

1 **Seasonal forecasting of groundwater levels in principal aquifers of the United**
2 **Kingdom**

3

4 J.D. Mackay^{a*}, C.R. Jackson^a, A. Brookshaw^b, A.A. Scaife^b, J. Cook^c and R.S. Ward^a.

5

6 ^a British Geological Survey, Environmental Science Centre, Keyworth, Nottingham NG12

7 5GG, UK

8 ^b Met Office, FitzRoy Road, Exeter, EX1 3PB, UK

9 ^c ISA Lille, Lille, France.

10 * Corresponding author: Tel.: +44 1159 363100; fax: +44 1159 363200; email:

11 joncka@bgs.ac.uk

12

13 **Abstract**

14 To date, the majority of hydrological forecasting studies have focussed on using medium-
15 range (3 to 15 days) weather forecasts to drive hydrological models and make predictions of
16 future river flows. With recent developments in seasonal (1 to 3 months) weather forecast
17 skill, such as those from the latest version of the UK Met Office global seasonal forecast
18 system (GloSea5), there is now an opportunity to use similar methodologies to forecast
19 groundwater levels in more slowly responding aquifers on seasonal timescales. This study
20 uses seasonal rainfall forecasts and a lumped groundwater model to simulate groundwater
21 levels at 21 locations in the United Kingdom up to three months into the future. The results
22 indicate that the forecasts have skill; outperforming a persistence forecast and
23 demonstrating reliability, resolution and discrimination. However, there is currently little to

24 gain from using seasonal rainfall forecasts over using site climatology for this type of
25 application. Furthermore, the forecasts are not able to capture extreme groundwater levels,
26 primarily because of inadequacies in the driving rainfall forecasts. The findings also show
27 that the origin of forecast skill, be it from the meteorological input, groundwater model or
28 initial condition, is site specific and related to the groundwater response characteristics to
29 rainfall and antecedent hydro-meteorological conditions.

30

31 **Keywords**

32 Seasonal forecasting; ensemble forecasting; groundwater level forecasting; AquiMod;
33 GloSea5.

34

35

36

37

38

39

40

41

42

43

44

45 **1. Introduction**

46 Often a cleaner and more reliable source of drinking water than surface reservoirs,
47 groundwater aquifers comprise the world's largest freshwater resource and provide
48 resilience to climate extremes which may increase in frequency with future climate change
49 (Alley et al., 2002; Mishra and Singh, 2010; Sukhija, 2008). Under prolonged dry climatic
50 conditions groundwater drought can develop, often characterised by significantly low
51 groundwater levels which persist for months to years (Lanen and Peters, 2000; Marsh et al.,
52 2007). This may lead to the drying up of significant water-bearing wells and the degradation
53 of ecologically important rivers and springs. Conversely, lasting wet conditions can induce
54 anomalously high groundwater levels resulting in persistent flooding, potentially at large
55 economic cost (Huntingford et al., 2014; Pinault et al., 2005; Upton and Jackson, 2011).
56 Proper management of these resources is vital to ensure their sustainability and to reduce
57 the risk and impacts from groundwater level extremes.

58
59 One possible way forward is to forecast future groundwater levels so that management
60 strategies can be employed in advance of likely future events. However, these approaches
61 generally require some insight into future weather patterns and an understanding of site-
62 specific hydrogeological characteristics that control the non-linear groundwater discharge
63 response to changes in groundwater levels (Eltahir and Yeh, 1999; Moore and Bell, 1999).
64 This paper attempts to do this by using state-of-the-art seasonal weather forecasts to drive
65 a series of groundwater models to forecast groundwater levels up to three months into the
66 future.

67

68 The majority of groundwater level forecasting studies have been conducted using black-box
69 modelling approaches (Jakeman et al., 2006) whereby an empirical relationship between
70 groundwater level time-series and one or more predictor variables is found using an
71 optimization algorithm (Sahu, 2003). Typically, meteorological covariates, including rainfall
72 and temperature, are used because these perturb groundwater recharge fluxes. Flow
73 through the unsaturated zone and saturated aquifer can slow the response of groundwater
74 level to rainfall events (Alley et al., 2002). Accordingly, a suitable characterisation of this
75 lagged response may be sufficient for forecasting future groundwater levels in aquifers,
76 given up-to-date weather data.

77

78 The most widely used method to characterise the lagged response of groundwater levels to
79 meteorological predictor variables is the non-parametric Artificial Neural Network (ANN), a
80 flexible tool that is able to implement multiple statistical models to replicate patterns in
81 time-series (Maier and Dandy, 2000). Daliakopoulos et al. (2005) used neural networks to
82 forecast monthly groundwater levels in a highly heterogeneous alluvial aquifer in Crete,
83 Greece. Trichakis et al. (2009) also used ANNs to forecast the change in hydraulic head in a
84 complex karstic limestone aquifer in Greece which proved to be accurate up to a 90-day
85 lead time. Taormina et al. (2012) forecast groundwater levels on an hourly time-step for a
86 flashy shallow coastal aquifer in the Venice lagoon and found that they could accurately
87 reproduce groundwater depths for several months ahead. These, along with other studies
88 that have used ANNs (Nourani et al., 2008; Sreekanth et al., 2009; Ying et al., 2014) all show
89 significant forecasting skill months into the future. However, there are two key limitations
90 with these approaches: i) not all aquifers exhibit a significant lagged response to antecedent

91 weather; and ii) to forecast more than one time-step ahead these studies used retrospective
92 observed meteorological predictor variables which would not be available ahead of time.

93

94 Tsanis et al. (2008) recognised the second issue and adapted the work of Daliakopoulos et
95 al. (2005) to include a precipitation projection model which, if used in combination with
96 seasonally averaged temperature data, could simulate groundwater levels up to 30 months
97 ahead, achieving a $R^2 > 0.9$. It should be noted, however, that it is likely that this high
98 correlation score largely reflects the model's ability to capture a downward groundwater
99 level trend induced by steady abstractions in the dry season. Even so, it does demonstrate
100 the possibility of using meteorological forecasts to extend the lead time of real-time
101 groundwater level projections.

102

103 Alternative black box methods such as support vector machines (Behzad et al., 2010;
104 Suryanarayana et al., 2014; Vapnik, 1999; Yoon et al., 2011) and wavelet decompositions
105 (Adamowski and Chan, 2011; Maheswaran and Khosa, 2013; Partal and Kişi, 2007) have also
106 been used for groundwater level forecasting in the past with promising levels of skill.

107 Mendicino et al. (2008) took a different approach by using a simple conceptual distributed
108 water balance model to derive average groundwater storage over the most southern
109 peninsular of Italy, the outputs of which were used to derive a groundwater drought index.

110 They found that due to the persistence of low groundwater levels in the summer months,
111 droughts could be forecast months prior to their occurrence based on model simulations of
112 the current groundwater storage.

113

114 While these studies have shown some skill, the relative infancy of groundwater level
115 forecasting science becomes apparent when compared to the abundance of studies
116 focussed on forecasting other hydrological variables such as river discharge for flood
117 forecasting (see Cloke and Pappenberger (2009) and Cuo et al. (2011) for two
118 comprehensive reviews of these applications). Here, forecasters are not afforded the luxury
119 of long response times to prior weather patterns. At the catchment scale, river flow
120 response time to rainfall is typically of the order of minutes to hours. As such, forecasters
121 drive their hydrological models with medium-range weather forecast products from
122 numerical weather prediction (NWP) centres, which typically offer lead times of 3 to 15
123 days. These extended lead times may allow water resource managers and contingency
124 planners to implement mitigation strategies in advance of extreme events. Of course, the
125 benefit of increased lead time comes at a cost; namely that these meteorological products
126 are inherently uncertain due to the non-linear, chaotic nature of the atmosphere (Lorenz,
127 1963). In response, river flow forecasters now adopt probabilistic methodologies that
128 incorporate this uncertainty rather than relying on a single deterministic forecast. A popular
129 approach that couples probability with determinism is ensemble forecasting (Lewis, 2005)
130 whereby a number of deterministic weather forecasts with differing initial conditions are
131 used to drive the hydrological model. If these realisations are assumed independent and of
132 the same random process, it is possible to assign probabilities to the occurrence or
133 exceedance of given flow thresholds. This probabilistic, ensemble-based approach provides
134 more consistent and skilful outlooks from which users can manage risks more effectively
135 (Addor et al., 2011; Buizza, 2008). One may also cascade other uncertainties, such as those
136 associated with the hydrological model parameterisation, through the forecasting system
137 (Beven, 2006; Pappenberger et al., 2005; Zappa et al., 2010; Zappa et al., 2011). A well

138 established approach for this is the Generalised Likelihood Uncertainty Estimation (GLUE)
139 method (Beven and Binley, 1992; Beven and Binley, 2013), whereby an informal likelihood
140 function is used to weight an ensemble of behavioural models. It should be noted, however,
141 that due to the computational burden, such approaches for real-time hydrological
142 forecasting applications are still not widely used today.

143

144 The response of groundwater levels to rainfall generally operate on longer time scales (days
145 to months) than river flows. As such, strategies to mitigate an imposing groundwater
146 drought, for example, can only be properly formulated with a good understanding of the
147 likely future groundwater levels over a similar time scale. Here, longer-range weather
148 forecasts on the scale of months would be required, like those produced by the latest
149 version of the Met Office global seasonal forecast system (GloSea5) which are now showing
150 increased skill up to a three month lead time (Scaife et al., 2014). To date, however, the
151 majority of seasonal forecasting studies have been undertaken by the river flow forecasting
152 community. Yossef et al. (2012) investigated the potential for forecasting monthly and
153 seasonal river flow extremes in 20 large river basins around the world by driving the global
154 hydrological model, PCR-GLOBWB (Sperna Weiland et al., 2010) with observed
155 meteorological forcing data. They found that they could capture observed flood and
156 drought events given skilful meteorological inputs. More recently, Svensson et al. (2015)
157 used GloSea5 seasonal rainfall forecasts to drive a 1 km resolution water balance model
158 (Bell et al., 2013) and forecast winter (December-January-February) river flows across the
159 UK. The forecasts correlated with observed winter river flows with a median correlation
160 score of 0.45. They also found a clear geographical contrast in the source of predictability
161 whereby the initial condition was the strongest source of predictability in the more

162 permeable, baseflow-dominated catchments of south-east England, and the skill was much
163 more dependent on the meteorological forcing data for the flashy catchments in the north-
164 west of Great Britain. The role of river flow response characteristics on seasonal forecast
165 skill was also found to be important for global seasonal river flow forecasting by Yossef et al.
166 (2013). Indeed, contrasting response characteristics to rainfall can also be found in
167 groundwater level time-series (e.g. see the work of Bloomfield and Marchant, 2013), and
168 these are likely to influence the sensitivity of groundwater level forecasts to the
169 meteorological forcing data.

170

171 To summarise, skilful forecasts of groundwater levels would provide useful information to
172 water resource managers and contingency planners which could help to mitigate hazards
173 such as groundwater flooding and drought, both of which can lead to social, economic and
174 environmental degradation. Experience gained from the river flow forecasting community
175 shows that skilful ensemble hydrological forecasts can be generated using driving data from
176 medium-range NWP models. However, because aquifers generally respond to prevailing
177 weather patterns over a number of months, the insight gained over a 15-day lead time may
178 be small. This has led most studies to rely on the lagged response of groundwater levels to
179 past weather patterns to make forecasts. However, it may be possible to extend the skilful
180 forecast lead time using seasonal weather forecast products to drive groundwater models,
181 an approach that is already showing some skill in river flow forecasting experiments.

