

1 **Supplementary Information 1: RCM/bias-correction information**

2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21

*Table S1. Snowmelt model parameters*

Parameter	Definition	Value	Units
lapse_rate	Temperature lapse rate (reduction in air temperature with increasing elevation)	0.0059	°C m <sup>-1</sup>
tsnow	Threshold temperature, below which precipitation is snow	1.0	°C
tmelt	Temperature threshold for snowmelt	0.0	°C
mfac	Melt factor	6.0	mm °C <sup>-1</sup> day <sup>-1</sup>
tdrel	Threshold temperature for drainage release	0.0	°C
k1	Storage time constants of two-outlet liquid water store	0.5	day <sup>-1</sup>
k2		0.9	
Scfac	Critical water retention capacity	0.18	-
snowfrac	Under-catch factor for gauged rainfall falling as snow (not applied to RCM data)	1.0	-

*Table S2. RCM variables used in the calculation of PET.*

Variable	Description	Units
huss	Specific humidity	-
rls	Radiation, net long wave	W m <sup>-2</sup>
rss	Radiation, net short wave	W m <sup>-2</sup>
sfcWind	Wind speed	m s <sup>-1</sup>
tas	Mean temperature	°C
psl	sea level air pressure (used to derive surface air pressure at gridbox altitude, psurf)	hPa
pr_bc	precipitation (after bias-correction)	mm day <sup>-1</sup>

## Supplementary info 2: PDM Calibration

### *Model configurations*

A full list of the options tested in the PDM model for each catchment is given in Table S3 and summarised below following the notation of Moore (2007).

- *Runoff generation and groundwater recharge.* In the “Full” version of the model, drainage to the groundwater store is described by a recharge time-constant and exponent ( $k_g$ ,  $b_g$ ) and, optionally, a soil tension storage capacity ( $S_t$ ). For the “Reduced” form of model, surface runoff is simply split so that a fixed fraction ( $\alpha$ ) enters the surface store while the remainder enters the groundwater store. In both cases the water absorption capacity of the soil is described by a Pareto distribution characterised by a shape parameter ( $b$ ), a maximum storage ( $c_{max}$ ), and optionally a minimum storage ( $c_{min}$ ). A “Classic” version of the Reduced model employs a rectangular distribution for the soil’s water absorption capacity ( $c_{min} = 0$ ,  $b = 1$ ).
- *Surface water routing.* The surface runoff component of total flow is related to the volume of water in the surface store using a time constant ( $k_1$ ) and exponent ( $m$ ). Exponents of  $m = 1, 2, 3$  are trialled, as is two linear ( $m = 1$ ) stores in series in a discretely-equivalent transfer function form.
- *Groundwater routing.* The baseflow component of total flow is related to the volume of water in the groundwater store using a time constant ( $k_b$ ) and exponent ( $m$ ). Values of  $m = 2, 3$  are trialled.
- *Groundwater (GW) extension.* A standard implementation of PDM conserves water throughout, albeit with the option of applying a multiplicative factor (*rainfac*) to the precipitation input. Conceptually, this factor might compensate for a lack of representativeness in the data used to estimate catchment precipitation. It may also serve to account for losses or gains of water affecting the catchment itself. Alternatively, functionality within the GW extension can be considered to address catchment water conservation issues. This extension, subject to data availability, allows modelling of underflows at the catchment outlet, external springs, pumped abstractions, and the incorporation of well level data. Under eFLaG, only the Spring Factor option (*springfac*) is invoked and repurposed to infer unknown net water exchanges affecting the catchment via the groundwater storage. It serves as a multiplicative factor representing either net losses from ( $1 \geq \text{springfac} > 0$ ), or net gains to ( $\text{springfac} < 0$ ), the baseflow.

### *Calibration process*

A three-stage calibration process was applied independently for each of the model configurations in Table S3, each starting from a number of different choices for the initial parameter. The design of this process was motivated by the desire to find a procedure that could be applied automatically across many disparate catchments without a tendency to either get blocked in local optimums or produce unphysical models.

The following three calibration stages were employed for all model configurations, except those employing the GW extension.

- *Stage 1.* Four or five dominant parameters were segregated according to whether they were judged to control mainly the slow response ( $c_{max}$ ,  $k_b$  for Reduced Models;  $k_b$ ,  $k_g$ ,  $b_g$  for Full Models) or the fast response ( $k_1$ ,  $\alpha$  for Reduced Models;  $k_1$ ,  $c_{max}$  for Full Models). Then, (i)

the slow parameters were calibrated to optimise the  $KGE'_{log}$ , (ii) the fast parameters were calibrated to optimise the  $KGE'_{sqrt}$ , and then (iii)  $rainfac$  was calibrated to achieve zero bias. These steps were iterated six times to achieve convergence.

- *Stage 2.* All parameters calibrated in Stage 1 are re-calibrated simultaneously to maximise the  $KGE'_{sqrt}$ . The  $rainfac$  was then recalibrated to achieve zero bias. This process was iterated three times to ensure convergence.
- *Stage 3.* Additional parameters controlling the distribution of the soil water absorption capacity ( $S_t$ ,  $b$  and  $c_{min}$  for the Full Model;  $b$  and  $c_{min}$  for the Reduced Model; none for the “Classic” model), and one parameter ( $b_e$ ) controlling the sensitivity of the conversion of Potential Evaporation (PE) to Actual Evaporation (AE) with available soil moisture, were each calibrated separately to optimise  $KGE'_{sqrt}$ . Stage 2 was then repeated.

