

Review comments from Bettina Schaepli

Comment #1.1. This manuscript proposes a new method to correct forecasted streamflows based on the forecast error of the previous time step. The method represents a modification of the commonly applied autoregressive correction. The paper is well written, the method and the results concisely presented and discussed. However, the presented results did not convince me that the new method really outperforms the reference method; this might easily be improved by showing more details of the performed tests.

Response: Thank you very much for the time and effort you have taken to carefully review this manuscript. We feel your comments are very valuable for us to improve the quality of the manuscript considerably. We appreciate your positive feedback. We have paid serious attention to your suggestions and have attempted to address all your concerns, especially on the model performance comparison. Please refer to specific responses to your comments below.

My suggestion for moderate revisions of this paper are:

General comment on used terms

Comment #1.2. I would carefully revise the used wording to clearly distinguish between “forecast” (prediction of the system state at a given moment in time) from the more general “prediction”. At the moment, the two terms are used interchangeably, which might sometimes be misleading, especially because the discussed streamflow correction only applies to forecasting.

Response: We have now used the word *forecast* throughout the paper, and have removed the word *prediction*. We occasionally use the word ‘simulation’ to differentiate instances where forecast rainfalls would never be used to force a rainfall runoff model, and make the distinction in the text, as follows:

“Our study is aimed at streamflow forecasting applications, so we preserve the distinction between observed and forecast forcings by referring to streamflows modelled with observed rainfall as *simulations* and those modelled with forecast rainfall as *forecasts*. As the forecast rainfall we use is observed rainfall, the terms *forecast* and *simulation* are interchangeable.” Lines 198-203

Intro

Comment #1.3. As far as I see, the Kavetski et al. 2003 reference does not discuss forecasting and thus also not updating procedures but parameter estimation. Please give here references for papers that actually use streamflow correction / updating in a forecast setting.

Response: Thanks for the suggestion. We now use Morawietz et al., 2011 in the revision as a reference for updating procedures used in the context of streamflow forecasting.

Comment #1.4. In the general discussion of streamflow prediction errors, you might want to add the recent reference by Pianosi, F., and Raso, L.: Dynamic modeling of predictive uncertainty by regression on absolute errors, Water Resources Research, 48, W03516, 10.1029/2011wr010603, 2012.

Response: Thanks. We have added Pianosi and Raso (2012) as a reference for heteroscedastic prediction errors.

Method, section 2 and 5

Comment #1.5. Eq. 2 as well as following eqs. does not show the involved parameters

Response: We have carefully revised all equations and notations. The corresponding new equation comes with the definitions of the parameters involved.

Comment #1.6. Eq. 4: which part of the equation does result in the “median value”? should be corrected;

Response: Thanks for reading the manuscript so closely. We provide the proof in Appendix A to show why the updated streamflow is the median.

Comment #1.7. Reference for max. likelihood formulations in eq. 6 , 7?

Response: we add Li et al. (2013) as a reference for the likelihood formulations.

Comment #1.8. In general, I think the superscript notation is not nice to read, why not use two different variable names and subscripts for the parameters?

Response: Thanks for the suggestion. We have carefully and thoroughly revised the notations to increase the readability. For example, we don't use complex superscripts any more.

Comment #1.9. I am not convinced by the current structure with section 5 presenting the new approach; instead of having an “idea-flow” paper structure (method - result 1 - new method - result 2), I would introduce the new method in section 3.

Response: This suggestion is really valuable to improve the presentation of this manuscript. We have followed the suggestion to change the structure of the manuscript.

Comment #1.10. Eq. 8: the same variable name is used for something new, to avoid, what is QM?

Response: We have updated the notations completely and avoided the duplication of variable names.

Comment #1.11. P. 6045 last line: word missing

Response: In the revision, the estimation of all three models is described in Section 2.2. We don't need this sentence any more.

Comment #1.12. P. 6046 first paragraph: would be useful before eq. 8

Response: We have re-worded the motivation/idea behind the RAR-Norm model and placed this paragraph before the definition as suggested by you.

Comment #1.13. Likelihood formulation of the new approach?

Response: We have added the likelihood formulation.

Case study

Comment #1.14. The GR4J model: do any specificities of the model influence the obtained results? (to be mentioned in results section?)

Response: There may be. GR4J may be more prone to fluctuations in base hydrological model performance than other models, as pointed out by our other reviewer. We have added the following discussion of this matter:

“We note that the poor performance of the hydrological model may be specific to the GR4J model, and many not occur in other hydrological models. Evin et al. (2014) estimated hydrological model and error model parameters jointly using GR4J and another hydrological model, HBV, for the three US catchments tested here. While they did not assess the performance of the base hydrological models, they found that HBV tended to perform more robustly when combined with different error models. It is possible that we may have achieved more stable base model performance had we used HBV or another hydrological model. We note, however, that our conclusions can probably be generalised to other hydrological models that do not offer robust base model performance under joint parameter estimation (e.g. GR4J). Because the RAR-Norm model essentially limits the range of updating that can be applied through the AR-Norm model, it will tend to rely more heavily on the base hydrological model, and therefore will tend to favour parameter sets that encourage good stand-alone performance of the base model. For those hydrological models that already produce robust base model performance under joint parameter estimation (perhaps HBV), RAR-Norm is unlikely to undermine this performance for the same reasons. We see some evidence of this in our experiments with GR4J: when the performance of the base hydrological model is already strong relative to the updated forecasts for the AR-Norm and AR-Raw models (e.g. the Tarwin, Mitta Mitta, or Guadalupe catchments), the RAR-Norm model base hydrological model also performs strongly.” (Lines 418-437)

Comment #1.15. P. 6041, line 21: “we then predict streamflow”: not clear here whether in prediction or in forecast mode

Response: We agree – see responses to comments 1.2 and 1.16. This sentence now reads:

“We then generate streamflow forecasts in that year (1999) with model parameters estimated from the remaining data.”

Comment #1.16. The use of “simulation” and “prediction” is confusing; I recommend using the term “forecast” for simulations with forecasted rainfall and the term “simulation” in the other case

Response: Thanks for this – we agree this makes things clearer. We have now changed the terms we use, as described in response to comment #1.2.

Results

Comment #1.17. I suggest a new results section presenting all the results

Response: Thank you for your valuable recommendation. We have followed your suggestion and made all results into two sections: Section 4). All methods are now described before results are documents.

Comment #1.18. In any case, the results of the new method require a separate section (now part of section 5 presenting also the method)

Response: We have followed your suggestion and presented all results in Section 4 (See Comment #1.17)

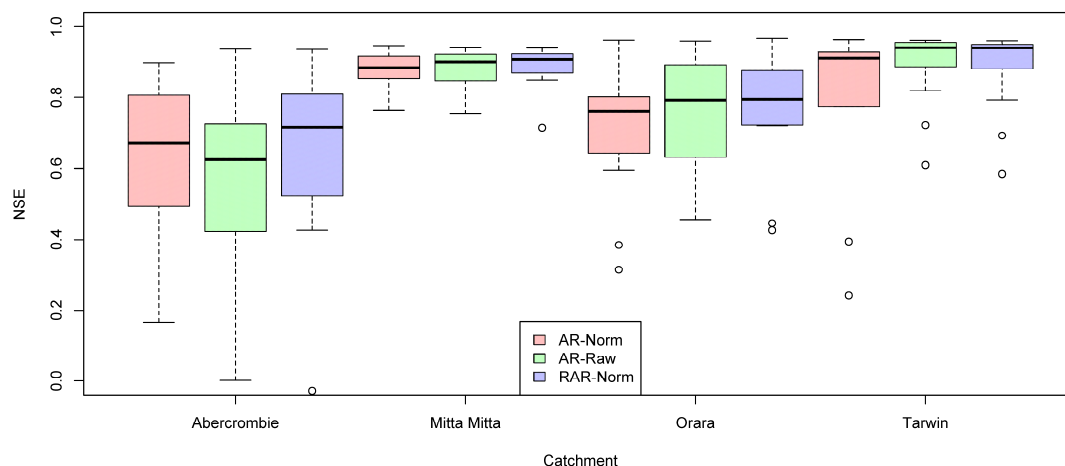
Comment #1.19. P. 6047, line 13: “notable better performance”: as far as I see, only 2 out of 4 cases show a slightly better performance; from fig. 7, the improvement of RAR over AR-raw is not evident

Response: We have made several additions to demonstrate the performance of our RAR-Norm model:

- 1) We now assess our model on three additional US catchments, and these confirm our results
- 2) We use a greater range of metrics and analyses, including assessing the performance of AR models on all instances where streams are i) rising and ii) receding. This demonstrates the general tendency of AR-Norm models to perform least well when flows are rising, as well as the general tendency of the AR-Raw model to perform least well when streamflows are receding. The RAR-Norm model tends to combine the best aspects of the AR-Norm and AR-Raw models.
- 3) We have chosen a different example that more clearly shows the problem of over-correction of receding flows by the AR-Raw model.
- 4) We have changed the structure of the paper, following your suggestions, to group all results together. In doing this, we are better able to present that the RAR-Norm model outperforms both other models (see also response to Comment 1.21)

Comment #1.20. The NSE results are aggregated, how do they look like for individual cross-validation experiments? Are the NSE samples really significantly better with RAR (different distributions with higher mean) or is this pure chance?

