible over-correction by the AR-Norm model,
 228 identified by the condition $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$. Figure 3 shows that about 10-25% of the
 229 AR-Norm updated forecasts have an error update that is larger than the forecast error
 230 on the previous day and therefore are susceptible to over-correction. The frequency of
 231 these instances varies somewhat from catchment to catchment. The RAR-Norm model
 232 identifies 10-30% of the forecasts as possible instances of problematic updating, and
 233 the AR-Norm model identifies a similar number of instances (slightly fewer – they are
 234 not identical because the parameters for each model are inferred independently).

235 Figure 4 presents a time-series for the Orara catchment that shows the instances
236 susceptible to over-correction for the AR-Norm model. These instances all occur
237 when the streamflow rises. The RAR-Norm model effectively rectifies the problem of
238 over-correction caused by the AR-Norm model. We note that there is nothing that
239 forces the instances susceptible to over-correction identified by the AR-Norm model
240 to be the same as those identified by the RAR-Norm models because the two models
241 are calibrated independently (and therefore base hydrological model simulations may
242 be different). However, the restriction defined in the RAR-Norm model is largely
243 applied to the instances where the AR-Norm model is susceptible to over-correction.

244 **4.2 Over-correction of forecasts as the hydrograph recedes**

245 The second adverse effect of conventional AR models is over-correction of forecasts
246 as streamflows recede. An example is presented in Figure 5 where the AR-Raw model
247 causes over-correction. Here, the base hydrological model over-estimates the receding
248 hydrograph on 05/10/1993. The magnitude of the error update given by the AR-Raw
249 model cannot adjust according to the value of the forecast. As a result, the AR-Raw
250 model updates the forecast on 06/10/1993 by a large amount, resulting in serious
251 under-estimation (the forecast streamflow is nearly zero), and an artificial distortion
252 of the hydrograph. (We note that we have seen this problem become much worse in
253 unpublished experiments of forecasts made for several time-steps into the future,
254 sometimes resulting in forecasts of zero flows during large floods.) In contrast, the
255 AR-Norm model performs better in this example, giving a smaller magnitude of error
256 update by recognising that the hydrograph is moving downward. It is generally true
257 that in applying the AR-Raw model, over-correction may occur when the streamflow
258 is receding. The RAR-Norm model produces updated streamflow similar to the AR-
259 Norm model when the hydrograph recedes rapidly and avoids the over-correction by
260 the AR-Raw model on 06/10/1993.

261 Figure 6 provides more examples of the over-correction caused by the AR-Raw model
262 from a longer time-series plot for the Abercrombie catchment. There are three clear
263 instances of over-correction, all occurring on the time step immediately after large
264 peaks in observed streamflows. The RAR-Norm model works better than the AR-Raw
265 model to avoid the three instances of over-correction for the Abercrombie catchment.
266 Overall, the RAR-Norm model takes a conservative position when streamflow
267 changes rapidly, either rising or falling. When streamflow changes rapidly, it is

268 difficult to anticipate the magnitude of forecast error. Accordingly the conventional
269 AR models are prone to over-correction in such instances.

270 **4.3 Poor stand-alone performance of the base hydrological model**

271 The third adverse effect with conventional AR error models is the stand-alone
272 performance of the base hydrological model (GR4J). As noted above, the parameters
273 of the base hydrological model are estimated jointly with each error model. For
274 streamflow forecasting, we expect to obtain a reasonably accurate forecast from the
275 base hydrological model followed by an updating procedure as an auxiliary means to
276 improve the forecast accuracy. At lead times of many time-steps (e.g., streamflow
277 forecasts generated from medium-range rainfall forecasts) the magnitude of AR error
278 updates becomes rapidly smaller (tending to zero), and thus the performance of the
279 base hydrological model is crucial for realistic forecasts at longer lead times. While
280 we investigate only forecasts at a lead time of one time step in this study, we aim to
281 develop methods that can be applied to forecasts at longer lead times. Further, if the
282 base hydrological model does not replicate important catchment processes realistically,
283 the performance of the hydrological model outside the calibration period may be less
284 robust.

285 Figure 7 presents the Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970)
286 calculated from the base hydrological model and the error models. When the AR-
287 Norm model is used, the forecasts from the base hydrological model are very poor for
288 the Orara catchment ($NSE < 0$). The scatter plot in Figure 8 shows a serious over-
289 estimation of the streamflow simulation for the Orara. When the AR-Norm model is
290 used, the base hydrological model greatly over-estimates discharge and the AR-Norm
291 model then attempts to correct this systematic over-estimation. This is also shown in
292 Figure 4 where the base hydrological model has a strong tendency to over-estimate
293 streamflows for a range of streamflow magnitudes. The base hydrological model with
294 the AR-Norm model also performs poorly for the Abercrombie catchment (Figure 7).
295 In this case, the base hydrological model tends to under-estimate streamflows (results
296 not shown). For the other three catchments, however, the base hydrological model
297 with the AR-Norm model performs reasonably well.

298 In general, the AR-Raw base hydrological model performs as well or better than the
299 AR-Norm base hydrological model. The AR-Raw base hydrological model is notably

300 better than the AR-Norm base hydrological model in the Abercrombie and Orara
301 catchments (Figure 7). This suggests that more robust performance can be expected of
302 base hydrological models when AR models are applied to raw errors.

303 The RAR-Norm model generally improves the performance of the AR-Norm base
304 hydrological model to a level similar to the AR-Raw base hydrological model (Figure
305 7). The improvement over the AR-Norm base hydrological model is especially
306 evident for the Orara (Figures 4 and 7) and Abercrombie catchments (Figures 7).

307 We note that for the AR-Norm models, the updated forecasts are not always better
308 than forecasts generated by the base hydrological models. For the Tarwin and
309 Guadalupe catchments, AR-Norm forecasts are not as good as the forecasts generated
310 by the AR-Norm base hydrological model. This points to a tendency to overfit the
311 parameters to the calibration period, resulting in the error model undermining the
312 performance of the base hydrological model under cross-validation. Such a lack of
313 robustness is highly undesirable in forecasting applications, where the hydrological
314 models should be able to operate in conditions that differ from those experienced
315 during calibration. Note that this problem also occurs in the RAR-Norm model
316 (Guadalupe) and in the AR-Raw model (Abercrombie, Guadalupe) but to a much
317 smaller degree.

318 In general, the updated forecasts from the RAR-Norm model show similar or better
319 forecast accuracy, as measured by NSE, than both the AR-Raw model and the AR-
320 Norm model (Figure 7). We note that the Orara catchment is an exception: here the
321 AR-Raw model shows slightly better performance than RAR-Norm model.
322 Conversely, the RAR-Norm model shows notably better performance than both the
323 AR-Norm and AR-Raw models in the Abercrombie and Guadalupe catchments. This
324 suggests the RAR-Norm model may work better in intermittently flowing catchments,
325 although further testing is required to establish that this is true for a greater range of
326 catchments.

327 **4.4 Further analyses**

328 We further evaluate the NSE of the three different error models calibrated when
329 streamflows are receding (i.e. $\tilde{Q}_t \leq \tilde{Q}_{t-1}$) and rising (i.e. $\tilde{Q}_t > \tilde{Q}_{t-1}$) (Table 2). For the
330 receding streamflows (constituting 70-85% of streamflows), the AR-Raw model leads
331 to the overall worst forecast accuracy because of the over-correction explained in

332 Section 4.1. This is especially evident for the Abercrombie catchment (and, to a lesser
333 degree, the Guadalupe catchment). The RAR-Norm model significantly outperforms
334 the other two models for the Abercrombie catchment and shares similar forecast
335 accuracy to the (strongly performing) AR-Norm model for the other catchments.
336 When streamflows are rising (which also includes streamflow peaks), the AR-Norm
337 model can cause over-correction and leads to the least accurate forecasts (in terms of
338 NSE), and the RAR-Norm model behaves similarly to the AR-Raw model, which
339 consistently provides the most accurate forecasts. (The only exception is the
340 Guadalupe River, where the AR-Raw model clearly outperforms the RAR-Norm
341 model when streamflows are rising. This is somewhat compensated for by the
342 markedly better performance the RAR-Norm model offers over the AR-Raw model
343 when streamflows are receding for this catchment, leading to better forecasts overall
344 (Figure 7).) We conclude that the AR-Norm model generally tends to perform least
345 well when streamflows recede, and that the AR-Raw model tends to perform least
346 well when streamflows rise. We also conclude that the RAR-Norm model tends to
347 combine the best elements of the AR-Norm and AR-Raw models, leading to the best
348 overall performance.