When the GW extension was employed, the three stages were modified so that (i), the Spring Factor was used to achieve zero bias, (ii), greater emphasis was placed on obtaining suitable ground- and soil-water storage parameters (beginning in Stage 1), and (iii) initial parameters were chosen that were more suitable for slowly responding catchments.

### *Calibrated model selection*

A calibrated PDM model is produced for each model configuration, each initial parameter choice, and at each of the three calibration stages, yielding a total of  $46 \times 3 = 138$  possible calibrations per catchment. Figure S1 shows the  $KGE'_{sqrt}$  values, colour coded according to the PDM model configuration, for each of these calibrations for catchment 2001 (Helmsdale at Kilphedir, North West Scotland) and catchment 39089 (Gade at Bury Mill, Hertfordshire and North London area). Any calibrations yielding extreme parameters, including those found to be storing excessive quantities of water, are automatically judged to be unphysical and are shown in black.

The Helmsdale catchment demonstrates several features that are typical across most catchments with low or medium Base Flow Index (BFI) (Figure S1a, BFI = 0.47). These are, (i) calibrations with different model configurations or different initial parameter choices often yield similar metric values, (ii) there are only a small number of calibrations that produce unphysical models and poorer metric values, and (iii) the use of the model configurations employing the GW extension do not generally produce the best model performances. In contrast, for catchments with very high BFI, such as the Gade catchment (Figure S1b, BFI = 0.89), the GW extension is often essential: other model configurations typically produce poor and variable metric values, or unphysical models.

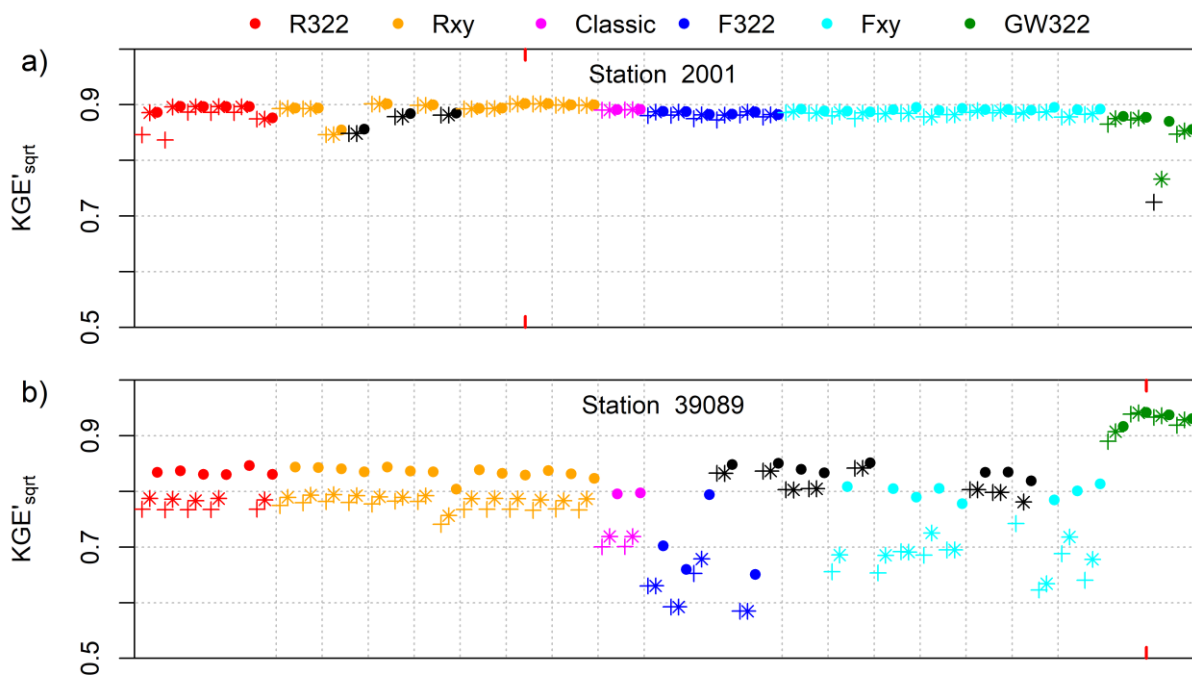
For final model selection, calibrations automatically deemed unphysical were excluded. The use of the GW extension at catchments with Base Flow Index (BFI) lower than 0.7, and at the Leven at Linnbrane (85001), were also judged to be unphysical and therefore also excluded. Final selection among the remaining candidates was based on the highest weighted sum of the modified Kling-Gupta Efficiency calculated on square root ( $KGE'_{sqrt}$ ) and log ( $KGE'_{log}$ ) transformed flows ( $m^3/s$  units), with weights of 0.8 and 0.2 respectively. This was motivated by the above mentioned equifinality of  $KGE'_{sqrt}$ . The  $KGE'_{log}$  metric is usually not recommended (Santos et al., 2018), but its limited use here, compared to a selection based fully on highest  $KGE'_{sqrt}$ , led to some improvement in low flow metrics (e.g., a mean reduction of 4.3% in  $Q95_{APE}$ ) with only minimal reductions in  $KGE'_{sqrt}$  (a mean reduction of 0.0012). The use of an alternative low flow metric in the selection process, rather than the  $KGE'_{log}$ , would be expected to produce similar results.

