

1 Deep learning for monthly rainfall-runoff modelling: a large-sample 2 comparison with conceptual models across Australia

3 Stephanie R. Clark¹, Julien Lerat², Jean-Michel Perraud², Peter Fitch²

4 ¹CSIRO, Environment, Sydney, NSW, Australia

5 ²CSIRO, Environment, Canberra, ACT, Australia

6 *Correspondence to:* Stephanie Clark (stephanie.clark@csiro.au)

7

8

9 **Abstract**

10 A deep learning model designed for time series predictions, the long short-term memory (LSTM) architecture is regularly
11 producing reliable results in local and regional rainfall-runoff applications around the world. Recent large-sample-hydrology
12 studies in North America and Europe have shown the LSTM to successfully match conceptual model performance at a daily
13 timestep over hundreds of catchments. Here we investigate how these models perform in producing monthly runoff
14 predictions in the relatively dry and variable conditions of the Australian continent. The monthly timestep matches historic
15 data availability and is also important for future water resources planning, however it provides significantly smaller training
16 data sets than daily time series. In this study, a continental-scale comparison of monthly deep learning (LSTM) predictions
17 to conceptual rainfall-runoff model (WAPABA) predictions is performed on almost 500 catchments across Australia with
18 performance results aggregated over a variety of catchment sizes, flow conditions, and hydrological record lengths. The study
19 period covers a wet phase followed by a prolonged drought, introducing challenges for making predictions outside of known
20 conditions – challenges that will intensify as climate change progresses. The results show that LSTMs matched or exceeded
21 WAPABA prediction performance for more than two-thirds of the study catchments; the largest performance gains of LSTM
22 versus WAPABA occurred in large catchments; the LSTM models struggled less to generalise than the WAPABA models
23 (eg. making predictions under new conditions); and catchments with few training observations due to the monthly timestep
24 did not demonstrate a clear benefit with either WAPABA or LSTM.

25 **Key words:** Hydrology and water resources, machine learning, deep learning, benchmarking to traditional models, neural
26 networks or ANN, LSTM

27

28 **Major points**

- 29
- 30 1. A deep learning model (single-layer LSTM) matched or exceeded performance of a WAPABA rainfall-runoff
model in 69% of study catchments.
 - 31 2. Monthly datasets contain enough information to train the LSTMs to this level.
 - 32 3. Generalisation to new conditions was found to improve with use of the LSTM, with implications for modelling
33 under climate change.

34 1. Introduction

35 With progressively variable climate conditions and the ever-increasing accessibility of hydrologic data, there comes the
36 opportunity to reconsider how available data is being used to efficiently predict streamflow runoff on a large scale.
37 Hydrological researchers are increasingly turning to emerging machine learning techniques such as deep learning to analyse
38 this increasing volume of data, due to the relative ease of extracting useful information from large datasets and producing
39 accurate predictions about future conditions without the need for detailed knowledge about the underlying physical systems.
40 Machine learning models have been shown to be capable of obtaining more information from hydrological datasets than is
41 abstracted with traditional models, due to their automatic feature engineering and ability to effectively capture high-
42 dimensional and long-term relationships ([Nearing et al., 2021](#), [Frame et al., 2021](#)). The continually evolving machine learning
43 field will continue to offer novel opportunities that can be harnessed for hydrological data analyses, and it is important to
44 understand how these methods relate to classical models. Here, a basic machine learning model is benchmarked against a
45 traditional conceptual model over a large sample of catchments as a step towards a general understanding of the use of deep
46 learning models as a tool for the task of monthly rainfall-runoff modelling in Australian catchments.

47 Deep learning models have been shown in many applications to provide accurate hydrological predictions and classifications
48 ([Kratzert et al., 2018](#), [Shen et al., 2021](#), [Reichstein et al., 2019](#), [Frame et al., 2022](#)). These models are particularly useful to
49 hydrological studies as they provide the potential to quickly add and remove predictors ([Shen, 2018](#)), scale to multiple
50 catchments ([Kratzert et al., 2019](#), [Lees et al., 2021](#)), automatically extract useful and abstract information from large datasets
51 ([Reichstein et al., 2019](#), [Shen, 2018](#)), make predictions in areas with little or no data ([Kratzert et al., 2019b](#), [Majeske et al.,](#)
52 [2022](#), [Ouma et al., 2022](#), [Choi et al., 2022](#)), and extrapolate proficiently to larger hydrologic events than are seen in the
53 training dataset ([Frame et al., 2022](#)), [Li et al., 2021](#), [Song et al., 2022](#)).

54 The long short-term memory network (LSTM, ([Hochreiter and Schmidhuber, 1997](#))), is a deep learning model that is gaining
55 popularity in hydrology for daily time series predictions at individual basins or groups of basins due to its ability to efficiently
56 and accurately produce predictions without requiring assumptions about the physical processes generating the data. The
57 LSTM is a type of recurrent neural network (RNN), an extension of the multilayer perceptron that is specifically designed for
58 use with time series data through its sequential consideration of input data. The LSTM further extends the RNN to incorporate
59 gates and memory cells, allowing for input data to be remembered over much longer time periods and for unimportant data
60 to be forgotten from the network. LSTMs make predictions by taking into account both the short and long temporal patterns
61 in a time series as well as incorporating information from exogenous predictors. The data-driven detection of intercomponent,
62 spatial and temporal relationships by these deep learning models can be of particular benefit when attempting to represent
63 systems in which the physical characteristics are not well defined and the intervariable relationships are complex.

64 The increasing popularity of the LSTM in hydrology is due to its ability to capture the short-term interactions between rainfall
65 and runoff, as well as the long-term patterns and interactions arising from longer-frequency drivers such as climate, catchment
66 characteristics, land use and changing anthropogenic activity. A growing number of publications are applying LSTMs to
67 hydrological simulations and comparing results to process-based or conceptual modelling results.

68 A gap exists in the literature concerning a comparison of LSTM models and conceptual models at a monthly time step over
69 a large sample of catchments. The conditions in which LSTMs or conceptual models may have an advantage for monthly
70 rainfall-runoff modelling, in a general sense, are not yet understood as most machine learning applications in hydrology are
71 individual-basin case studies ([Papacharalampous et al., 2019](#)) at a daily timestep or higher frequency (eg. [Li et al., 2021](#),
72 [Yokoo et al., 2022](#)). Though the LSTM has successfully matched conceptual model performance in some large-sample-
73 hydrology studies at daily timesteps (eg. In the USA ([Kratzert et al., 2019b](#)) and the UK ([Lees et al., 2021](#))) it is yet unknown

74 how these models compare to conceptual models for monthly runoff predictions in relatively dry conditions such as those
75 characterised by Australian catchments.

76 Monthly hydrological models are important tools for water resources assessments as hydrologic data has historically been
77 recorded at a monthly or longer frequency based on the schedule of manually-collected measurements. Furthermore, the
78 monthly timestep is often the most practical for water resources planning with many decisions requiring only monthly
79 streamflow predictions. With their simpler structure, fewer parameters and lower data requirements compared to daily models
80 ([Hughes, 1995](#), [Mouelhi et al., 2006](#)), monthly models are also useful tools to investigate uncertainty in rainfall-runoff model
81 structure ([Huard and Mailhot, 2008](#)) and to support probabilistic seasonal streamflow forecasting systems ([Bennett et al.,
82 2017](#)). Due to data availability, models designed to run on monthly timesteps can be used across much larger areas, informing
83 important large-scale water resources decision-making. For these reasons, generalisable models at monthly timesteps are
84 vital. However, the monthly timestep is traditionally a difficult one to model as it requires extracting both short and long-
85 term hydrologic processes ([Machado et al., 2011](#)). In a machine-learning context, the monthly time step differs significantly
86 from the daily time step as it drastically reduces the size of the data set available for model training (by a factor of 30). As
87 the convergence of machine learning algorithms typically improves with larger data sets, a central research question of this
88 paper is to explore the capacity of the LSTM algorithm to cope with the reduced amount of input data imposed by the monthly
89 time step.

90 LSTMs have been used to model the rainfall-runoff relationship at a monthly time step in a limited number of localised
91 studies, showing potential for this application on a broader scale. [Ouma et al. \(2022\)](#) used monthly aggregated data due to
92 low data availability in three scarcely-gauged basins the Nzoia River basin, Kenya. [Majeske et al. \(2022\)](#) trained LSTMs
93 with spatially- and temporally-limited data for three sub-basins of the Ohio River Basin, claiming the daily timestep was
94 superfluous and cumbersome in some conditions. [Lee et al. \(2020\)](#) found the LSTM adept at preserving long-term memory
95 in monthly streamflow at a single station on the Colorado River over a 97-year study without any weakening of the short-
96 term memory structure. [Yuan et al. \(2018\)](#) used a novel method for parameter calibration in an LSTM for monthly rainfall-
97 runoff estimation at a single station on the Astor River basin in northern Pakistan. [Song et al. \(2022\)](#) found the LSTM better
98 reproduced observed monthly runoff and simulated extreme runoff events than a physically-based model at five discharge
99 stations in the Yeongsan River basin in South Korea.

100 Large-sample hydrologic studies that assess methods on a large number of catchments are being increasingly called for in
101 the field of hydrology ([Papacharalampous et al., 2019](#), [Mathevet et al., 2020](#), [Gupta et al., 2014](#)). [Papacharalampous et al.
102 \(2019\)](#) compared the performance of a number of statistical and machine learning methods (no LSTM) on 2000 generated
103 timeseries and over 400 real-world river discharge timeseries and determined that the machine learning and stochastic
104 methods provided similar forecasting results. [Mathevet et al. \(2020\)](#) compared daily conceptual model performance (no
105 machine learning) for runoff prediction in over 2000 watersheds, determining that performance depended more on catchment
106 and climate characteristics than on model structure. [Kratzert et al. \(2018\)](#) found individual daily-scale LSTMs were able to
107 predict runoff with accuracies comparable to a baseline hydrological model for over 200 differently complex catchments.
108 ([Kratzert et al., 2019b](#)) found a global LSTM trained on over 500 basins in the United States with daily data produced better
109 individual catchment runoff predictions than conceptual and physically-based models calibrated on each catchment
110 individually. ([Lees et al., 2021](#)) produced a global LSTM to model almost 700 catchments in Great Britain, finding that this
111 model outperformed a suite of benchmark conceptual models, showing particular robustness in arid catchments and
112 catchments where the water balance does not close. ([Jin et al., 2022](#)) compared machine learning daily rainfall-runoff models
113 to process-based models for over 50 catchments in the Yellow River Basin in China. ([Frame et al., 2021](#)) found that a global

114 LSTM with climate forcing data performed similarly or outperformed a process-based model on over 500 US catchments,
115 and that in catchments where hydrologic conditions are not well understood the LSTM was a better choice.

116 This study aims to determine the ability of a simple machine learning model (a single-layer LSTM) to match or exceed the
117 performance of a conceptual monthly rainfall-runoff model (the WAPABA model ([Wang et al., 2011](#))) for predicting runoff,
118 using inputs derived from easily accessible climate variables. The goal here is not to maximise LSTM performance to cutting-
119 edge machine learning standards, rather to ascertain the minimum performance level that a non-expert user might expect to
120 obtain from basic usage of an LSTM with the input data regularly used in a conceptual model. A frequently heard reason for
121 hydrological researchers not engaging machine learning approaches is the small data size associated with individual
122 catchment time series, and it is of interest to examine the lower limits of data availability required to fit an LSTM with
123 individual catchment monthly data sets.

124 A comparison is made on almost 500 basins across Australia, representing a wide variety of catchment types, hydro-climate
125 conditions, and with differing amounts of historical data. The prediction performance of the LSTM machine learning models
126 is compared to the WAPABA conceptual models for each individual catchment. The proportion of catchments in which the
127 runoff prediction performance of the conceptual model is met or exceeded by the machine learning model is determined.
128 Conditions under which the machine learning models or the conceptual models may have an advantage are investigated, such
129 as catchment size, flow level, and length of historical record. The central questions of this study are:

- 130 1) In general, do LSTMs match conceptual model prediction performance on Australian catchments?
- 131 2) Is the reduced number of data points due to the monthly time step an issue for training an LSTM?
- 132 3) Under what conditions is the LSTM of particular benefit or drawback? (eg. catchment size, flow level, amount of
133 training data, etc.)

134 The results of this large-sample analysis of LSTM performance over the Australian continent will assist in understanding
135 whether LSTMs are a justifiable alternative to conceptual models for monthly rainfall-runoff prediction in Australia and
136 similar environments, including if monthly data sets are sufficient to produce accurate predictions with the LSTM. Building
137 on the results of this study, further benefits of deep learning could be harnessed through the creation of larger-scale models
138 that encompass climatic, hydrologic and anthropogenic patterns spanning multiple catchments, allowing for the sharing of
139 information under similar conditions and the potential transfer of knowledge between data-rich and data-scarce regions, or
140 models that blend conceptual models into the machine learning network structure.

141

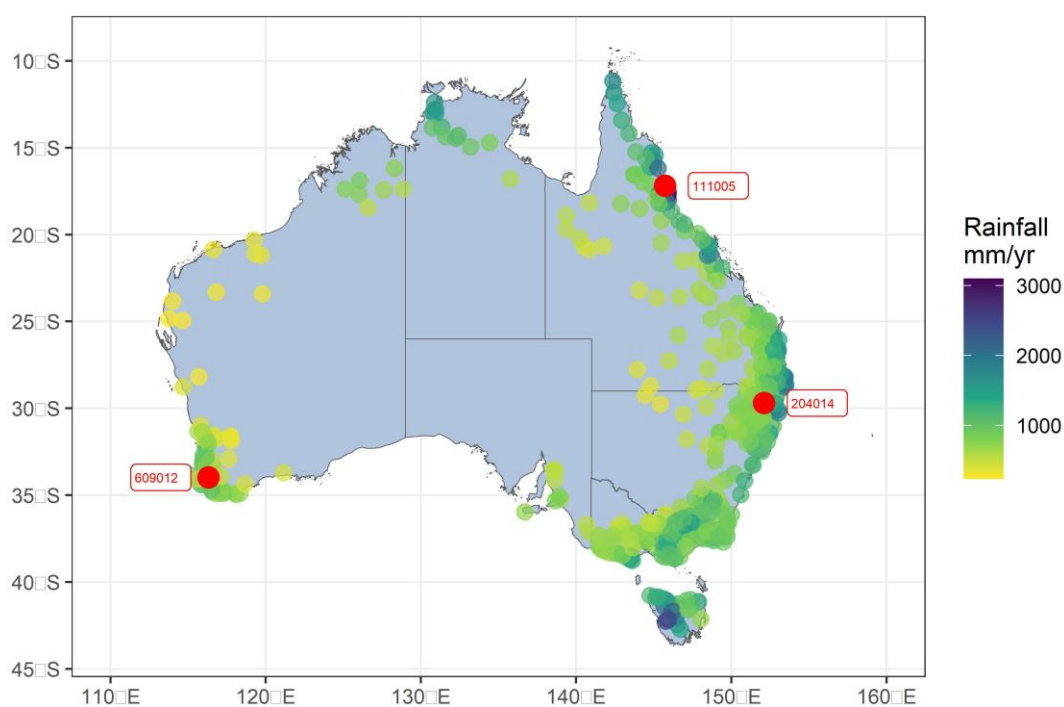
142 **2. Data and Methods**

143 **2.1. Data**

144 The catchment and climate data used in this study are from a dataset curated by [Lerat et al. \(2020\)](#) comprising a selection of
145 basins across Australia. The dataset spans all main climate regions of the continent, providing data from a variety of rainfall,
146 aridity and runoff regimes, as described in Table 1. Catchments where some data were marked as suspicious (e.g. high flow
147 data with large uncertainties, inconsistencies, suspected errors) or with more than 30% missing data were excluded. This left
148 496 catchments in the study, with locations as shown in Figure 1. The area of the individual catchments ranges from
149 approximately 5 km² to 120,000 km².