349 We have shown that over-corrections can lead to inaccurate deterministic forecasts,
350 and we now discuss the consequences for the probabilistic predictions given by each
351 of the error models. We assess probabilistic forecast skill with skill scores derived
352 from two probabilistic verification measures: the Continuous Rank Probability Score
353 (CRPS) and the Root Mean Square Error in Probability (RMSEP) (denoted by
354 CRPS_SS and RMSEP_SS, respectively) (Wang and Robertson, 2011). Both skill
355 scores are calculated with respect to a reference forecast. The reference forecast is
356 generated by resampling historical streamflows: for a forecast issued for a given
357 month/year (e.g. February 1999), we randomly draw a sample of 1000 daily
358 streamflows that occurred in that month (e.g. February) from other years with
359 replacement (e.g. years other than 1999). Table 3 compares these two skill scores
360 calculated for the all catchments. The RAR-Norm model performs best across the
361 range of skill scores and catchments, attaining the highest CRPS_SS in 4 of the 7
362 catchments and the highest RMSEP_SS in 4 of 7 catchments. Even where RAR-Norm
363 was not the best performed model, it performs very similarly to the best performing
364 model in all cases. Interestingly, the AR-Raw model tends to outperform the AR-

365 Norm model in CRPS_SS while the reverse is true for RMSEP_SS. The CRPS tests
366 how appropriate the spread of uncertainty is for each probabilistic forecast, while
367 RMSEP puts little weight on this. The results suggest that while the median forecasts
368 of AR-Norm tends to be slightly more accurate than those of the AR-Raw model, the
369 forecast uncertainty is represented slightly better by the AR-Raw model.

370 To better understand how reliably the forecast uncertainty is quantified by each model,
371 we produce Probability Integral Transform (PIT) uniform probability plots (Wang and
372 Robertson, 2011) in Figure 9. There are two main points to draw from these plots.
373 First, the curves are very similar for all error models (a partial exception is the San
374 Marcos catchment, where the AR-Raw model is slightly closer to the one-to-one line
375 than the other models). This demonstrates that in general the models produce similarly
376 reliable uncertainty distributions. Second, all models show an inverted S-shaped curve,
377 which indicates that the uncertainty ranges are too wide. This underconfidence is a
378 result of using a Gaussian distribution to characterise the error. The Gaussian
379 distribution is not flexible enough to represent the high degree of kurtosis in the
380 distribution of the residuals after error updating (partly because the errors become
381 very small after updating). We are presently experimenting with other distributions in
382 order to address this issue, and will seek to publish this work in future. For the
383 purposes of the present study, we conclude that the three error models are similarly
384 reliable.

385 **5. Discussion and conclusions**

386 For streamflow forecasting, rainfall-runoff models are often augmented with an
387 updating procedure that corrects the forecast using information from recent simulation
388 errors. The most popular updating approach uses autoregressive (AR) models that
389 exploit the “memory” in model errors. AR models may be applied to raw errors
390 directly or to normalised errors.

391 We demonstrate three adverse effects of AR error updating procedures on seven
392 catchments. The first adverse effect is possible over-correction on the rising limb of
393 the hydrograph. The AR-Norm model can exhibit the tendency to over-correct the
394 peaks or on the rise of a hydrograph, because error updating can be (overly) amplified
395 by the back-transformation. The second adverse effect is the tendency to over-correct
396 receding hydrographs. This tendency is most prevalent in the AR-Raw model, which

397 can fail to recognise that a large error update may not be appropriate for small
398 streamflows.

399 The third adverse effect is that the stand-alone performance of base hydrological
400 model can be poor when the parameters of the rainfall-runoff model and the error
401 model are jointly estimated. We show that poor base hydrological model performance
402 is particularly prevalent in the AR-Norm model. The poor performance appears to
403 occur in catchments with highly skewed streamflow observations (the intermittent
404 Abercrombie River, and the Orara River, a catchment in a subtropical climate). For
405 example, in the Orara River, the base hydrological model tends to greatly over-
406 estimate streamflows, and then relies on the error updating to correct the over-
407 estimates. This is not desirable in real-time forecasting applications for two major
408 reasons. First, modern streamflow forecasting systems often extend forecast lead-
409 times with rainfall forecast information (Bennett et al., 2014). The magnitude of AR
410 updating decays with lead time, and forecasts at longer lead times rely heavily on the
411 performance of the base hydrological model. Second, hydrological models are
412 designed to simulate various components of natural systems, such as baseflow
413 processes or overland flow. In theory, simulating these processes correctly will allow
414 the model to perform well for climate conditions that may substantially differ from
415 those experienced during the parameter estimation period. If the hydrological model
416 parameters do not reflect the natural processes for a given catchment, the hydrological
417 model may be much less robust outside the parameter estimation period.

418 We note that the poor performance of the hydrological model may be specific to the
419 GR4J model, and may not occur in other hydrological models. Evin et al. (2014)
420 estimated hydrological model and error model parameters jointly using GR4J and
421 another hydrological model, HBV, for the three US catchments tested here. While
422 they did not assess the performance of the base hydrological models, they found that
423 HBV tended to perform more robustly when combined with different error models. It
424 is possible that we may have achieved more stable base model performance had we
425 used HBV or another hydrological model. We note, however, that our conclusions can
426 probably be generalised to other hydrological models that do not offer robust base
427 model performance under joint parameter estimation (e.g. GR4J). Because the RAR-
428 Norm model limits the range of updating that can be applied, it will tend to rely more
429 heavily on the base hydrological model, and therefore will tend to favour parameter

430 sets that encourage good stand-alone performance of the base model. For those
431 hydrological models that already produce robust base model performance under joint
432 parameter estimation (perhaps HBV), RAR-Norm is unlikely to undermine this
433 performance for the same reasons. We see some evidence of this in our experiments
434 with GR4J: when the performance of the base hydrological model is already strong
435 relative to the updated forecasts for the AR-Norm and AR-Raw models (e.g. the
436 Tarwin, Mitta Mitta, or Guadalupe catchments), the RAR-Norm model base
437 hydrological model also performs strongly.

438 The tendency of the AR-Norm model to over-correct rising streamflows is probably
439 generic. In particular, transformations other than the log-sinh transformation may still
440 lead to over-correction at the peak of hydrograph. The proof in Appendix A shows
441 that if a transformation satisfies some conditions (first derivate is positive and second
442 derivate is negative), it will tend to correct more for higher forecast streamflows and
443 can cause the problem of over-correction. The conditions given by Appendix A are
444 generally true for many other transformations used for data normalisation and
445 variance stabilisation in hydrological applications, such as logarithm transformation
446 or the Box-Cox transformation with the power parameter less than 1.

447 We use joint parameter inference to calibrate hydrological model and error model
448 parameters, in order to address the true nature of underlying model errors. Inferring
449 parameters of the error model and the base hydrological model independently – i.e.,
450 first inferring parameters of the base hydrological model, holding these constant and
451 then inferring the error model parameters - relies on simplified and often invalid error
452 assumptions (it assumes independent, homoscedastic and Gaussian errors), but
453 nonetheless could be a pragmatic alternative to the joint parameter inference to reduce
454 computational demands. The over-correction of conventional AR models is
455 independent of the parameter inference, whether the error and base hydrological
456 model parameters are inferred jointly or independently.