Response: We show box-plots of the NSE values from each cross-validation period for the four Australian catchments below (note that for display purposes we have limited the vertical axis to [0, 1] and this means some of the outliers are not displayed). These support our contention that the RAR-Norm model generally leads to better performing forecasts in two ways: 1) the means of the box plots are similar or higher than the next best performing model and 2) The distributions NSE scores of RAR-norm tend to be as narrow or narrower than the next best performing model, indicating that the performance of RAR-norm model is more robust under cross-validation.



Box-plots of NSE values for each cross-validation period for Australian catchments. Dark lines, mean values; boxes, interquartile range; whiskers, [0.1, 0.9] intervals; points, outliers.

Discussion

Comment #1.21. I am not convinced that the paper shows that the new method leads to a more robust performance of the base hydrological model. This should be shown in a more convincing way by presenting some more detailed results of all the simulations.

Response: We believe the additional metrics, analyses and catchments we have added (see response to comment 1.19) add weight to our conclusion that the RAR-Norm model is an improvement over the conventional AR-Raw and AR-Norm models we test. We summarise this in the conclusion as follows:

“The RAR-Norm model is a modification of the AR-Norm: in most instances it operates as the AR-Norm model, but in instances of possible over-correction it relies on the error in untransformed streamflows at the previous time step. That is, RAR-Norm is essentially a more conservative error model than AR-Norm: in situations where streamflows change rapidly, it opts to update with whichever error (transformed or untransformed) is smaller. This forces greater reliance on the base hydrological model to simulate streamflows accurately, leading to more robust performance in the base hydrological model. The RAR-Norm model clearly outperforms the AR-Norm model in both the updated and base model forecasts, as well as ameliorating the problem of over-correcting rising streamflows. The RAR-Norm model’s advantage over the AR-Raw model is less clear: both the base hydrological model and the updated forecasts produced by the AR-Raw model perform similarly to (or sometimes slightly better than) the RAR-Norm model. However, the RAR-Norm model clearly addresses the problem of over-correcting receding streamflows that occurs in the AR-Raw model. As we show, this type of over-correction can seriously distort event hydrographs, and cause forecasts of near zero flows reasonably substantial flows are observed. While these instances are not very common, the failure in the forecast is a serious one. As we note earlier, the over-correction of receding flows is likely to be exacerbated when producing forecasts at lead times of more than one time step. Accordingly, we contend that the RAR-Norm model is preferable to both AR-Norm and AR-Raw models for streamflow forecasting applications.” (Lines 458-479)

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

29

30 **1. Introduction**

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

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

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

60 There is no agreement on whether it is better to apply an AR model to normalised or
61 raw errors. Recent work by Evin et al. (2013) found that an AR model applied to raw

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

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

82 **2. Autoregressive error models**

83 **2.1 Formulations**

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

$$89 \quad f(x) = b^{-1} \log\{\sinh(a + bx)\} \quad (1)$$

90 is applied to stabilise variance and normalise data.

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

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

$$94 \quad Z_t = \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t, \quad (2)$$

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

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

$$97 \quad Z_t = f\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\} + \varepsilon_t. \quad (3)$$

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

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

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

$$112 \quad \tilde{Q}_t^* = f^{-1}\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\}, \quad (4)$$

113 and

$$114 \quad \tilde{Q}_t^* = \tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1}), \quad (5)$$

115 where $f^{-1}(x)$ is the inverse of log-sinh transformation (or back-transformation). The
 116 magnitude of the error update by the AR-Raw model, $\tilde{Q}_t^* - \tilde{Q}_t$, is dependent only on
 117 the difference between Q_{t-1} and \tilde{Q}_{t-1} . In contrast, the magnitude of the error update

118 by the AR-Norm model is dependent not only on the difference between Q_{t-1} and
 119 \tilde{Q}_{t-1} , but also on \tilde{Q}_t . Put differently, the AR-Norm model uses errors calculated in
 120 the transformed domain, and this means that the error in the original domain can be
 121 amplified (or reduced) by the back-transformation (Equation (4)). The AR-Raw model
 122 uses errors calculated in the original domain and no back-transformation is used in
 123 calculating \tilde{Q}_t^* (Equation (5)), meaning that the error in the original domain cannot be
 124 amplified (or reduced). In Appendix B, we show that the AR-Norm model gives
 125 greater error updates for larger values of \tilde{Q}_t .

126 We will demonstrate in Section 4 that the AR-Norm and AR-Raw models can
 127 sometimes cause over-correction of forecasts. Motivated to overcome the potential for
 128 over-correction, we introduce a modification of the AR-Norm model, called the
 129 restricted AR-Norm model (referred to as *RAR-Norm*). A condition
 130 $|\tilde{Q}_t^* - \tilde{Q}_t| \leq |Q_{t-1} - \tilde{Q}_{t-1}|$ is used to limit the correction amount to not exceeding the error
 131 in the last time step in absolute value. The updated streamflow is given by

$$132 \quad \tilde{Q}_t^* = \begin{cases} f^{-1} \{ \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) \} & \text{if } D_t \leq |Q_{t-1} - \tilde{Q}_{t-1}| \\ \tilde{Q}_t + (Q_{t-1} - \tilde{Q}_{t-1}) & \text{otherwise.} \end{cases} \quad (6)$$

133 where

$$134 \quad D_t = \left| f^{-1} \{ \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) \} - \tilde{Q}_t \right|. \quad (7)$$

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

$$136 \quad Z_t = \begin{cases} \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t & \text{if } D_t \leq |Q_{t-1} - \tilde{Q}_{t-1}| \\ f(\tilde{Q}_t + Q_{t-1} - \tilde{Q}_{t-1}) + \varepsilon_t & \text{otherwise.} \end{cases} \quad (8)$$

137 2.2 Estimation

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

142 (a) for AR-Norm

$$143 \quad L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_t J_{Z_t \rightarrow Q_t} \phi(Z_t | \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}), \sigma^2), \quad (9)$$

144 (b) for AR-Raw

$$145 \quad L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_t J_{Z_t \rightarrow Q_t} \phi(Z_t | f\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\}, \sigma^2) \quad ,$$

$$146 \quad (10)$$

147 (c) for RAR-Norm

$$148 \quad L(\theta_E, \theta_H) = \prod_t P(Q_t | \tilde{Q}_t, \tilde{Q}_{t-1}; \theta_E, \theta_H) = \prod_{t: D_t \leq |\tilde{Q}_{t-1} - \tilde{Q}_t|} J_{Z_t \rightarrow Q_t} \phi(Z_t | \tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}), \sigma^2)$$

$$149 \quad + \prod_{t: D_t > |\tilde{Q}_{t-1} - \tilde{Q}_t|} J_{Z_t \rightarrow Q_t} \phi(Z_t | f\{\tilde{Q}_t + \rho(Q_{t-1} - \tilde{Q}_{t-1})\}, \sigma^2), \quad (11)$$

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

157 3. Data

158 We use daily data from four Australian catchments and three catchments from the
 159 United States (US; Figure 1, Table 1). Australian streamflow data are taken from the
 160 Catchment Water Yield Estimation Tool (CWYET) dataset (Vaze et al., 2011).
 161 Australian rainfall and potential evaporation data are derived from the Australian
 162 Water Availability Project (AWAP) dataset (Jones et al., 2009). All data for the US
 163 catchments come from the Model Intercomparison Experiment (MOPEX) dataset
 164 (Duan et al., 2006). The selected US catchments are amongst the 12 catchments used
 165 by Evin et al. (2014) to compare joint and postprocessor approaches to estimate
 166 hydrological uncertainty, and allows us to compare results with that study (the other
 167 catchments used by Evin et al. (2014) are influenced by snowmelt, which is not
 168 considered in the hydrological model used in this study). The Abercrombie River and

169 the Guadalupe River intermittently experience periods of very low (to zero) flow,
170 while the other rivers flow perennially (Table 1). Such dry catchments are challenging
171 for hydrological simulations and error modelling. All catchments have high-quality
172 streamflow records with very few missing data.