150

Variable	Min	Q25	Median	Q75	Max
Catchment area (km ²)	4	180	449	1,456	119,000
Mean rainfall (mm/y)	237	691	887	1130	3097
Mean PET (mm/y)	918	1280	1500	1755	2321
Mean runoff (mm/y)	0.5	46	130	275	2213
Aridity index rainfall/PET (-)	0.11	0.44	0.61	0.81	2.61
Daily rainfall skewness (-)	2.4	4.8	5.9	7.4	16.7
Runoff coeff. Runoff/rainfall (-)	0.001	0.069	0.150	0.255	0.902
% zero flows in daily series	0.0	0.0	3.4	23.7	74.0



152

153 **Figure 1: Locations of the 496 study catchments, coloured by mean annual rainfall. The three labelled catchments, which will be**
 154 **used as examples during the study, represent a wet catchment (111005 in Northern Queensland), a temperate catchment (204014**
 155 **in New South Wales), and a dry catchment (609012 in Western Australia).**

156

157 Observed runoff data were collected from the Bureau of Meteorology's Water Data online portal
 158 (<http://www.bom.gov.au/waterdata>), rainfall and temperature data are from the Bureau of Meteorology's AWAP archive
 159 (Jones et al., 2009), and potential evapotranspiration data was computed by the Penman equation as part of the AWRA-L
 160 landscape model developed jointly by CSIRO and the Bureau of Meteorology (Frost et al., 2018). Rainfall, temperature and
 161 evapotranspiration are averaged from daily grids (5x5km) over each of the catchments.

162 The runoff records begin between January 1950 and September 1982, and end between October 2016 and June 2020. The
 163 number of runoff observations per catchment ranges from 425 to 846 with a median dataset size of 613 observations. The
 164 rainfall and potential evapotranspiration data cover the period from 1911 to 2020 continuously. The resulting dataset consists
 165 of a set of 496 time series ranging from 37 to 70 years in length, with a median record length of 51 years.

166

167 **2.1.1. Training and testing data split**

168 The data set for each catchment is split into two portions for modelling – in machine learning these are referred to as ‘training’
169 and ‘testing’ sets, corresponding to the traditional ‘calibration’ and ‘validation’ sets used in hydrologic modelling. The
170 training data set runs from January 1950 (or the start of the station’s record, if later) to December 1995 for all catchments.
171 The testing data set begins in January 1996 for all catchments and ends in July 2020 (or at the end of the station’s record, if
172 sooner). This split is chosen to divide the streamflow records into two relatively even periods, but also to distinguish an early
173 wet period from a testing period characterised by the Millennium Drought over south-eastern and eastern Australia ([Van Dijk
174 et al., 2013](#)). WAPABA and LSTM models were trained and evaluated using the same data splits, giving identical duration
175 and dataset size.

176 When split into training and testing sets at the beginning of January 1996, between 38% and 72% of the data from each
177 catchment becomes the training set. The length of the training data record for individual catchments ranges from 14 to 47
178 years, with the smallest data set used for training containing 172 observations. Typically in machine learning, a portion of
179 the training data is held back to be used during the model fitting process to monitoring for over-fitting and to signal early
180 stopping of training if necessary. Since the training data sets in this study are already small by machine learning standards,
181 this has not been done as it would reduce the number of training observations significantly, as well as lead to a smaller
182 training dataset than used in the WAPABA models. A sensitivity test has been performed to justify this choice, and it was
183 found that training the LSTMs with 20% of the training data reserved for this task (ie. with the data split into training: 64%,
184 validation: 16%, testing: 20%) produced no apparent benefit in prediction performance.

185 **2.2. Models**

186 **2.2.1. Deep learning time series models (LSTMs)**

187 The long short-term memory network, LSTM ([Hochreiter and Schmidhuber, 1997](#)), is an updated recurrent neural network
188 (RNN) specifically designed for deep learning with time series data. The inclusion of gates and memory cells increases the
189 length of time series the LSTM is able to process; three gates (input, output and forget gates) regulate the flow of information
190 into and out of the memory cell, determining which information from the past is to be retained and which can be forgotten.
191 In this way, each member of the LSTM output becomes a function of the relevant input at previous timesteps.

192 The LSTM network consists of an input layer, one or more hidden layers, and an output layer. The layers are connected by a
193 set of updatable weights, with the same weights applying to all timesteps of the data. Memory cells shadow each node on the
194 hidden layer, retaining important information over long time periods. Each node of the input layer represents a variable of
195 the input data set. Observations are fed into the network along with a pre-specified number of predictor values from previous
196 timesteps (known as the lookback length, or lag) which are cycled sequentially through the network. Network weights are
197 updated by backpropagating the gradient of the error between the modelled and observed outputs. For detailed information
198 on the mathematical functioning of the LSTM, see ([Goodfellow et al., 2016](#)) and ([Kratzert et al., 2018](#)).

199 In this study, a separate LSTM is trained for each catchment. Input to the LSTMs are monthly averaged measurements of:
200 rainfall depth (P), potential evapotranspiration (E), average maximum daily temperature over the month, and net monthly
201 (effective) rainfall (P^*) computed for month t by summing daily effective rainfall, as shown here:

$$P_t^* = \sum_{d=0}^{d=\text{days}(t)} \max(0, P_d - E_d) \quad 1$$

202 Standard scaling of the input data is performed per catchment as follows:

$$\tilde{X}_t = \frac{X_t - \mu_x}{\sigma_x} \quad 2$$

203 where X_t is an input variable for month t , μ_x is its mean and σ_x its standard deviation over the training period. The target
 204 variable for LSTM training is monthly average runoff. Observed runoff values are scaled by taking the square root and then
 205 transforming to the range [-1,1] per catchment, as follows:

$$Y_t = 2 \frac{\sqrt{Q_t} - Y_0}{Y_1 - Y_0} - 1 \quad 3$$

206 where Q_t is the observed runoff for month t , and Y_0 and Y_1 are the minimum and maximum square root transformed flow
 207 over the training period, respectively. The square root transform is chosen to be conceptually consistent with the objective
 208 function of the WAPABA model calibration (as described below, mean absolute error of the square roots of flows). Note that
 209 the same scaling constants ($\mu_x, \sigma_x, Y_0, Y_1$) used during LSTM training are also applied to LSTM inputs and targets for the
 210 testing period. Using scaling constants only derived from the training data ensures that the training process is not
 211 incorporating any information from the testing data set.

212 The loss function used for training the LSTM is the mean absolute error (MAE) performed on the transformed runoff, as
 213 follows:

$$L = \sum_t |Y_t - \hat{Y}_t| \quad 4$$

214 where \hat{Y}_t is the output of the network for month t and Y_t is the transformed runoff for the same month.

215 Hyperparameters, or parameters controlling the LSTM training algorithm, were selected after a grid search (over 1016
 216 separate runs) on a randomly selected catchment (14207) with a good length data record and tested on a small additional
 217 subset of catchments. As the purpose of this study was not to optimise catchment-specific predictions results, a more
 218 comprehensive hyperparameter search by catchment was deemed unnecessary. The hyperparameter space searched was:
 219 initial learning rate δ_0 (1e-3 to 1e-4), sequence (lookback or lag) length (6, 9, 12, 15, 18, 21, 24 months) and number of
 220 hidden nodes (10, 20, 30, 40, 50, 60). The hyperparameter set that performed the best predictions over the training period
 221 was selected for use in all LSTMs: 10 nodes on a single hidden layer, run with a sequence length 6 months, and an initial
 222 learning rate δ_0 of 0.0001. Subsequent to this hyperparameter search, the effect of raising the initial learning rate for faster
 223 convergence while using input and recurrent dropout to prevent overfitting was investigated on all catchments. Empirically,
 224 and counter to our intuition, this never improved training performance and so the initial learning rate δ_0 of 0.0001 was
 225 retained. The learning rate was allowed to vary during training with a patience of 3 epochs without improvement before
 226 multiplying by a factor of 0.2 to obtain a new learning rate. The dataset was divided into 400 steps-per-epoch for training;
 227 data was sent through the model in batches with a weight update after each (an epoch, or iteration, is concluded when the
 228 entire dataset has been run through the model once). The LSTM training was implemented using a gradient descent algorithm
 229 run for a maximum of 100 epochs. Training was set to stop early if the training error failed to decrease over 5 consecutive
 230 epochs. The LSTMs were implemented with Tensorflow in Python, using numeric seeds to ensure reproducible outcomes.

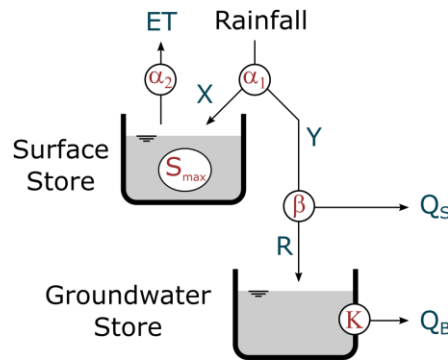
231 **2.2.2. WAPABA rainfall-runoff models**

232 The WAPABA model is a conceptual monthly rainfall-runoff model introduced by [Wang et al. \(2011\)](#). The model is an
 233 evolution of the Budyko framework proposed by [Zhang et al. \(2008\)](#) where water fluxes are partitioned using parameterised
 234 curves. The model uses two inputs, mean monthly rainfall and potential evapotranspiration, and operates in five stages. First,
 235 input rainfall is split between effective rainfall that will eventually leave the catchment, and catchment consumption that
 236 replenishes soil moisture and evaporates. Second, catchment consumption is portioned between soil moisture replenishment
 237 and actual evapotranspiration. Third, effective rainfall is partitioned between surface water (fast) and groundwater (slow)
 238 stores. Fourth, the groundwater store is drained to provide a baseflow contribution. Fifth, the surface water and baseflow are
 239 added to obtain the final simulated runoff for the month. The model has five parameters described in Table 2 [which interact](#)
 240 [as depicted in](#) Figure 2.

241 Table 2 WAPABA model parameters

Name	Description	Unit	Minimum	Maximum
alpha1	Exponent of the catchment consumption/effective rainfall curve	Dimensionless	1.0	10.0
alpha2	Exponent of the soil moisture storage/evapotranspiration curve	Dimensionless	1.0	10.0
Beta	Partition between groundwater recharge and surface runoff	Dimensionless	0.0	1.0
Smax	Maximum water-holding capacity of soil store	mm	5.0	6000.0
Inverse K	Inverse of groundwater store time constant	1/day	0.000274	1.0

242



243

244 **Figure 2: WAPABA model schematic**

245 A separate WAPABA model is run for each study catchment. The WAPABA models were trained (calibrated) and tested
 246 (validated) over the same periods as the LSTMs: 1950 to 1995 inclusive for training, and 1996 to June 2020 for testing. The
 247 model was calibrated with a warm-up period of 2 years to avoid possible bias associated with initial values. WAPABA
 248 parameters were optimized over the training period using the Shuffle Complex Evolution algorithm ([Duan et al., 1993](#)) with
 249 the Swift software package ([Perraud et al., 2015](#)). The objective function used for the WAPABA models is the same as the
 250 one used for LSTM, i.e. the mean absolute error (MAE) on the square root of runoff (see Equation 4).

251 2.3. Performance evaluation

252 Predictions from the conceptual (WAPABA) and machine learning (LSTM) models for all catchments are compared to
253 observed runoff, assessing each model's predictive capabilities on the set of catchments. Runoff prediction performance is
254 reported here using the following metrics.

255 The Nash Sutcliffe Efficiency (NSE, [Nash and Sutcliffe, 1970](#)) is the most often used performance metric in hydrology. It
256 can be considered a normalised form of mean squared error (MSE) and is defined as:

$$NSE = 1 - \frac{\sum_t (Q_{obs}^t - Q_{mod}^t)^2}{\sum_t (Q_{obs}^t - \mu_{obs})^2} = 1 - \frac{E}{V} \quad 5$$

257 where Q_{obs}^t and Q_{mod}^t are the observed and modelled discharges for month t , respectively, and μ_{obs} is the average observed
258 discharge over the training or testing period. The ratio of the sum of squared errors, $E = \sum_t (Q_{obs}^t - Q_{mod}^t)^2$, to the variance,
259 $V = \sum_t (Q_{obs}^t - \mu_{obs})^2$, is subtracted from a maximum score of 1. An NSE closer to 1 indicates better predictive capability
260 of the model, and an NSE less than 0 indicates the model mean squared error is larger than the observation variance.

261 The NSE metric alone cannot provide an accurate description of model performance due to its focus on high flow regime
262 ([Schaeffli and Gupta, 2007](#)). The reciprocal NSE focuses the error metric on low flows ([Pushpalatha et al., 2012](#)) by comparing
263 the reciprocals of the observed and modelled flows. It is calculated as:

$$RecipNSE = 1 - \frac{\sum_t \left(\frac{1}{(Q_{obs}^t + 1)} - \frac{1}{(Q_{mod}^t + 1)} \right)^2}{\sum_t \left(\frac{1}{(Q_{obs}^t + 1)} - \frac{1}{(\mu_{obs} + 1)} \right)^2} \quad 6$$

264 The Kling-Gupta efficiency (KGE, [Gupta et al., 2009](#)) provides an alternative to metrics based on sum of squared error
265 such as the two previous ones, by equally weighting measures of bias of the mean, variability, and correlation into a single
266 metric as follows:

$$KGE = 1 - \sqrt{\left(1 - \frac{\mu_{sim}}{\mu_{obs}}\right)^2 + \left(1 - \frac{\sigma_{sim}}{\sigma_{obs}}\right)^2 + (1 - \rho)^2} \quad 7$$

267 where μ_x and σ_x are the mean and the standard deviation and ρ is the Pearson correlation coefficient between the simulated
268 and observed data.