457 In order to mitigate the adverse effects of conventional AR updating procedures, we
458 introduce a new updating procedure called the RAR-Norm model. The RAR-Norm
459 model is a modification of the AR-Norm model: in most instances it operates as the
460 AR-Norm model, but in instances of possible over-correction it relies on the error in
461 untransformed streamflows at the previous time step. That is, RAR-Norm is
462 essentially a more conservative error model than AR-Norm: in situations where

463 streamflows change rapidly, it opts to update with whichever error (transformed or
 464 untransformed) is smaller. This forces greater reliance on the base hydrological model
 465 to simulate streamflows accurately, leading to more robust performance in the base
 466 hydrological model. The RAR-Norm model clearly outperforms the AR-Norm model
 467 in both the updated and base model forecasts, as well as ameliorating the problem of
 468 over-correcting rising streamflows. The RAR-Norm model's advantage over the AR-
 469 Raw model is less clear: both the base hydrological model and the updated forecasts
 470 produced by the AR-Raw model perform similarly to (or sometimes slightly better
 471 than) the RAR-Norm model. However, the RAR-Norm model clearly addresses the
 472 problem of over-correcting receding streamflows that occurs in the AR-Raw model.
 473 As we show, this type of over-correction can seriously distort event hydrographs, and
 474 cause forecasts of near zero streamflows when reasonably substantial streamflows are
 475 observed. While these instances are not very common, the failure in the forecast is a
 476 serious one. As we note earlier, the over-correction of receding streamflows is likely
 477 to be exacerbated when producing forecasts at lead times of more than one time step.
 478 Accordingly, we contend that the RAR-Norm model is preferable to both AR-Norm
 479 and AR-Raw models for streamflow forecasting applications.

480 **Appendix A**

481 For brevity we only show the case of the AR-Norm model; analogous arguments can
 482 be used to prove the cases of the AR-Raw and RAR-Norm models. The streamflow
 483 ensemble forecast Q_t given by the AR-Norm model defined by (1) can be written as

$$484 \quad Q_t = \max \left[f^{-1} \left\{ \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t \right\}, 0 \right]. \quad (\text{A1})$$

485 where negative values after the back-transformation are assigned zero values. Because
 486 we assume that ε_t is a standard normal random variable, to show that \tilde{Q}_t^* is the
 487 median of Q_t we need only show that $P(Q_t \leq \tilde{Q}_t^*) = 0.5$, which can be proved as
 488 follows:

$$489 \quad \begin{aligned} P(Q_t \leq \tilde{Q}_t^*) &= P \left(\max \left[f^{-1} \left\{ \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t \right\}, 0 \right] \leq \tilde{Q}_t^* \right) \\ &= P \left(f^{-1} \left\{ \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t \right\} \leq \tilde{Q}_t^* \text{ and } 0 \leq \tilde{Q}_t^* \right). \end{aligned} \quad (\text{A2})$$

490 Because \tilde{Q}_t^* always has a non-negative value, we have

$$\begin{aligned}
491 \quad P(Q_t \leq \tilde{Q}_t^*) &= P\left(f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t\} \leq f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\}\right) \\
&= P(\varepsilon_t \leq 0) = 0.5
\end{aligned} \tag{A3}$$

492 Appendix B

493 We will show analytically that the AR-Norm model gives a larger magnitude of the
494 error update for a higher forecast streamflow.

495 Firstly, we will show that the first derivate of the log-sinh transform f defined by (3)
496 is positive and the second derivate is negative (i.e. $f'(x) > 0$ and $f''(x) < 0$) for any
497 $b > 0$ and any x . Following some simple manipulation, we have

$$498 \quad f'(x) = \frac{\cosh(a+bx)}{\sinh(a+bx)} > 0 \quad \text{and} \quad f''(x) = \frac{-b}{\sinh^2(a+bx)} < 0 \tag{B1}$$

499 Using the differentiation of inverse functions, we find the first and second derivatives of
500 the inverse transform f^{-1}

$$501 \quad [f^{-1}]'(x) = \frac{1}{f'\{f^{-1}(x)\}} > 0 \quad \text{and} \quad [f^{-1}]''(x) = \frac{-f''\{f^{-1}(x)\}}{[f'\{f^{-1}(x)\}]^3} > 0, \tag{B2}$$

502 for any $b > 0$ and any x .

503 Next, we will derive the difference in magnitudes of the error update between low and
504 high forecast streamflows. For the sake of notation simplicity, we rewrite $q = \tilde{Z}_t$ and
505 $u = \rho(Z_{t-1} - \tilde{Z}_{t-1})$ and assume that $u > 0$. Using Equation (4), the updated streamflow
506 can be written as $\tilde{Q}_t^* = f^{-1}(q+u)$. The magnitude of the error update can be written as

$$507 \quad |\tilde{Q}_t^* - \tilde{Q}_t| = |f^{-1}(q+u) - f^{-1}(q)| = \begin{cases} \int_0^u [f^{-1}]'(x+q) dx & \text{if } u > 0 \\ 0 & \\ \int_u^0 [f^{-1}]'(x+q) dx & \text{otherwise.} \end{cases} \tag{B3}$$

508 Suppose that we have two forecast streamflows $\tilde{Q}_{t,1} \leq \tilde{Q}_{t,2}$ and denote the normalised
509 forecast streamflow by $q_1 = \tilde{Z}_{t,1}$ and $q_2 = \tilde{Z}_{t,2}$ and the updated streamflow by $\tilde{Q}_{t,1}^*$ and
510 $\tilde{Q}_{t,2}^*$. Because f is an increasing function, we have $q_1 \leq q_2$. The difference in the
511 magnitude of the error update between $\tilde{Q}_{t,1}$ and $\tilde{Q}_{t,2}$ can be derived as

$$\begin{aligned}
512 \quad |\tilde{Q}_{t,1} - \tilde{Q}_{t,1}^*| - |\tilde{Q}_{t,2} - \tilde{Q}_{t,2}^*| &= \begin{cases} \int_0^u \left\{ [f^{-1}]'(x+q_1) - [f^{-1}]'(x+q_2) \right\} dx & \text{if } u > 0 \\ \int_u^0 \left\{ [f^{-1}]'(x+q_1) - [f^{-1}]'(x+q_2) \right\} dx & \text{otherwise.} \end{cases} \quad (B4)
\end{aligned}$$

513 From (A2), we have shown that $[f^{-1}]'$ is a positive increasing function and this
514 ensures that $[f^{-1}]'(x+q_1) - [f^{-1}]'(x+q_2) \leq 0$. Finally we have

$$515 \quad |\tilde{Q}_{t,1} - \tilde{Q}_{t,1}^*| \leq |\tilde{Q}_{t,2} - \tilde{Q}_{t,2}^*|. \quad (B5)$$

516 Therefore, the error update at larger forecast streamflows is always larger than the
517 error update at lower forecast streamflows.

518 **Acknowledgments**

519 This work is part of the WIRADA (Water Information Research and Development
520 Alliance) streamflow forecasting project funded under CSIRO Water for a Healthy
521 Country Flagship. We would like to thank Durga Shrestha (CSIRO) for valuable
522 suggestions that led to substantial strengthening of the manuscript. We would like to
523 thank two reviewers, Bettina Schaepli and Mark Thyer, for their careful reviews and
524 valuable recommendations, which have improved the quality of this manuscript
525 considerably.

526 **Table of Tables**

527 Table 1: Catchment characteristics.

528 Table 2: Comparison of the NSE calculated at (a) the receding limb and (b) the rising
529 limb of the hydrograph for three different error models.530 Table 3: Comparison of the skill scores based on CRPS and RMSEP (denoted by
531 CRPS_SS and RMSEP_SS) for three different error models.532 **Table of Figures**

533 Figure 1: Map of US (top) and Australian (bottom) catchments.

534 Figure 2: An example of over-correction caused by the AR-Norm model in the Mitta
535 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,
536 without error updating). Solid lines: forecasts with error updating.537 Figure 3: The fraction of instances where $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$ (i.e., instances where over-
538 correction may occur in the AR-Norm model and where error updating is restricted in
539 the RAR-Norm model) for the AR-Norm and RAR-Norm models for Australian
540 catchments.

541 Figure 4: Forecast streamflows for the Orara catchment for an example 1-year period.

542 Top panel shows streamflows forecast with AR-Norm model, bottom panel shows
543 streamflows forecast with the RAR-Norm model. Dashed lines: forecasts from the
544 base hydrological model (i.e., without error updating). Solid lines: forecasts with error
545 updating. Tick marks in the x-axis denote the instance of updating where546 $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$.547 Figure 5: An example of over-correction caused by the AR-Raw model in the Mitta
548 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,
549 without error updating). Solid lines: forecasts with error updating.550 Figure 6: Forecast streamflows for the Abercrombie catchment for the period between
551 01/08/1997 and 15/09/1997. Top panel shows streamflows forecast with AR-Raw
552 model, bottom panel shows streamflows forecast with the RAR-Norm model. Dashed
553 lines: forecasts from the base hydrological model (i.e., without error updating). Solid
554 lines: forecasts with error updating. Gray shading denotes instances of over-correction
555 caused by the AR-Raw model.556 Figure 7: NSE of streamflows forecast with the AR-Norm, AR-Raw and RAR-Norm
557 models (colours). Performance of the corresponding base hydrological models is
558 shown by hatched blocks.559 Figure 8: Comparison of the observed streamflows (Q_t) and forecast streamflows (\tilde{Q}_t),
560 as forecast: 1) with the base hydrological model (circles), and 2) with the base
561 hydrological model and error updating models (dots) for the Orara catchment.562 Figure 9: PIT-uniform probability plots. Curves on the diagonal indicate perfectly
563 reliable forecasts.