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

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

195 To demonstrate the problems of over-correction of errors in updating and poor stand-
196 alone performance of the base hydrological model, we consider only streamflow
197 forecasts for one time step ahead. We will consider longer lead times in future work.
198 Forecasts are generated using observed rainfall (i.e., a ‘perfect’ rainfall forecast) as
199 input. In streamflow forecasting, forecasts may be generated from rainfall information
200 that comes from a different source (e.g., a numerical weather prediction model). Our
201 study is aimed at streamflow forecasting applications, so we preserve the distinction

202 between observed and forecast forcings by referring to streamflows modelled with
 203 observed rainfall as *simulations* and those modelled with forecast rainfall as *forecasts*.
 204 As the forecast rainfall we use is observed rainfall, the terms *forecast* and *simulation*
 205 are interchangeable.

206 4. Results

207 4.1 Over-correction of forecasts as the hydrograph rises

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

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

229 Figure 3 shows instances of possible over-correction by the AR-Norm model,
 230 identified by the condition $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$. Figure 3 shows that about 10-25% of the
 231 AR-Norm updated forecasts have an error update that is larger than the forecast error
 232 on the previous day and therefore are susceptible to over-correction. The frequency of
 233 these instances varies somewhat from catchment to catchment. The RAR-Norm model

234 identifies 10-30% of the forecasts as possible instances of problematic updating, and
235 the AR-Norm model identifies a similar number of instances (slightly fewer – they are
236 not identical because the parameters for each model are inferred independently).

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

246 **4.2 Over-correction of forecasts as the hydrograph recedes**

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

263 Figure 6 provides more examples of the over-correction caused by the AR-Raw model
264 from a longer time-series plot for the Abercrombie catchment. There are three clear
265 instances of over-correction, all occurring on the time step immediately after large

266 peaks in observed streamflows. The RAR-Norm works better than the AR-Raw model
267 to avoid the three instances of over-correction for the Abercrombie catchment.
268 Overall, the RAR-Norm model takes a conservative position when streamflow
269 changes rapidly, either rising or falling. When streamflow changes rapidly, it is
270 difficult to anticipate the magnitude of forecast error. Accordingly the conventional
271 AR models are prone to over-correction in such instances.

272 **4.3 Poor stand-alone performance of the base hydrological model**

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

287 Figure 7 presents the Nash-Sutcliffe efficiency (NSE) (Nash and Sutcliffe, 1970)
288 calculated from the base hydrological model and the error models. When the AR-
289 Norm model is used, the forecasts from the base hydrological model are very poor for
290 the Orara catchment ($NSE < 0$). The scatter plot in Figure 8 shows a serious over-
291 estimation of the streamflow simulation for the Orara. When the AR-Norm model is
292 used, the base hydrological model greatly over-estimates discharge and the AR-Norm
293 model then attempts to correct this systematic over-estimation. This is also shown in
294 Figure 4 where the base hydrological model has a strong tendency to over-estimate
295 streamflows for a range of streamflow magnitudes. The base hydrological model with
296 the AR-Norm model also performs poorly for the Abercrombie catchment (Figure 7).
297 In this case, the base hydrological model tends to under-estimate streamflows (results

298 not shown). For the other three catchments, however, the base hydrological model
299 with the AR-Norm model performs reasonably well.

300 In general, the AR-Raw base hydrological model performs as well or better than the
301 AR-Norm base hydrological model. The AR-Raw base hydrological model is notably
302 better than the AR-Norm base hydrological model in the Abercrombie and Orara
303 catchments (Figure 7). This suggests that more robust performance can be expected of
304 base hydrological models when AR models are applied to raw errors.

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

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

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

330 4.4 Further analyses

331 We further evaluate the NSE of the three different error models calibrated when
332 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
333 receding streamflows (constituting 70-85% of streamflows), the AR-Raw model leads
334 to the overall worst forecast accuracy because of the over-correction explained in
335 Section 4.1. This is especially evident for the Abercrombie catchment (and, to a lesser
336 degree, the Guadalupe catchment). The RAR-Norm model significantly outperforms
337 the other two models for the Abercrombie catchment and shares similar forecast
338 accuracy to the (strongly performing) AR-Norm model for the other catchments.
339 When streamflows are rising (which also includes streamflow peaks), the AR-Norm
340 model can cause over-correction and leads to the least accurate forecasts (in terms of
341 NSE), and the RAR-Norm model behaves similarly to the AR-Raw model, which
342 consistently provides the most accurate forecasts. (The only exception is the
343 Guadalupe River, where the AR-Raw model clearly outperforms the RAR-Norm
344 model when streamflows are rising. This is somewhat compensated for by the
345 markedly better performance the RAR-Norm model offers over the AR-Raw model
346 when streamflows are receding for this catchment, leading to better forecasts overall
347 (Figure 7).) We conclude that the AR-Norm model generally tends to perform least
348 well when streamflows recede, and that the AR-Raw model tends to perform least
349 well when streamflows rise. We also conclude that the RAR-Norm model tends to
350 combine the best elements of the AR-Norm and AR-Raw models, leading to the best
351 overall performance.

352 We have shown that over-corrections can lead to inaccurate deterministic forecasts,
353 and we now discuss the consequences for the probabilistic predictions given by each
354 of the error models. We assess probabilistic forecast skill with skill scores derived
355 from two probabilistic verification measures: the Continuous Rank Probability Score
356 (CRPS) and the Root Mean Square Error in Probability (RMSEP) (denoted by
357 CRPS_SS and RMSEP_SS, respectively) (Wang and Robertson, 2011). Both skill
358 scores are calculated with respect to a reference forecast. The reference forecast is
359 generated by resampling historical streamflows: for a forecast issued for a given
360 month/year (e.g. February 1999), we randomly draw a sample of 1000 daily
361 streamflows that occurred in that month (e.g. February) from other years with
362 replacement (e.g. years other than 1999). Table 3 compares these two skill scores

363 calculated for the all catchments. The RAR-Norm model performs best across the
364 range of skill scores and catchments, attaining the highest CRPS_SS in 4 of the 7
365 catchments and the highest RMSEP_SS in 4 of 7 catchments. Even where RAR-Norm
366 was not the best performed model, it performs very similarly to the best performing
367 model in all cases. Interestingly, the AR-Raw model tends to outperform the AR-
368 Norm model in CRPS_SS while the reverse is true for RMSEP_SS. The CRPS tests
369 how appropriate the spread of uncertainty is for each probabilistic forecast, while
370 RMSEP puts little weight on this. The results suggest that while the median forecasts
371 of AR-Norm tends to be slightly more accurate than those of the AR-Raw model, the
372 forecast uncertainty is represented slightly better by the AR-Raw model.

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

388 **5. Discussion and conclusions**

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

394 We demonstrate three adverse effects of AR error updating procedures on seven
395 catchments. The first adverse effect is possible over-correction on the rising limb of
396 the hydrograph. The AR-Norm model can exhibit the tendency to over-correct the
397 peaks or on the rise of a hydrograph, because error updating can be (overly) amplified
398 by the back-transformation. The second adverse effect is the tendency to over-correct
399 receding hydrographs. This tendency is most prevalent in the AR-Raw model, which
400 can fail to recognise that a large error update may not be appropriate for small
401 streamflows.

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

421 We note that the poor performance of the hydrological model may be specific to the
422 GR4J model, and many not occur in other hydrological models. Evin et al. (2014)
423 estimated hydrological model and error model parameters jointly using GR4J and
424 another hydrological model, HBV, for the three US catchments tested here. While
425 they did not assess the performance of the base hydrological models, they found that
426 HBV tended to perform more robustly when combined with different error models. It

427 is possible that we may have achieved more stable base model performance had we
428 used HBV or another hydrological model. We note, however, that our conclusions can
429 probably be generalised to other hydrological models that do not offer robust base
430 model performance under joint parameter estimation (e.g. GR4J). Because the RAR-
431 Norm model essentially limits the range of updating that can be applied through the
432 AR-Norm model, it will tend to rely more heavily on the base hydrological model,
433 and therefore will tend to favour parameter sets that encourage good stand-alone
434 performance of the base model. For those hydrological models that already produce
435 robust base model performance under joint parameter estimation (perhaps HBV),
436 RAR-Norm is unlikely to undermine this performance for the same reasons. We see
437 some evidence of this in our experiments with GR4J: when the performance of the
438 base hydrological model is already strong relative to the updated forecasts for the AR-
439 Norm and AR-Raw models (e.g. the Tarwin, Mitta Mitta, or Guadalupe catchments),
440 the RAR-Norm model base hydrological model also performs strongly.