269 Finally, Bias is a measure of consistent under-forecasting or over-forecasting of the mean, defined as:

$$Bias = \frac{\mu_{sim} - \mu_{obs}}{\mu_{obs}} \quad 8$$

270 Comparison of performance metrics between catchments using normalised indexes

271 When comparing metrics across model types and catchments, a normalised difference in NSE values is used. The NSE metric
272 can reach into large negative values in dry catchments when the variance of the observations is very small compared to the
273 model errors ([Mathevet et al., 2006](#)), as can be seen from Equation 5. Differences between large negative values of NSE have
274 a much smaller implication than the same absolute difference between values of NSE closer to 1. To allow for a comparison
275 between the WAPABA and LSTM models at catchments of various aridities, the normalised difference in NSE is calculated
276 following [Lerat et al. \(2012\)](#):

$$Diff_NSE_{norm} = \frac{NSE_2 - NSE_1}{(1 - NSE_1) + (1 - NSE_2)} = \frac{NSE_2 - NSE_1}{2 - (NSE_1 + NSE_2)} \quad 9$$

277 where NSE_1 and NSE_2 are the NSE values corresponding to the two models to be compared. Substituting in $NSE = 1 - \frac{E}{V}$
 278 from Equation 5 into Equation 9, the normalised difference in NSE can be seen to represent a percentage difference in the
 279 sum of squared errors between the two models being compared:

$$Diff_NSE_{norm} = \frac{NSE_2 - NSE_1}{2 - (NSE_1 + NSE_2)} = \frac{E_1 - E_2}{E_1 + E_2} \quad 10$$

280 A similar formula is applied to reciprocal NSE and KGE. The normalised difference between the bias for two models is
 281 calculated as:

$$Diff_Bias_{norm} = \frac{|Bias_1| - |Bias_2|}{|Bias_1| + |Bias_2|} \quad 11$$

282 To simplify the comparison of model results across the large number of catchments, model performances at each catchment
 283 are classified as similar if the normalised difference between WAPABA and LSTM metrics lies within +/- 0.05 at that
 284 catchment, following [Lerat et al. \(2020\)](#). Therefore in this paper, a ‘similar’ NSE denotes that the sum of squared errors of
 285 the WAPABA and LSTM models at an individual catchment differ by no more than 5%. For differences greater than this,
 286 the catchments are classified by the model type producing the higher metric. The selection of the threshold of 0.05 was based
 287 on the recommendations of ([Lerat et al., 2020](#)) and the authors’ experience relative to the use of the NSE, KGE and bias
 288 metrics.

289

290 **3. Results**

291 For each of the study catchments, a WAPABA model and an LSTM model have been trained using monthly data over the
 292 training period, and the prediction performance of the models are evaluated here on monthly data from the testing period
 293 (data unseen by the model during training) using the metrics described above. A general comparison of WAPABA and LSTM
 294 prediction performance is first made over all catchments with a continental-scale analysis of the performance metrics, to
 295 determine:

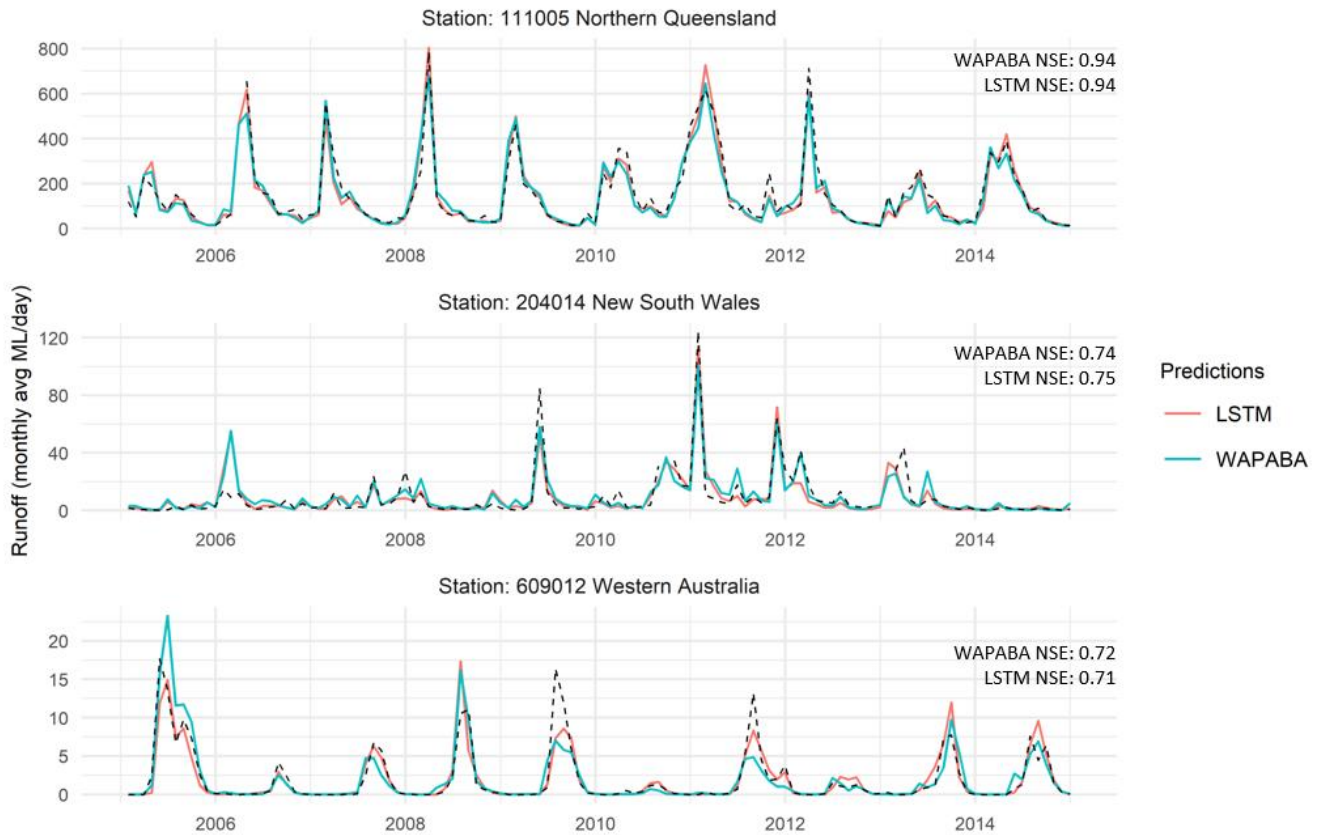
- 296 1) the proportion of overall catchments for which the WAPABAs or the LSTMs produced better predictions,
 297 and
- 298 2) differences at individual catchments in WAPABA versus LSTM prediction performance.

299 A comparison of model performance is then made in relation to various catchment and time series characteristics (eg.
 300 catchment size, flow level, record length), to determine if an association exists between these properties and the relative
 301 performance of the conceptual and machine learning models.

302 **Example prediction results**

303 As a sample of the modelling output, Figure 3 shows the WAPABA and LSTM runoff predictions along with the
 304 corresponding observed runoff for the three stations highlighted in Figure 1 (over the testing period). These hydrographs are
 305 representative of a wet catchment in Northern Queensland (Mulgrave River at the Fisheries, 111005), a temperate catchment

306 in NSW (Mann River at Mitchell, 204014), and a dry, intermittent catchment in Western Australia (Blackwood River at
 307 Winnejup, 609012). NSE values of each of the predictions are noted. The WAPABA and LSTM predictions both match the
 308 observed data reasonably well in the three catchments. The performance of the models, in particular for the Blackwood River
 309 at Winnejup is remarkable because of the difficulty in modelling dry, intermittent catchments (Wang et al., 2020). The next
 310 sections provide a more detailed assessment of the performance over all catchments using quantitative metrics.



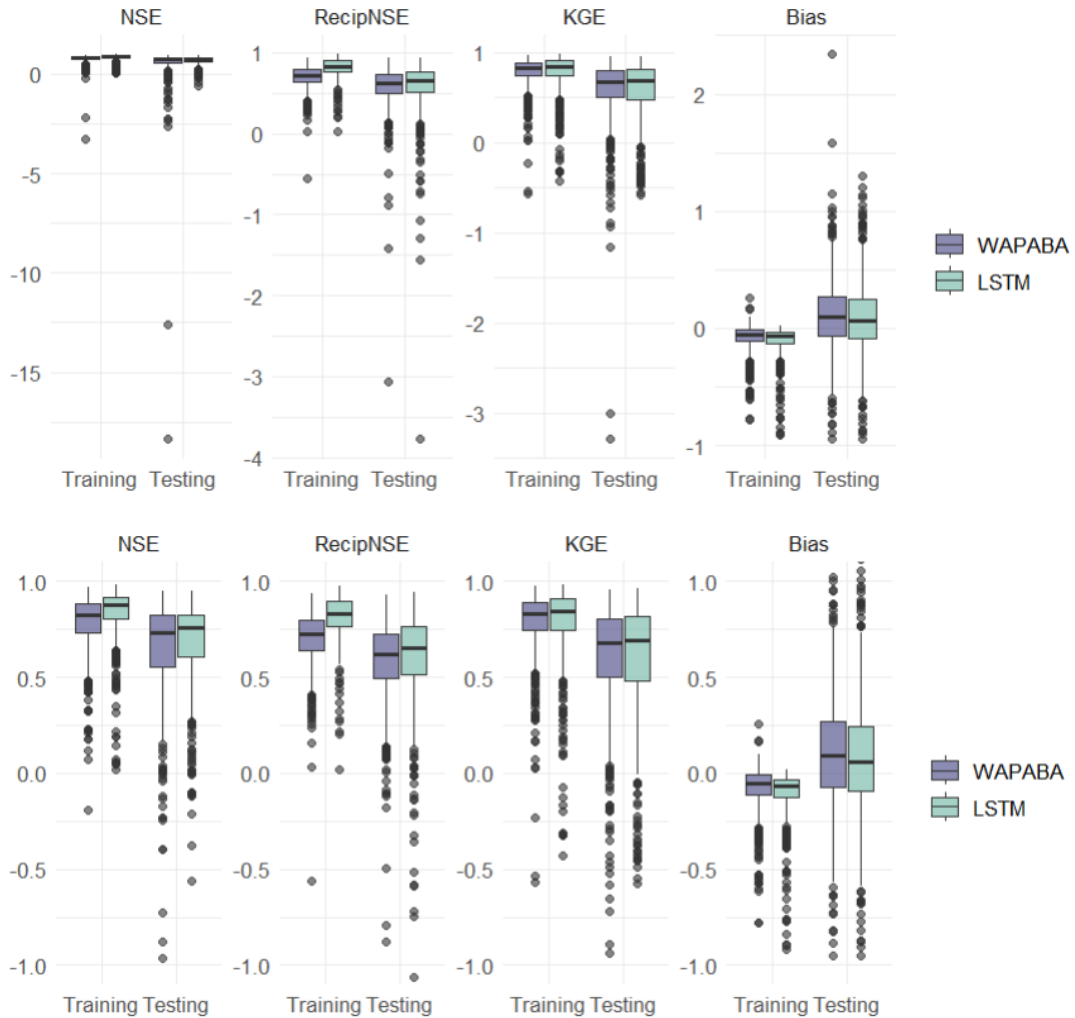
311
 312 **Figure 3: Observed data (black dashed line) and predicted runoff (by WAPABA and LSTM models) over the testing period for**
 313 **the Mulgrave River at the Fisheries (111005), Mann River at Mitchell (204014) and the Blackwood River at Winnejup (609012).**
 314 **Catchment locations are shown on Figure 1.**

315
 316 **Large-sample performance summary**

317 The general runoff prediction performance of WAPABA and LSTM models on a continent-wide basis is summarized in
 318 Figure 4. From the models run for each catchment, metrics are determined on the training portion (calibration) and testing
 319 portion (validation) separately and gathered here in boxplots. Median and quartiles of NSE, reciprocal NSE, KGE and Bias
 320 over all catchments are shown for each model type, with each data point representing an individual catchment. All data is
 321 shown on the top panel, and due to a few large (negative) outliers the same figure is shown with a restricted y-axis for
 322 visualization purposes on the lower panel. Higher values of the first three metrics (NSE, reciprocal NSE and KGE) indicate
 323 a better match of predicted runoff with observed runoff, whereas lower values of Bias indicate better prediction results.

324 Figure 4 shows that across the set of study catchments the median values of NSE, Reciprocal NSE, and KGE are slightly
 325 higher for LSTM than for WAPABA during both the training and testing phases. Bias has a slightly lower median for the
 326 LSTM. As expected, both model types perform better during the training phase than the testing phase for all metrics. The
 327 difference between WAPABA and LSTM performance is relatively large during the training period but similar during testing,

328 indicating perhaps a higher tendency towards overfitting by the ML models than traditional modellers would be expecting.
 329 The interquartile ranges increase from training to testing (longer boxes during testing), indicating a greater spread of
 330 performance results when the models are run on data not seen during the training phase. Over all catchments, the median
 331 NSE is: 0.74 with the WAPABA models and 0.76 with the LSTM models (on testing data). See Table 3 for median values
 332 of these metrics.



333

334

335 **Figure 4: Performance metrics summary for the set of 496 catchments (zoomed in on lower panel, excluding outliers < -1). Median**
 336 **values of LSTM performance metrics are slightly higher than WAPABA for NSE, Reciprocal NSE and KGE (higher indicates**
 337 **better performance), and slightly lower for Bias (lower Bias is preferable). For all four metrics on both models, the longer testing**
 338 **boxes indicate more spread in performance results when predicting on new data.**

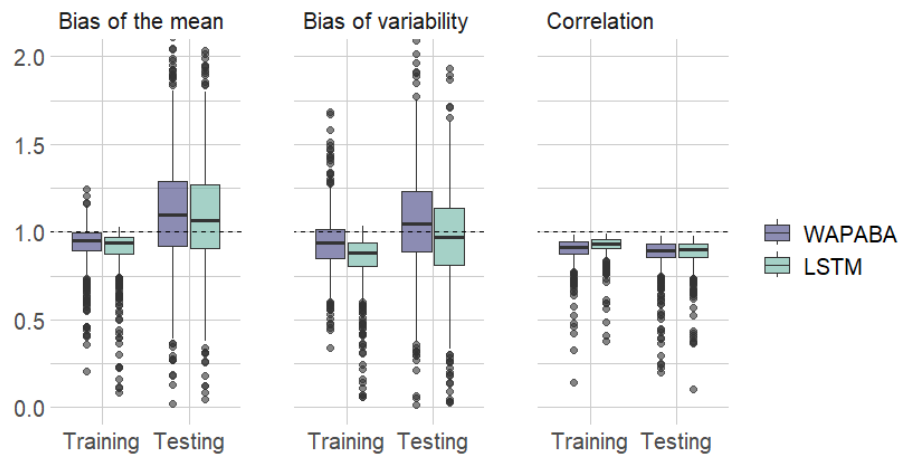
339 Aggregated performance metrics may mask performance variability within certain aspects of the time series ([Mathevet et](#)
 340 [al., 2020](#)). The KGE has the benefit of being easily decomposed into three components for further error analysis: bias of the
 341 mean (ratio of mean of simulations to mean of observations), bias of variability (ratio of standard deviation of simulations
 342 to standard deviation of the observations), and correlation (matching of the timing and shape of the time series to the
 343 observations).