Part of the quality control for the PDM model was the examination of RCM flows (simrcm) for each catchment. For the Misbourne at Little Missenden (catchment 39127, BFI = 0.96), this revealed unphysically smooth multi-decade recessions beginning in the immediate future for some RCM ensemble members. The cause of this behaviour was found to be a very large minimum point soil water capacity combined with the use of the Reduced Model with no soil drainage to groundwater: a combination that inhibited runoff generation when, due to climate change, the soil water store became depleted. Because of this, the best performing Full Model (F-GW322) was chosen as offering a better hydrological representation of the catchment and more realistic predictions under climate change. This highlights the possibility of unphysical calibrations achieving good metric values against historical river flows ( $KGE'_{sqr t} = 0.927$  was achieved for Misbourne), while producing unphysical results when climate change pushes hydrological conditions outside of their historical regime. This possibility is expected to be more associated with high BFI catchments as these will have greater sensitivity to the longer-term average trends in the weather that become apparent under climate change. By using multiple disparate hydrological models in the eFLaG project (PDM, GR4J, GR6J and G2G), over-reliance on a single model can be avoided.

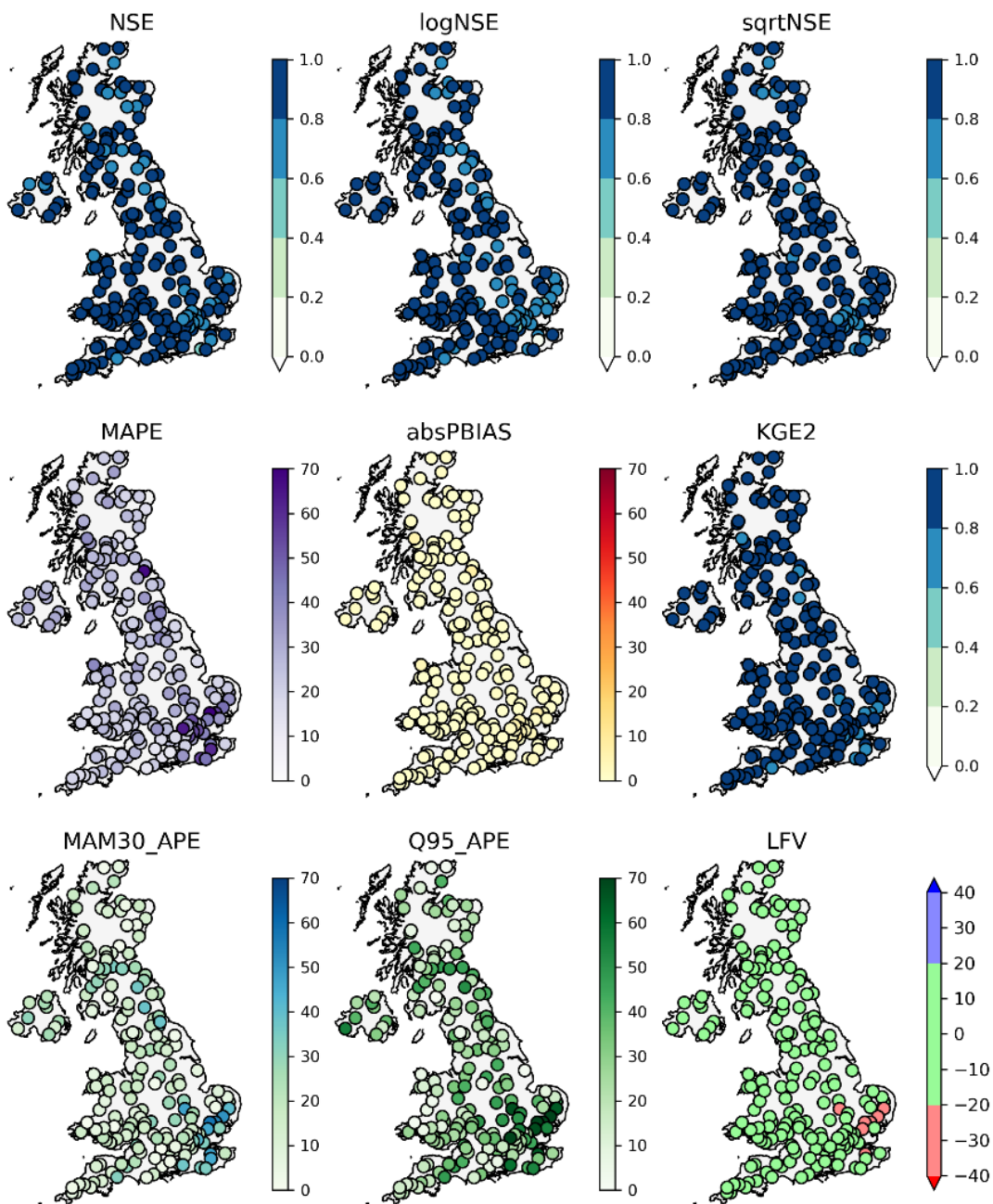
The number of times each model was selected for one of the 200 catchments is listed in Table S3.

**Table S3.** Full list of PDM models trialled for each catchment, and the number of catchments for which each model was selected. For the “Surface routing exponent” column, “22” indicates use of two linear reservoirs in series.