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

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

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

483 **Appendix A**

484 For simplicity we only show the case of the AR-Norm model and analogues
 485 arguments can be used to prove the cases of the AR-Raw and RAR-Norm models.
 486 The streamflow ensemble forecast Q_t given by the AR-Norm model defined by (1)
 487 can be written as

$$488 \quad Q_t = \max \left[f^{-1} \left\{ \tilde{Z}_t + \rho (Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t \right\}, 0 \right]. \quad (\text{A1})$$

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

$$492 \quad P(Q_t \leq \tilde{Q}_t^*) = P\left(\max\left[f^{-1}\left\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t\right\}, 0\right] \leq \tilde{Q}_t^*\right) \quad (\text{A2})$$

$$= P\left(f^{-1}\left\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t\right\} \leq \tilde{Q}_t^* \text{ and } 0 \leq \tilde{Q}_t^*\right)$$

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

$$494 \quad P(Q_t \leq \tilde{Q}_t^*) = P\left(f^{-1}\left\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1}) + \varepsilon_t\right\} \leq f^{-1}\left\{\tilde{Z}_t + \rho(Z_{t-1} - \tilde{Z}_{t-1})\right\}\right) . \quad (\text{A3})$$

$$= P(\varepsilon_t \leq 0) = 0.5$$

495 Appendix B

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

498 Firstly, we will show that the first derivate of the log-sinh transform f defined by (3)
 499 is positive and the second derivate is negative (i.e. $f'(x) > 0$ and $f''(x) < 0$) for any
 500 $b > 0$ and any x . Following some simple manipulation, we have

$$501 \quad f'(x) = \frac{\cosh(a+bx)}{\sinh(a+bx)} > 0 \quad \text{and} \quad f''(x) = \frac{-b}{\sinh^2(a+bx)} < 0 \quad (\text{B1})$$

502 Using the differentiation of inverse functions, we find the first and second derivates of
 503 the inverse transform f^{-1}

$$504 \quad \left[f^{-1}\right]'(x) = \frac{1}{f'\{f^{-1}(x)\}} > 0 \quad \text{and} \quad \left[f^{-1}\right]''(x) = \frac{-f''\{f^{-1}(x)\}}{\left[f'\{f^{-1}(x)\}\right]^3} > 0, \quad (\text{B2})$$

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

506 Next, we will derive the difference of magnitudes of the error update between low and
 507 high forecast streamflows. For the sake of notation simplicity, we rewrite $q = \tilde{Z}_t$ and
 508 $u = \rho(Z_{t-1} - \tilde{Z}_{t-1})$ and assume that $u > 0$. Using Equation (4), the updated streamflow
 509 can be written as $\tilde{Q}_t^* = f^{-1}(q+u)$. The magnitude of the error update can be written as

$$510 \quad |\tilde{Q}_t^* - \tilde{Q}_t| = |f^{-1}(q+u) - f^{-1}(q)| = \begin{cases} \int_0^u [f^{-1}]'(x+q) dx & \text{if } u > 0 \\ \int_u^0 [f^{-1}]'(x+q) dx & \text{otherwise.} \end{cases} \quad (\text{B3})$$

511 Suppose that we have two forecast streamflows $\tilde{Q}_{t,1} \leq \tilde{Q}_{t,2}$ and denote the normalised
 512 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
 513 $\tilde{Q}_{t,2}^*$. Because f is an increasing function, we have $q_1 \leq q_2$. The difference in the
 514 magnitude of the error update between $\tilde{Q}_{t,1}$ and $\tilde{Q}_{t,2}$ can be derived as

$$515 \quad |\tilde{Q}_{t,1} - \tilde{Q}_{t,1}^*| - |\tilde{Q}_{t,2} - \tilde{Q}_{t,2}^*| = \begin{cases} \int_0^u \left\{ [f^{-1}]'(x+q_1) - [f^{-1}]'(x+q_2) \right\} dx & \text{if } u > 0 \\ \int_u^0 \left\{ [f^{-1}]'(x+q_1) - [f^{-1}]'(x+q_2) \right\} dx & \text{otherwise.} \end{cases} \quad (\text{B4})$$

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

$$518 \quad |\tilde{Q}_{t,1} - \tilde{Q}_{t,1}^*| \leq |\tilde{Q}_{t,2} - \tilde{Q}_{t,2}^*|. \quad (\text{B5})$$

519 Therefore, the error update at larger forecast streamflows is always larger than error
 520 update at lower forecast streamflows.

521 Acknowledgments

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

528 **Table of Tables**

529 Table 1: Catchment characteristics.

530 Table 2: Comparison of the NSE calculated at (a) the receding limb and (b) the rising
531 limb of the hydrograph for three different error models.532 Table 3: Comparison of the skill scores based on CRPS and RMSEP (denoted by
533 CRPS_SS and RMSEP_SS) for three different error models.534 **Table of Figures**

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

536 Figure 2: An example of over-correction caused by the AR-Norm model in the Mitta
537 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,
538 without error updating). Solid lines: forecasts with error updating.539 Figure 3: The fraction of instances where $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$ (i.e., instances where over-
540 correction may occur in the AR-Norm model and where error updating is restricted in
541 the RAR-Norm model) for the AR-Norm and RAR-Norm models for Australian
542 catchments.543 Figure 4: Forecast streamflows for the Orara catchment for an example 1-year period.
544 Top panel shows streamflows forecast with AR-Norm model, bottom panel shows
545 streamflows forecast with the RAR-Norm model. Dashed lines: forecasts from the
546 base hydrological model (i.e., without error updating). Solid lines: forecasts with error
547 updating. Tick marks in the x-axis denote the instance of updating where
548 $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$.549 Figure 5: An example of over-correction caused by the AR-Raw model in the Mitta
550 Mitta catchment. Dashed lines: forecasts from the base hydrological model (i.e.,
551 without error updating). Solid lines: forecasts with error updating.552 Figure 6: Forecast streamflows for the Abercrombie catchment for the period between
553 01/08/1997 and 15/09/1997. Top panel shows streamflows forecast with AR-Raw
554 model, bottom panel shows streamflows forecast with the RAR-Norm model. Dashed
555 lines: forecasts from the base hydrological model (i.e., without error updating). Solid
556 lines: forecasts with error updating. Gray shading denotes instances of over-correction
557 caused by the AR-Raw model.

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

561 Figure 8: Comparison of the observed streamflows (Q_t) and forecast streamflows
562 (\tilde{Q}_t), as forecast: 1) with the base hydrological model (circles), and 2) with the base
563 hydrological model and error updating models (dots) for the Orara catchment.

564 Figure 9: PIT-uniform probability plots. Curves on the diagonal indicate perfectly
565 reliable forecasts.

566

567

568 **References**

569 Bates, B. C., and Campbell, E. P.: A Markov chain Monte Carlo scheme for parameter
570 estimation and inference in conceptual rainfall-runoff modeling, *Water Resour Res*,
571 37, 937-947, 10.1029/2000wr900363, 2001.

572 Bennett, J. C., Robertson, D. E., Shrestha, D. L., Wang, Q. J., Enever, D.,
573 Hapuarachchi, P., and Tuteja, N. K.: A System for Continuous Hydrological
574 Ensemble Forecasting (SCHEF) to lead times of 9 days, *J Hydrol*,
575 10.1016/j.jhydrol.2014.08.010, 2014.

576 Del Giudice, D., Honti, M., Scheidegger, A., Albert, C., Reichert, P., and
577 Rieckermann, J.: Improving uncertainty estimation in urban hydrological modeling by
578 statistically describing bias, *Hydrol. Earth Syst. Sci.* , 17, 4209-4225, 10.5194/hess-
579 17-4209-2013, 2013.

580 Duan, Q., Schaake, J., Andreassian, V., Franks, S., Goteti, G., Gupta, H. V., Gusev, Y.
581 M., Habets, F., Hall, A., Hay, L., Hogue, T., Huang, M., Leavesley, G., Liang, X.,
582 Nasonova, O. N., Noilhan, J., Oudin, L., Sorooshian, S., Wagener, T., and Wood, E.
583 F.: Model Parameter Estimation Experiment (MOPEX): An overview of science
584 strategy and major results from the second and third workshops, *J Hydrol*, 320, 3-17,
585 10.1016/j.jhydrol.2005.07.031, 2006.

586 Duan, Q. Y., Sorooshian, S., and Gupta, V. K.: Optimal Use of the Sce-Ua Global
587 Optimization Method for Calibrating Watershed Models, *J Hydrol*, 158, 265-284,
588 10.1016/0022-1694(94)90057-4, 1994.

589 Engeland, K., Renard, B., Steinsland, I., and Kolberg, S.: Evaluation of statistical
590 models for forecast errors from the HBV model, *J Hydrol*, 384, 142-155,
591 10.1016/j.jhydrol.2010.01.018, 2010.