182

183 This paper presents a new probabilistic groundwater level forecasting approach that utilises
184 state-of-the-art GloSea5 multi-member seasonal forecasts of rainfall produced by the UK
185 Met Office to drive a series of groundwater models up to three months into the future. A

186 parsimonious lumped conceptual groundwater model, *AquiMod* (Mackay et al., 2014),
187 which simulates groundwater levels at observation boreholes has been used. The models
188 have been calibrated to simulate groundwater level time-series at 21 locations across the
189 UK and in different aquifers with contrasting hydrogeological properties and response
190 characteristics to rainfall. The skill of the groundwater level forecasts is evaluated over the
191 four UK seasons using a 14-year sequence of *GloSea5* rainfall reforecasts. For comparison,
192 reforecasts using rainfall climatology and observed rainfall have also been evaluated.
193 Consideration of the catchment response characteristics and their influence on forecast skill
194 are also made. From these analyses, this study seeks to provide a first evaluation of the
195 potential for national, real-time seasonal groundwater level forecasting.

196 **2. Methodology**

197 **2.1. Study catchments**

198 In total, 21 groundwater catchments, each with an observation borehole and associated
199 groundwater level record were selected for this study from a database of 181 groundwater
200 level time-series held in the National Groundwater Level Archive (Marsh and Hannaford,
201 2008). They were selected because: i) they are situated in unconfined aquifers for which the
202 *AquiMod* groundwater model is best suited; ii) they are located away from any significant
203 groundwater abstractions; and iii) they have continuous monthly groundwater level records
204 that cover the operational 14-year *GloSea5* reforecast period from March 1996 to February
205 2010 (MacLachlan et al., 2014) and at least 15 years of data prior to this for model
206 calibration. The boreholes penetrate into some of the UK's principal aquifers including the
207 Cretaceous Chalk and Lower Greensand, the Jurassic and Magnesian Limestone and the

208 Permo-Triassic Sandstone (Figure 1). Between 16 and 34 years of continuous groundwater
209 level data were available for model calibration.

210

211 Figure 2 shows the raw groundwater level time-series for four of the observation boreholes.
212 Also included are the groundwater level auto-correlation plots and the rainfall-groundwater
213 level cross-correlation plots. It can be seen that groundwater level fluctuations contrast
214 between the catchments. For example, Ashton Farm shows a sinusoidal pattern with
215 relatively uniform amplitude while the New Red Lion borehole shows more variable
216 amplitude with multiple winter peaks. The West Dean cross-correlation plot shows the
217 highest correlation between groundwater and rainfall at a lag of zero, indicating a very rapid
218 and flashy response. This is in contrast to the smooth Heathlanes hydrograph which exhibits
219 relatively small seasonal variability, but more pronounced inter-annual fluctuations. The
220 auto-correlation and cross-correlation plots for this site indicate significant persistence in
221 levels and a much longer response time to rainfall. Also note that because this borehole is in
222 a high storage Sandstone aquifer, annual groundwater levels typically fluctuate by only 0.5
223 m. In contrast, the water table at New Red Lion in the low porosity Jurassic Limestone
224 aquifer can vary by as much as 20 m in one year.

225 **2.2. Aquimod**

226 Aquimod takes monthly rainfall and potential evapotranspiration (PET) driving data and
227 uses conceptual hydrological equations to simulate the downward movement of water
228 through the soil and unsaturated zone and the lateral flow and subsequent discharge of
229 groundwater through the saturated zone (Figure 3). A soil module divides rainfall between
230 evapotranspiration, runoff and soil drainage. The soil drainage is attenuated through the

231 unsaturated zone using a Weibull distribution transfer function, before reaching the
232 saturated zone as groundwater recharge. Discharge from the saturated zone is calculated
233 using a Darcy flux equation. The reader is referred to Mackay et al. (2014) for a more
234 comprehensive description of the underlying theory and model code.

235

236 The Aquimod code was chosen for this study because it was designed specifically for
237 simulating groundwater levels at observation boreholes. It includes in built Monte Carlo
238 parameter sampling, has a small computational burden and also allows the user to
239 incorporate different saturated zone model structures with variable levels of complexity.
240 Mackay et al. (2014) showed that this model can efficiently capture the non-linear
241 groundwater level dynamics in a range of hydrogeological settings. They also showed that a
242 two or three layer aquifer representation is generally most efficient and these structures
243 have been adopted in this study (Figure 3).

244 **2.3. Model calibration**

245 The models were driven with observed monthly rainfall and PET data and calibrated against
246 observed groundwater levels prior to the reforecast period. Rainfall data were obtained
247 from the national 5 km gridded dataset held by the UK Met Office National Climate
248 Information Centre (Perry et al., 2009). This is comprised of rain gauge data interpolated
249 onto a regular grid using inverse-distance weighting. PET data were extracted from the Met
250 Office Rainfall and Evapotranspiration Calculation System (MORECS) dataset (Field, 1983)
251 which uses synoptic station data in conjunction with a modified version of the Penman-
252 Monteith equation to determine the monthly average PET rate on a 40 km grid of over the
253 UK (Monteith and Unsworth, 2008).

254

255 Following the methodology of Mackay et al. (2014), eight of the possible 16 model
256 parameters were fixed based on known catchment characteristics while the remaining were
257 used as calibration parameters (Table 1). A Monte Carlo procedure was used to randomly
258 select 10^6 unique parameter sets from a user-defined parameter space for each model
259 structure. Here we considered the uncertainty in model structure and parameter selection
260 by adopting the GLUE methodology and using the well established Nash-Sutcliffe efficiency
261 (NSE) score (Bennett et al., 2013; Nash and Sutcliffe, 1970) as the informal likelihood
262 measure. Only those models that exceeded a NSE score of 0.5 were deemed behavioural.
263 Those that did not achieve this were assigned a likelihood of zero.

264

265 Between 1780 and 2470 behavioural models were obtained for the 21 study catchments
266 achieving a maximum efficiency (NSE_{max}) between 0.71 and 0.94 and a containment ratio
267 (CR) (Xiong and O'Connor, 2008), which specifies the percentage of observations captured
268 within specified upper and lower prediction bounds, between 65.3 and 97.8% when using
269 the GLUE 95% confidence interval (Table 2). Figure 4 shows the observed and simulated
270 groundwater levels with the GLUE 95% prediction bounds for Bussels, where AquiMod
271 achieved the highest NSE_{max} , and Therfield Rectory, where AquiMod scored the lowest
272 NSE_{max} . It can be seen that the GLUE prediction bounds for Bussels contain almost all of the
273 observations and the best model closely replicates the timing and seasonality in the
274 hydrograph including the pronounced 1976 drought period. For Therfield Rectory, AquiMod
275 captures the timing and seasonality of the hydrograph and most of the observations during
276 the drought of 1973. However, it fails to capture the rapid recession in 1992, and some of

277 the peak levels observed in 1961, 1979 and 1988. These deficiencies are considered in more
278 detail in the discussion.

279 **2.4. Reforecast climate data**

280 Monthly rainfall inputs for the 14-year reforecast period were taken from the GloSea5
281 model. These comprised four reforecasts per year representing the four seasons: i) spring
282 March-April-May (MAM); ii) summer June-July-August (JJA); iii) autumn September-October-
283 November (SON); iv) winter December-January-February (DJF). Each consisted of an
284 ensemble of one, two and three month ahead rainfall, averaged over the UK. The GloSea5
285 winter and summer forecasts were made up of 24 ensemble members while the spring and
286 autumn forecasts were made up of 12 ensemble members. All were downscaled to the
287 catchment scale using linear models defined by ordinary least squares regression between
288 observed catchment rainfall and observed UK average rainfall. Figure 5a and Figure 5b show
289 the relationship between seasonal UK average and seasonal catchment rainfall for the
290 Ashton Farm and New Red Lion observation boreholes. The fitted linear regression models
291 are shown by the solid black lines. It can be seen that for Ashton Farm, a linear
292 approximation of the scale relationship is satisfactory, giving an R^2 score of 0.51, while for
293 New Red Lion, this approximation is less adequate, where the model only explains 31% of
294 the variance. In general, however, the linear regression models demonstrated a good fit,
295 with a mean R^2 score of 0.46 across the study catchments. These models were then used to
296 downscale the GloSea5 forecasts of UK average rainfall for each catchment. The downscaled
297 GloSea5 seasonal rainfall forecasts for Ashton Farm showed the most skill, where the
298 ensemble mean seasonal rainfall correlated with the observed catchment rainfall with an R^2
299 of 0.44 (Figure 5c). In contrast, the downscaled GloSea5 seasonal rainfall forecasts for New

300 Red Lion showed negligible correlation with the observed catchment rainfall (Figure 5d).

301 Overall, the skill of the downscaled rainfall forecasts was low with a mean R^2 of 0.19 across
302 the study catchments.

303

304 For each seasonal reforecast at a given location, the population of behavioural models were
305 run for two years using observed rainfall and PET data to initialise the soil and unsaturated
306 zone in the models. Their initial groundwater levels were fixed to the latest observation. The
307 models were then run for a further three months using the rainfall and PET data described
308 above, producing an ensemble of predictions with $n*m$ members, where n is the number of
309 behavioural models, and m is the rainfall ensemble size. The predicted groundwater level
310 probability density function was then constructed using the predefined GLUE likelihoods
311 and assuming equal probability of occurrence for each rainfall ensemble member.

312 **2.5. Skill analysis**

313 When evaluating forecast skill, it is often useful to establish categorical events for which the
314 observed and forecast frequencies can be compared. Here, three categorical events were
315 chosen for each catchment; below, near and above normal groundwater levels, defined by
316 monthly terciles from the observed groundwater level data. Jolliffe and Stephenson (2012)
317 detail a vast number of forecast verification metrics. We have chosen to use four
318 quantitative metrics which assess different aspects of forecast skill for a given categorical
319 event including:

320

- 321 1. **Frequency bias:** The ratio of the total number of forecast occurrences to the total
322 number of observed events. Here, the forecast event was defined as that which had
323 the highest forecast probability.
- 324 2. **Reliability:** The consistency between the forecast probabilities and the observed
325 relative frequencies. Here, a negatively oriented reliability score derived from the
326 decomposition of the brier score (Murphy, 1973) has been used.
- 327 3. **Relative operating characteristic (ROC) score:** This measures the capacity to
328 correctly discriminate between the occurrence and non-occurrence of an event. A
329 value greater than 0.5 indicates that the hit rate exceeds the false alarm rate.
- 330 4. **Continuous ranked probability score (CRPS):** Calculated as the integrated square
331 difference between the cumulative distributions of the forecasts and observations.
332 This is a probabilistic generalisation of the mean absolute error.

333

334 We chose to convert the CRPS into a skill score, the CRPSS, by comparing the groundwater
335 level forecasts to a reference persistence forecast. A persistence-type benchmark was
336 deemed the most rigorous test given that hydrogeological memory can serve as a potential
337 source of skill. We evaluated three different benchmarks against historical observed
338 groundwater levels including i) persisting the latest observed groundwater level; ii)
339 perturbing the latest observed groundwater level using the monthly mean changes in
340 groundwater levels taken from historical data; and iii) persisting the percentile location of
341 the initial groundwater level in the distribution of historical groundwater levels for that
342 month over the following three months (i.e. if the initial condition was near-normal, the
343 forecast for the subsequent three months would remain in this category). We found that

344 the third approach was the best, consistently outperforming the other two benchmarks and
345 so this was deemed the most rigorous test for forecast skill.

346

347 To complement the benchmark tests, the groundwater models have also been driven with
348 two other meteorological inputs including: i) an unskilful rainfall ensemble made up of re-
349 sampled observed catchment data; and ii) a best case deterministic rainfall input using
350 observed data.