344 Table 3: Median values of metrics over the set of catchments (n=496)

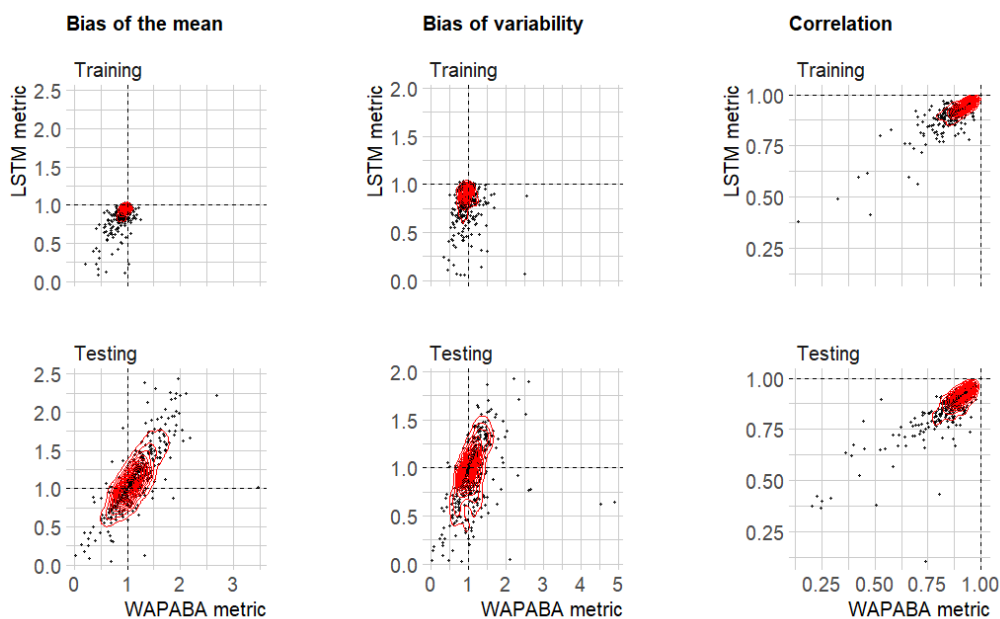
	NSE	Reciprocal NSE	KGE	Bias	Bias of the mean	Bias of variability	Correlation
WAPABA	0.74	0.62	0.68	0.09	1.10	1.05	0.90
LSTM	0.76	0.65	0.70	0.06	1.07	0.97	0.90

345 In Figure 5 and Table 3, model performance is assessed with respect to each component of the KGE metric. Boxplots of the
 346 decomposed KGE components are shown by model type and training/testing period. During testing, the medians of bias of
 347 the mean and standard deviation are above zero for WAPABA and greater for WAPABA than LSTM. This indicates that
 348 mean streamflow and streamflow variability tend to be overestimated more by the WAPABA models compared to the
 349 LSTMs. The LSTM median bias of variability is below zero, and therefore streamflow variability is more prone to
 350 underestimation. For bias of the mean and standard deviation, the depth of the boxplots increases from training to testing,
 351 indicating the bias values from individual catchments are more diverse during the testing period.

352 The scatterplots in the lower part of Figure 5 compare the KGE components at individual catchments for the WAPABA
 353 and LSTM models (each dot represents a catchment), separately for training and testing portions of the data. Most values of
 354 bias of the mean (left column) are between 0 and 1 during training (underestimating) yet during testing values extend
 355 beyond 2, indicating the mean flow in many catchments is overestimated by both model types on the testing data. The
 356 observable correlation in testing period bias of the mean between WAPABA and LSTM indicates that this error is not
 357 specific to model type. Correlation between simulations and observed data is similar for both model types and remains
 358 relatively constant between training and testing periods (right column).



359



360

361 **Figure 5: KGE decomposition into three components: bias of the mean, bias of variability, and correlation. Each dot represents an**
 362 **individual catchment (large outliers have been omitted for visualization purposes.) The mean flow and variability (left and middle**

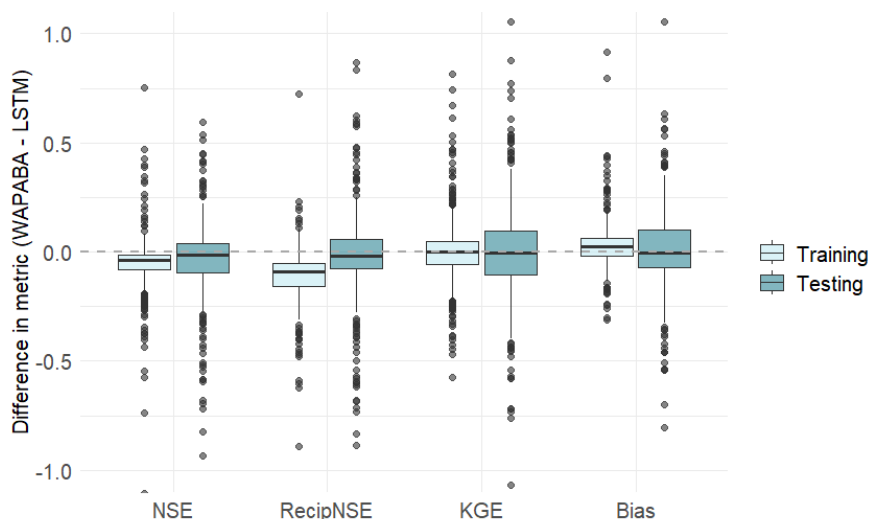
363 columns) tend to be underestimated during training and both under- and overestimated during testing by both model types. The
364 correlation (right column) remains similar during training and testing.

365

366 Performance differences at individual catchments

367 The differences between WAPABA and LSTM performance at each catchment (eg. $NSE_i = NSE_{i,WAPABA} - NSE_{i,LSTM}$ for
368 catchment i) are summarised in Figure 6. Values above zero indicate higher metrics obtained by WAPABA, and values below
369 zero indicate higher metrics obtained by the LSTM model at a specific catchment.

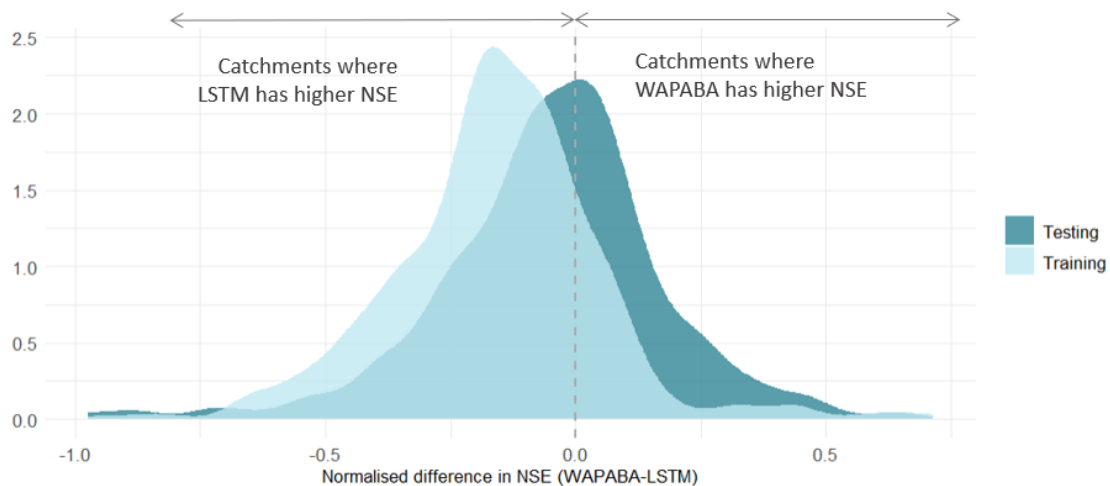
370 The boxplots indicate that median differences in WAPABA and LSTM prediction performance at each catchment (measured
371 by NSE, Reciprocal NSE, KGE and Bias on the testing data) are very close to zero. However, there are outliers (black dots)
372 representing large performance differences between WAPABA and LSTM models, both positive and negative. These indicate
373 that each model provides advantages for predicting runoff in certain catchments. In this figure the boxplots are restricted to
374 the range [-1,1] for visualisation purposes. A version of this figure including the large outliers is provided in Figure A1 of
375 the Appendix.



376

377 **Figure 6: Difference in the metrics (WAPABA – LSTM) for each catchment. A positive value indicates WAPABA has a higher**
378 **metric for that catchment, and a negative value indicates LSTM has a higher metric. The median difference in each metric lies**
379 **close to zero for the testing portion of the dataset, signifying overall similarity in catchment-specific metrics between model types.**
380 **Large negative outliers have been excluded from this figure for visualisation purposes, but are included in the reproduction in the**
381 **Appendix.**

382 This data set represents a range of catchments across Australia, some being characterised by highly arid conditions. To enable
383 comparisons between these diverse catchments, the impact of large negative NSE values which can occur at very dry
384 catchments is minimised by calculating the normalised differences in NSE between the WAPABA and LSTM predictions at
385 each catchment, as per Equation 9. The normalised differences fall into the range [-1,1], facilitating comparison. This
386 distribution is shown in Figure 7 for the 496 catchments. The portion of the distribution lying to the right of the vertical
387 dashed line corresponds to catchments with better prediction by WAPABA and catchments to the left have better prediction
388 by LSTM. The x-axis corresponds to percentage differences between the sum of squared errors of the two model types (ie. -
389 0.5 indicates a 50% performance gain by LSTM and 0.5 indicates a 50% performance gain by WAPABA).



390

391

392

393

394

Figure 7: Distribution of normalized differences between WAPABA and LSTM prediction performance at individual catchments (measured by NSE). The values on the x-axis represent percentage/100 difference in sum of squared errors between WAPABA and LSTM at the same catchment (ie 0.5 -> 50% difference in sum of squared errors). The catchments under the curve on the right of the dashed line have better predictions by the WAPABA model and on the left by the LSTM model.

395

396

397

398

399

In Figure 7, it can be seen that during the training period the majority of catchments are to the left of the line indicating better prediction by LSTM, and in the testing period there is a more even split. The median normalised difference in NSE across the 496 catchments over the training period is -0.15 (mean -0.16) and -0.04 (mean -0.05) during the testing period. This equates to a median 15% performance advantage by LSTM versus WAPABA during training and 4% during testing based on sum of squared errors.

400

401

402

This figure suggests that in general there is little overall advantage of either the WAPABA or LSTM models when predicting on unseen data across the whole sample of catchments. However, the width of the distribution indicates that both the WAPABA and LSTM models have advantages at certain individual catchments, which will be explored in the next section.

403

404

405

406

Figure 8 quantifies the proportion of catchments with similar or better prediction performance by either WAPABA or LSTM (on the testing data). ‘Similarity’ is defined here as an absolute normalised difference in NSE of less than 0.05 between WAPABA and LSTM predictions, meaning the sum of squared errors of the WAPABA and LSTM models at an individual catchment differ by no more than 5%.

407

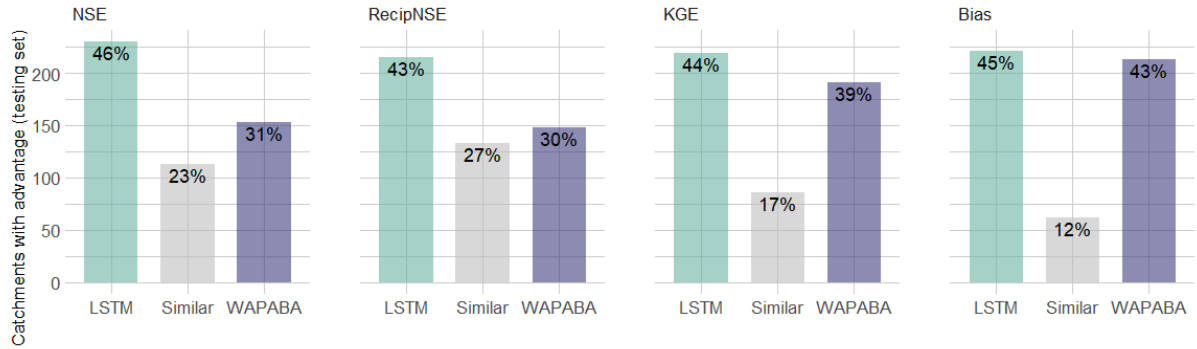
408

409

410

411

The LSTM models produce similar or higher NSE values for 69% of the catchments when tested on data not seen during the training process (and 89% of the catchments during training, not shown). It can also be seen that 70% of catchments have similar or higher reciprocal NSE (focusing on low flow predictions) with the LSTM, 61% have similar or higher KGE with the LSTM (higher being preferable), and 57% have similar or lower Bias (lower being preferable) with the LSTM model compared to WAPABA on the same catchment.



413

414 **Figure 8: Percentage of catchments with similar or better performance metrics on the testing portion of the data (note better Bias**
 415 **is lower, all others is higher). For catchments in the ‘similar’ category, the sum of squared errors of the WAPABA and LSTM**
 416 **predictions differ by less than 5%. The LSTM model produces predictions with similar or higher NSE values compared to the**
 417 **WAPABA predictions for 69% of the catchments.**

418

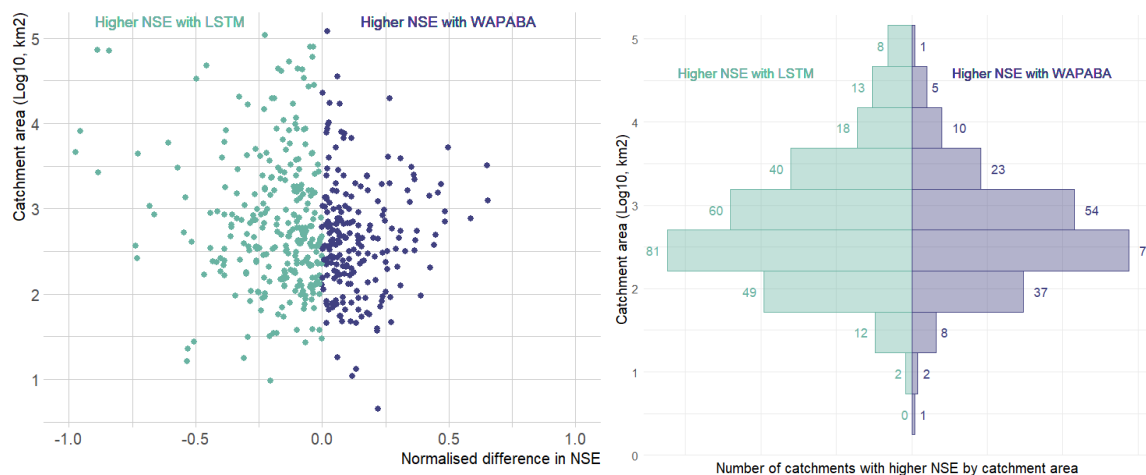
419 **Prediction performance comparison by catchment or time series characteristics**

420 In this section, it is investigated if the abilities of WAPABA and LSTM to accurately predict runoff at individual catchments
 421 vary based on attributes such as catchment area, flow level and length of historical record.