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

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

- 597 accounting for error autocorrelation and heteroscedasticity, *Water Resour Res*, 50,
598 2350-2375, 10.1002/2013WR014185, 2014.
- 599 Jones, D. A., Wang, W., and Fawcett, R.: High-quality spatial climate data-sets for
600 Australia, *Australian Meteorological and Oceanographic Journal*, 58, 233-248, 2009.
- 601 Kavetski, D., Franks, S. W., and Kuczera, G.: Confronting Input Uncertainty in
602 Environmental Modelling, in: *Calibration of Watershed Models*, edited by: Duan, Q.,
603 Gupta, H. V., Sorooshian, S., Rousseau, A. N., and Turcotte, R., American
604 Geophysical Union, Washington D.C., 49-68, 2003.
- 605 Li, M., Wang, Q. J., and Bennett, J.: Accounting for seasonal dependence in
606 hydrological model errors and prediction uncertainty, *Water Resour Res*, 49, 5913-
607 5929, 10.1002/wrcr.20445, 2013.
- 608 Morawietz, M., Xu, C. Y., and Gottschalk, L.: Reliability of autoregressive error
609 models as post-processors for probabilistic streamflow forecasts, *Adv. Geosci.*, 29,
610 109-118, 10.5194/adgeo-29-109-2011, 2011.
- 611 Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models
612 part I — A discussion of principles, *J Hydrol*, 10, 282-290, 10.1016/0022-
613 1694(70)90255-6, 1970.
- 614 Perrin, C., Michel, C., and Andreassian, V.: Improvement of a parsimonious model
615 for streamflow simulation, *J Hydrol*, 279, 275-289, 10.1016/S0022-1694(03)00225-7,
616 2003.
- 617 Pianosi, F., and Raso, L.: Dynamic modeling of predictive uncertainty by regression
618 on absolute errors, *Water Resour Res*, 48, W03516, 10.1029/2011wr010603, 2012.
- 619 Salamon, P., and Feyen, L.: Disentangling uncertainties in distributed hydrological
620 modeling using multiplicative error models and sequential data assimilation, *Water*
621 *Resour Res*, 46, W12501, 10.1029/2009wr009022, 2010.
- 622 Schaeffli, B., Talamba, D. B., and Musy, A.: Quantifying hydrological modeling errors
623 through a mixture of normal distributions, *J Hydrol*, 332, 303-315,
624 10.1016/j.jhydrol.2006.07.005, 2007.

- 625 Schoups, G., and Vrugt, J. A.: A formal likelihood function for parameter and
626 predictive inference of hydrologic models with correlated, heteroscedastic, and non-
627 Gaussian errors, *Water Resour Res*, 46, W10531, 10.1029/2009wr008933, 2010.
- 628 Seo, D. J., Herr, H. D., and Schaake, J. C.: A statistical post-processor for accounting
629 of hydrologic uncertainty in short-range ensemble streamflow prediction, *Hydrol.*
630 *Earth Syst. Sci. Discuss.*, 3, 1987-2035, 10.5194/hessd-3-1987-2006, 2006.
- 631 Shrestha, D. L., and Solomatine, D. P.: Data - driven approaches for estimating
632 uncertainty in rainfall - runoff modelling, *International Journal of River Basin*
633 *Management*, 6, 109-122, 10.1080/15715124.2008.9635341, 2008.
- 634 Thyer, M., Kuczera, G., and Wang, Q. J.: Quantifying parameter uncertainty in
635 stochastic models using the Box-Cox transformation, *J Hydrol*, 265, 246-257,
636 10.1016/S0022-1694(02)00113-0, 2002.
- 637 Vaze, J., Perraud, J. M., Teng, J., Chiew, F. H. S., Wang, B., and Yang, Z.: Catchment
638 Water Yield Estimation Tools (CWYET), the 34th World Congress of the
639 International Association for Hydro- Environment Research and Engineering: 33rd
640 Hydrology and Water Resources Symposium and 10th Conference on Hydraulics in
641 Water Engineering, Brisbane, 2011.
- 642 Wang, Q. J., and Robertson, D. E.: Multisite probabilistic forecasting of seasonal
643 flows for streams with zero value occurrences, *Water Resour Res*, 47, W02546,
644 10.1029/2010WR009333, 2011.
- 645 Wang, Q. J., Shrestha, D. L., Robertson, D. E., and Pokhrel, P.: A log-sinh
646 transformation for data normalization and variance stabilization, *Water Resour Res*,
647 48, W05514, 10.1029/2011WR010973, 2012.
- 648 World Meteorological Organization: Simulated real-time intercomparison of
649 hydrological models, World Meteorological Organization, Geneva, Switzerland, 1992.
- 650 Xiong, L. H., and O'Connor, K. M.: Comparison of four updating models for real-time
651 river flow forecasting, *Hydrolog Sci J*, 47, 621-639, 10.1080/02626660209492964,
652 2002.

653 Xu, C. Y.: Statistical analysis of parameters and residuals of a conceptual water
654 balance model - Methodology and case study, *Water Resour Manag*, 15, 75-92,
655 10.1023/A:1012559608269, 2001.

656

657 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%

658

659

660 Table 2: Comparison of the NSE calculated at (a) the receding limb and (b) the rising
 661 limb of the hydrograph for three different error models.
 662

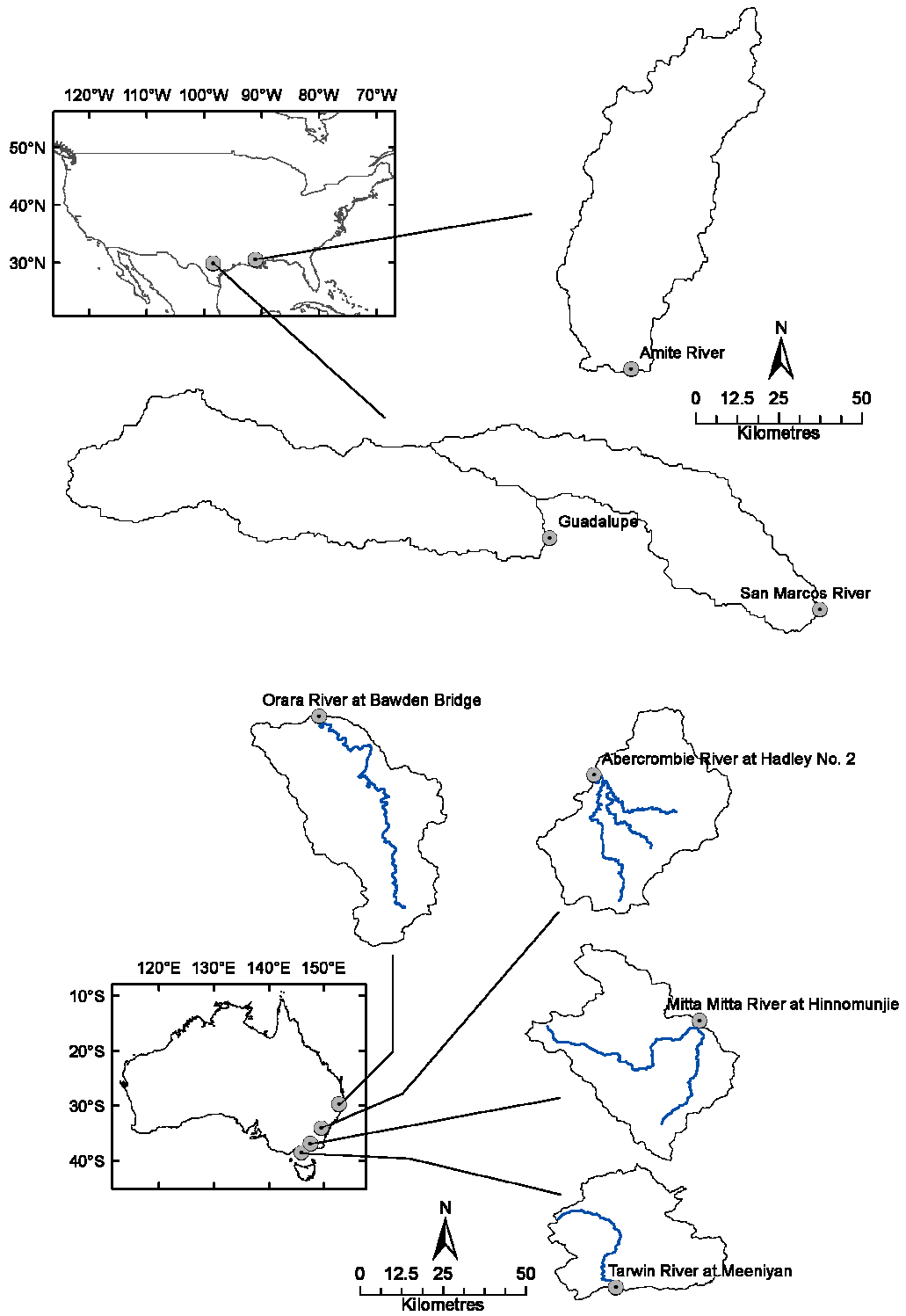
	(a) $\tilde{Q}_t \leq \tilde{Q}_{t-1}$				(b) $\tilde{Q}_t > \tilde{Q}_{t-1}$			
	Proportion of flows	AR- Norm	AR- Raw	RAR- Norm	Proportion of flows	AR- Norm	AR- Raw	RAR- Norm
Abercrombie	82%	0.11	-0.41	0.52	19%	0.58	0.66	0.65
Mitta Mitta	82%	0.95	0.91	0.95	18%	0.81	0.86	0.86
Orara	85%	0.94	0.91	0.95	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