351 **3. Results**

352 It is known that groundwater levels respond to rainfall differently between the catchments.
353 It is therefore likely that the models will also respond differently. This is examined in the
354 first part of the results by undertaking a sensitivity analysis of the models. The results from
355 this are used to organise the models into a number of response type groups. Note here, and
356 in the text that follows, the term model refers to the population of behavioural models for a
357 given catchment rather than a single model realisation. The remainder of this section
358 analyses the skill of the forecasts for each of the response groups, first by using the skill
359 metrics outlined above and then by analysing a selection of forecast time-series plots.

360 **3.1. Groundwater level response to rainfall**

361 It is postulated that because of the contrasting response characterises to rainfall across the
362 catchments, the calibrated models will exhibit different sensitivities to rainfall over the
363 three month forecast horizon. Understanding these sensitivities is important because they
364 influence the added value of using seasonal rainfall forecasts to simulate future
365 groundwater levels.

366

367 A relative measure of sensitivity to rainfall has been derived for each of the calibrated
368 models. To do this, each model was spun-up using observed rainfall and PET and then run
369 for three months using six arbitrary synthetic rainfall inputs ranging from 0 to 5 mm d⁻¹. This
370 process was repeated using each of the months in the reforecast sequence as the initial
371 condition. The sensitivity was then calculated for each month as the range of the six
372 groundwater level forecasts, normalised with respect to the model specific yield. This
373 normalisation step accounts for the different storage properties of each model to allow for
374 easier inter-model comparison.

375

376 Figure 6a shows how the model sensitivity to rainfall changes over the reforecast period for
377 one, two and three month simulations for the Rockley observation borehole in the Chalk
378 aquifer. As would be expected, the sensitivity increases with lead time as the influence of
379 the initial condition diminishes, but there is also a seasonal cycle with peak sensitivity during
380 the winter and considerably reduced sensitivity in the summer months. Given that the
381 climate data for the forecasts are fixed, these variations are a result of perturbations in the
382 initial conditions only. This can be explained by the initial soil moisture deficit (SMD)
383 conditions in the model (Figure 6b) which generally develop in the warmer summer months
384 and must be satisfied before recharge (Figure 6d) is initiated. In the winter months the SMD
385 is small and so small changes in the rainfall input can significantly perturb the modelled
386 groundwater level. Despite this, the sensitivity can increase as the soil moisture deficit
387 develops (for example see year 2003 boxed in Figure 6). This behaviour can be attributed to
388 the initial groundwater level condition, which shows to be receding, and the quadratic
389 groundwater discharge response to a unit rise in groundwater head. In other words, as the

390 groundwater level recedes, the discharge response to an influx of recharge is smaller, and so
391 the sensitivity increases.

392

393 Similar seasonal fluctuations in sensitivity were observed for all of the study catchments,
394 but the magnitude varied substantially. The reason for this is likely to be multifaceted, but it
395 can be attributed primarily to the different model response times to rainfall. It is possible to
396 investigate this by considering the calibrated unsaturated zone Weibull distribution transfer
397 function in Aquimod which spreads the flux of water from the soil zone to the water table
398 over a number of months. This transfer function can be evaluated at lags covering the three
399 month forecast horizon to define a model response characteristic, P , which specifies the
400 percentage of modelled effective rainfall that reaches the water table over this period.

401 Figure 7a shows that this value ranges between 20 – 95% and that the relationship between
402 P and the derived model sensitivities can be approximated with an exponential curve ($R^2 =$
403 0.79) that shows that as P increases, the model sensitivity to rainfall also increases. The
404 permeability of each catchment is also likely to influence the sensitivity to rainfall. Indeed, a
405 closer fit is obtained if the model sensitivity is normalised by the catchment baseflow index
406 (BFI) (Figure 7b), taken from Marsh and Hannaford (2008), which defines the proportion of
407 effective rainfall that contributes to groundwater flow. Furthermore, the calibrated model
408 sensitivities also correlate well with an independent inference of the response time to
409 rainfall for each catchment estimated by the peak lead lag correlation (CC_{max}) between
410 observed rainfall and de-seasonalised groundwater levels (Figure 7c), obtaining an R^2 of
411 0.76.

412

413 These findings demonstrate that the more pronounced the AquMod lagging mechanism,
414 the less sensitive the three month simulations are to the choice of rainfall input. A similar
415 relationship between the sensitivity and an independent estimation of the peak
416 groundwater level response time to rainfall for each borehole, further indicates that the
417 catchment response time has a clear influence on the sensitivity, and therefore is also likely
418 to influence the skill of the forecasts. Consequently, the catchments have been split into
419 three equally sized groups representing slowly responding ($3 \leq CC_{\max} \leq 10$), moderately
420 responding ($1 \leq CC_{\max} \leq 2$) and quickly responding ($0 \leq CC_{\max} \leq 1$) groundwater catchments.
421 These are indicated in Figure 7a-c by the circles, squares and triangles respectively and the
422 analyses in the subsequent sections are conducted using this grouping.

423 3.2. Skill metrics

424 For the purpose of this skill analysis the reforecasts have been subdivided into 36 different
425 assessment groups for which an independent assessment of skill has been conducted. These
426 groups comprise the three categorical events, the four seasons and the three groundwater
427 response groups. Figure 8 shows the four skill measurements for all of the assessment
428 groups using the three different climate inputs.

429

430 The frequency bias ranges between 0.61 and 0.5 (Figure 8a). In the summer (JJA), there is a
431 consistent under forecasting of below normal levels which is mainly offset by a positive
432 frequency bias for near normal events. The winter (DJF) forecasts show the opposite
433 pattern, under forecasting above normal events and over forecasting below normal events.
434 The fact that groundwater levels tend to peak in the winter and trough in the summer
435 indicates that there is a tendency for the forecasts to miss the groundwater level extremes.

436 Indeed, on average, the above and below normal events show negative frequency biases of
437 -0.04 and -0.11 respectively while the near normal event category shows a positive
438 frequency bias of 0.13. This deficiency cannot be attributed to the driving rainfall data as all
439 assessment groups demonstrate that they are insensitive to this, except for the quickly
440 responding catchments in winter and autumn (SON) where using the best case observed
441 climate reduces the bias by approximately half. This insensitivity is also apparent in the
442 other skill metrics, indicating that the skill or lack of it stems more from the model than the
443 rainfall input in most situations.

444

445 Generally, the forecasts are more reliable when predicting above and below normal events
446 than near normal events, especially during winter and autumn and during the summer for
447 the quickly responding catchments (Figure 8b). Figure 9 shows the reliability diagrams for
448 the quickly responding catchments in winter. It can be seen that for the above and below
449 normal events the reliability curves follow the line of perfect reliability closely indicating
450 good consistency between the forecast probabilities and observed relative frequencies. In
451 contrast, there is a tendency for the forecasts to predict closer to base rate probabilities
452 (0.33) for the near normal events as indicated by the flat reliability curves which imply a lack
453 of forecast resolution. This is reflected in the ROC scores (Figure 8c) which are smaller on
454 average for the near normal events indicating that the forecasts are less efficient at
455 discriminating these events. Even so, all of the ROC scores obtained were greater than 0.5
456 showing that the number of hits exceeded the number of false alarms. The forecasts were
457 also able to discriminate below normal events with an average ROC score of 0.87 using the
458 downscaled GloSea climate which is particularly encouraging.

459

460 The ROC scores also demonstrate a clear relationship with the catchment response times
461 where the less sensitive, slowly responding catchments have greater discrimination capacity
462 than the quickly responding catchments. However, this again appears to be an artefact of
463 the model skill rather than to do with the sensitivity to the rainfall input. Even so, the use of
464 observed climate consistently improves the discrimination capacity of the forecasts,
465 particularly for the quickly responding catchments where improvements of up to 0.14 are
466 shown.

467

468 From the 36 assessment groups, 35 return a positive CRPSS when using the climatology
469 rainfall input (Figure 8d). This indicates that even climatology yields forecasts that are a
470 better predictor of groundwater levels than a persistence forecast. Slightly fewer (31) of the
471 groups return a positive CRPSS using the observed rainfall and only 30 when using the
472 downscaled GloSea data. All suggest that the forecasts consistently outperform the
473 persistence approach.

474 **3.3. Time-series analysis**

475 Finally, the forecasts have been evaluated over three time periods which contain important
476 historical events including: i) the onset and persistence of below normal levels in 1996 and
477 1997, a period where many parts of the UK experienced groundwater drought; ii) the
478 subsequent transition back to normal levels in 1997 and 1998, broadly associated with the
479 end of the drought; and (iii) the onset and peak of above normal levels in the winter of
480 2000/2001, a period where many boreholes recorded their highest ever levels and where
481 there was widespread groundwater flooding, particularly in the Chalk of south-east England.
482 Figure 10 shows the number of catchments in each response category that successfully

483 forecast each event using the three different climate inputs. A forecast was only deemed a
484 success if all of the observed groundwater levels were contained within the forecast
485 uncertainty bounds as defined by the limits of the ensemble. In addition, Figure 11 displays
486 time-series plots for several of the study catchments over these events which have been
487 used to compare the observed groundwater levels (black dots) against the ensemble mean
488 forecasts (thick dashed lines) using the GloSea5 and observed rainfall inputs. The
489 uncertainty bounds (thin dashed lines) are also shown.

490

491 The forecasts were least effective at capturing the high levels of winter 2000/2001 when
492 using the downscaled GloSea and climatology rainfall, but showed significant improvements
493 when driven with observed data. This is demonstrated in Figure 11a where the observed
494 initial groundwater rise (time steps one to three) at the quickly responding New Red Lion
495 borehole, is only replicated by the ensemble mean forecast when using the observed rainfall
496 input. The forecast using the downscaled GloSea rainfall does not capture this due to
497 underestimating the seasonal rainfall by almost 130 mm. Note that the GloSea forecast is
498 able to capture the peak groundwater level at the fourth time step. This is partly because
499 this corresponds to the one month ahead forecast, and therefore the model was initialised
500 at the above normal levels from the previous time step.

501

502 For the below normal levels of 1996 and 1997, the choice of climate has less impact on the
503 success rate which is demonstrated by the New Red Lion reforecast in Figure 11b. It can be
504 seen here that regardless of the rainfall input, the ensemble mean overestimates the
505 groundwater levels and the uncertainty bounds do not capture the gradual recession of the
506 hydrograph and even extend into the above normal range between time steps five and six.

507 This insensitivity could be explained by the large soil moisture deficit that would likely
508 develop over this period. Therefore, the forecasts are much more reliant on the skill of the
509 model, which in this case does not capture the groundwater discharge and subsequent
510 hydrograph recession adequately.

511

512 Some catchments, such as the quickly responding Bussels catchment (Figure 11c) did
513 demonstrate significant sensitivity to the rainfall input during this drought period. Here, it
514 can be seen that when using the observed rainfall, the ensemble mean follows the
515 persistent low groundwater levels closely, but when using the downscaled GloSea rainfall,
516 the ensemble mean forecast actually predicts a sharp rise in groundwater level almost back
517 to normal conditions (time steps seven to nine) due to the downscaled GloSea forecast
518 overestimating rainfall by 100 mm for this period.

519

520 The highest overall success rates using the downscaled GloSea inputs were recorded for the
521 return to normal levels in 1997 and 1998. For the moderately responding Rockley
522 observation borehole (Figure 11d), it can be seen that the two ensemble mean forecasts
523 using the GloSea and observed rainfall inputs are similar. Furthermore, both capture all of
524 the observations in their uncertainty bounds which was observed for most of the
525 catchments for this period.