564

565

566 **References**

567 Bates, B. C., and Campbell, E. P.: A Markov chain Monte Carlo scheme for parameter
568 estimation and inference in conceptual rainfall-runoff modeling, *Water Resour Res*,
569 37, 937-947, 10.1029/2000wr900363, 2001.

570 Bennett, J. C., Robertson, D. E., Shrestha, D. L., Wang, Q. J., Enever, D.,
571 Hapuarachchi, P., and Tuteja, N. K.: A System for Continuous Hydrological
572 Ensemble Forecasting (SCHEF) to lead times of 9 days, *J Hydrol*,
573 10.1016/j.jhydrol.2014.08.010, 2014.

574 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P., and
575 Rieckermann, J.: Improving uncertainty estimation in urban hydrological modeling by
576 statistically describing bias, *Hydrol. Earth Syst. Sci.* , 17, 4209-4225, 10.5194/hess-
577 17-4209-2013, 2013.

578 Duan, Q., Schaake, J., Andreassian, V., Franks, S., Goteti, G., Gupta, H. V., Gusev, Y.
579 M., Habets, F., Hall, A., Hay, L., Hogue, T., Huang, M., Leavesley, G., Liang, X.,
580 Nasonova, O. N., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T., and Wood, E.
581 F.: Model Parameter Estimation Experiment (MOPEX): An overview of science
582 strategy and major results from the second and third workshops, *J Hydrol*, 320, 3-17,
583 10.1016/j.jhydrol.2005.07.031, 2006.

584 Duan, Q. Y., Sorooshian, S., and Gupta, V. K.: Optimal Use of the Sce-Ua Global
585 Optimization Method for Calibrating Watershed Models, *J Hydrol*, 158, 265-284,
586 10.1016/0022-1694(94)90057-4, 1994.

587 Engeland, K., Renard, B., Steinsland, I., and Kolberg, S.: Evaluation of statistical
588 models for forecast errors from the HBV model, *J Hydrol*, 384, 142-155,
589 10.1016/j.jhydrol.2010.01.018, 2010.

590 Evin, G., Kavetski, D., Thyer, M., and Kuczera, G.: Pitfalls and improvements in the
591 joint inference of heteroscedasticity and autocorrelation in hydrological model
592 calibration, *Water Resour Res*, 49, 4518-4524, 10.1002/wrcr.20284, 2013.

593 Evin, G., Thyer, M., Kavetski, D., McInerney, D., and Kuczera, G.: Comparison of
594 joint versus postprocessor approaches for hydrological uncertainty estimation

- 595 accounting for error autocorrelation and heteroscedasticity, *Water Resour Res*, 50,
596 2350-2375, 10.1002/2013WR014185, 2014.
- 597 Jones, D. A., Wang, W., and Fawcett, R.: High-quality spatial climate data-sets for
598 Australia, *Australian Meteorological and Oceanographic Journal*, 58, 233-248, 2009.
- 599 Kavetski, D., Franks, S. W., and Kuczera, G.: Confronting Input Uncertainty in
600 Environmental Modelling, in: *Calibration of Watershed Models*, edited by: Duan, Q.,
601 Gupta, H. V., Sorooshian, S., Rousseau, A. N., and Turcotte, R., American
602 Geophysical Union, Washington D.C., 49-68, 2003.
- 603 Li, M., Wang, Q. J., and Bennett, J.: Accounting for seasonal dependence in
604 hydrological model errors and prediction uncertainty, *Water Resour Res*, 49, 5913-
605 5929, 10.1002/wrcr.20445, 2013.
- 606 Morawietz, M., Xu, C. Y., and Gottschalk, L.: Reliability of autoregressive error
607 models as post-processors for probabilistic streamflow forecasts, *Adv. Geosci.*, 29,
608 109-118, 10.5194/adgeo-29-109-2011, 2011.
- 609 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models
610 part I — A discussion of principles, *J Hydrol*, 10, 282-290, 10.1016/0022-
611 1694(70)90255-6, 1970.
- 612 Perrin, C., Michel, C., and Andreassian, V.: Improvement of a parsimonious model
613 for streamflow simulation, *J Hydrol*, 279, 275-289, 10.1016/S0022-1694(03)00225-7,
614 2003.
- 615 Pianosi, F., and Raso, L.: Dynamic modeling of predictive uncertainty by regression
616 on absolute errors, *Water Resour Res*, 48, W03516, 10.1029/2011wr010603, 2012.
- 617 Salamon, P., and Feyen, L.: Disentangling uncertainties in distributed hydrological
618 modeling using multiplicative error models and sequential data assimilation, *Water*
619 *Resour Res*, 46, W12501, 10.1029/2009wr009022, 2010.
- 620 Schaeffli, B., Talamba, D. B., and Musy, A.: Quantifying hydrological modeling errors
621 through a mixture of normal distributions, *J Hydrol*, 332, 303-315,
622 10.1016/j.jhydrol.2006.07.005, 2007.

- 623 Schoups, G., and Vrugt, J. A.: A formal likelihood function for parameter and
624 predictive inference of hydrologic models with correlated, heteroscedastic, and non-
625 Gaussian errors, *Water Resour Res*, 46, W10531, 10.1029/2009wr008933, 2010.
- 626 Seo, D. J., Herr, H. D., and Schaake, J. C.: A statistical post-processor for accounting
627 of hydrologic uncertainty in short-range ensemble streamflow prediction, *Hydrol.*
628 *Earth Syst. Sci. Discuss.*, 3, 1987-2035, 10.5194/hessd-3-1987-2006, 2006.
- 629 Shrestha, D. L., and Solomatine, D. P.: Data - driven approaches for estimating
630 uncertainty in rainfall - runoff modelling, *International Journal of River Basin*
631 *Management*, 6, 109-122, 10.1080/15715124.2008.9635341, 2008.
- 632 Thyer, M., Kuczera, G., and Wang, Q. J.: Quantifying parameter uncertainty in
633 stochastic models using the Box-Cox transformation, *J Hydrol*, 265, 246-257,
634 10.1016/S0022-1694(02)00113-0, 2002.
- 635 Vaze, J., Perraud, J. M., Teng, J., Chiew, F. H. S., Wang, B., and Yang, Z.: Catchment
636 Water Yield Estimation Tools (CWYET), the 34th World Congress of the
637 International Association for Hydro- Environment Research and Engineering: 33rd
638 Hydrology and Water Resources Symposium and 10th Conference on Hydraulics in
639 Water Engineering, Brisbane, 2011.
- 640 Wang, Q. J., and Robertson, D. E.: Multisite probabilistic forecasting of seasonal
641 flows for streams with zero value occurrences, *Water Resour Res*, 47, W02546,
642 10.1029/2010WR009333, 2011.
- 643 Wang, Q. J., Shrestha, D. L., Robertson, D. E., and Pokhrel, P.: A log-sinh
644 transformation for data normalization and variance stabilization, *Water Resour Res*,
645 48, W05514, 10.1029/2011WR010973, 2012.
- 646 World Meteorological Organization: Simulated real-time intercomparison of
647 hydrological models, World Meteorological Organization, Geneva, Switzerland, 1992.
- 648 Xiong, L. H., and O'Connor, K. M.: Comparison of four updating models for real-time
649 river flow forecasting, *Hydrolog Sci J*, 47, 621-639, 10.1080/02626660209492964,
650 2002.

651 Xu, C. Y.: Statistical analysis of parameters and residuals of a conceptual water
652 balance model - Methodology and case study, *Water Resour Manag*, 15, 75-92,
653 10.1023/A:1012559608269, 2001.