422 *Catchment size*

423 Figure 9 shows the association of prediction performance with catchment area. The left panel shows the catchment area
 424 compared to the normalised difference in NSE between LSTM and WAPABA prediction performance for each catchment.
 425 Data points are coloured according to the model that produced the better prediction for that catchment. This figure indicates
 426 the largest performance gains of LSTM versus WAPABA occurred in large catchments (points furthest to the left are found
 427 in the upper portion of the plot). Splitting the catchments into quintiles by area, the results can be analysed for the largest
 428 20% of catchments. Of these catchments, over three-quarters (78%) had similar or better runoff predictions with the LSTM
 429 (with similarity defined as less than 5% difference in sum of squared errors compared to WAPABA predictions). In this top
 430 quintile of catchments by area, those with higher NSE values from the LSTM show a greater average advantage (average
 431 24% lower sum of squared errors, maximum 97% lower), than those with better WAPABA predictions (average 15% lower
 432 sum of squared errors, maximum 65% lower).

433 The mirrored histogram in the right panel of Figure 8 shows catchments stratified into bins by area (log base 10), coloured
 434 and counted by the model type that produced the better runoff prediction at each catchment. The LSTM models produced
 435 higher NSEs for a greater number of catchments than the WAPABA models in all of the bins, except the lowest bin (where
 436 n=1).



437

438 **Figure 9: Model performance by catchment size. Left panel: Each data point represents the normalized difference in prediction**
 439 **performance at an individual catchment, arranged by catchment size. The spread of data points in the top left quadrant indicates**
 440 **that in large catchments the performance gain of LSTM over WAPABA can exceed 90% in terms of sum of squared errors. Right**
 441 **panel: count of catchments in each size bin that have better performance with each model.**

442

443 *Flow level*

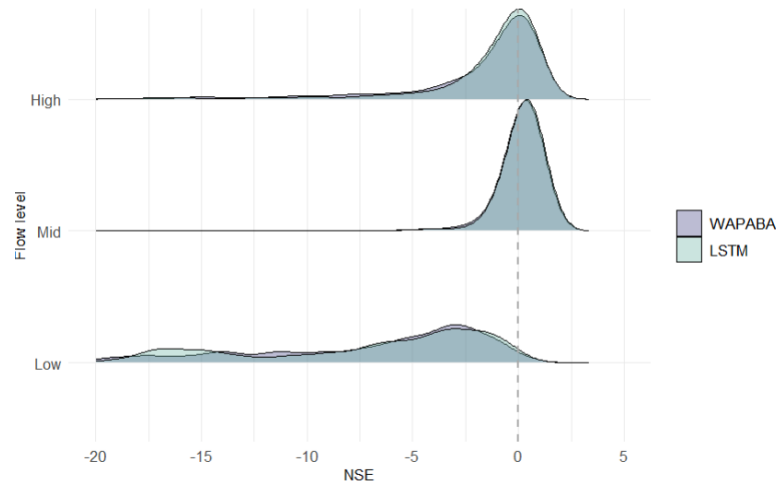
444 Model performance is compared for high, medium and low flow portions of the time series. For each station, each observation
 445 is categorised based on its flow level. High flows are defined here as the top 5% of flow values and low flows as the lower
 446 10% of flows at each station (calculated excluding zeros) over all observed data during the study period. The training and
 447 testing portions of the time series over all the catchments have different distributions of flow levels, as listed in Table 4.
 448 During the testing portion of the study period, conditions are dryer with more no-flow and low-flow observations, and fewer
 449 medium- and high-flow observations than during training.

450 Table 4: Distribution of flow levels during training and testing

Flow level	Training observations (n)	Testing observations (n)
No flow	18,728	21,690
Low	11,967	14,668
Medium	127,584	96,089
High	9,192	4,203

451 For comparison purposes in this section, the raw observed and modelled flow data are standardized by station based on the
 452 mean and standard deviation of all observations at that station during the study period. The observed mean is subtracted from
 453 each value before dividing by the standard deviation of the observations, allowing for basins with a range of flow volumes
 454 to be compared.

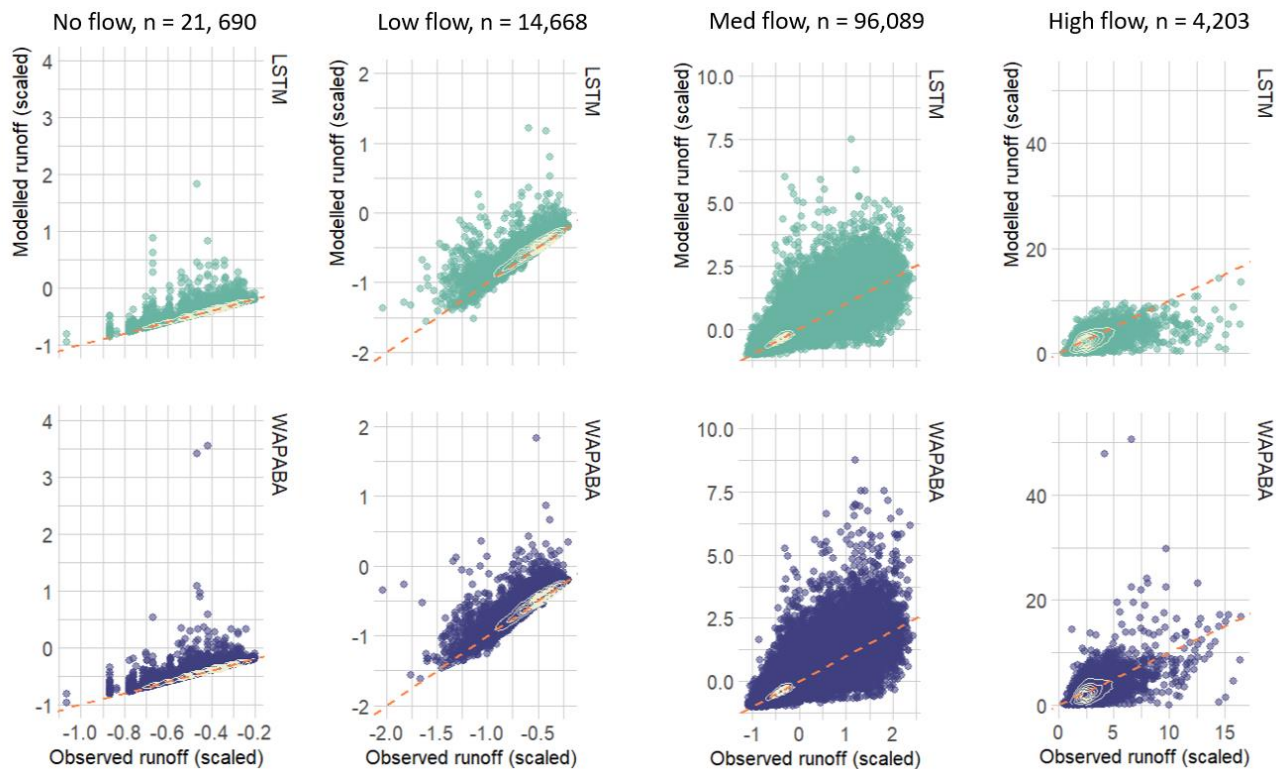
455 Figure 10 shows that when NSE is calculated separately for the low, medium and high flow measurements at each catchment,
 456 both model types have similar NSE distributions. Medium flows are better predicted (NSE peak closer to 1) than high flows,
 457 and low flows appear to be poorly represented by both WAPABA and the LSTM.



458

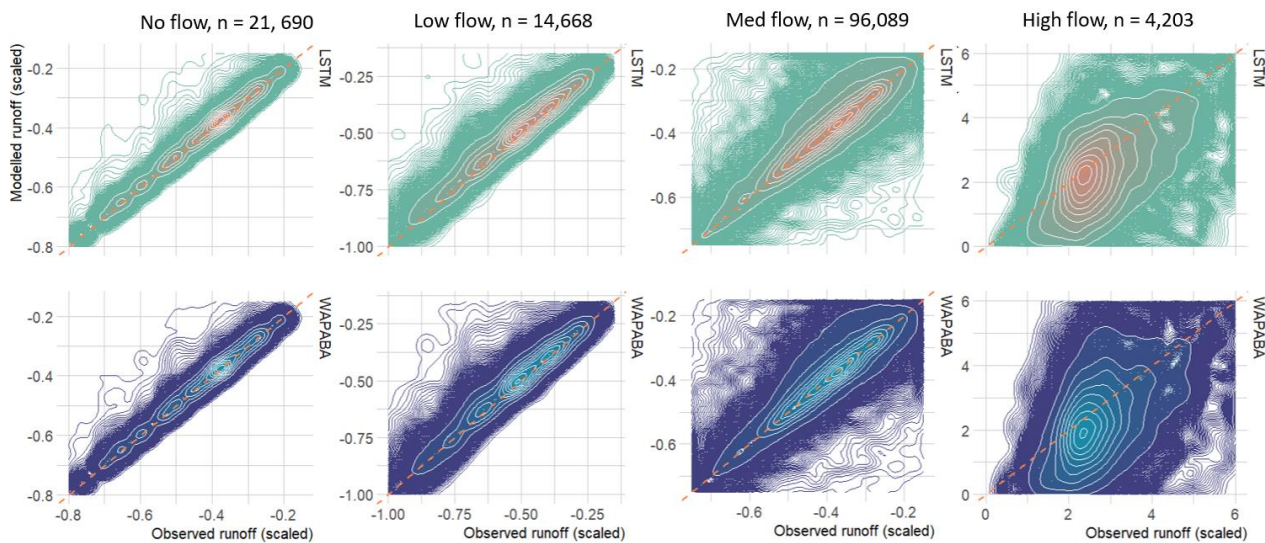
459 **Figure 10: NSE distributions calculated separately by flow level over all catchments. Both model types have similar distributions**
 460 **of NSE by flow. Medium flows are best represented, followed by high and then low flows.**

461 Figure 11 compares the standardized modelled flow to the standardized observations for all testing observations at all stations.
 462 Kernel density contours split the data into 10 density regions on each plot and a 1:1 line is added to aid interpretation. The
 463 lower panel focuses on the regions of highest density for each subset of flows. Note that the standardization procedure used
 464 in this section leads to standardized ‘no-flow’ data points that do not fall exactly on zero in the plot even though the raw flow
 465 values at these points are zero. For no-flows and low flows (left two panels), the densest portions of the observation/prediction
 466 clouds are closely aligned along the 1:1 line indicating similar predictions obtained with both models. The magnitude of the
 467 outliers (beyond the outermost contour) is greatest above the 1:1 line indicating that prediction errors for no-flows and low
 468 flows are dominated by overestimations. For medium flow levels, the contours again follow the 1:1 line. The contours tend
 469 to expand upwards as flow size increases, indicating a tendency towards more overestimation with higher flows. The shape
 470 of the contours is similar for both models. On the upper panel it can be seen that the edges of the data cloud expand upwards
 471 and outwards as the flows increase. The medium flow prediction errors with largest magnitude tend to be overestimations,
 472 with the WAPABAs producing greater overestimations than the LSTMs on the higher flows (still in this medium-flow
 473 subset). For high flows (on the far right panel), the majority tend to be underestimated by both LSTM and WAPABA (central
 474 density located below the 1:1 line), though there is a difference in the outliers – most of the larger errors in LSTM high flow
 475 predictions are underestimations, whereas the high-magnitude WAPABA errors are both over- and underestimations of high
 476 flows.



477

478



479

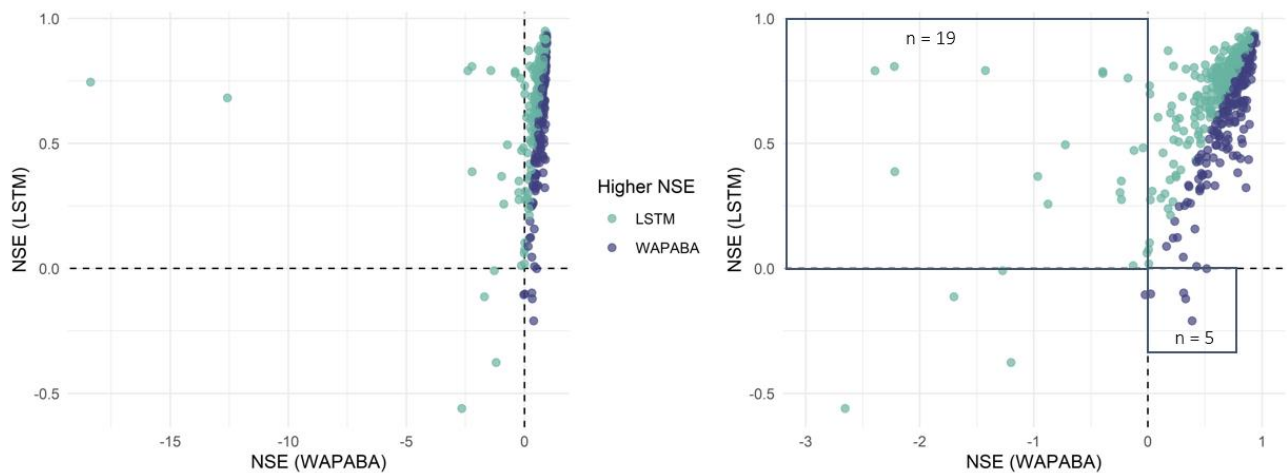
480 **Figure 11: Prediction performance related to flow level. Upper panel: Observed vs. modelled flow pairs (normalized data) at all**
 481 **stations, separated into no-flow, low, medium and high flows (testing data only). The densest portion of the data cloud is identified**
 482 **with density contours. Note that the data have been standardized based on observed mean and standard deviation leading to non-**
 483 **zero values in the ‘no-flow’ category. Lower panel: Comparison of density distributions of the data, zoomed in on the kernel density**
 484 **contours. In general, the largest errors on medium flows tend to be overestimations (by both models) and on high flows tend to be**
 485 **underestimations (by both models) or overestimations (by WAPABA).**

486

487 *Poorly predicted catchments*

488 Figure 12 compares the NSEs for WAPABA and LSTM runoff predictions by catchment. Each dot represents an individual
 489 catchment, coloured according to the model with higher NSE at that catchment. The top left quadrant contains catchments

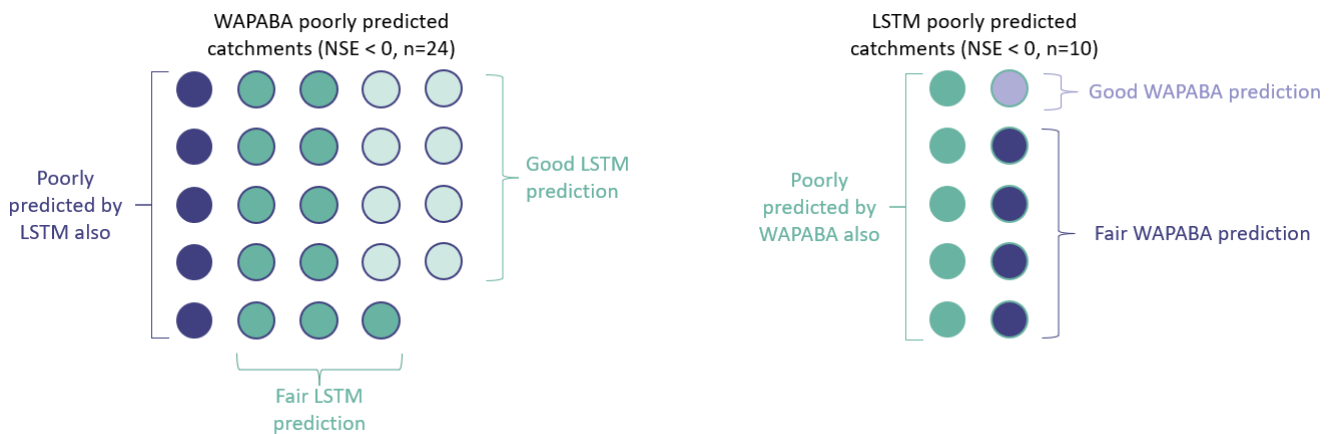
490 where $NSE_{WAPABA} < 0$ and $NSE_{LSTM} > 0$ ($n=19$), and the lower right quadrant contains catchments where $NSE_{LSTM} < 0$ and
 491 $NSE_{WAPABA} > 0$ ($n=5$).