526 **4. Discussion**

527 This study has demonstrated that skilful seasonal forecasts of groundwater levels at
528 observation boreholes can be generated by using seasonal weather forecasts to drive
529 parsimonious conceptual groundwater models. The forecasts were proficient at

530 discriminating between below, near and above normal future groundwater levels and they
531 consistently outperformed a reference persistence forecast system. They also demonstrated
532 good reliability, particularly for the seasonal forecasts of spring groundwater levels. These
533 positive attributes have also been demonstrated for the quickly responding catchments,
534 indicating that the skill can extend beyond the peak response time of these groundwater
535 systems.

536

537 The skill of the forecasts originates from a combination of the driving climate data, the
538 groundwater models and the initial groundwater level condition. For those catchments
539 where groundwater levels respond more slowly to rainfall, the groundwater models and the
540 initial conditions have a stronger influence on the forecast skill than the rainfall input.

541 However, there is no clear indication that the sensitivity to the rainfall input directly affects
542 the forecast skill. Rather, the relationship between groundwater level response time to
543 rainfall and forecast skill appears to be primarily controlled by the groundwater model
544 efficiency. Indeed, when conducting this work, we could find no apparent correlation
545 between skill and geographical location like, for example, the work of Svensson et al. (2015).

546 However, we suggest that with a larger sample size of boreholes, and by evaluating the
547 forecast skill at longer lead times, where meteorological driving data plays a more crucial
548 role in the forecast skill, such relationships may become more apparent. Indeed, while all of
549 the response groups demonstrated forecast skill, it remains to be seen at what lead time
550 this skill diminishes.

551

552 The origin of forecast skill also changes as a response to antecedent hydro-meteorological
553 conditions. Here, it was found that when a large soil moisture deficit is developed during the

554 model spin up and initialisation, the subsequent forecasts are less sensitive to the rainfall
555 input. As such, we noted that for the summer forecasts, the skill derives mainly from the
556 groundwater models and their internal hydrogeological memory. This has potential
557 implications because some of the models have shown deficiencies, such as poor
558 representation of the hydrograph recession, which materialised as forecast errors. Some of
559 these deficiencies are likely to result from imperfect model calibration, errors in the
560 meteorological input data and observed groundwater levels, or from inadequacies of the
561 model structure and parameters. In this study, no account of input error was made, but we
562 did acknowledge some of the model uncertainties by using an equifinality of acceptable
563 model structures and parameter sets. Of course, this approach in itself may also propagate
564 forecasting errors. For example, the choice of the NSE as a measure of model likelihood was
565 subjective, and as with any objective function, is subject to undesirable properties that are
566 likely to manifest themselves as modelling errors (Smith et al., 2008). There is also evidence
567 that model appropriateness depends strongly on hydro-climatic conditions (Herman et al.,
568 2013) and that it may be beneficial to develop better suited limits of acceptability which can
569 be relaxed dynamically as a mean to implicitly account for input errors (Liu et al., 2009).

570

571 In contrast to the forecasts issued following dry conditions, during winter, when the soil is
572 generally more saturated, the forecasts are more sensitive to the driving rainfall data, and
573 as such the meteorological forecasts play a more crucial role in the skill of the groundwater
574 level forecasts. The winter forecasts using the downscaled GloSea and climatology rainfall
575 inputs both consistently outperformed the persistence approach, although it should be
576 noted that using the downscaled GloSea5 rainfall data showed no significant improvement
577 over using the site climatology inputs, and in some cases showed to be a worse rainfall

578 predictor. This is perhaps not surprising given that UK rainfall has complex spatio-temporal
579 signatures that make deriving robust downscaling transformations difficult. Certainly, the
580 linear downscaling model employed showed to be inadequate for some sites, and improving
581 this should be a high priority for improving site-specific hydrological forecasts like these.
582 However, it may be possible to improve this using more sophisticated non-linear
583 downscaling and post processing techniques which have shown to be effective for medium-
584 range ensemble streamflow forecasts (Verkade et al., 2013). Further data assimilation could
585 also provide enhancements in skill through dynamic updating of state variables and forecast
586 errors, although to date there is limited evidence that this is useful for seasonal or
587 groundwater level forecasting applications (Liu et al., 2012). There are also other seasonal
588 weather forecasting models which could be used for these types of applications, such as the
589 System 4 from the European Centre for Medium-range Weather Forecasts (ECMWF) which
590 has shown “marginally useful” degrees of reliability over Northern Europe (Weisheimer and
591 Palmer, 2014).

592

593 It is important to note that the interpretation of skill in this study is primarily based on
594 analysis of the verification metrics and by comparing the forecasts to the benchmark results.
595 Pappenberger et al. (2015) compared a range of benchmarks for medium-range river flow
596 forecasting and they note that the best benchmarks are the ones that are hardest to beat.
597 While considerable effort was made to select appropriate benchmarks and avoid reporting
598 “naïve” skill, it should be noted that the persistence benchmark used is less skilful for those
599 boreholes that exhibit significant inter-annual groundwater level fluctuations, and so there
600 is likely to be positive bias in the CRPSS reported for the more slowly responding
601 catchments. Certainly, an equivalent thorough examination of benchmark performance to

602 that of Pappenberger et al. (2015) is also needed for seasonal groundwater level
603 forecasting.
604
605 It is also important to consider that, the verification metrics used in this study only give an
606 average indication of the forecast's ability to reliably discriminate between the occurrence
607 of below, near and above normal levels over the 14-year reforecast sequence. When looking
608 at the extreme 2000/2001 high groundwater level event specifically, only two of the 21
609 groundwater level forecasts were able to capture it within their uncertainty bounds when
610 using the downscaled GloSea and climatology inputs. For the 1996/1997 drought period, the
611 timing of the return to normal conditions could only be predicted when using observed
612 rainfall data. This is an important issue, as it is arguably extreme events like these that, if
613 foreseeable, would provide the most economic, environmental and societal benefit. That of
614 course is not to say that these forecasts are not useful; on the contrary the Environment
615 Agency in England, for example, routinely use measures of aquifer levels relative to normal
616 conditions to inform agricultural communities about future prospects for spray irrigation
617 and this approach can be used to help aid decision making processes for these needs. It
618 does however mean that if we wish to forecast the initiation or end of extreme events on a
619 seasonal time scale at the catchment or borehole resolution, then further enhancements in
620 the skill and the use of seasonal rainfall forecasts are required.

621 **5. Conclusions**

622 Using seasonal weather forecasts to drive 21 conceptual groundwater models, this study
623 has shown that skilful seasonal forecasts of groundwater levels at observation boreholes
624 can be generated up to three months into the future. Site-specific groundwater level

625 response characteristics to rainfall result in contrasting sensitivities to the driving rainfall
626 input across the study catchments. These sensitivities have also shown to be strongly
627 controlled by prevailing weather conditions, where dry conditions tend to result in forecasts
628 that are strongly controlled by the groundwater model, and wet conditions result in
629 forecasts that are much more reliant on good driving rainfall data. This has important
630 implications for where the skill or lack of it derives from, and more importantly, where
631 future improvements can be made. There are clearly issues with correctly forecasting
632 extreme groundwater levels which are primarily due to lack of skill in the driving rainfall
633 data. In particular it is recommended that future work should focus these aspects:

- 634 1. Investigate the best practice for data assimilation, downscaling and post processing
635 of seasonal weather forecasts for hydrological forecast applications.
- 636 2. Compare the use of different seasonal forecast products such as those produced by
637 the ECMWF System 4 model.
- 638 3. Examine the maximum skilful forecast lead time for different aquifers in relation to
639 their response characteristics to rainfall.

640 **6. Acknowledgements**

641 Mackay, Jackson and Ward were supported by core science funds of the Natural
642 Environment Research Council (NERC) British Geological Survey's (BGS) Groundwater
643 Science and Environmental Modelling Directorates. Scaife was supported by the Joint
644 DECC/Defra Met Office Hadley Centre Climate Programme (GA01101) and Brookshaw was
645 supported by the UK Public Weather Service research program. The groundwater level data
646 for this study were taken from the NERC BGS National Groundwater Level Archive. The
647 climate data were made available by the NERC Centre for Ecology and Hydrology, and the

648 UK Met-Office. Mackay, Jackson and Ward publish with the permission of the Executive
649 Director of the British Geological Survey.

650 7. References

- 651 Adamowski, J. & Chan, H. F. 2011. A wavelet neural network conjunction
652 model for groundwater level forecasting. *Journal of Hydrology*, 407, 28-
653 40.
- 654 Addor, N., Jaun, S., Fundel, F. & Zappa, M. 2011. An operational hydrological
655 ensemble prediction system for the city of Zurich (Switzerland): skill,
656 case studies and scenarios. *Hydrol. Earth Syst. Sci.*, 15, 2327-2347.
- 657 Alley, W. M., Healy, R. W., LaBaugh, J. W. & Reilly, T. E. 2002. Flow and Storage
658 in Groundwater Systems. *Science*, 296, 1985-1990.
- 659 Behzad, M., Asghari, K. & Coppola, E. 2010. Comparative Study of SVMs and
660 ANNs in Aquifer Water Level Prediction. *Journal of Computing in Civil
661 Engineering*, 24, 408-413.
- 662 Bell, V. A., Davies, H. N., Kay, A. L., Marsh, T. J., Brookshaw, A. & Jenkins, A.
663 2013. Developing a large-scale water-balance approach to seasonal
664 forecasting: application to the 2012 drought in Britain. *Hydrological
665 Processes*, 27, 3003-3012.
- 666 Bennett, N. D., Croke, B. F. W., Guariso, G., Guillaume, J. H. A., Hamilton, S. H.,
667 Jakeman, A. J., Marsili-Libelli, S., Newham, L. T. H., Norton, J. P., Perrin,
668 C., Pierce, S. A., Robson, B., Seppelt, R., Voinov, A. A., Fath, B. D. &
669 Andreassian, V. 2013. Characterising performance of environmental
670 models. *Environmental Modelling & Software*, 40, 1-20.
- 671 Beven, K. 2006. A manifesto for the equifinality thesis. *Journal of Hydrology*,
672 320, 18-36.
- 673 Beven, K. & Binley, A. 1992. The future of distributed models: Model
674 calibration and uncertainty prediction. *Hydrological Processes*, 6, 279-
675 298.
- 676 Beven, K. & Binley, A. 2013. GLUE: 20 years on. *Hydrological Processes*, n/a-
677 n/a.
- 678 Bloomfield, J. P. & Marchant, B. P. 2013. Analysis of groundwater drought
679 building on the standardised precipitation index approach. *Hydrol. Earth
680 Syst. Sci.*, 17, 4769-4787.
- 681 Boorman, D. B., Hollis, J. M. & Lilly, A. 1995. Report No. 126 Hydrology of soil
682 types: a hydrologically-based classification of the soils of the United
683 Kingdom. Institute of Hydrology, Wallingford, UK.