Recharge form	Model code	Groundwater routing exponent	Surface routing exponent	Initial parameter choices	Final selections (out of 200)
Reduced	R322	3	22	6	25
Reduced	R33	3	3	2	44
Reduced	R32	3	2	2	34
Reduced	R31	3	1	2	3
Reduced (Classic)	C31	3	1	2	5
Reduced	R222	2	22	2	4
Reduced	R23	2	3	2	9
Reduced	R22	2	2	2	16
Reduced	R21	2	1	2	2
Full	F322	3	22	6	5
Full	F33	3	3	2	2
Full	F32	3	2	2	6
Full	F31	3	1	2	8
Full	F222	2	22	2	5
Full	F23	2	3	2	4
Full	F22	2	2	2	7
Full	F21	2	1	2	2
Reduced GW extension	R-GW322	3	22	2	15 (out of 26)
Full GW extension	F-GW322	3	22	2	4 (out of 26)

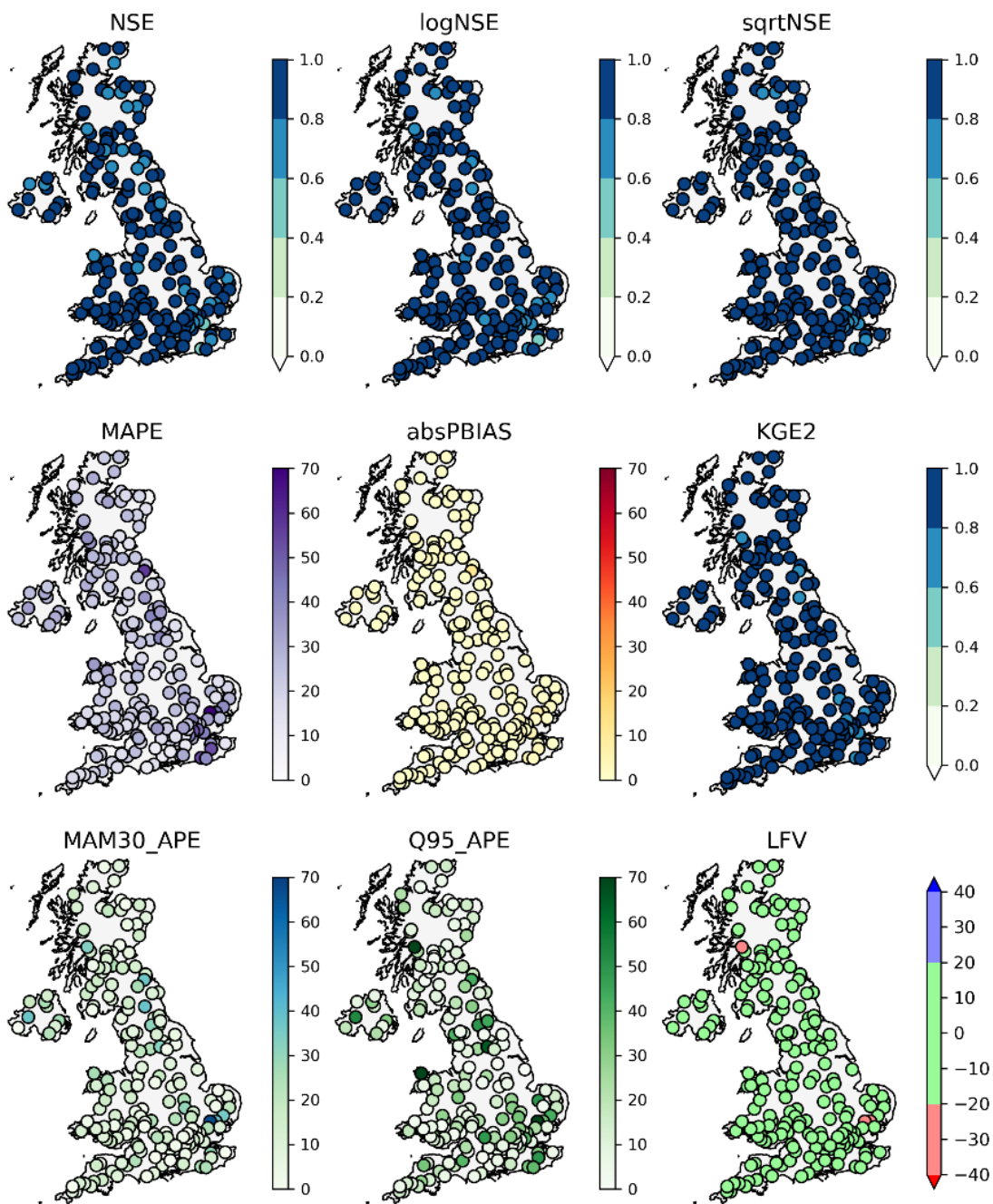


**Figure S1.** The  $KGE'_{sqrt}$  of modelled river flow for each model configuration, calibration stage, and initial parameter choice. Crosses, asterisks and circles indicate performance at calibration stages 1, 2, and 3 respectively. Different colours and dashed lines are used to separate different model configurations. Black is used to show calibrations that resulted in unphysical model parameters. Red ticks on the upper and lower x-axis indicate the final model selection. Catchments are: (a) Helmsdale at Kilphedir (2001), BFI = 0.47, and (b) Gade at Bury Mill (39089), BFI = 0.89.



153

154 **Figure S2: Performance results for GR4J (for metrics see Table 3 in main paper)**