663

664 Table 3: Comparison of the skill scores based on CRPS and RMSEP (denoted by
 665 CRPS_SS and RMSEP_SS) for three different error models.
 666

	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

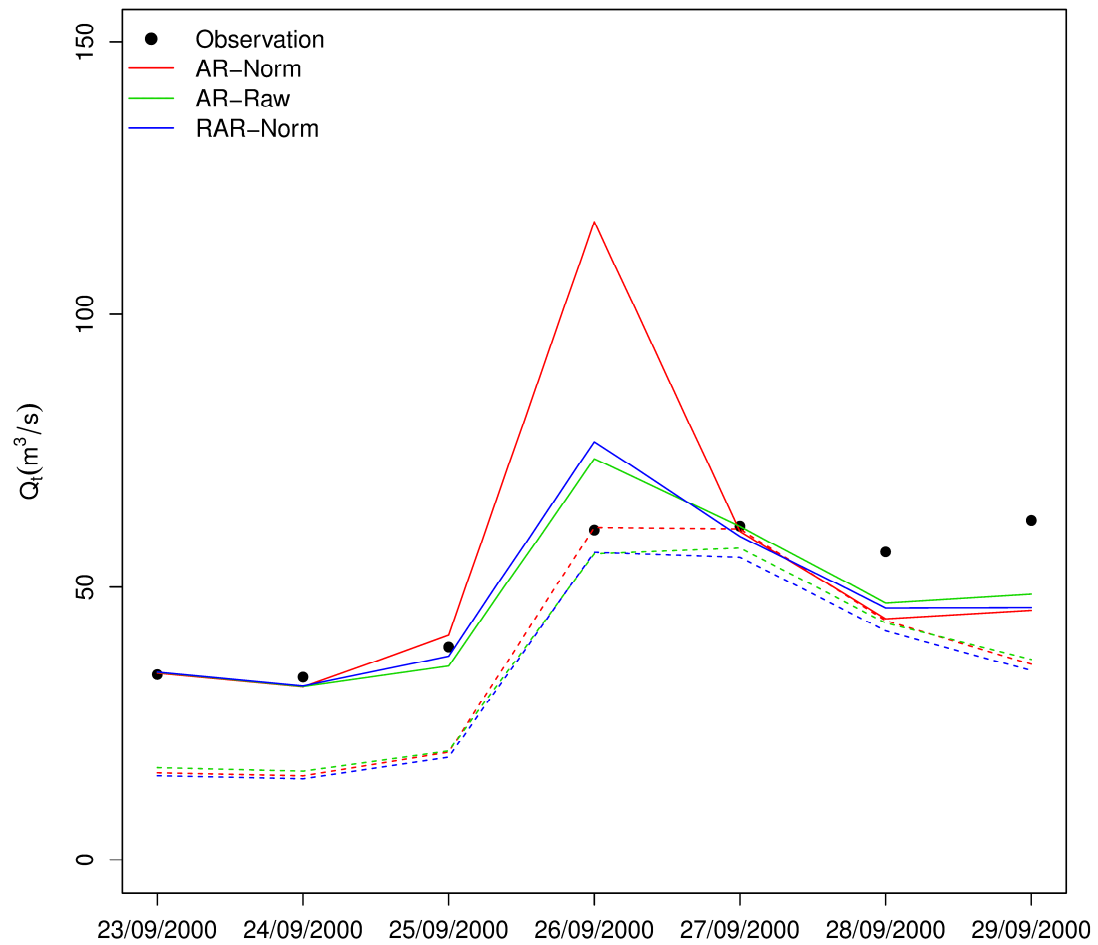
667



668

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

670



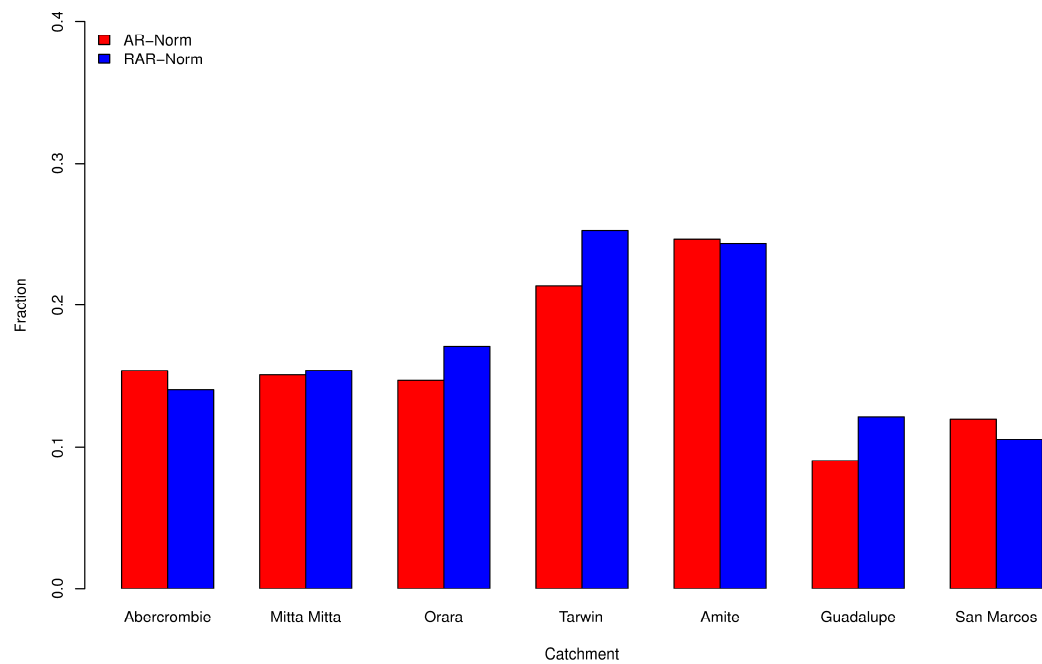
671

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

675

676

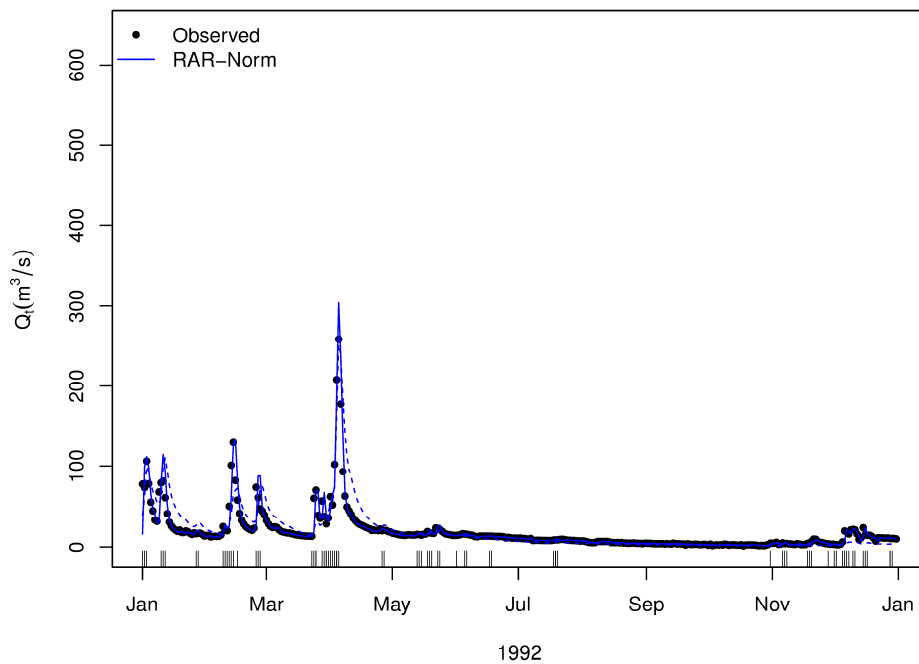
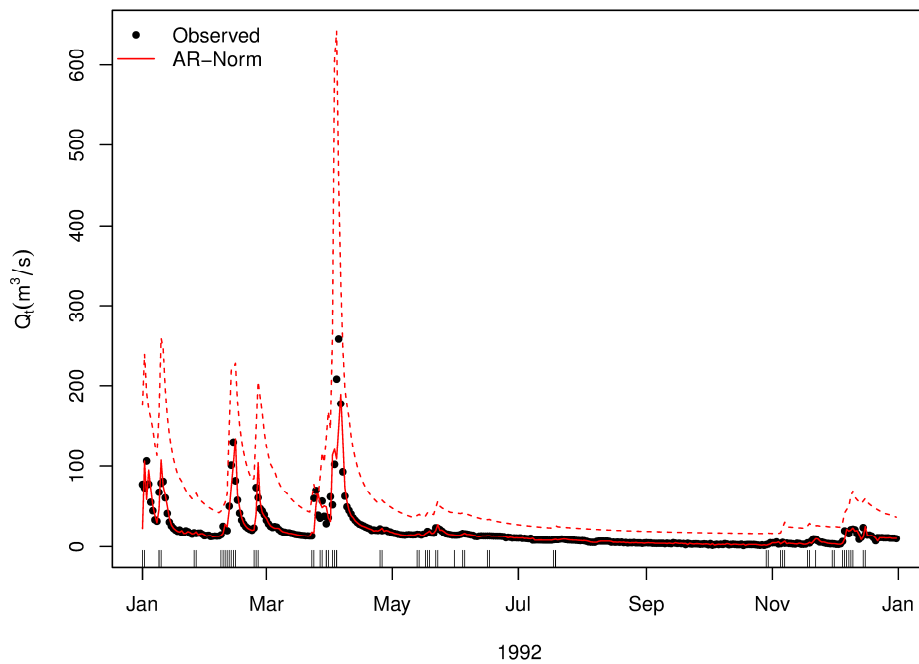
677