- 684 Buizza, R. 2008. The value of probabilistic prediction. *Atmospheric Science*
685 *Letters*, 9, 36-42.
- 686 Cloke, H. L. & Pappenberger, F. 2009. Ensemble flood forecasting: A review.
687 *Journal of Hydrology*, 375, 613-626.
- 688 Cuo, L., Pagano, T. C. & Wang, Q. J. 2011. A Review of Quantitative
689 Precipitation Forecasts and Their Use in Short- to Medium-Range
690 Streamflow Forecasting. *Journal of Hydrometeorology*, 12, 713-728.
- 691 Daliakopoulos, I. N., Coulibaly, P. & Tsanis, I. K. 2005. Groundwater level
692 forecasting using artificial neural networks. *Journal of Hydrology*, 309,
693 229-240.
- 694 Eltahir, E. A. B. & Yeh, P. J. F. 1999. On the asymmetric response of aquifer
695 water level to floods and droughts in Illinois. *Water Resources Research*,
696 35, 1199-1217.
- 697 Field, M. 1983. The meteorological office rainfall and evaporation calculation
698 system — MORECS. *Agricultural Water Management*, 6, 297-306.
- 699 Herman, J. D., Reed, P. M. & Wagener, T. 2013. Time-varying sensitivity
700 analysis clarifies the effects of watershed model formulation on model
701 behavior. *Water Resources Research*, 49, 1400-1414.
- 702 Huntingford, C., Marsh, T., Scaife, A. A., Kendon, E. J., Hannaford, J., Kay, A. L.,
703 Lockwood, M., Prudhomme, C., Reynard, N. S., Parry, S., Lowe, J. A.,
704 Screen, J. A., Ward, H. C., Roberts, M., Stott, P. A., Bell, V. A., Bailey, M.,
705 Jenkins, A., Legg, T., Otto, F. E. L., Massey, N., Schaller, N., Slingo, J. &
706 Allen, M. R. 2014. Potential influences on the United Kingdom's floods of
707 winter 2013/14. *Nature Clim. Change*, 4, 769-777.
- 708 Jakeman, A. J., Letcher, R. A. & Norton, J. P. 2006. Ten iterative steps in
709 development and evaluation of environmental models. *Environmental*
710 *Modelling & Software*, 21, 602-614.
- 711 Jolliffe, I. T. & Stephenson, D. B. 2012. Forecast Verification: A Practitioner's
712 Guide in Atmospheric Science. Wiley, Chichester, UK.
- 713 Lanen, H. A. J. & Peters, E. 2000. Definition, Effects and Assessment of
714 Groundwater Droughts. In: VOGT, J. & SOMMA, F. (eds.) *Drought and*
715 *Drought Mitigation in Europe*. Springer Netherlands.
- 716 Lewis, J. M. 2005. Roots of Ensemble Forecasting. *Monthly Weather Review*,
717 133, 1865-1885.
- 718 Liu, Y., Freer, J., Beven, K. & Matgen, P. 2009. Towards a limits of acceptability
719 approach to the calibration of hydrological models: Extending
720 observation error. *Journal of Hydrology*, 367, 93-103.
- 721 Liu, Y., Weerts, A. H., Clark, M., Hendricks Franssen, H. J., Kumar, S.,
722 Moradkhani, H., Seo, D. J., Schwanenberg, D., Smith, P., van Dijk, A. I. J.
723 M., van Velzen, N., He, M., Lee, H., Noh, S. J., Rakovec, O. & Restrepo, P.

- 724 2012. Advancing data assimilation in operational hydrologic forecasting:
725 progresses, challenges, and emerging opportunities. *Hydrol. Earth Syst.*
726 *Sci.*, 16, 3863-3887.
- 727 Lorenz, E. N. 1963. Deterministic Nonperiodic Flow. *Journal of the Atmospheric*
728 *Sciences*, 20, 130-141.
- 729 Mackay, J. D., Jackson, C. R. & Wang, L. 2014. A lumped conceptual model to
730 simulate groundwater level time-series. *Environmental Modelling &*
731 *Software*, 61, 229-245.
- 732 MacLachlan, C., Arribas, A., Peterson, K. A., Maidens, A., Fereday, D., Scaife, A.
733 A., Gordon, M., Vellinga, M., Williams, A., Comer, R. E., Camp, J., Xavier,
734 P. & Madec, G. 2014. Global Seasonal forecast system version 5
735 (GloSea5): a high-resolution seasonal forecast system. *Quarterly Journal*
736 *of the Royal Meteorological Society*, n/a-n/a.
- 737 Maheswaran, R. & Khosa, R. 2013. Long term forecasting of groundwater levels
738 with evidence of non-stationary and nonlinear characteristics.
739 *Computers & Geosciences*, 52, 422-436.
- 740 Maier, H. R. & Dandy, G. C. 2000. Neural networks for the prediction and
741 forecasting of water resources variables: a review of modelling issues
742 and applications. *Environmental Modelling & Software*, 15, 101-124.
- 743 Marsh, T., Cole, G. & Wilby, R. 2007. Major droughts in England and Wales,
744 1800–2006. *Weather*, 62, 87-93.
- 745 Marsh, T. J. & Hannaford, J. 2008. UK Hydrometric Register. Centre for Ecology
746 and Hydrology, Wallingford, UK.
- 747 Mendicino, G., Senatore, A. & Versace, P. 2008. A Groundwater Resource Index
748 (GRI) for drought monitoring and forecasting in a mediterranean climate.
749 *Journal of Hydrology*, 357, 282-302.
- 750 Mishra, A. K. & Singh, V. P. 2010. A review of drought concepts. *Journal of*
751 *Hydrology*, 391, 202-216.
- 752 Monteith, J. L. & Unsworth, M. H. 2008. Principles of Environmental Physics:
753 Third Edition. Elsevier, London, UK.
- 754 Moore, R. J. & Bell, V. A. 1999. Incorporation of groundwater losses and well
755 level data in rainfall-runoff models illustrated using the PDM. *Hydrol.*
756 *Earth Syst. Sci.*, 6, 25-38.
- 757 Murphy, A. H. 1973. A New Vector Partition of the Probability Score. *Journal of*
758 *Applied Meteorology*, 12, 595-600.
- 759 Nash, J. E. & Sutcliffe, J. V. 1970. River flow forecasting through conceptual
760 models part I — A discussion of principles. *Journal of Hydrology*, 10, 282-
761 290.

- 762 Nourani, V., Mogaddam, A. A. & Nadiri, A. O. 2008. An ANN-based model for
763 spatiotemporal groundwater level forecasting. *Hydrological Processes*,
764 22, 5054-5066.
- 765 Pappenberger, F., Beven, K. J., Hunter, N. M., Bates, P. D., Gouweleeuw, B. T.,
766 Thielen, J. & de Roo, A. P. J. 2005. Cascading model uncertainty from
767 medium range weather forecasts (10 days) through a rainfall-runoff
768 model to flood inundation predictions within the European Flood
769 Forecasting System (EFFS). *Hydrol. Earth Syst. Sci.*, 9, 381-393.
- 770 Pappenberger, F., Ramos, M. H., Cloke, H. L., Wetterhall, F., Alfieri, L., Bogner,
771 K., Mueller, A. & Salamon, P. 2015. How do I know if my forecasts are
772 better? Using benchmarks in hydrological ensemble prediction. *Journal*
773 *of Hydrology*, 522, 697-713.
- 774 Partal, T. & Kisi, Ö. 2007. Wavelet and neuro-fuzzy conjunction model for
775 precipitation forecasting. *Journal of Hydrology*, 342, 199-212.
- 776 Perry, M., Hollis, D. & Elms, M. 2009. The Generation of Daily Gridded Datasets
777 of Temperature and Rainfall for the UK. Exeter, UK.
- 778 Pinault, J. L., Amraoui, N. & Golaz, C. 2005. Groundwater-induced flooding in
779 macropore-dominated hydrological system in the context of climate
780 changes. *Water Resources Research*, 41, W05001.
- 781 Sahu, B. K. 2003. Time Series Modelling in Earth Sciences. A.A. Balkema,
782 Netherlands.
- 783 Scaife, A. A., Arribas, A., Blockley, E., Brookshaw, A., Clark, R. T., Dunstone, N.,
784 Eade, R., Fereday, D., Folland, C. K., Gordon, M., Hermanson, L., Knight, J.,
785 R., Lea, D. J., MacLachlan, C., Maidens, A., Martin, M., Peterson, A. K.,
786 Smith, D., Vellinga, M., Wallace, E., Waters, J. & Williams, A. 2014.
787 Skillful long-range prediction of European and North American winters.
788 *Geophysical Research Letters*, 41, 2014GL059637.
- 789 Smith, P., Beven, K. J. & Tawn, J. A. 2008. Informal likelihood measures in
790 model assessment: Theoretic development and investigation. *Advances*
791 *in Water Resources*, 31, 1087-1100.
- 792 Sperna Weiland, F. C., van Beek, L. P. H., Kwadijk, J. C. J. & Bierkens, M. F. P.
793 2010. The ability of a GCM-forced hydrological model to reproduce
794 global discharge variability. *Hydrol. Earth Syst. Sci.*, 14, 1595-1621.
- 795 Sreekanth, P. D., Geethanjali, N., Sreedevi, P. D., Ahmend, S., Ravi Kumar, N. &
796 Kamala Jayanthi, P. D. 2009. Forecasting groundwater level using
797 artificial neural networks *Current Science*, 96, 933-939.
- 798 Sukhija, B. S. 2008. Adaptation to climate change: strategies for sustaining
799 groundwater resources during droughts. *Geological Society, London*,
800 *Special Publications*, 288, 169-181.

- 801 Suryanarayana, C., Sudheer, C., Mahmood, V. & Panigrahi, B. K. 2014. An
802 integrated wavelet-support vector machine for groundwater level
803 prediction in Visakhapatnam, India. *Neurocomputing*, 145, 324-335.
- 804 Svensson, C., Brookshaw, A., Scaife, A. A., Bell, V. A., Mackay, J. D., Jackson, C.
805 R., Hannaford, J., Davies, H. N., Arribas, A. & Stanley, S. 2015. Long-range
806 forecasts of UK winter hydrology. *Environmental Research Letters*, 10,
807 064006.
- 808 Taormina, R., Chau, K.-w. & Sethi, R. 2012. Artificial neural network simulation
809 of hourly groundwater levels in a coastal aquifer system of the Venice
810 lagoon. *Engineering Applications of Artificial Intelligence*, 25, 1670-1676.
- 811 Trichakis, I. C., Nikolos, I. K. & Karatzas, G. P. 2009. Optimal selection of
812 artificial neural network parameters for the prediction of a karstic
813 aquifer's response. *Hydrological Processes*, 23, 2956-2969.
- 814 Tsanis, I. K., Coulibaly, P. & Daliakopoulos, I. N. 2008. Improving groundwater
815 level forecasting with a feedforward neural network and linearly
816 regressed projected precipitation. *Journal of Hydroinformatics*, 10, 317-
817 330.
- 818 Upton, K. A. & Jackson, C. R. 2011. Simulation of the spatio-temporal extent of
819 groundwater flooding using statistical methods of hydrograph
820 classification and lumped parameter models. *Hydrological Processes*, 25,
821 1949-1963.
- 822 Vapnik, V. N. 1999. An overview of statistical learning theory. *Neural Networks*,
823 *IEEE Transactions on*, 10, 988-999.
- 824 Verkade, J. S., Brown, J. D., Reggiani, P. & Weerts, A. H. 2013. Post-processing
825 ECMWF precipitation and temperature ensemble reforecasts for
826 operational hydrologic forecasting at various spatial scales. *Journal of*
827 *Hydrology*, 501, 73-91.
- 828 Weisheimer, A. & Palmer, T. N. 2014. On the reliability of seasonal climate
829 forecasts.
- 830 Xiong, L. & O'Connor, K. M. 2008. An empirical method to improve the
831 prediction limits of the GLUE methodology in rainfall-runoff modeling.
832 *Journal of Hydrology*, 349, 115-124.
- 833 Ying, Z., Wenxi, L., Haibo, C. & Jiannan, L. 2014. Comparison of three
834 forecasting models for groundwater levels: a case study in the semiarid
835 area of west Jilin Province, China. *Journal of Water Supply*, 63, 671-683.
- 836 Yoon, H., Jun, S.-C., Hyun, Y., Bae, G.-O. & Lee, K.-K. 2011. A comparative study
837 of artificial neural networks and support vector machines for predicting
838 groundwater levels in a coastal aquifer. *Journal of Hydrology*, 396, 128-
839 138.