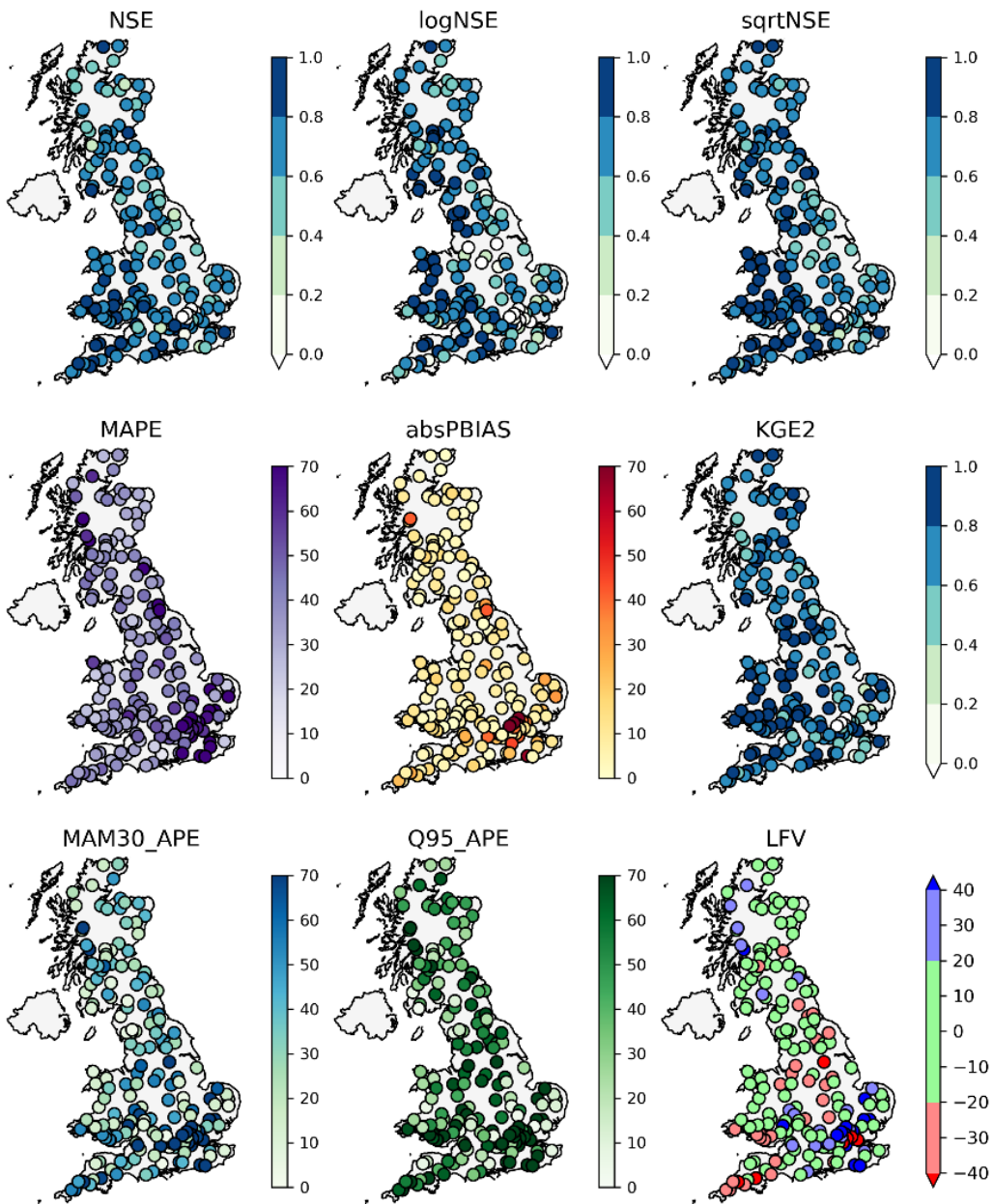


155

156

157 **Figure S3: Performance results for GR6J**



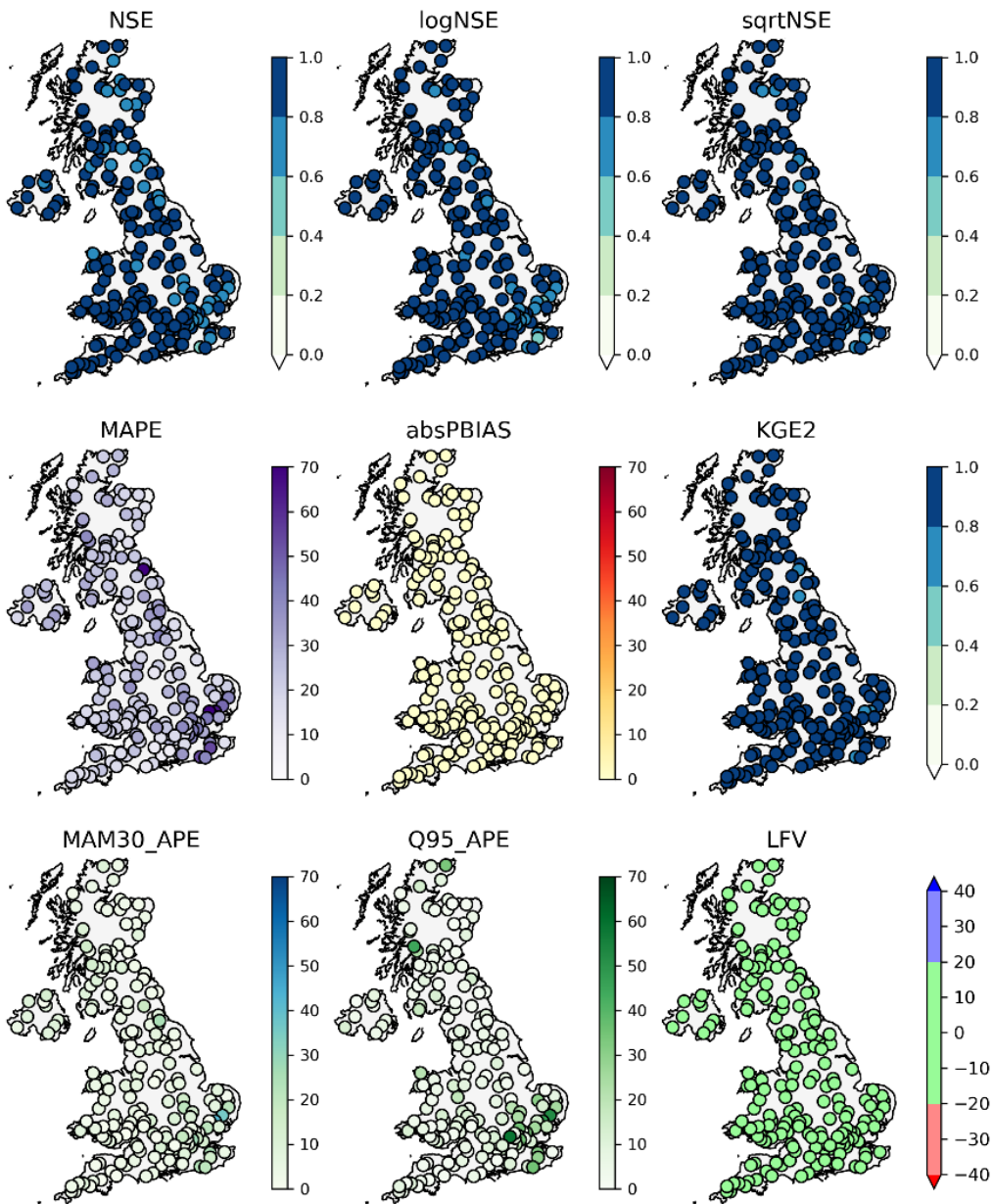


158