492
 493 **Figure 12: Comparison of NSEs on testing data - each data point represents the WAPABA and LSTM values of NSE for a single**
 494 **catchment, coloured by the model which provides the best prediction at that catchment. On the right panel, two far-left outliers**
 495 **have been removed to enable better viewing of the other datapoints. Catchments in the upper left quadrant are those in which**
 496 **runoff is poorly predicted by WAPABA ($NSE < 0$) and better predictions ($NSE > 0$) are obtained with LSTM. The lower right**
 497 **quadrant correspondingly shows catchments in which the NSE values from LSTM are below 0 and WAPABA has better**
 498 **predictions ($NSE > 0$).**

499 WAPABA and LSTM predictions at each catchment are classified into poor ($NSE < 0$), fair ($0 \leq NSE \leq 0.5$) or good (NSE
 500 > 0.5) categories. In this set of catchments, the runoff at 5 catchments is poorly predicted ($NSE < 0$) by both model types
 501 (lower left quadrant of Figure 12). All other catchments are better represented by one model or the other, with either
 502 WAPABA or LSTM producing predictions with NSEs above 0.

503 For the 5% ($n=24$) of overall catchments that are poorly represented by WAPABA ($NSE < 0$), runoff predictions at 23 of
 504 these catchments (96%) are improved with use of the LSTM. In fact, one-third ($n=8$) of these have ‘good’ predictions by the
 505 LSTM ($NSE > 0.5$). Conversely, for the 2% of catchments ($n=10$) that are poorly represented by the LSTM, 60% are improved
 506 with use of WAPABA, and one-tenth ($n=1$) have ‘good’ predictions by WAPABA (in this catchment the LSTM prediction
 507 is on the border of poor and fair ($NSE=0.001$)). Figure 13 depicts the number of catchments poorly represented by each
 508 model and how these specific catchments are represented by the alternate model. For half of the catchments with poor LSTM
 509 predictions, WAPABA does poorly as well; whereas in 79% of the catchments with poor WAPABA predictions, fair or good
 510 predictions were obtained with the LSTM.



511
 512 **Figure 13: Number of catchments with poor runoff predictions by each model type. Colouring indicates the prediction results from**
 513 **the alternate model type. One-third of WAPABA poorly predicted catchments have good predictions with the LSTM. One-tenth**
 514 **of LSTM poorly predicted catchments have good predictions with the WAPABA. Results are denoted as poor ($NSE < 0$), fair ($0 \leq$
 515 **$NSE \leq 0.5$), or good ($NSE > 0.5$).****

516

517 ***Generalising to changing conditions***

518 The ability of a model to generalise outside of the conditions encountered during training is important, especially in the
519 context of a changing climate. A model that is able to make predictions on unseen (testing) data to a comparable performance
520 level as on the training data will provide confidence in making predictions into the future when external conditions are not
521 expected to remain constant. In this data set it is known that conditions differ between the training and testing data, with
522 wetter climate conditions during the training period and a dryer testing period.

523 It was found that 2% (n=11) of WAPABA models struggled with generalising outside of the training period, with ‘good’
524 (NSE > 0.5) runoff predictions during training but ‘very poor’ predictions (NSE < -0.5) during the testing period. The testing
525 predictions for all of these catchments were improved by use of the LSTM, and at 4 of these catchments ‘good’ predictions
526 (NSE > 0.5) were obtained with the LSTMs. Conversely, one LSTM model produced ‘good’ training runoff predictions and
527 ‘very poor’ testing predictions. This catchment was one of the 11 that also had poor generalisation (and ‘very poor’
528 predictions) with the WAPABA.

529

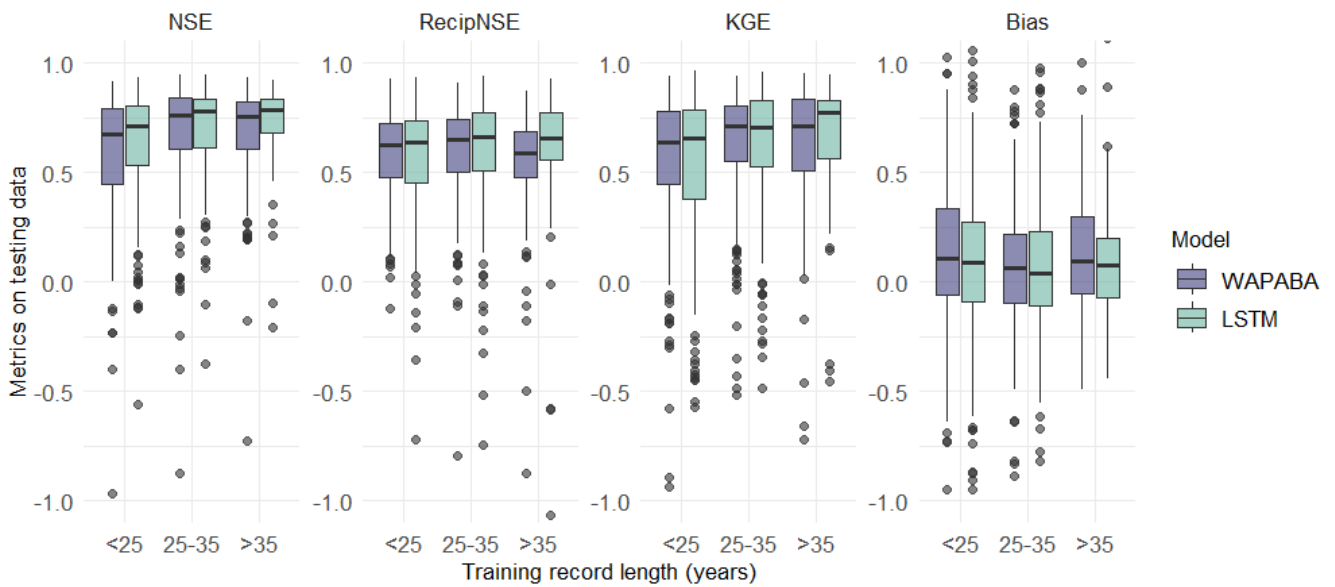
530 ***Historical record length and data set size***

531 The performance of each model type is compared to the length of historical records available at each station. Training data
532 length has been categorized here as 14-25 years (38% of stations), 25-35 years (40%), and 35-47 years (23%).

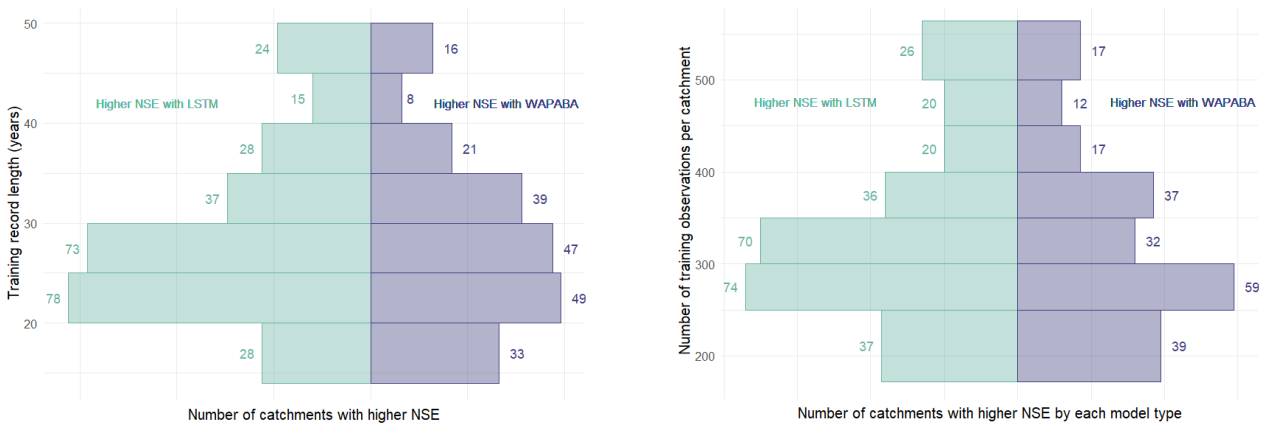
533 Figure 14 (top panel) shows prediction performance varying slightly with record length (for visualisation purposes, this figure
534 is shown without large negative outliers – the figure including outliers is provided in Figure A2 of the Appendix). Stations
535 with medium record length tend to have slightly better predictions according to the four metrics than those with shorter
536 records. The performance levels tend to even out as record lengths increase beyond 35 years, and there is even a slight decline
537 in the WAPABA reciprocal NSE.

538 Considering catchments individually, the median normalised difference in NSE between WAPABA and LSTM predictions
539 (on testing data) is just slightly below zero for all record lengths: -0.03 (<25 years of record), -0.04 (25-35 years), -0.04 (>35
540 years). This indicates that, in each of the short, medium and long record length categories, at least half of the individual
541 catchments have higher NSEs with the LSTMs.

542 The mirrored histogram in the lower left panel of Figure 14 quantifies the number of catchments within 5-year bins of record
543 length in which runoff is better predicted by the LSTM or by the WAPABA. In six of the eight bins, the majority of
544 catchments are better represented by the LSTM.



545



546

547 **Figure 14: Effect of record length and training data size on prediction performance for each model type. Upper panel: Medians of**
 548 **the NSE and KGE on testing data increase with record length for both WAPABA and LSTM predictions (large negative outliers**
 549 **have been excluded for visualization purposes, but are included in the corresponding figure in the Appendix). Lower left panel:**
 550 **Advantage of each model in 5-year increments of record length based on NSE values. Lower right panel: Advantage of each model**
 551 **based on number of training observations.**

552 Comparing performance based on the number of years of record does not take into account the actual size of the data sets,
 553 since measurement frequency differs at each station. Catchments in this study have between 172 and 564 training data
 554 observations (425-846 including testing data). The lower right panel of Figure 14 shows the number of catchments best
 555 modelled by the WAPABA or LSTM model (determined by higher NSE on the testing data) in relation to the number of
 556 training observations. Median NSE values of both the WAPABA and LSTM predictions increased with increasing number
 557 of training data points (not shown). Of particular note is that runoff at catchments with the smallest data sets (less than 250
 558 training data points) were similarly well predicted by both LSTM (median NSE = 0.67) and WAPABA (median NSE = 0.66).

559

560 4. Discussion

561 When considered over the entire study set of catchments, machine learning models were found to match conceptual model
 562 performance for the majority of catchments. The median NSE of runoff predictions was 0.74 with the WAPABA models and
 563 0.76 with the LSTMs, and the medians of other metrics were similarly aligned. At individual catchments, LSTM runoff
 564 prediction performance was similar to or exceeded WAPABA performance in 69% of the catchments in this study (based on
 565 the NSE metric). The median differences in metrics (NSE, Reciprocal NSE, KGE and Bias) between the model types at

566 individual catchments were close to zero, though the range of differences was wide in both directions suggesting many
567 catchments had noticeable prediction advantages with either the WAPABA or LSTM models.

568 Medium flows were similarly well-represented by both model types, with less accurate predictions for high flows and worse
569 again for low flows. Both WAPABA and LSTM models tend to overestimate low flows, while high flows are noticeably
570 underestimated by LSTM and both over- and underestimated by WAPABA. Across all flow levels, the mean flow is
571 prevalently overestimated during testing for both model types, though slightly more so by WAPABA (higher bias of the
572 mean). This overestimation is expected as the testing period in this study is drier than the training period and it is common
573 to have an overestimation of mean during dry periods ([Vaze et al., 2010](#)). Variability of streamflow tends towards
574 overestimation by WAPABA and underestimation by LSTM.

575 Larger catchments were found to have the potential for greater prediction improvements with the LSTM. This finding
576 supports the work of ([Fluet-Chouinard et al., 2022](#)), who found that deep learning methods compete especially well with
577 traditional models in larger non-regulated rivers where the influence of time lags is significant.

578 Though it is known that, in general, machine learning models benefit from large amounts of training data, it is often not
579 possible to provide large hydrological data sets. In this comparison, shorter training record lengths were not found to affect
580 one model type more than the other; the catchments with the smallest training data sets (less than 250 observations) did not
581 show a distinct prediction advantage with either WAPABA or LSTM (median NSEs of 0.66 and 0.67 respectively).

582 In past studies, traditional models have been found to struggle to make accurate runoff predictions under shifting meteorologic
583 data ([Saft et al., 2016](#)). This is an issue that researchers have noted deep learning models may have the potential to overcome
584 ([Li et al., 2021](#), [Wi and Steinschneider, 2022](#)). In this study, the variation in differences in prediction performance at
585 individual catchments is more evident during the testing portion than the training portion of the time series, implying that the
586 WAPABA and LSTM models may each have advantages or drawbacks for generalising to unseen data on various catchments.
587 It was found that in catchments where the WAPABA models provide good runoff predictions during training but struggle to
588 make accurate predictions on new data, the LSTM provides improved predictions in all cases (ie., for those with testing NSE
589 < 0 with WAPABA, all but one had NSE > 0 with the LSTM). In the opposite case, where the LSTM produced substantially
590 poorer predictions on testing data than training data, these predictions were not outdone by WAPABA. This improvement by
591 the LSTM in predicting beyond conditions experienced during training will become progressively important as climate
592 change continues.

593 Aside from scientific considerations, another important advantage of developing rainfall-runoff models using a machine
594 learning software framework is to easily share them among users and to benefit from software optimisation provided by well-
595 established frameworks such as Tensorflow, Keras, or Pytorch. Better benchmark datasets and centralised repositories will
596 be the key to advancement of machine learning in hydrology ([Nearing et al., 2021](#), [Shen et al., 2021](#)). Initiatives are being
597 made to grow reusable software for applying machine learning in hydrology and to benchmark these against other approaches
598 ([Abbas et al., 2022](#)) and ([Kratzert et al., 2022](#)).

599 **Metrics and models**

600 Certain caveats are acknowledged regarding the metrics and models used here. It is possible that the use of individual metrics
601 to compare predictions along the entire length of the time series may mask any variability in model performance that occurs
602 in subperiods of the time series ([Clark et al., 2021](#), [Mathevet et al., 2020](#)). These limitations were partially addressed by

603 comparing high, medium and low flow periods separately, though there are many other subdivisions of the time series that
604 have not been included in the scope of this study.