678

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

683



684

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

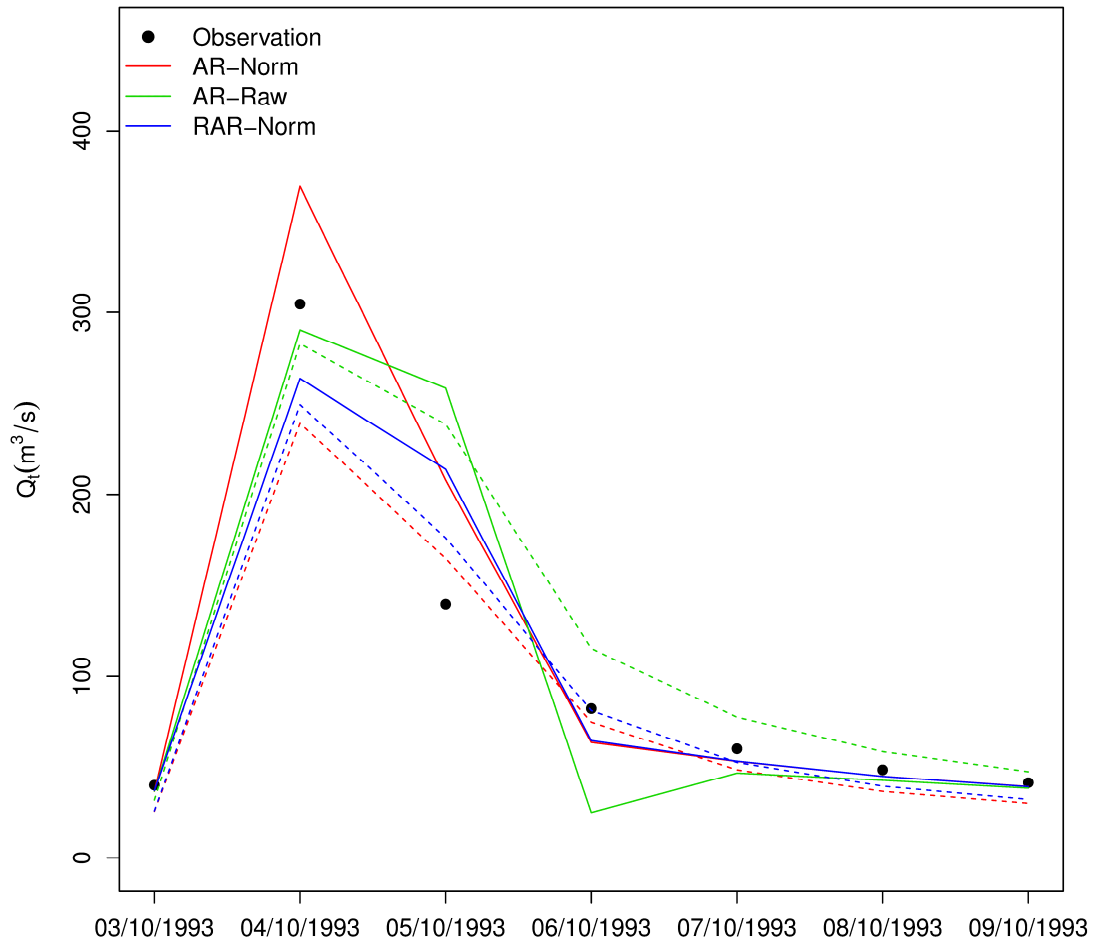
686 Top panel shows streamflows forecast with AR-Norm model, bottom panel shows

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

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

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

690 $D_t > |Q_{t-1} - \tilde{Q}_{t-1}|$.