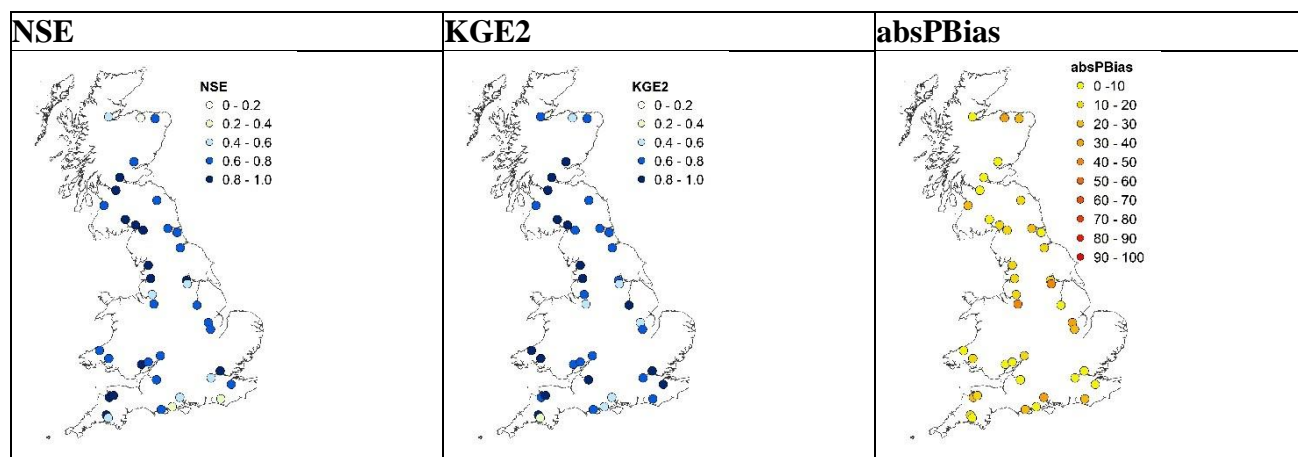
159

160 **Figure S4: Performance results for G2G**





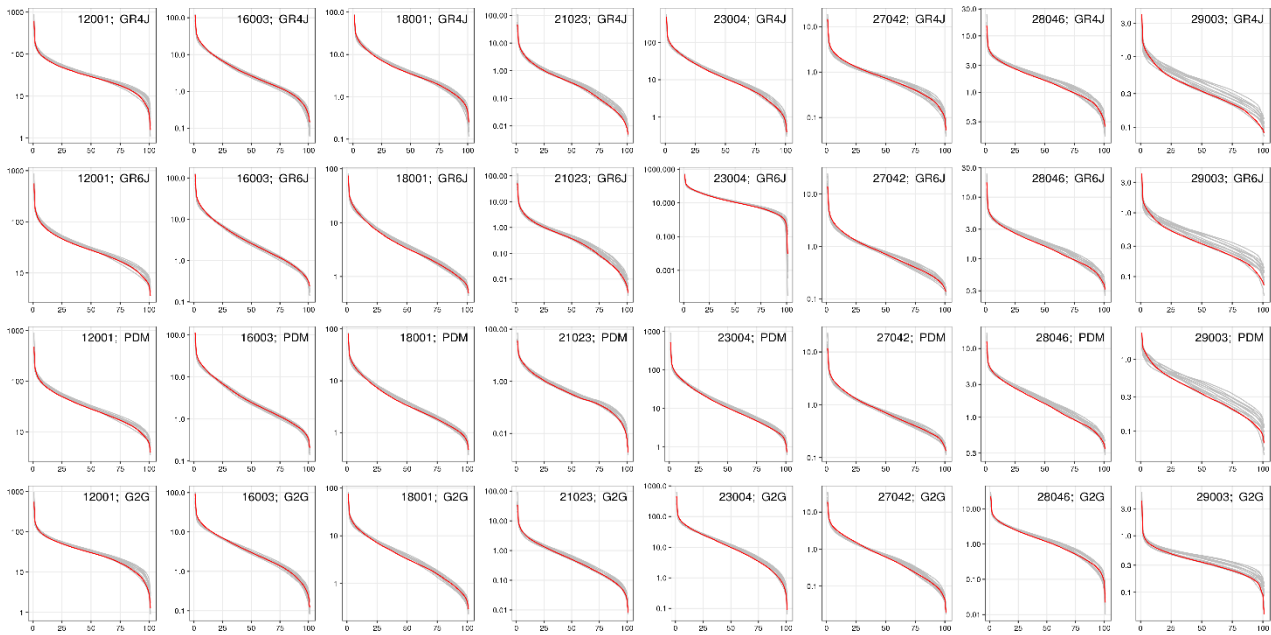
**Figure S5: Performance results for PDM**