605 WAPABA is only one example of a conceptual rainfall-runoff model. There are others that could have been chosen for this
606 analysis, though fewer are suitable for comparisons at a monthly time step than would be the case at the daily time step.
607 Model comparisons in [Wang et al. \(2011\)](#), [Bennett et al. \(2017\)](#) and the subsequent body of work with WAPABA in Australia
608 have established WAPABA as a reasonable benchmark against which to assess the machine learning model performance.

609 Though this study has focused on comparing the LSTM model to the WAPABA, readers may wonder if the more traditional
610 feed-forward neural network (FFNN) may suffice in producing as good results. The FFNN has been used in hydrology for
611 many years to model the relationship between climatic predictors and hydrological responses and many researchers are
612 familiar with this basic neural network structure. However, the FFNN is a static network and does not consider the sequential
613 nature of the input data. Though the six months of lagged predictor variables could be input as separate variables, this requires
614 an increase in the complexity of the training space and is not likely to be the optimal choice for time series data as the
615 cumulative impact of the predictor sequences may not be captured. Many studies have already considered the comparison of
616 FFNNs to LSTMs for rainfall-runoff modelling and have determined the LSTM to provide superior runoff predictions (eg.
617 [Rahimzad et al., 2021](#)). As an experiment, the FFNN has been run on this set of 496 catchments and added to the comparison
618 of overall model performance, shown in Figure A3 of the Appendix. It can be seen that the FFNN leads to lower NSE, KGE,
619 Reciprocal NSE, bias of the mean, bias of variability and correlation values, and therefore provides less accurate estimations
620 of runoff than both the WAPABA and the LSTM. For this reason, the FFNN has not been included in the bulk of this study.

621

622 **Future research directions**

623 Future work may entail an expansion of the architecture and complexity of the LSTM models used here, to determine what
624 advantages could be gained from the use of more sophisticated model setups. This may involve the development of hybrid
625 models blending existing conceptual models with LSTMs, the production of a global LSTM incorporating all of the time
626 series, or a type of transfer learning where a model trained on data from all catchments is fine-tuned on a catchment-by-
627 catchment basis, as in [Kratzert et al. \(2019\)](#).

628 A simple LSTM has been used in this study, with a single layer and no catchment-specific hyperparameter tuning. Through
629 appropriate tuning of the models' architecture and hyperparameters for each catchment, more accurate results could be
630 expected. For example, it is known that the performance of data-driven runoff models is heavily dependent on the amount of
631 lagged data that is used as input ([Jin et al., 2022](#)). In this study, a lag of 6 months has been used for all of the catchments, and
632 as such, only temporal patterns of up to 6 months are captured by the LSTMs used in this paper. Varying the length of lag on
633 a catchment-specific basis may lead to better performance.

634 Opportunities also exist for multiple time series analyses on this set of basins to capture patterns in hydrologic behaviour that
635 surpass the catchment scale. With multiple time series analysis one might expect to see greater benefits in the use of machine
636 learning over traditional hydrologic models, since these large-scale studies present obstacles to traditional modelling due to
637 their greater input data and parameter requirements to accurately describe the physical properties of the catchments ([Nearing
638 et al., 2021](#)). Deep learning models have been found to produce better predictions when trained on multiple rather than
639 individual basins ([Nearing et al., 2021](#)), and it has been noted that the training of LSTMs on large diverse sets of watersheds
640 may help improve the realism of hydrologic projections under climate change ([Wi and Steinschneider, 2022](#)).

641 Another consideration may be hybrid modelling frameworks, which combine aspects of conceptual models with machine
642 learning models. These have the potential to draw benefits from both types of models to produce more interpretable and
643 possibly more physically realistic predictions. By leveraging the particular strengths of each model type, the limitations
644 inherent in each may be reduced. For example, Okkan et al., (2021) embedded machine learning models into the internal
645 structure of a conceptual model, calibrating both the host and source models simultaneously, and found the product
646 outperformed each model individually. Li et al., (2023) replaced a set of internal modules of a physical model with embedded
647 neural networks, leading to improved interpretability as well as predictions that are comparable to pure deep learning (LSTM)
648 predictions. The authors found that replacing any of the internal modules improved performance of the process-based model.
649 In the Australian context, Kapoor et al., (2023) studied the use of deep learning components in the form of LSTMs and
650 convolutional neural networks to represent subprocesses in the GR4J rainfall-runoff conceptual model for a set of over 200
651 basins. It was found the hybrid models outperformed the conceptual model as well as the deep learning models when used
652 separately, and provided improved interpretability, better generalisation and an improvement in prediction performance in
653 arid catchments. In this case of this study, the soil moisture and groundwater recharge outputs derived from the WAPABA
654 model would likely be useful as additional predictors for the LSTM model.

655 The question of catchment-specific circumstances under which the LSTM may provide an advantage to monthly rainfall-
656 runoff modelling has been broached in an elementary fashion here, and a more sophisticated investigation would be warranted
657 in further studies. Investigation of multi-dimensional patterns of catchment or climate characteristics that may be associated
658 with differences in predictive performance between the model types could lead to a greater understanding of the value that
659 LSTMs could add to hydrologic modelling.

660 **5. Conclusion**

661 A continental-scale comparison of conceptual (WAPABA) and machine learning (LSTM) model predictions has been made
662 for monthly rainfall-runoff modelling on 496 diverse catchments across Australia. This large-sample analysis of monthly-
663 timescale models aggregates performance results over a variety of catchment types, flow conditions, and hydrological record
664 lengths.

665 The following conclusions have been found:

- 666 1. The LSTM models match or exceed WAPABA prediction performance at a monthly scale for the majority of
667 catchments (69%) in this study.
- 668 2. Both the WAPABA and LSTM models have advantages at certain individual catchments; whilst the median
669 difference in prediction performance is near zero, the distribution spreads in both directions.
- 670 3. Larger catchments were found to have the potential for greater prediction improvements with the LSTM .
- 671 4. Mean streamflow and streamflow variability tend to be overestimated more by the WAPABA models than the
672 LSTMs.
- 673 5. Both model types predict medium flows better than high or low flows. The majority of high flows were
674 underestimated by both models, however WAPABA also had some tendency towards large over-estimations of high
675 flows that wasn't seen with the LSTMs.
- 676 6. Generalisation to new conditions is found to improve with use of the LSTM. In this data set the testing period was
677 significantly drier than the training period, with implications for making predictions in the context of climate change.
678 At catchments in which WAPABA produced good predictions on the training data but very poor predictions on the
679 testing data, the testing predictions were universally improved with use of the LSTM; the opposite case was not

680 observed (ie., in the one catchment with poor generalisation by the LSTM, this was not improved upon by the
681 WAPABA).
682 7. Catchments with the smallest training data sets (< 250 observations) were similarly well predicted by both model
683 types.

684 It has been shown that similar performance to traditional models is able to be reached despite the LSTM being fit using
685 limited data on single catchments and a basic model setup. With refinement of the LSTM model architecture and
686 hyperparameter tuning specific to each catchment, it may be possible to increase the proportion of catchments for which the
687 LSTM provides good prediction performance. Other benefits may be realised by combining multiple catchments within a
688 global model to capture patterns that transcend catchment boundaries, incorporating hybrid modelling techniques or
689 transferring knowledge from data-rich catchments to data-poor catchments within Australia or from international source
690 catchments.

691

692 **Author contributions**

693 PF and JMP designed the experiment with conceptual inputs from JL and SC. PF and JMP developed the LSTM model code
694 and performed the simulations, as JL performed the WAPABA simulations. SC conducted the comparison and prepared the
695 manuscript with contributions from all co-authors.

696

697 **Competing interests**

698 The authors declare that they have no conflict of interest.

699

700 **Acknowledgments**

701 The authors would like to thank the CSIRO Digital Water and Landscapes initiative for their support and for the funding of
702 this project.

703

704 **Data and code availability**

705 All data used in this paper are accessible through the website of the Australian Bureau of Meteorology. Rainfall and potential
706 evapotranspiration can be downloaded from the Australian Water Outlook portal at the following address:
707 <https://awo.bom.gov.au/>. Streamflow can be downloaded from the Water Data Online portal at the following address:
708 <http://www.bom.gov.au/waterdata/>. Catchment characteristics (e.g. area) can be obtained from the Geofabric dataset
709 available at the following address: <http://www.bom.gov.au/water/geofabric/>. The deep learning source code used in this paper
710 is available at: <https://csiro-hydroinformatics.github.io/monthly-lstm-runoff/> including an overview and instructions for
711 retrieving the source code and setting up batch calibrations on a Linux cluster. The code is made available under a CSIRO
712 open-source software license for research purposes.

References

- 715 ABBAS, A., BOITHIAS, L., PACHEPSKY, Y., KIM, K., CHUN, J. A. & CHO, K. H. 2022. AI4Water v1. 0: an open-source python
716 package for modeling hydrological time series using data-driven methods. *Geoscientific Model Development*, 15,
717 3021-3039.
- 718 BENNETT, J. C., WANG, Q. J., ROBERTSON, D. E., SCHEPEN, A., LI, M. & MICHAEL, K. 2017. Assessment of an ensemble
719 seasonal streamflow forecasting system for Australia. *Hydrology and Earth System Sciences*, 21, 6007-6030.
- 720 CHOI, J., LEE, J. & KIM, S. 2022. Utilization of the Long Short-Term Memory network for predicting streamflow in
721 ungauged basins in Korea. *Ecological Engineering*, 182, 106699.
- 722 CLARK, M. P., VOGEL, R. M., LAMONTAGNE, J. R., MIZUKAMI, N., KNOBEN, W. J., TANG, G., GHARARI, S., FREER, J. E.,
723 WHITFIELD, P. H. & SHOOK, K. R. 2021. The abuse of popular performance metrics in hydrologic modeling. *Water*
724 *Resources Research*, 57, e2020WR029001.
- 725 DUAN, Q., GUPTA, V. K. & SOROOSHIAN, S. 1993. Shuffled complex evolution approach for effective and efficient global
726 minimization. *Journal of optimization theory and applications*, 76, 501-521.
- 727 FLUET-CHOUINARD, E., AEBERHARD, W., SZEKELY, E., ZAPPA, M., BOGNER, K., SENEVIRATNE, S. & GUDMUNDSSON, L.
728 Machine learning-derived predictions of river flow across Switzerland. EGU General Assembly, 2022 Vienna,
729 Austria. Copernicus
- 730 FRAME, J. M., KRATZERT, F., KLOTZ, D., GAUCH, M., SHELEV, G., GILON, O., QUALLS, L. M., GUPTA, H. V. & NEARING, G. S.
731 2022. Deep learning rainfall–runoff predictions of extreme events. *Hydrology and Earth System Sciences*, 26,
732 3377-3392.
- 733 FRAME, J. M., KRATZERT, F., RANEY, A., RAHMAN, M., SALAS, F. R. & NEARING, G. S. 2021. Post-Processing the National
734 Water Model with Long Short-Term Memory Networks for Streamflow Predictions and Model Diagnostics.
735 *JAWRA Journal of the American Water Resources Association*, 57, 885-905.
- 736 FROST, A., RAMCHURN, A. & SMITH, A. 2018. The Australian Landscape Water Balance Model. *Bureau of Meteorology:*
737 *Melbourne, Australia*.
- 738 GOODFELLOW, I., BENGIO, Y., COURVILLE, A. & BENGIO, Y. 2016. *Deep learning*, MIT press Cambridge.
- 739 GUPTA, H. V., KLING, H., YILMAZ, K. K. & MARTINEZ, G. F. 2009. Decomposition of the mean squared error and NSE
740 performance criteria: Implications for improving hydrological modelling. *Journal of Hydrology*, 377, 80-91.
- 741 GUPTA, H. V., PERRIN, C., BLÖSCHL, G., MONTANARI, A., KUMAR, R., CLARK, M. & ANDRÉASSIAN, V. 2014. Large-sample
742 hydrology: a need to balance depth with breadth. *Hydrology and Earth System Sciences*, 18, 463-477.
- 743 HOCHREITER, S. & SCHMIDHUBER, J. 1997. Long short-term memory. *Neural Computation*, 9, 1735-1780.
- 744 HUARD, D. & MAILHOT, A. 2008. Calibration of hydrological model GR2M using Bayesian uncertainty analysis. *Water*
745 *Resources Research*, 44.
- 746 HUGHES, D. 1995. Monthly rainfall-runoff models applied to arid and semiarid catchments for water resource estimation
747 purposes. *Hydrological Sciences Journal*, 40, 751-769.
- 748 JIN, J., ZHANG, Y., HAO, Z., XIA, R., YANG, W., YIN, H. & ZHANG, X. 2022. Benchmarking data-driven rainfall-runoff
749 modeling across 54 catchments in the Yellow River Basin: Overfitting, calibration length, dry frequency. *Journal*
750 *of Hydrology: Regional Studies*, 42, 101119.
- 751 JONES, D. A., WANG, W. & FAWCETT, R. 2009. High-quality spatial climate data-sets for Australia. *Australian*
752 *Meteorological and Oceanographic Journal*, 58, 233.
- 753 KAPOOR, A., et al. 2023. DeepGR4J: A deep learning hybridization approach for conceptual rainfall-runoff
754 modelling. *Environmental Modelling & Software*, 169: 105831.
- 755 KRATZERT, F., GAUCH, M., NEARING, G. & KLOTZ, D. 2022. NeuralHydrology---A Python library for Deep Learning research
756 in hydrology. *Journal of Open Source Software*, 7, 4050.
- 757 KRATZERT, F., KLOTZ, D., BRENNER, C., SCHULZ, K. & HERRNEGGER, M. 2018. Rainfall–runoff modelling using Long Short-
758 Term Memory (LSTM) networks. *Hydrology and Earth System Sciences*, 22, 6005-6022.
- 759 KRATZERT, F., KLOTZ, D., SHALEV, G., KLAMBAUER, G., HOCHREITER, S., and NEARING, G. 2019. Towards learning
760 universal, regional, and local hydrological behaviors via machine learning applied to large-sample datasets,
761 *Hydrology and Earth System Sciences*, 23, 5089–5110
- 762 KRATZERT, F., KLOTZ, D., HERRNEGGER, M., SAMPSON, A. K., HOCHREITER, S. & NEARING, G. S. 2019b. Toward improved
763 predictions in ungauged basins: Exploiting the power of machine learning. *Water Resources Research*, 55, 11344-
764 11354.
- 765 LEE, T., SHIN, J.-Y., KIM, J.-S. & SINGH, V. P. 2020. Stochastic simulation on reproducing long-term memory of
766 hydroclimatological variables using deep learning model. *Journal of Hydrology*, 582, 124540.
- 767 LEES, T., BUECHEL, M., ANDERSON, B., SLATER, L., REECE, S., COXON, G. & DADSON, S. J. 2021. Benchmarking data-driven
768 rainfall–runoff models in Great Britain: a comparison of long short-term memory (LSTM)-based models with four
769 lumped conceptual models. *Hydrology and Earth System Sciences*, 25, 5517-5534.