- 840 Yossef, N. C., van Beek, L. P. H., Kwadijk, J. C. J. & Bierkens, M. F. P. 2012.
841 Assessment of the potential forecasting skill of a global hydrological
842 model in reproducing the occurrence of monthly flow extremes. *Hydrol.*
843 *Earth Syst. Sci.*, 16, 4233-4246.
- 844 Yossef, N. C., Winsemius, H., Weerts, A., van Beek, R. & Bierkens, M. F. P. 2013.
845 Skill of a global seasonal streamflow forecasting system, relative roles of
846 initial conditions and meteorological forcing. *Water Resources Research*,
847 49, 4687-4699.
- 848 Zappa, M., Beven, K. J., Bruen, M., Cofiño, A. S., Kok, K., Martin, E., Nurmi, P.,
849 Orfila, B., Roulin, E., Schröter, K., Seed, A., Szturc, J., Vehviläinen, B.,
850 Germann, U. & Rossa, A. 2010. Propagation of uncertainty from
851 observing systems and NWP into hydrological models: COST-731
852 Working Group 2. *Atmospheric Science Letters*, 11, 83-91.
- 853 Zappa, M., Jaun, S., Germann, U., Walser, A. & Fundel, F. 2011. Superposition
854 of three sources of uncertainties in operational flood forecasting chains.
855 *Atmospheric Research*, 100, 246-262.
- 856
857
858
859
860
861
862
863
864
865
866
867
868
869
870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

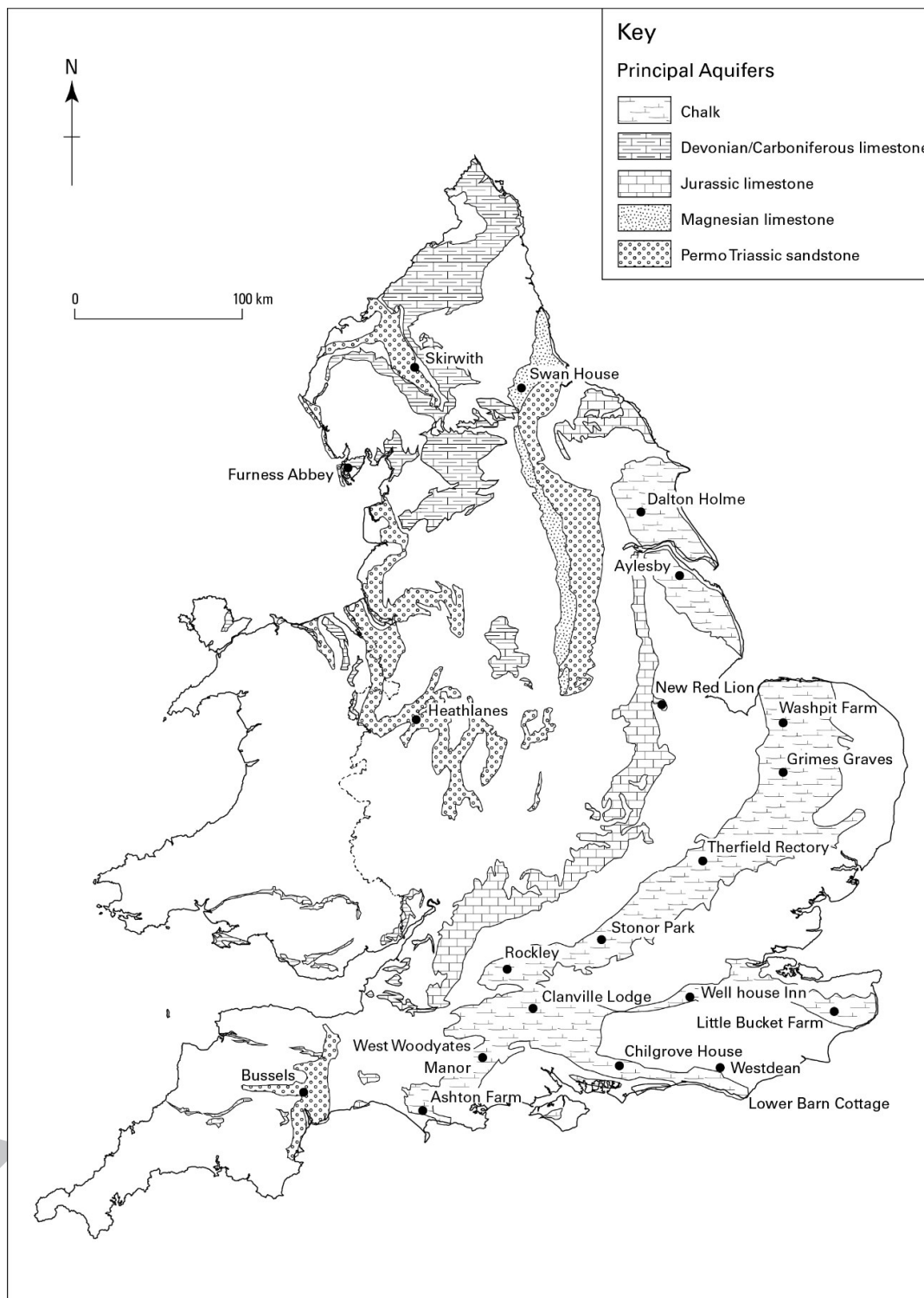
891

892

ACCEPTED MANUSCRIPT

893

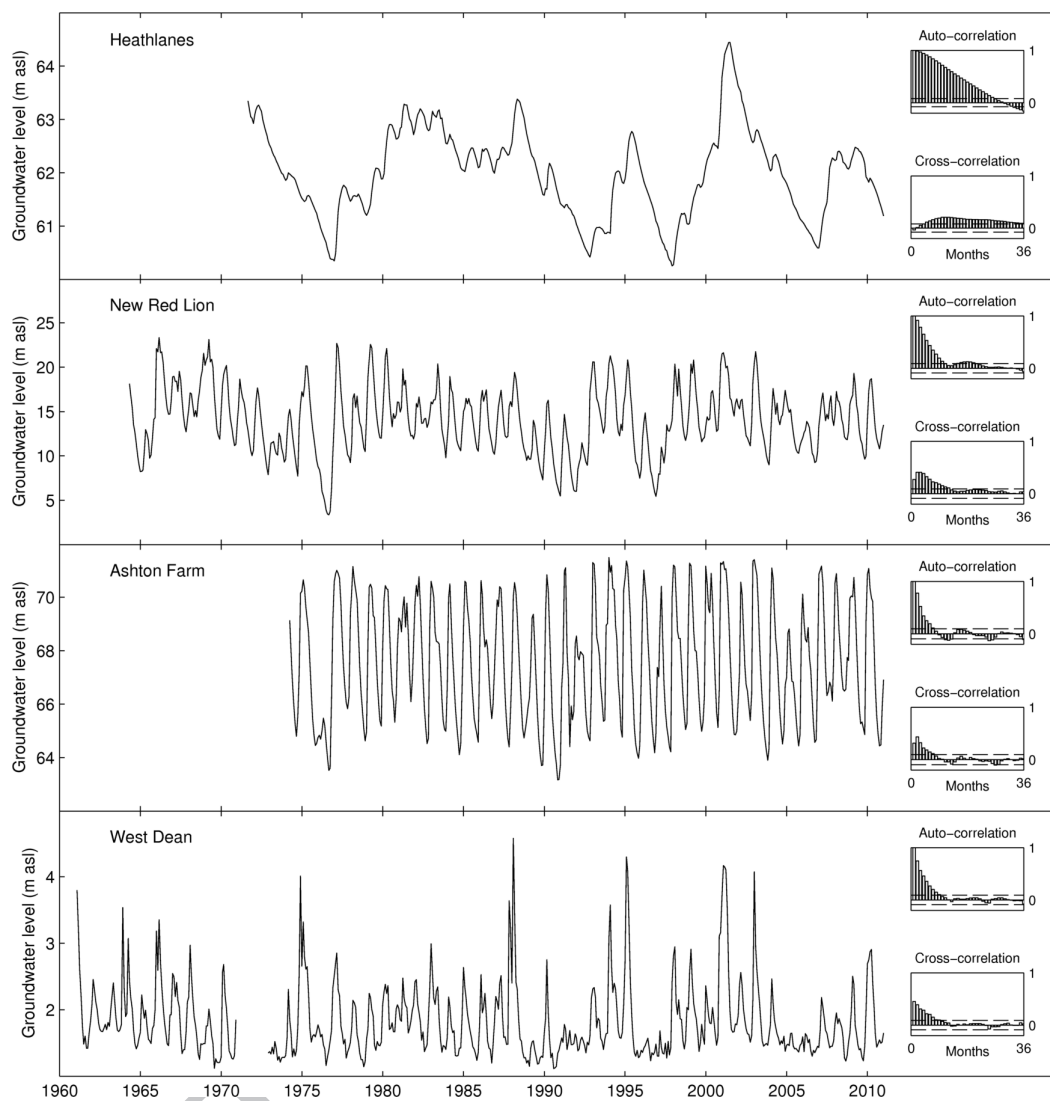
1. Figures



894

895 *Figure 1: Observation borehole locations across the principal aquifers of the UK.*

896



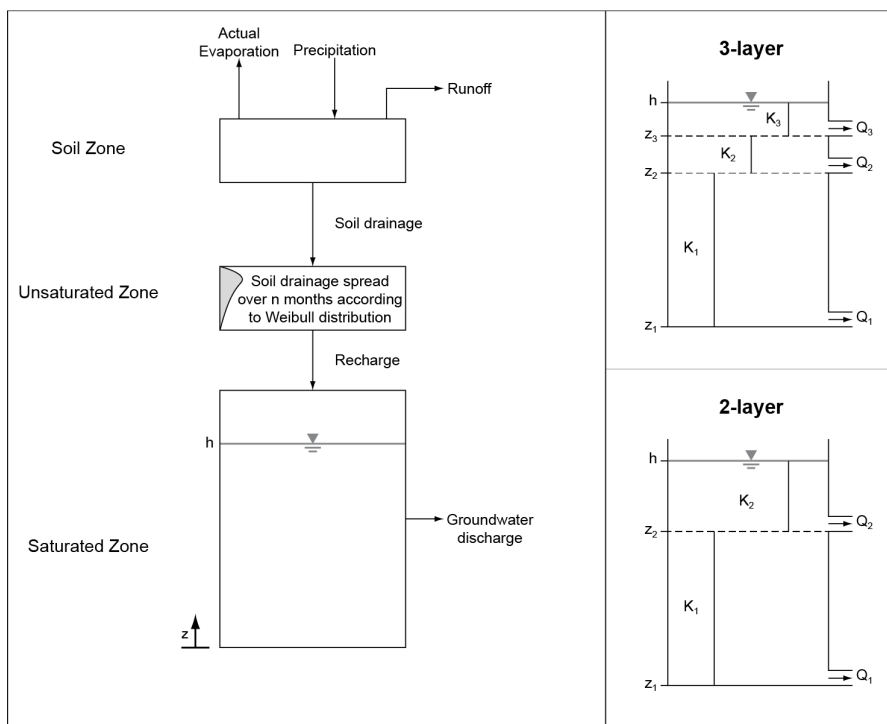
897

898 *Figure 2: Groundwater level time-series with groundwater level auto-correlation and rainfall-groundwater level*899 *cross-correlation plots. Note that the vertical scales vary across the plots.*