654

655 Table 1: Catchment characteristics.

Name	Country	Gauge Site	Area (km ²)	Rainfall (mm/yr)	Streamflow (mm/yr)	Runoff coefficient	Zero flows
Abercrombie	Aus	Abercrombie River at Hadley no. 2	1447	783	63	0.08	14.4%
Mitta Mitta	Aus	Mitta Mitta River at Hinnomunjie	1527	1283	261	0.20	0
Orara	Aus	Orara River at Bawden Bridge	1868	1176	243	0.21	0.6%
Tarwin	Aus	Tarwin River at Meeniyan	1066	1042	202	0.19	0
Amite	US	07378500	3315	1575	554	0.35	0
Guadalupe	US	08167500	3406	772	104	0.13	1.7%
San Marcos	US	08172000	2170	844	165	0.20	0%

656

657

658 Table 2: Comparison of the NSE calculated at (a) the receding limb and (b) the rising
 659 limb of the hydrograph for three different error models.
 660

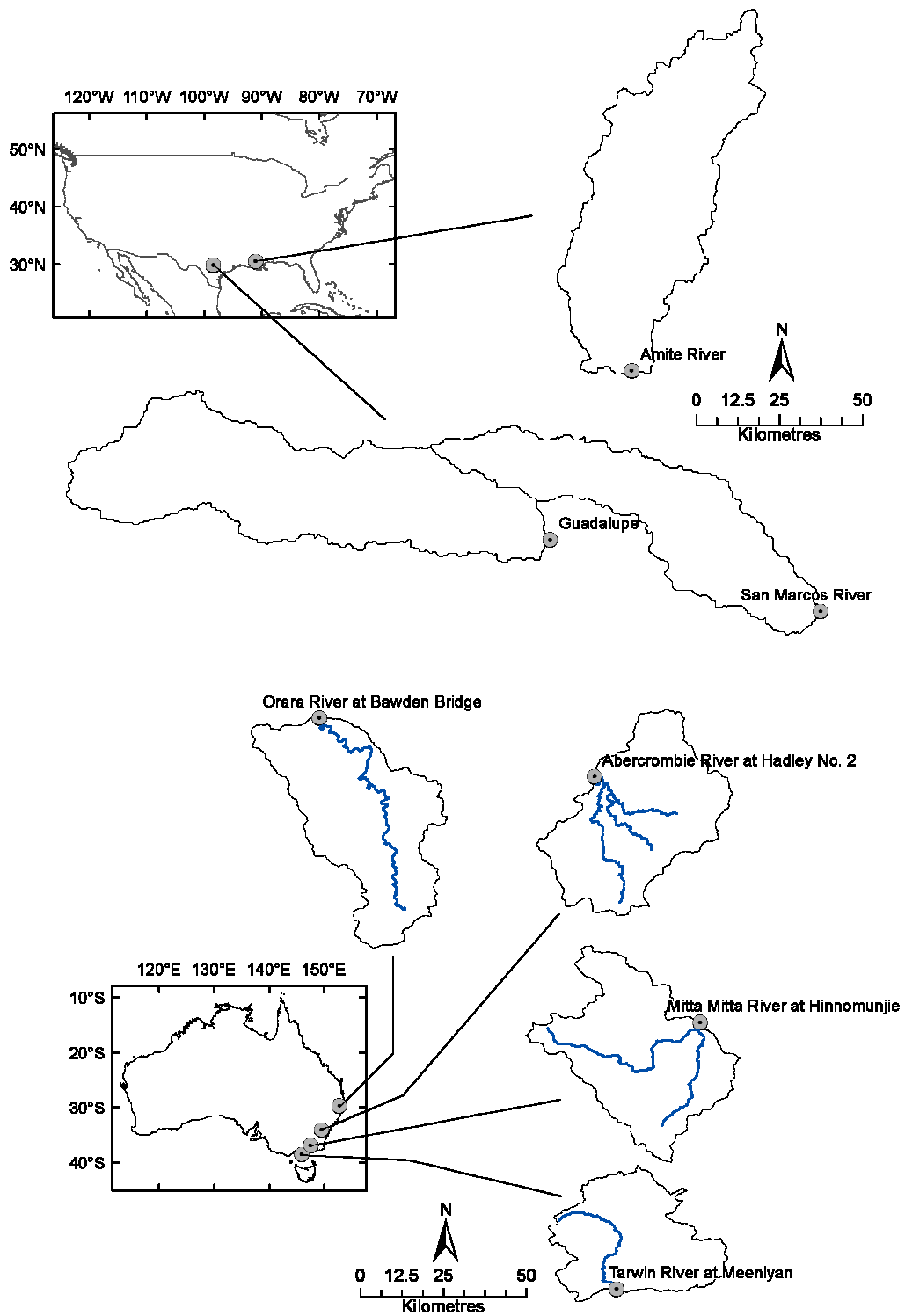
	(a) $\tilde{Q}_t \leq \tilde{Q}_{t-1}$				(b) $\tilde{Q}_t > \tilde{Q}_{t-1}$			
	Proportion of flows	AR- Norm	AR- Raw	RAR- Norm	Proportion of flows	AR- Norm	AR- Raw	RAR- Norm
Abercrombie	82%	0.11	-0.41	0.52	19%	0.58	0.66	0.65
Mitta Mitta	82%	0.95	0.91	0.95	18%	0.81	0.86	0.86
Orara	85%	0.94	0.91	0.95	15%	0.86	0.86	0.83
Tarwin	71%	0.90	0.91	0.90	29%	0.18	0.77	0.76
Amite	69%	0.76	0.82	0.84	31%	0.82	0.82	0.85
Guadalupe	83%	0.75	0.35	0.77	15%	0.24	0.55	0.45
San Marcos	82%	0.80	0.66	0.80	17%	0.63	0.64	0.64

661

662 Table 3: Comparison of the skill scores based on CRPS and RMSEP (denoted by
 663 CRPS_SS and RMSEP_SS) for three different error models.
 664

	CRPS_SS (%)			RMSEP_SS (%)		
	AR-Norm	AR-Raw	RAR-Norm	AR-Norm	AR-Raw	RAR-Norm
Abercrombie	64.1	62.3	66.3	75.1	73.7	74.7
Mitta Mitta	80.3	79.7	80.7	84.1	83.2	84.0
Orara	74.0	75.7	75.5	81.7	80.7	81.4
Tarwin	74.9	79.3	78.8	86.1	85.1	86.1
Amite	67.5	68.3	69.5	71.0	70.9	71.2
Guadalupe	57.4	60.9	59.8	76.3	75.2	77.2
San Marcos	68.8	66.0	68.9	73.9	73.9	74.3