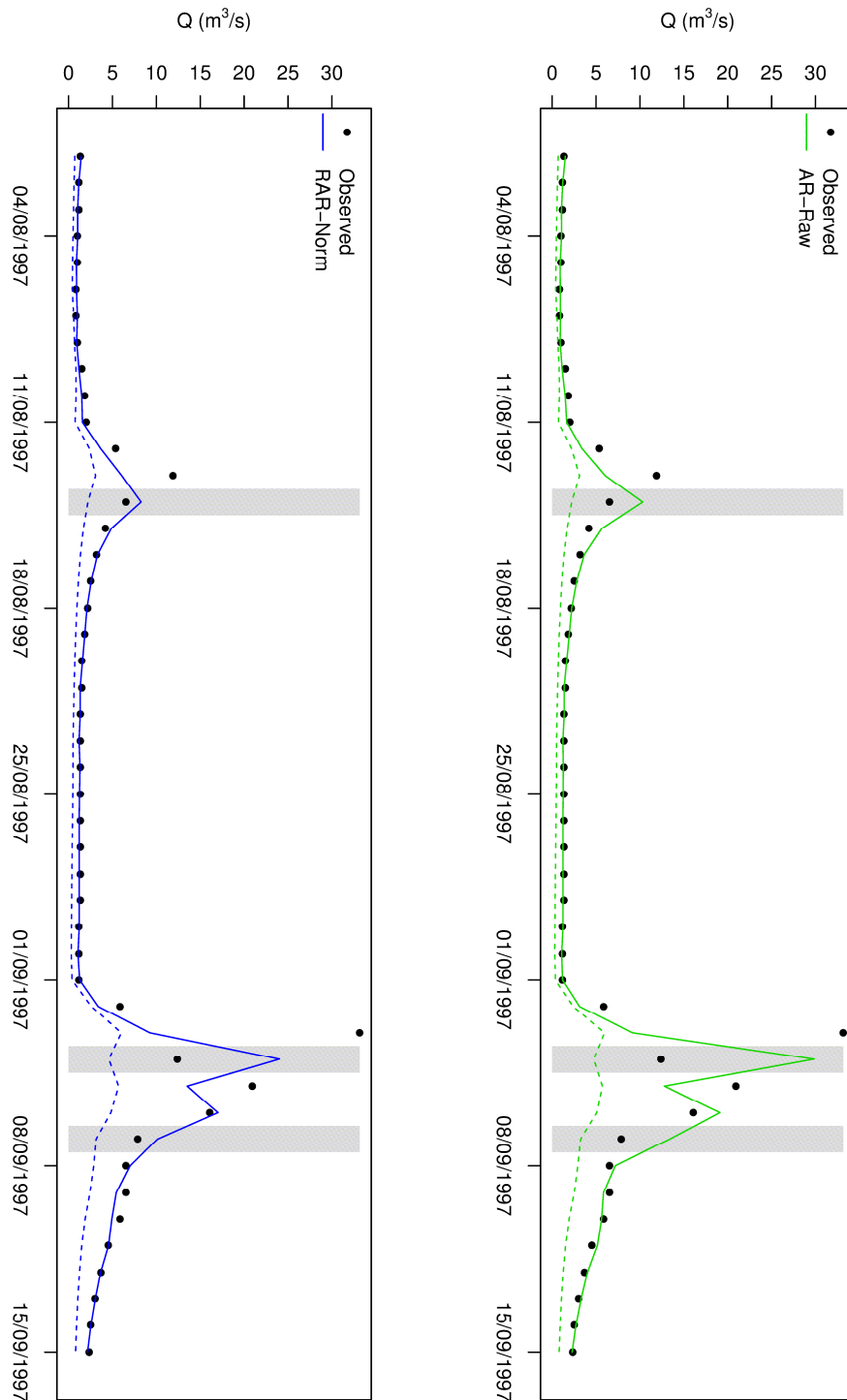
691

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

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

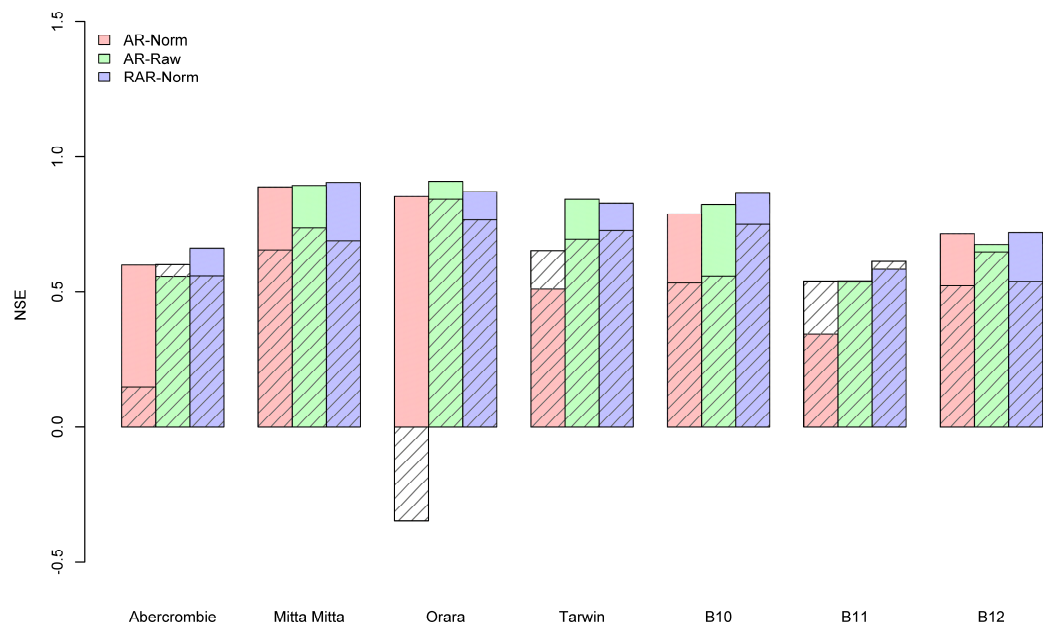
694 without error updating). Solid lines: forecasts with error updating.

695



696

697 Figure 6: Forecast streamflows for the Abercrombie catchment for the period between
 698 01/08/1997 and 15/09/1997. Top panel shows streamflows forecast with AR-Raw
 699 model, bottom panel shows streamflows forecast with the RAR-Norm model. Dashed
 700 lines: forecasts from the base hydrological model (i.e., without error updating). Solid
 701 lines: forecasts with error updating. Gray shading denotes instances of over-correction
 702 caused by the AR-Raw model.

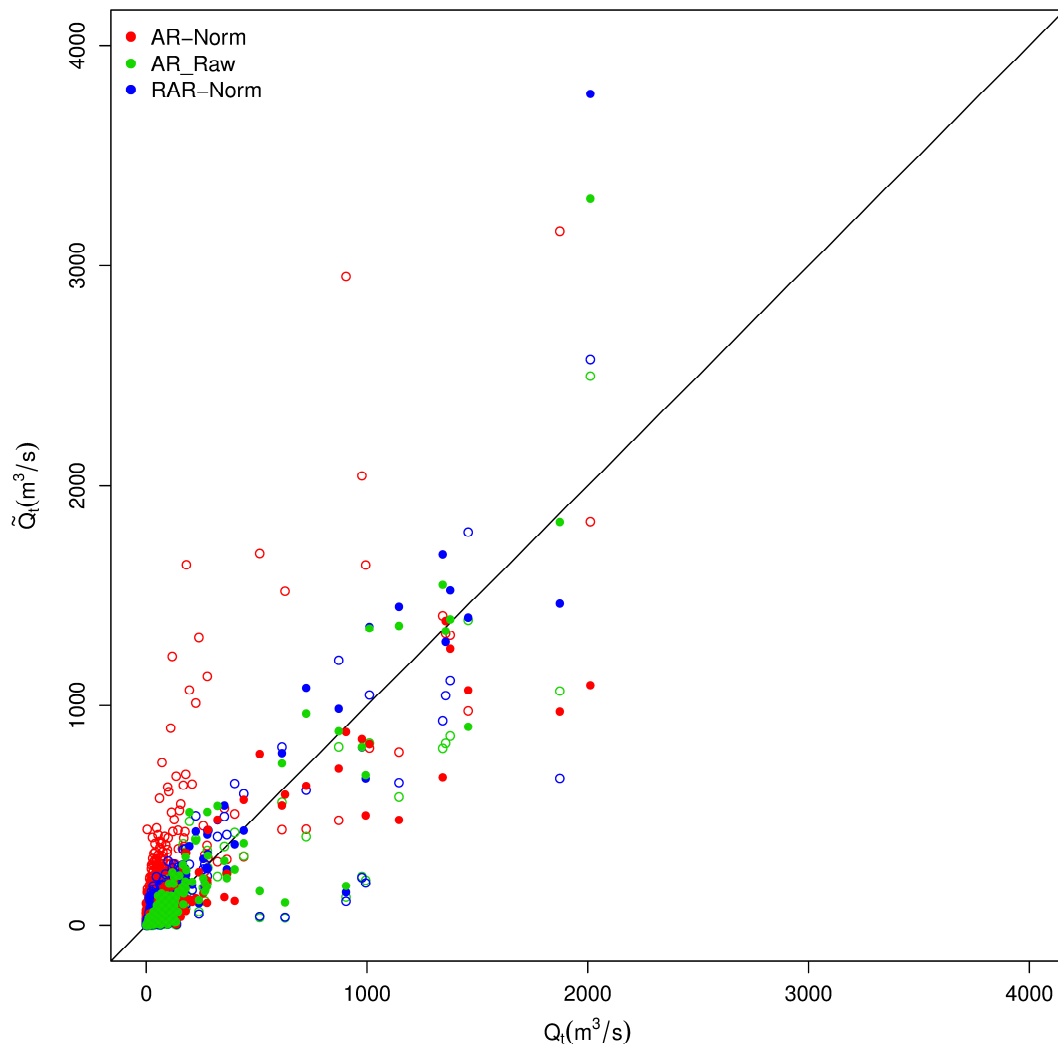


703

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

707

708

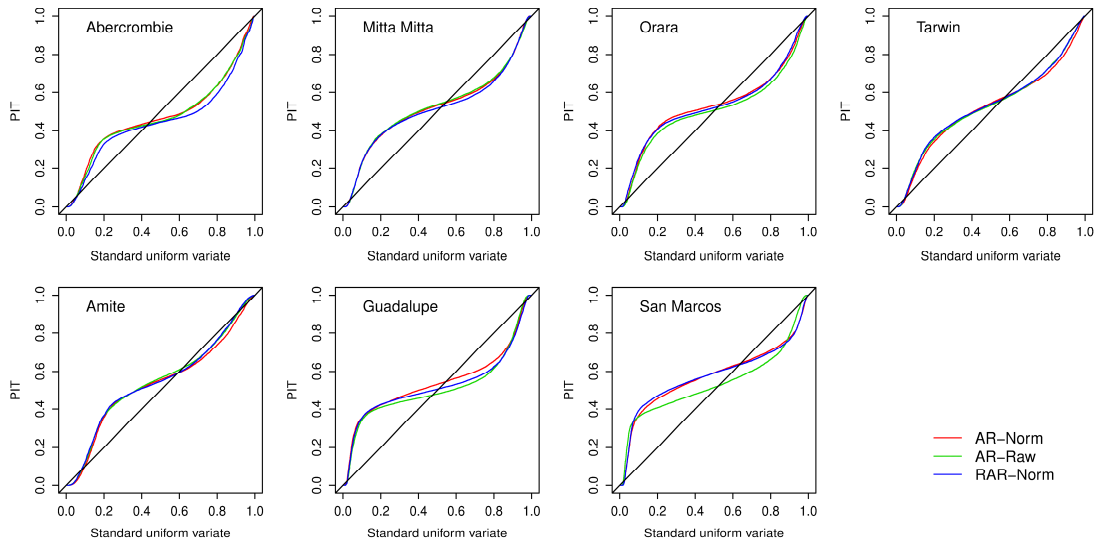


709

710 Figure 8: Comparison of the observed streamflows (Q_t) and forecast streamflows711 (\tilde{Q}_t), as forecast: 1) with the base hydrological model (circles), and 2) with the base

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

713



714

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

716 reliable forecasts.