770 LERAT, J., ANDRÉASSIAN, V., PERRIN, C., VAZE, J., PERRAUD, J.-M., RIBSTEIN, P. & LOUMAGNE, C. 2012. Do internal flow
771 measurements improve the calibration of rainfall-runoff models? *Water Resources Research*, 48.

772 LERAT, J., THYER, M., MCINERNEY, D., KAVETSKI, D., WOLDEMESKEL, F., PICKETT-HEAPS, C., SHIN, D. & FEIKEMA, P. 2020.
773 A robust approach for calibrating a daily rainfall-runoff model to monthly streamflow data. *Journal of Hydrology*,
774 591, 125129.

775 LI, B., et al. 2023 "Enhancing process-based hydrological models with embedded neural networks: A hybrid
776 approach. *Journal of Hydrology*, 625: 130107.

777 LI, W., KIAGHADI, A. & DAWSON, C. 2021. High temporal resolution rainfall-runoff modeling using long-short-term-
778 memory (LSTM) networks. *Neural Computing and Applications*, 33, 1261-1278.

779 MACHADO, F., MINE, M., KAVISKI, E. & FILL, H. 2011. Monthly rainfall-runoff modelling using artificial neural networks.
780 *Hydrological Sciences Journal—Journal des Sciences Hydrologiques*, 56, 349-361.

781 MAJESKE, N., ZHANG, X., SABAJ, M., GONG, L., ZHU, C. & AZAD, A. 2022. Inductive predictions of hydrologic events using a
782 Long Short-Term Memory network and the Soil and Water Assessment Tool. *Environmental Modelling &*
783 *Software*, 152, 105400.

784 MATHEVET, T., GUPTA, H., PERRIN, C., ANDRÉASSIAN, V. & LE MOINE, N. 2020. Assessing the performance and robustness
785 of two conceptual rainfall-runoff models on a worldwide sample of watersheds. *Journal of Hydrology*, 585,
786 124698.

787 MATHEVET, T., MICHEL, C., ANDRÉASSIAN, V. & PERRIN, C. 2006. A bounded version of the Nash-Sutcliffe criterion for
788 better model assessment on large sets of basins. *IAHS PUBLICATION*, 307, 211.

789 MOUELHI, S., MICHEL, C., PERRIN, C. & ANDRÉASSIAN, V. 2006. Stepwise development of a two-parameter monthly water
790 balance model. *Journal of Hydrology*, 318, 200-214.

791 NASH, J. E. & SUTCLIFFE, J. V. 1970. River flow forecasting through conceptual models part I—A discussion of principles.
792 *Journal of Hydrology*, 10, 282-290.

793 NEARING, G. S., KRATZERT, F., SAMPSON, A. K., PELISSIER, C. S., KLOTZ, D., FRAME, J. M., PRIETO, C. & GUPTA, H. V. 2021.
794 What role does hydrological science play in the age of machine learning? *Water Resources Research*, 57,
795 e2020WR028091.

796 OKKAN, U., ERSOY, Z. B., KUMANLIOGLU, A. A., & FISTIKOGLU, O. (2021). Embedding machine learning techniques into a
797 conceptual model to improve monthly runoff simulation: A nested hybrid rainfall-runoff modeling. *Journal of*
798 *Hydrology*, 598, 126433.

799 OUMA, Y. O., CHERUYOT, R. & WACHERA, A. N. 2022. Rainfall and runoff time-series trend analysis using LSTM recurrent
800 neural network and wavelet neural network with satellite-based meteorological data: case study of Nzoia
801 hydrologic basin. *Complex & Intelligent Systems*, 8, 213-236.

802 PAPACHARALAMPOUS, G., TYRALIS, H. & KOUTSOYIANNIS, D. 2019. Comparison of stochastic and machine learning
803 methods for multi-step ahead forecasting of hydrological processes. *Stochastic Environmental Research and Risk*
804 *Assessment*, 33, 481-514.

805 PERRAUD, J.-M., BRIDGART, R., BENNETT, J. C. & ROBERTSON, D. SWIFT2: High performance software for short-medium
806 term ensemble streamflow forecasting research and operations. 21st International Congress on Modelling and
807 Simulation, 2015. 2458-2464.

808 PUSHPALATHA, R., PERRIN, C., LE MOINE, N. & ANDRÉASSIAN, V. 2012. A review of efficiency criteria suitable for
809 evaluating low-flow simulations. *Journal of Hydrology*, 420, 171-182.

810 RAHIMZAD, M., MOGHADDAM N., A., ZOLFONOON, H. et al. 2021. Performance Comparison of an LSTM-based Deep
811 Learning Model versus Conventional Machine Learning Algorithms for Streamflow Forecasting. *Water Resources*
812 *Management* 35, 4167–4187.

813 REICHSTEIN, M., CAMPS-VALLS, G., STEVENS, B., JUNG, M., DENZLER, J. & CARVALHAIS, N. 2019. Deep learning and
814 process understanding for data-driven Earth system science. *Nature*, 566, 195-204.

815 SAFT, M., PEEL, M. C., WESTERN, A. W., PERRAUD, J. M. & ZHANG, L. 2016. Bias in streamflow projections due to climate-
816 induced shifts in catchment response. *Geophysical Research Letters*, 43, 1574-1581.

817 SCHAEFLI, B. & GUPTA, H. V. 2007. Do Nash values have value? *Hydrological processes*, 21, 2075-2080.

818 SHEN, C. 2018. A transdisciplinary review of deep learning research and its relevance for water resources scientists. *Water*
819 *Resources Research*, 54, 8558-8593.

820 SHEN, C., CHEN, X. & LALOY, E. 2021. Broadening the use of machine learning in hydrology. *Frontiers Media SA*.

821 SONG, Y. H., CHUNG, E.-S. & SHAHID, S. 2022. Differences in extremes and uncertainties in future runoff simulations using
822 SWAT and LSTM for SSP scenarios. *Science of The Total Environment*, 156162.

823 VAN DIJK, A. I., BECK, H. E., CROSBIE, R. S., DE JEU, R. A., LIU, Y. Y., PODGER, G. M., TIMBAL, B. & VINEY, N. R. 2013. The
824 Millennium Drought in southeast Australia (2001–2009): Natural and human causes and implications for water
825 resources, ecosystems, economy, and society. *Water Resources Research*, 49, 1040-1057.

826 VAZE, J., POST, D., CHIEW, F., PERRAUD, J.-M., VINEY, N. & TENG, J. 2010. Climate non-stationarity—validity of calibrated
827 rainfall-runoff models for use in climate change studies. *Journal of Hydrology*, 394, 447-457.

828 WANG, Q., PAGANO, T., ZHOU, S., HAPUARACHCHI, H., ZHANG, L. & ROBERTSON, D. 2011. Monthly versus daily water
829 balance models in simulating monthly runoff. *Journal of Hydrology*, 404, 166-175.

830 WANG, Q. J., BENNETT, J. C., ROBERTSON, D. E. & LI, M. 2020. A data censoring approach for predictive error modeling of
831 flow in ephemeral rivers. *Water Resources Research*, 56, e2019WR026128.

832 WI, S. & STEINSCHNEIDER, S. 2022. Assessing the physical realism of deep learning hydrologic model projections under
833 climate change. *Water Resources Research*, e2022WR032123.

834 YOKOO, K., ISHIDA, K., ERCAN, A., TU, T., NAGASATO, T., KIYAMA, M. & AMAGASAKI, M. 2022. Capabilities of deep
835 learning models on learning physical relationships: Case of rainfall-runoff modeling with LSTM. *Science of The*
836 *Total Environment*, 802, 149876.

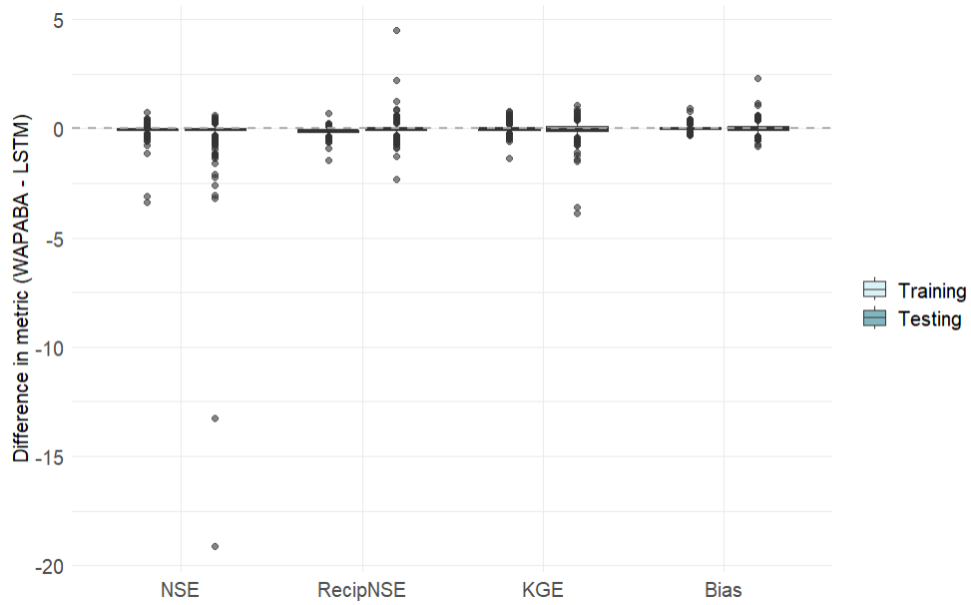
837 YUAN, X., CHEN, C., LEI, X., YUAN, Y. & MUHAMMAD ADNAN, R. 2018. Monthly runoff forecasting based on LSTM–ALO
838 model. *Stochastic Environmental Research and Risk Assessment*, 32, 2199-2212.

839 ZHANG, L., POTTER, N., HICKEL, K., ZHANG, Y. & SHAO, Q. 2008. Water balance modeling over variable time scales based
840 on the Budyko framework–Model development and testing. *Journal of Hydrology*, 360, 117-131.

841

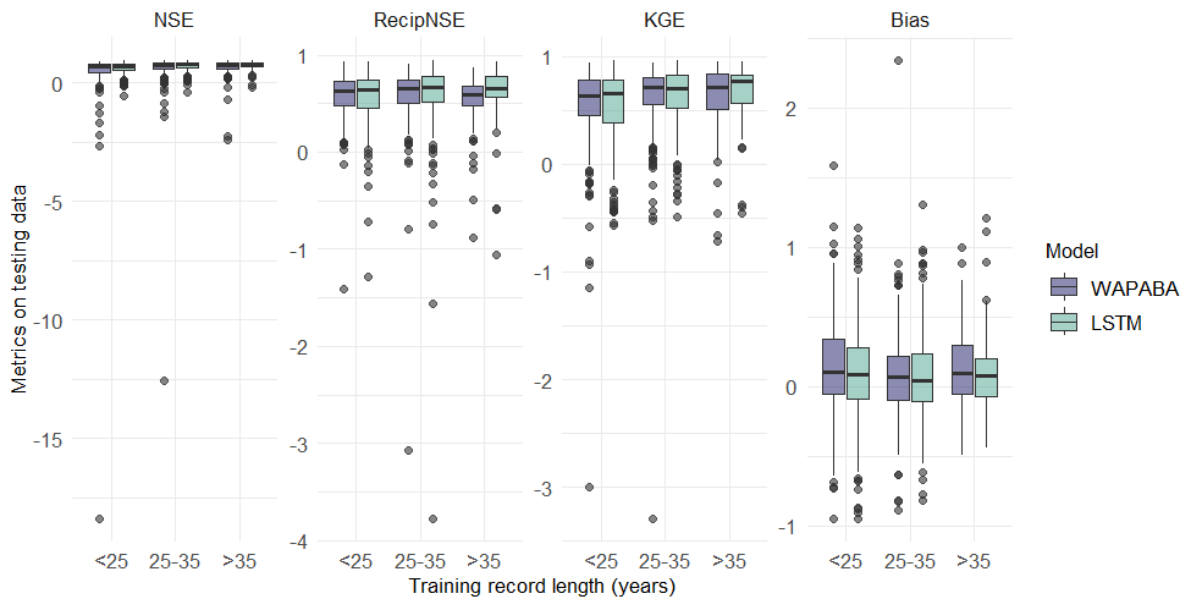
842 **Appendix**

843 Figures A1 and A2 are reproductions of figures in the report in which large outliers detract from a decent visualisation of the
 844 bulk of the data points. Here the entire data set is included, whereas the corresponding figures in the report are shown without
 845 the large outliers.



846

847 **Figure A1: Difference in the metrics (WAPABA – LSTM) for each catchment. A reproduction of Figure 15 that includes outliers.**



848

849 **Figure A2: Effect of record length and training data size on prediction performance for each model type. A reproduction of Figure**
 850 **13 that includes outliers.**

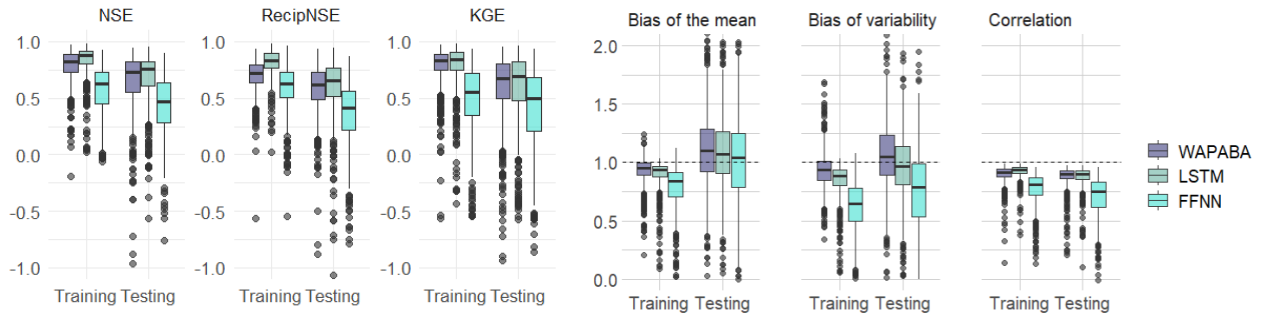
851

852 **Feed-forward Neural Network**

853 To investigate the use of a very simple neural network, the FFNN was run for the 496 catchments. Input variables were the
 854 same as for the LSTM and WAPABA, however 6 months of historical values were included with each training observation.

855 A grid search on 5 random catchments was conducted to select learning rate and batch size. Out of a search space of [8, 16,
856 32] for batchsize and [0.1, 0.01, 0.001, 0.0001] for learning rate, a batch size of 16 and learning rate of 0.01 were chosen.

857 Figure A3 includes the FFNN model results in the comparison with LSTM and WAPABA results. The FFNN values are
858 lower than WABAPA and LSTM indicating poorer runoff predictions over this set of catchments.



859

860 **Figure A3: FFNN metrics compared with WAPABA and LSTM metrics, corresponding to Figures 4 and 5.**

861