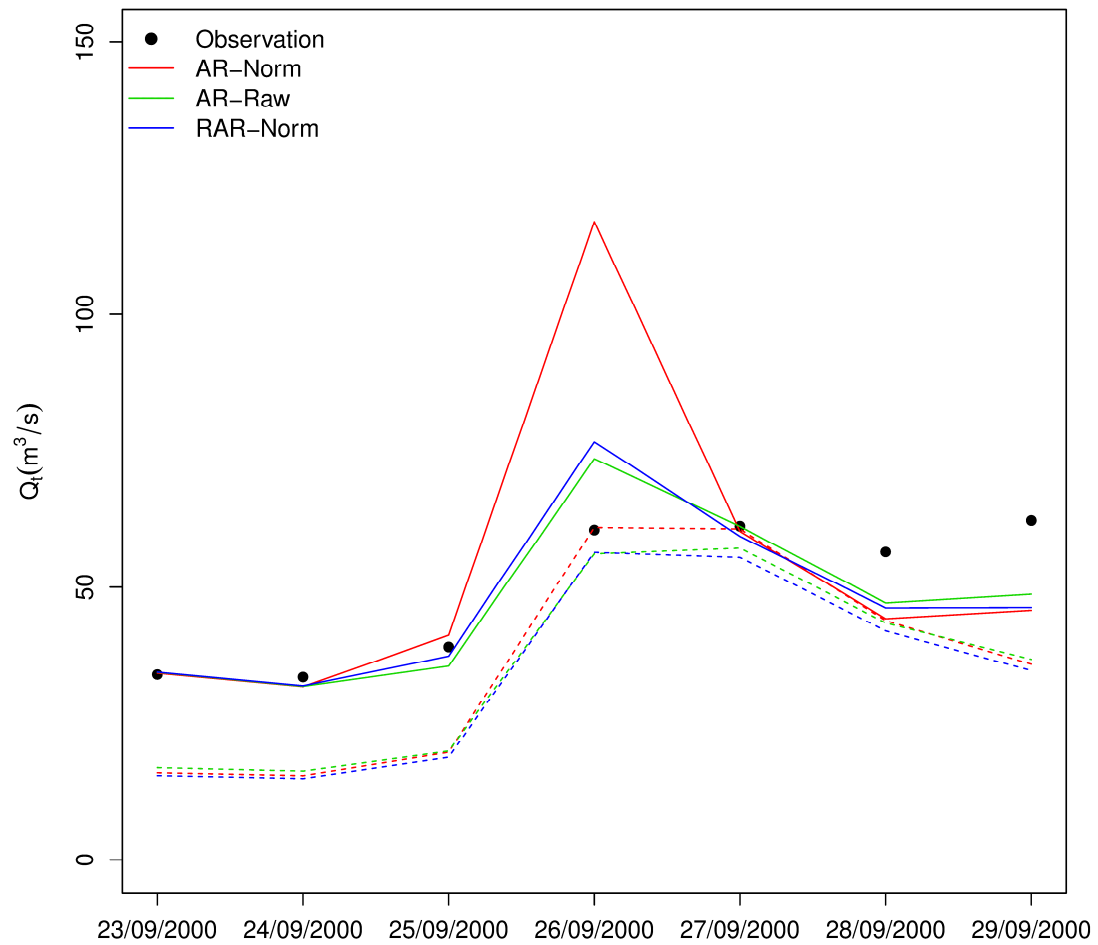
665



666

667 Figure 1: Map of US (top) and Australian (bottom) catchments.

668



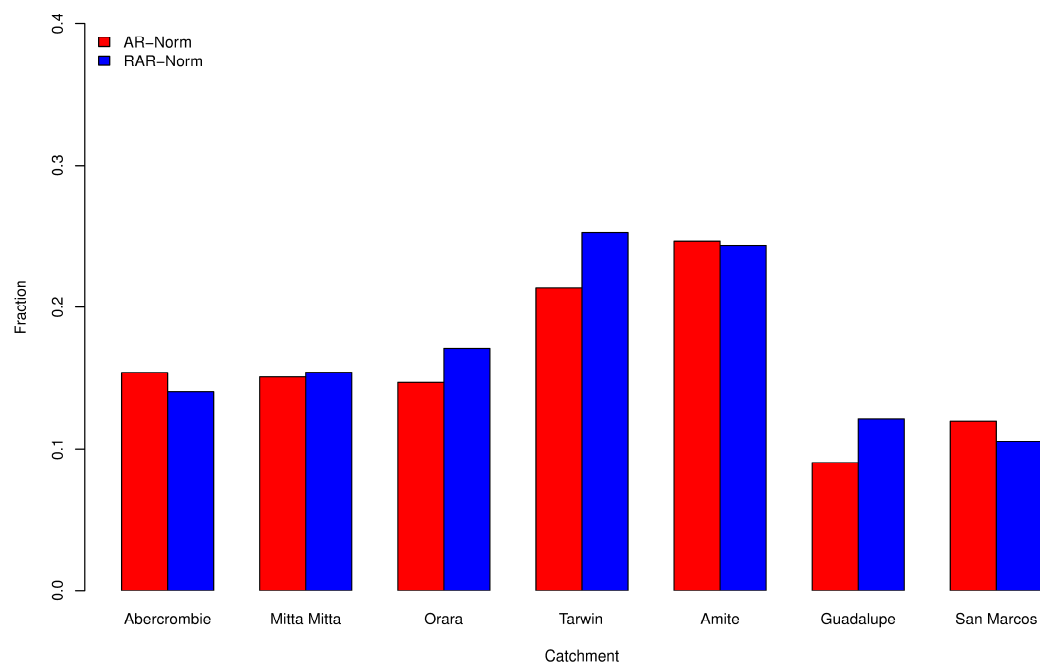
669

670 Figure 2: An example of over-correction caused by the AR-Norm model in the Mitta
 671 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,
 672 without error updating). Solid lines: forecasts with error updating.

673

674

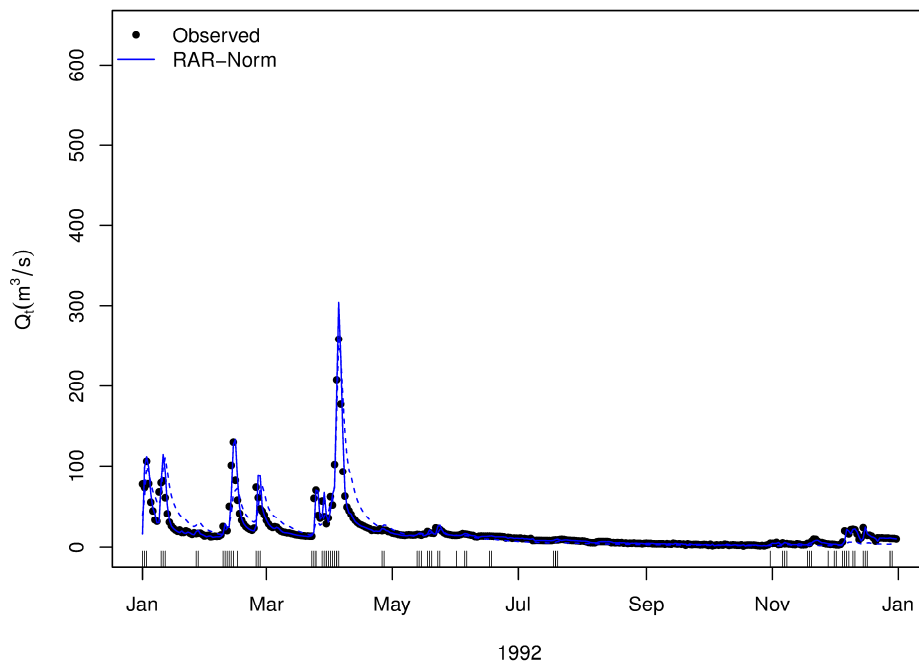
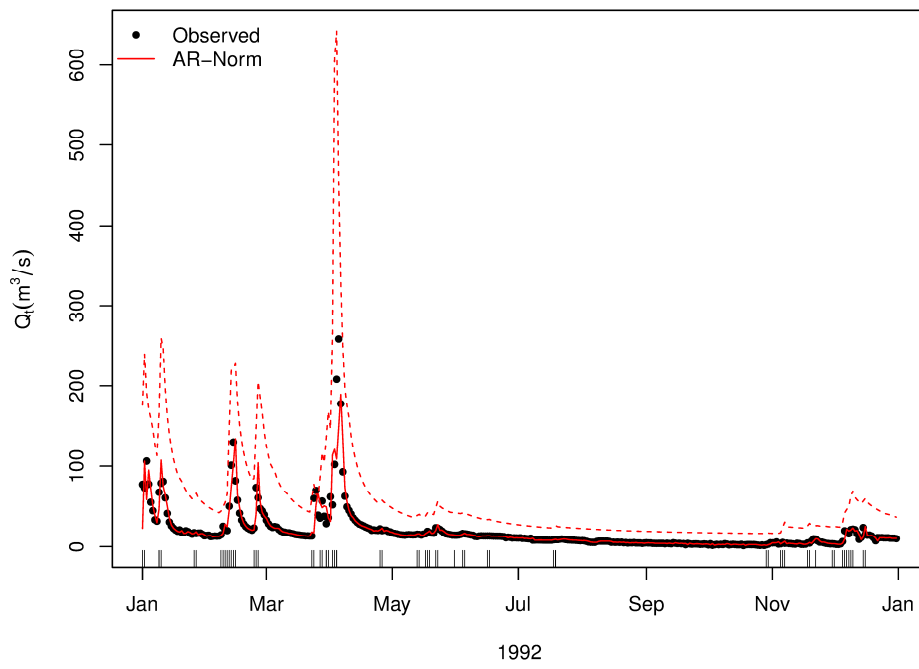
675



676

677 Figure 3: The fraction of instances where $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$ (i.e., instances where over-
 678 correction may occur in the AR-Norm model and where error updating is restricted in
 679 the RAR-Norm model) for the AR-Norm and RAR-Norm models for Australian
 680 catchments.