**Figure S6: Performance results for the distributed recharge model (ZOODRM)**

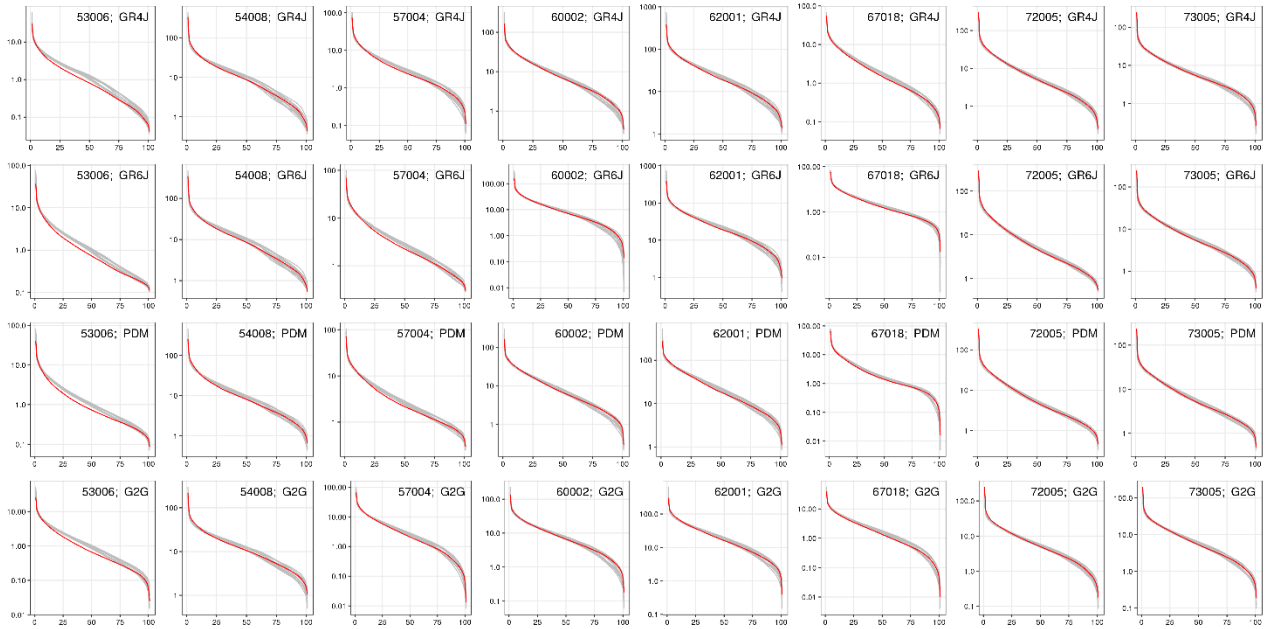
195

196 **Supplementary Information 4: River Flow and Groundwater level Duration Curves**



197

198 *Figure S7 -- Flow duration curves (FDCs) comparing the baseline flow regime in the 12 RCM*  
199 *ensemble members (simrcm, grey lines) to model observations (simobs, red line), 1989-2018. FDCs*  
200 *are featured for four hydrological models (GR4J, GR6J, PDM, G2G; rows) and eight catchments in*  
201 *eastern Scotland and north-east England (12001 Scottish Dee, 16003 Ruchill Water, 18001 Allan*  
202 *Water, 21023 Leet Water, 23004 South Tyne, 27042 Yorkshire Dove, 28046 Derbyshire Dove, 29003*  
203 *Lud; columns). The y-axis represents river flows (cumecs) on a logarithmic scale.*



204

205 *Figure S8 -- Flow duration curves (FDCs) comparing the baseline flow regime in the 12 RCM*  
206 *ensemble members (grey lines) to model observations (red line), 1989-2018. FDCs are featured for*  
207 *four hydrological models (GR4J, GR6J, PDM, G2G; rows) and eight catchments in Wales and north-*  
208 *west England (53006 Bristol Frome, 54008 Teme, 57004 Cynon, 60002 Cothi, 62001 Teifi, 67018*

Welsh Dee, 72005 Lune, 73005 Kent; columns). The y-axis represents river flows (cumecs) on a logarithmic scale.