900

901

902



903

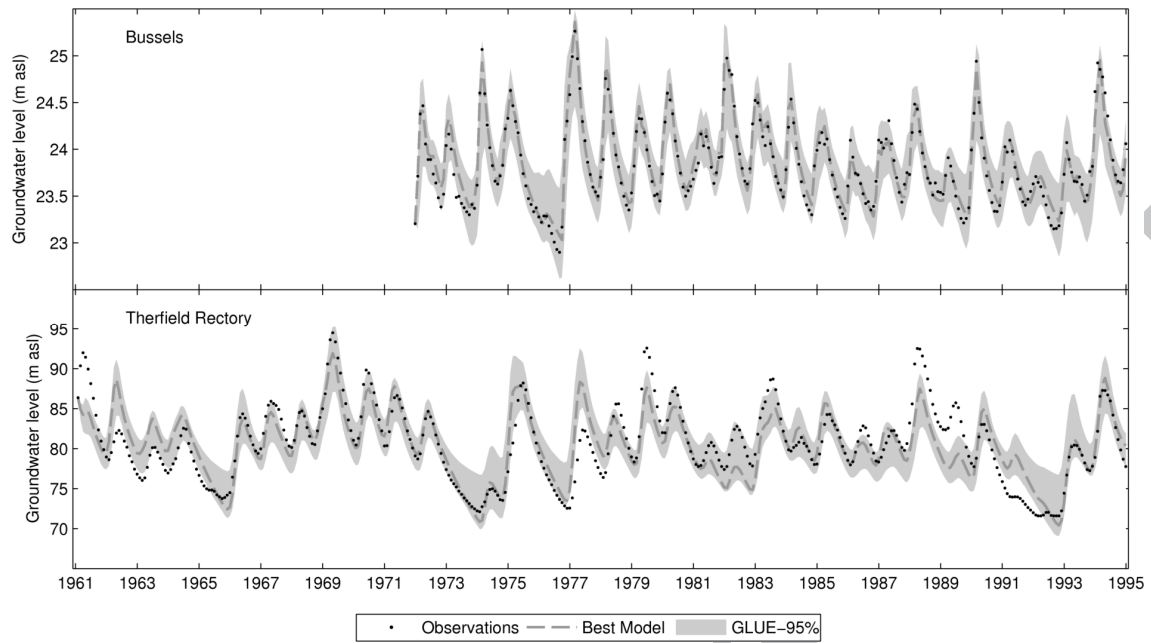
904 *Figure 3: Schematic of generalised AquiMod model structure (left) and different saturated zone component*905 *structures used in this study (right) after Mackay et al. (2014).*

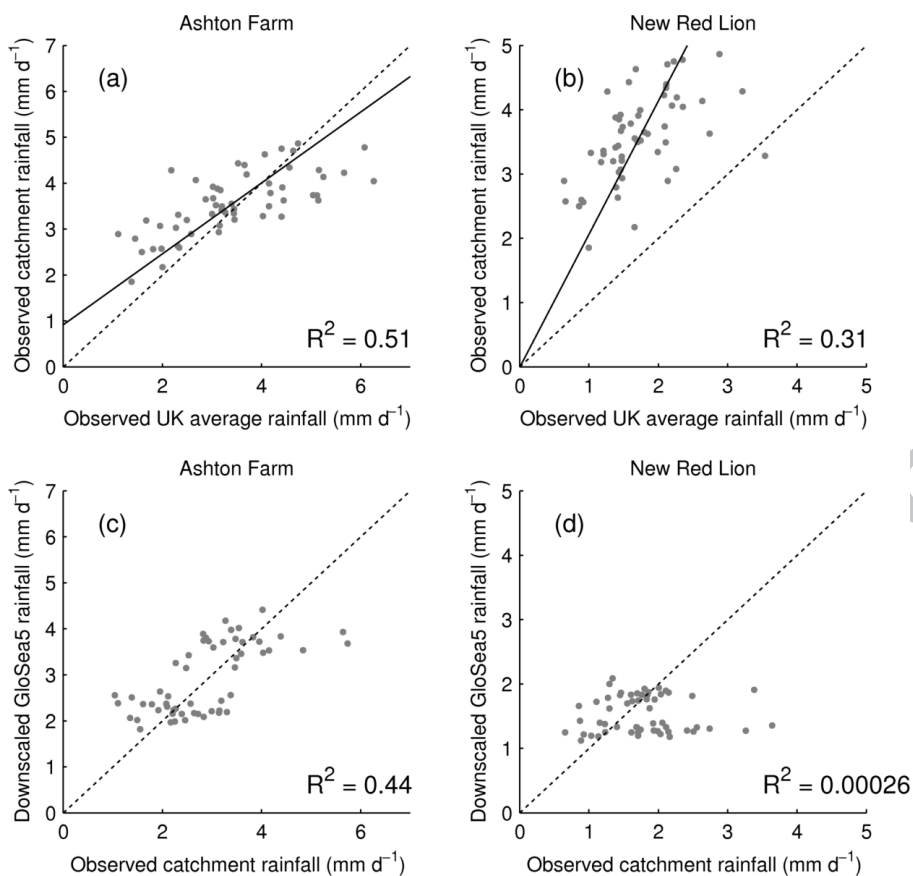
906

907

908

909





914

915

Figure 5: Linear regression models (solid black lines) fitted to downscale seasonal rainfall from UK average to

916

catchment scale for the Ashton Farm (a) and New Red Lion (b) observation boreholes. The resulting correlation

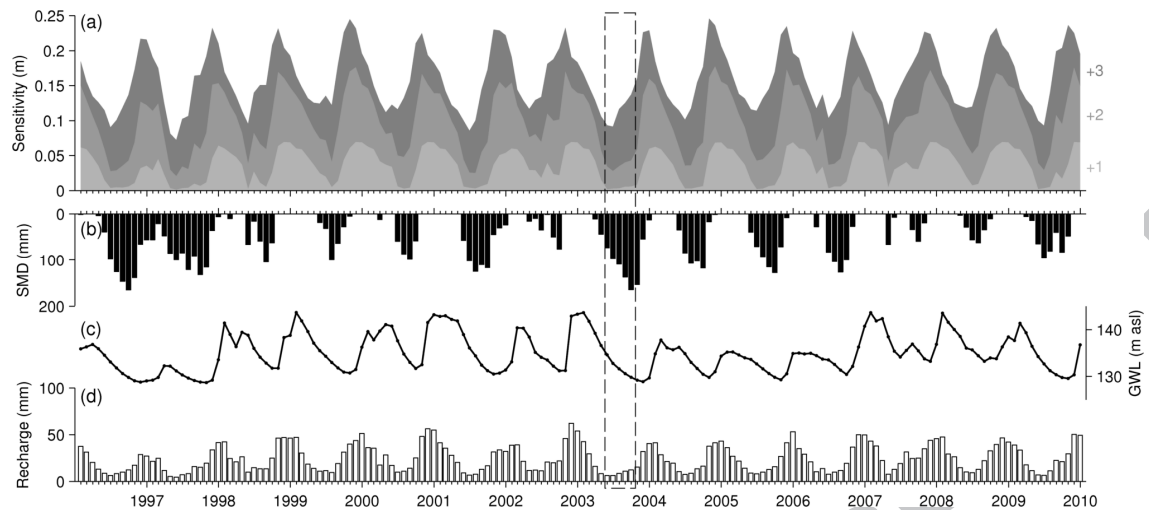
917

between the downscaled GloSea5 rainfall forecasts and the observed catchment rainfall is also shown for the

918

Ashton Farm (c) and New Red Lion (d) observation boreholes.

919

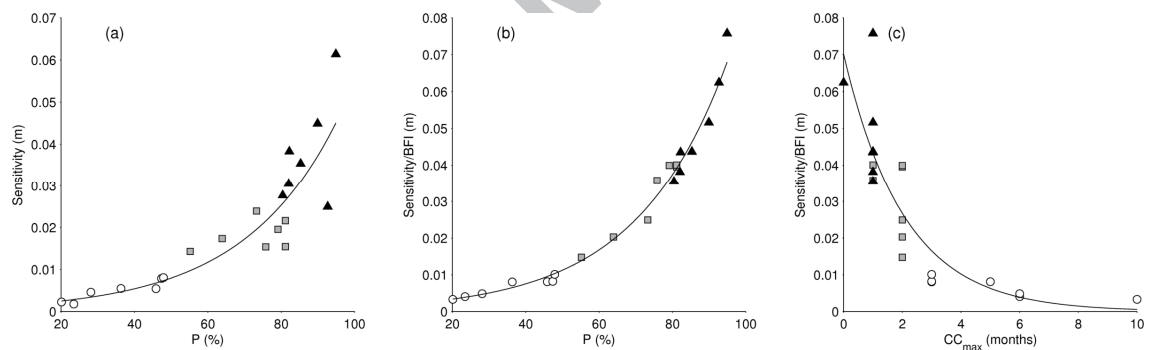


920

921 *Figure 6: Calculated monthly sensitivity to climate inputs for one, two and three month forecasts (a) ; Soil*
 922 *moisture deficit initial condition (b); Groundwater level initial condition (c); and mean monthly simulated*
 923 *recharge (d).*

924

925

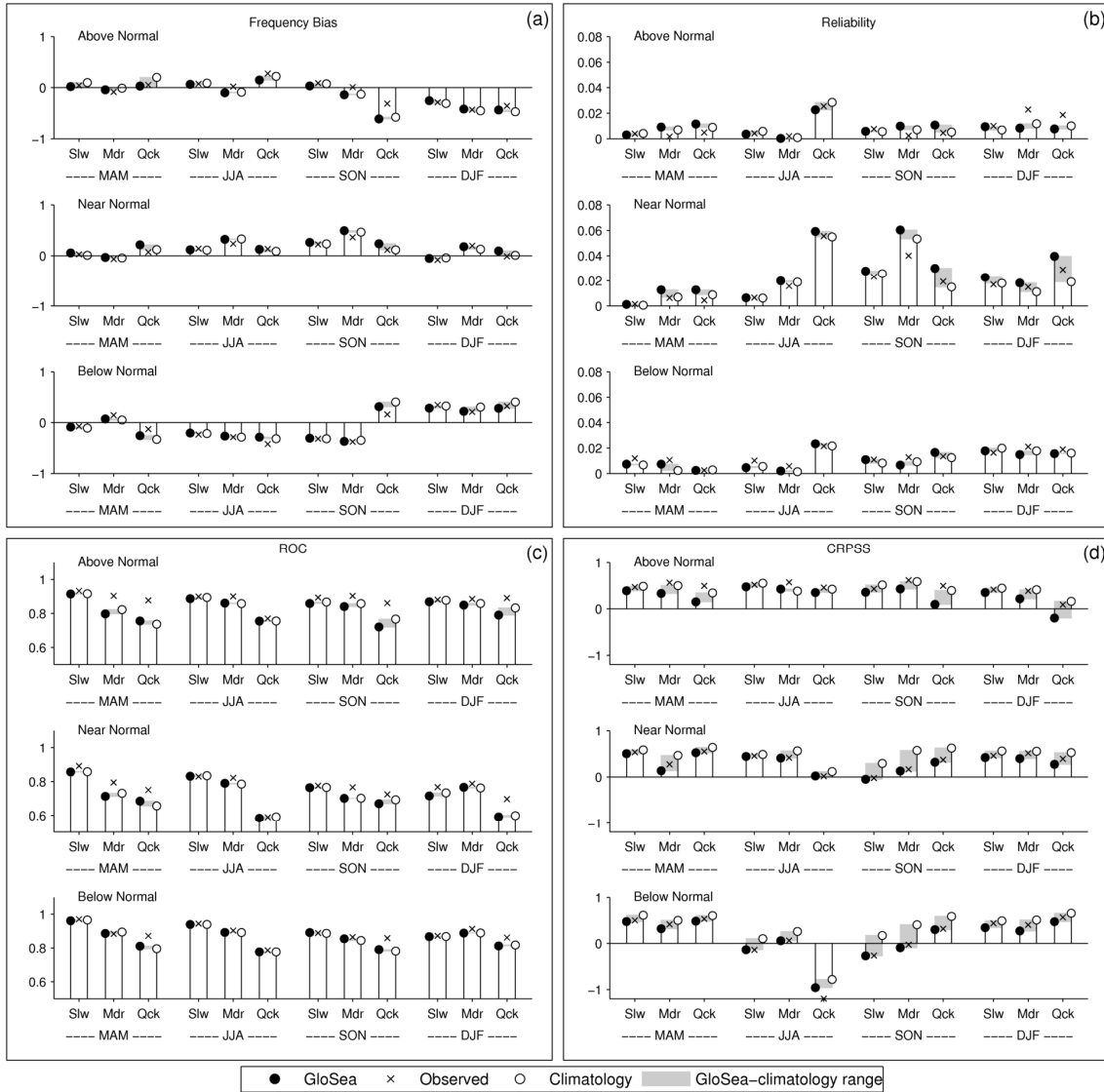


926

927 *Figure 7: Model response characteristic, P , against the derived model sensitivity (a); P against the sensitivity*
 928 *normalised with respect to the BFI (b); and the peak lead lag correlation between observed rainfall and de-*
 929 *seasonalised groundwater levels, CC_{max} , against the sensitivity normalised with respect to the BFI (c) for the 21*
 930 *catchment models. All data points are arranged into slowly (circles), moderately (squares) and quickly*
 931 *(triangles) responding catchments.*