681



682

683 Figure 4: Forecast streamflows for the Orara catchment for an example 1-year period.

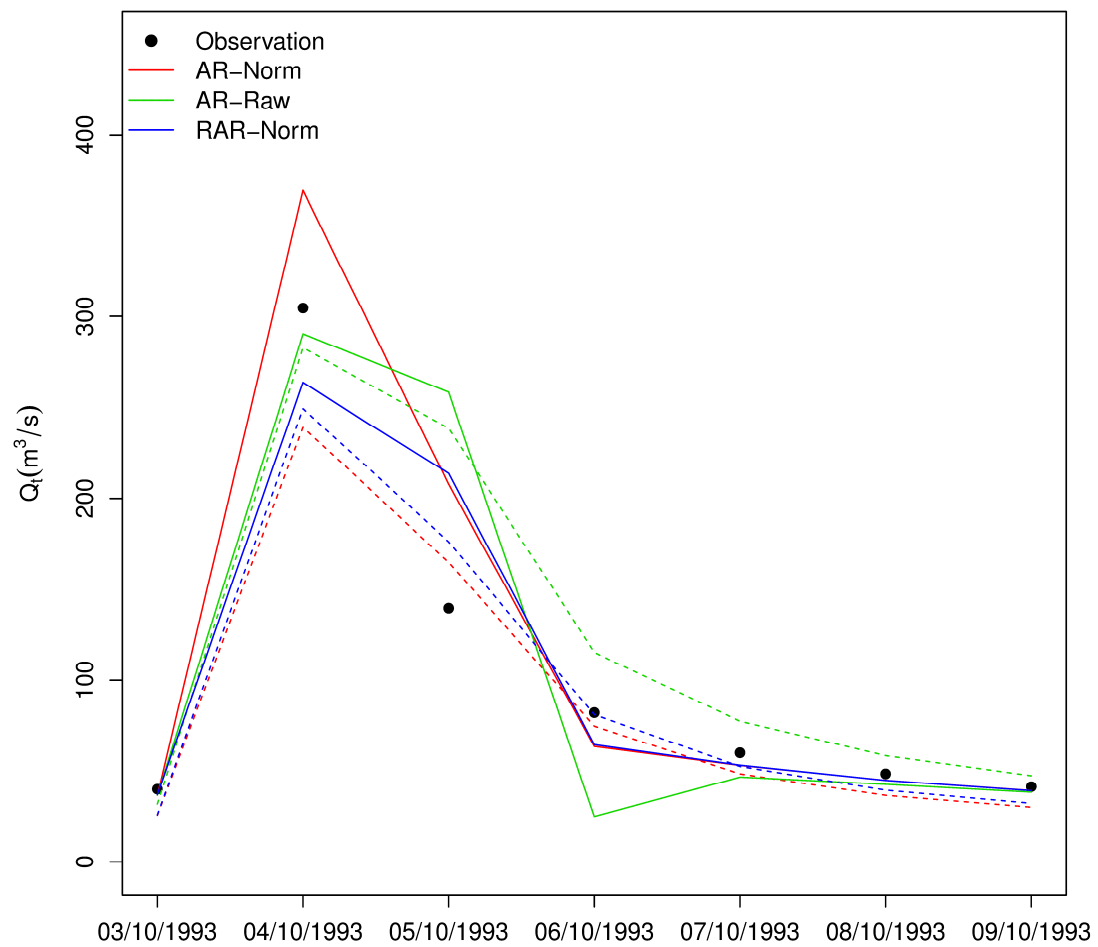
684 Top panel shows streamflows forecast with AR-Norm model, bottom panel shows

685 streamflows forecast with the RAR-Norm model. Dashed lines: forecasts from the

686 base hydrological model (i.e., without error updating). Solid lines: forecasts with error

687 updating. Tick marks in the x-axis denote the instance of updating where

688 $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$.



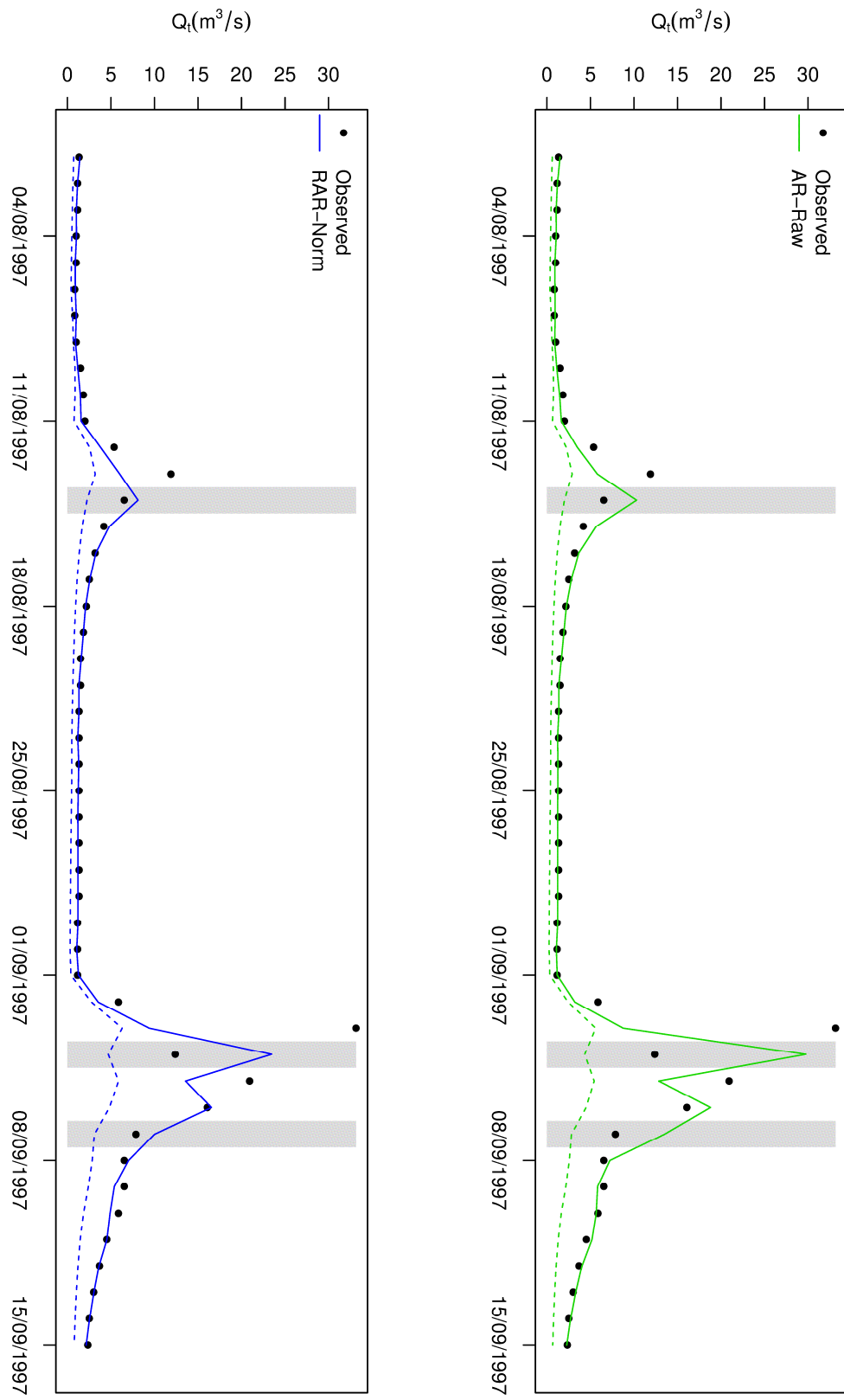
689

690 Figure 5: An example of over-correction caused by the AR-Raw model in the Mitta

691 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,

692 without error updating). Solid lines: forecasts with error updating.