932

933

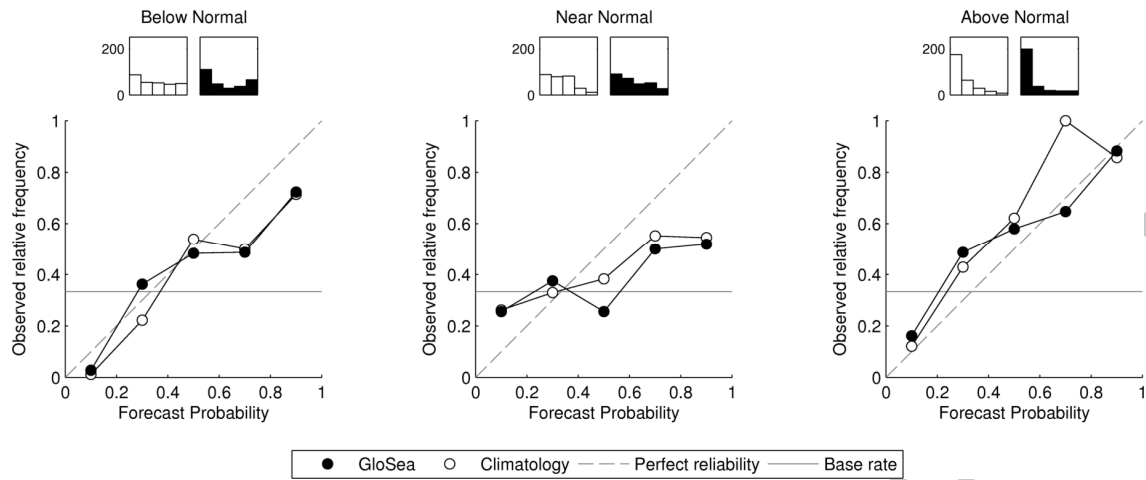


934

935 *Figure 8: Frequency bias (a), reliability (b), ROC (c), and CRPSS (d) metrics calculated from the reforecasts.*

936

937



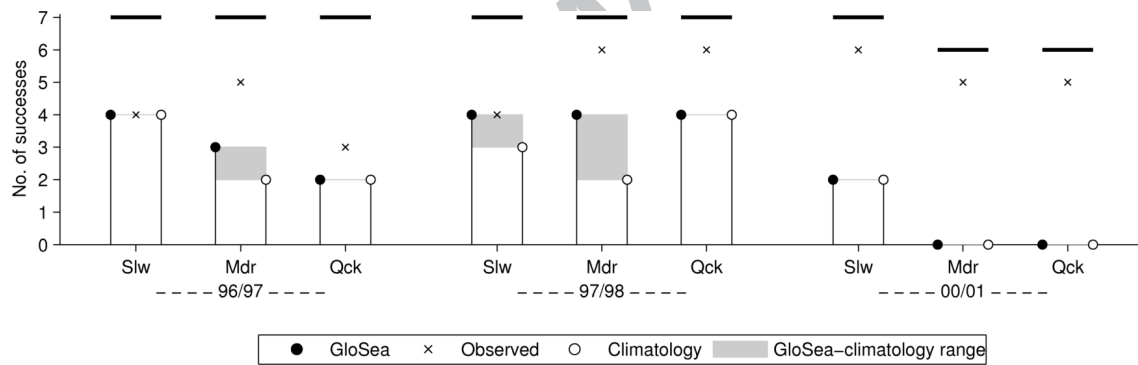
938

939 *Figure 9: Reliability diagrams for the quickly responding catchments. The histograms denote the sample sizes*

940 *for each point on the reliability curves.*

941

942

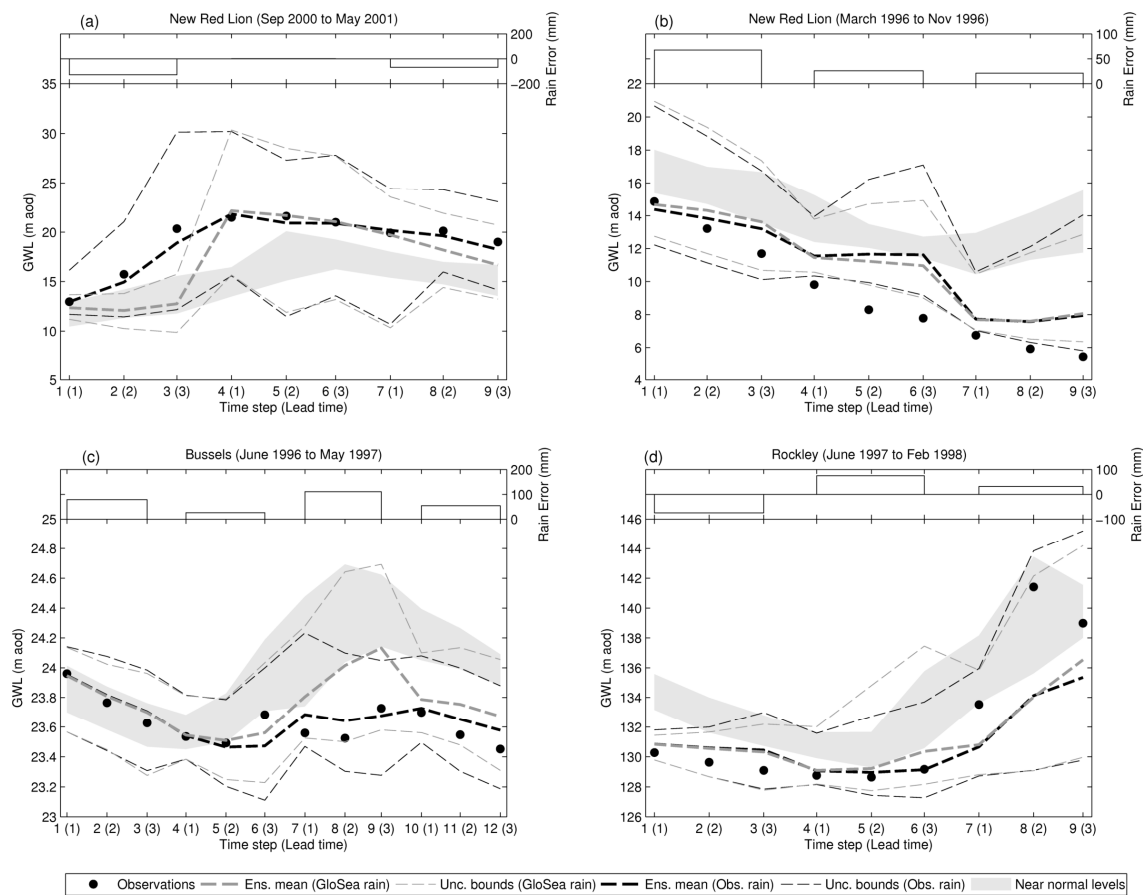


943

944 *Figure 10: Number of successful forecasts for three events. The solid black lines indicate the total number of*

945 *catchments with available observation data.*

946



947

948 *Figure 11: Comparison of the reforecasts using downscaled GloSea and observed rainfall inputs for different*
 949 *time periods.*

950

951

952

953

954

955

956

957

958 **2. Tables**959 *Table 1: List of AquiMod model parameters and calibration ranges.*

Module	Parameter (units)	Description	Typical calibration range
Soil	Δx (km)	Representative aquifer length	Fixed as distance between observation borehole and river discharging groundwater.
	BFI (-)	Baseflow index	Taken from Marsh and Hannaford (2008).
	FC (-)	Field capacity of the soil	Taken from Boorman et al. (1995).
	WP (-)	Wilting point of the soil	Taken from Boorman et al. (1995).
	Zr (mm)	Maximum rooting depth of vegetation	100 – 3000
	p (-)	Depletion factor of vegetation	0 – 1
Unsaturated Zone	n (-)	Maximum number of time-steps taken for soil drainage to reach the groundwater	Set based on cross-correlation analysis between rainfall and groundwater levels.
	k (-)	Weibull shape parameter	1 – 7
	λ (-)	Weibull scale parameter	1 – 12
Saturated Zone	K_i ($m d^{-1}$)	Hydraulic conductivity for layer i	0.01 – 100
	S (%)	Aquifer storage coefficient	0.1 – 20
	Z_i (m asl)	Outlet elevation for layer i	Deep outlet set to the known bottom elevation of aquifer.

Remaining outlet elevations

set after preliminary

calibration runs.

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981 Table 2: List of 21 observation boreholes with the number of behavioural models (n), the efficiency of the most
 982 efficient model (NSE_{max}) and the containment ratio using the GLUE 95% confidence bounds (CR).

Observation borehole	Aquifer	n	NSE_{max}	CR
Ashton Farm	Chalk	2155	0.89	94.4
Aylesby	Chalk	2470	0.82	96.9
Chilgrove House	Chalk	2125	0.91	97.8
Clanville Lodge	Chalk	2025	0.84	89.0
Dalton Holme	Chalk	2000	0.81	82.6
Grimes Graves	Chalk	1960	0.86	88.9
Little Bucket Farm	Chalk	2305	0.90	85.7
Rockley	Chalk	1835	0.88	94.1
Stonor Park	Chalk	2430	0.78	65.3
Therfield Rectory	Chalk	1915	0.71	68.9
Washpit Farm	Chalk	1910	0.91	96.3
Well House Inn	Chalk	1850	0.73	68.1
West Dean	Chalk	2210	0.83	92.2
West Woodyates Manor	Chalk	1780	0.86	84.8
New Red Lion	Jurassic Limestone	2155	0.74	77.0
Lower Barn Cottage	Lower Greensand	2120	0.81	79.5
Swan House	Magnesian Limestone	1960	0.86	89.6
Bussels	Permo-Triassic Sandstone	2090	0.94	97.5
Furness Abbey	Permo-Triassic Sandstone	2055	0.75	72.7
Heathlanes	Permo-Triassic Sandstone	2095	0.87	87.9
Skirwith	Permo-Triassic Sandstone	2390	0.83	87.6

983

984

985

986

987 **Highlights**

988 • We forecast groundwater levels 3 months into the future for 21 boreholes in the UK.

989 • We use GloSea5 seasonal rainfall forecasts to drive a conceptual groundwater
990 model.

991 • The forecasts consistently show more skill than a persistence forecasting approach.

992 • The forecasts are not able to capture extreme groundwater level events.

993 • Sensitivity to (skill derived from) rainfall forecasts is highly site specific.

994

ACCEPTED MANUSCRIPT