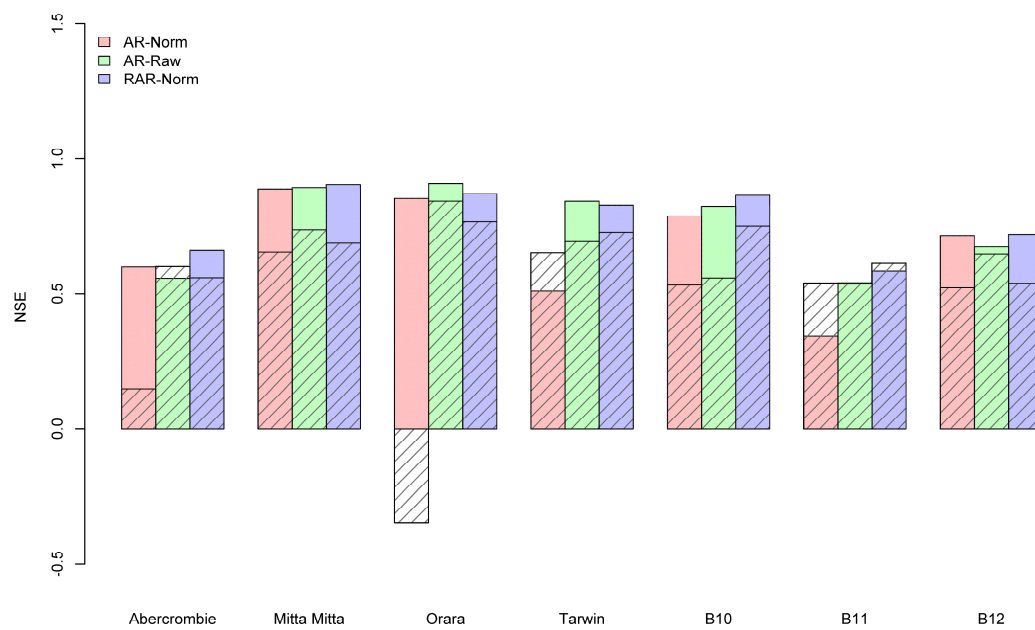
693



694

695 Figure 6: Forecast streamflows for the Abercrombie catchment for the period between
 696 01/08/1997 and 15/09/1997. Top panel shows streamflows forecast with AR-Raw
 697 model, bottom panel shows streamflows forecast with the RAR-Norm model. Dashed

698 lines: forecasts from the base hydrological model (i.e., without error updating). Solid
699 lines: forecasts with error updating. Gray shading denotes instances of over-correction
700 caused by the AR-Raw model.

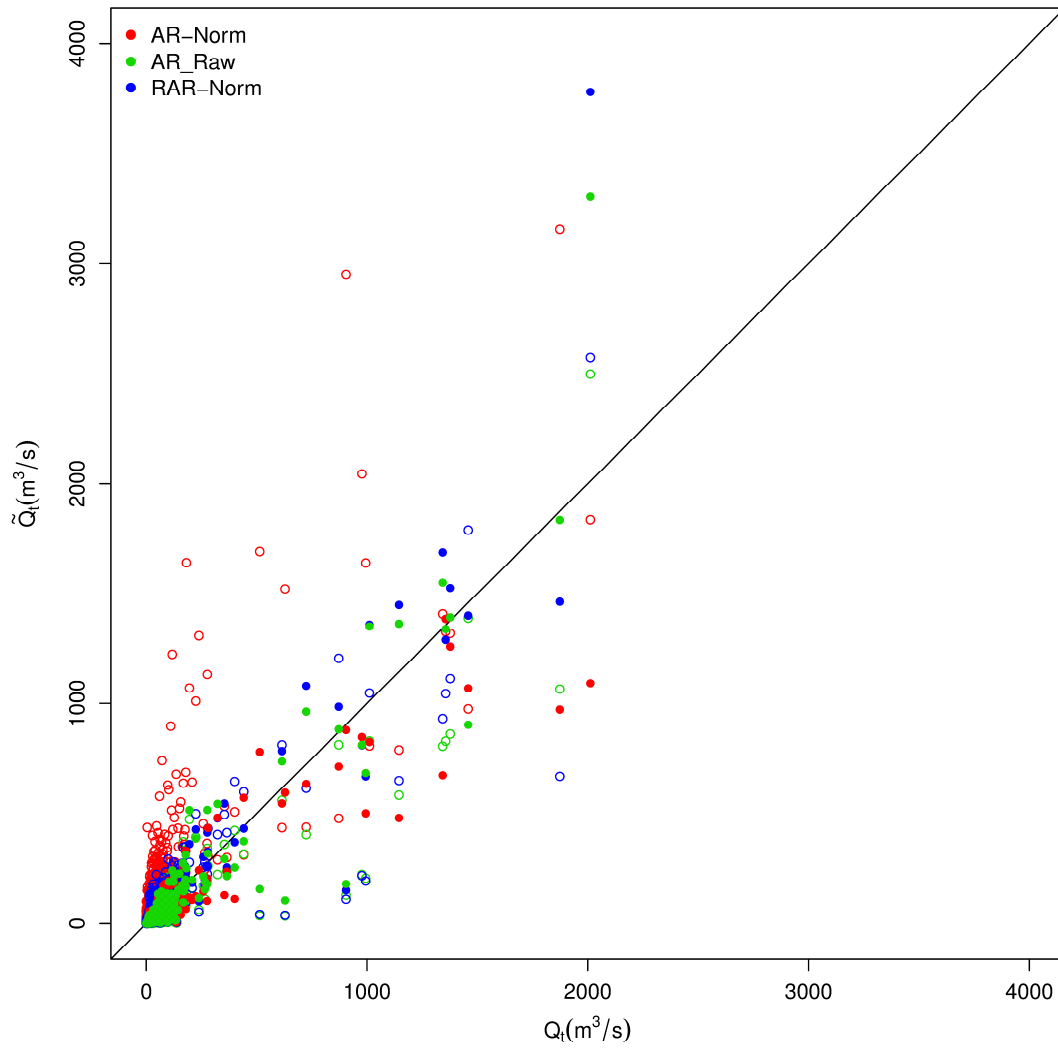


701

702 Figure 7: NSE of streamflows forecast with the AR-Norm, AR-Raw and RAR-Norm
 703 models (colours). Performance of the corresponding base hydrological models is
 704 shown by hatched blocks.

705

706



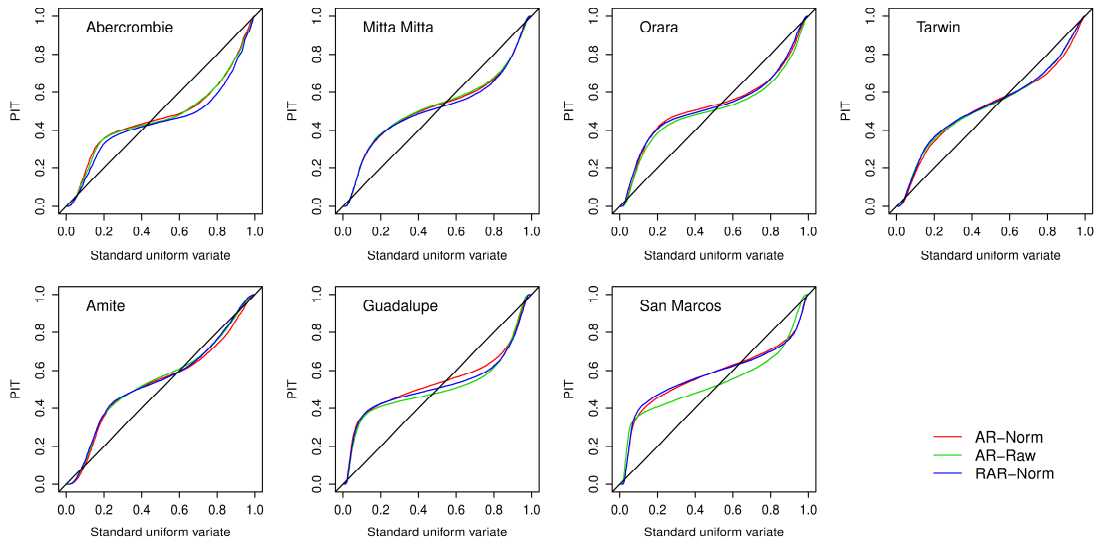
707

708 Figure 8: Comparison of the observed streamflows (Q_t) and forecast streamflows (\tilde{Q}_t),

709 as forecast: 1) with the base hydrological model (circles), and 2) with the base

710 hydrological model and error updating models (dots) for the Orara catchment.

711



712

713 Figure 9: PIT-uniform probability plots. Curves on the diagonal indicate perfectly

714 reliable forecasts.