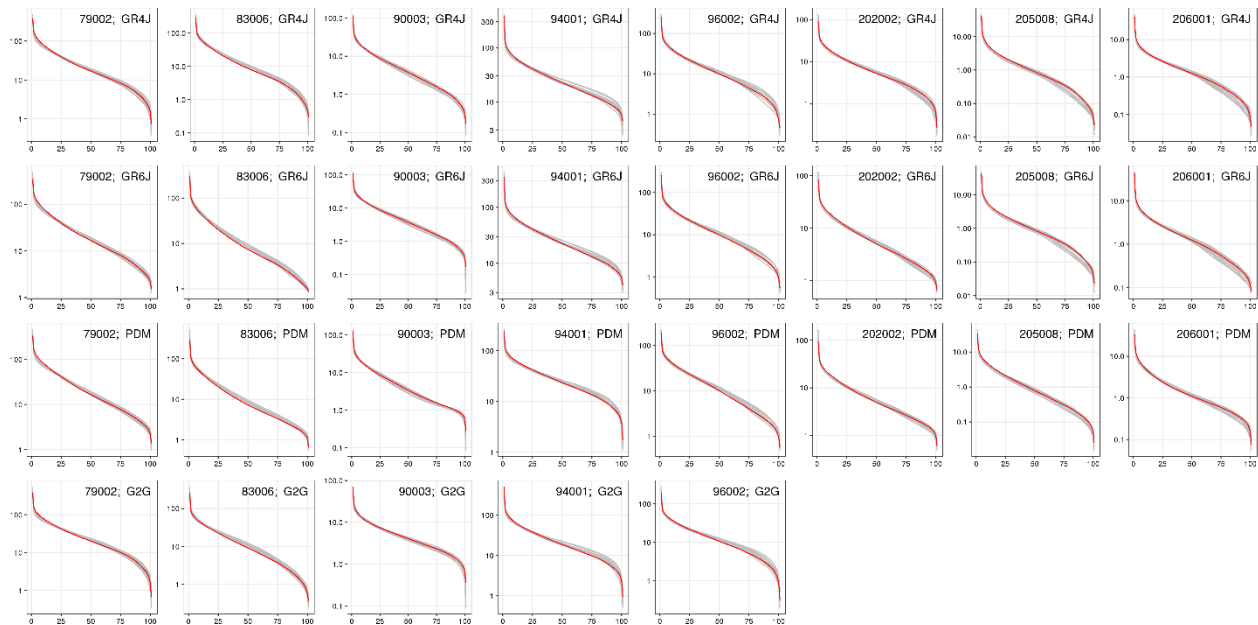
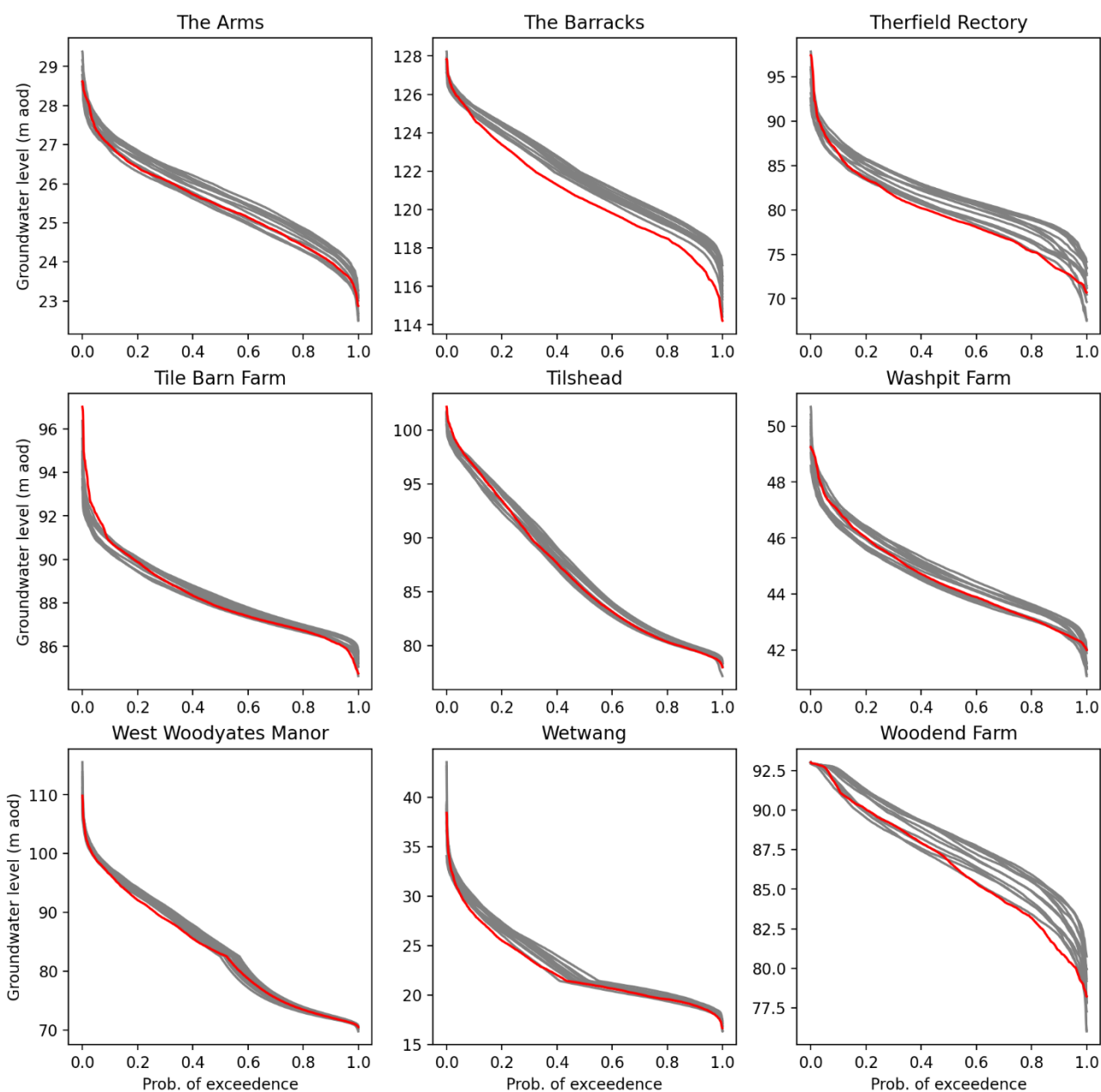
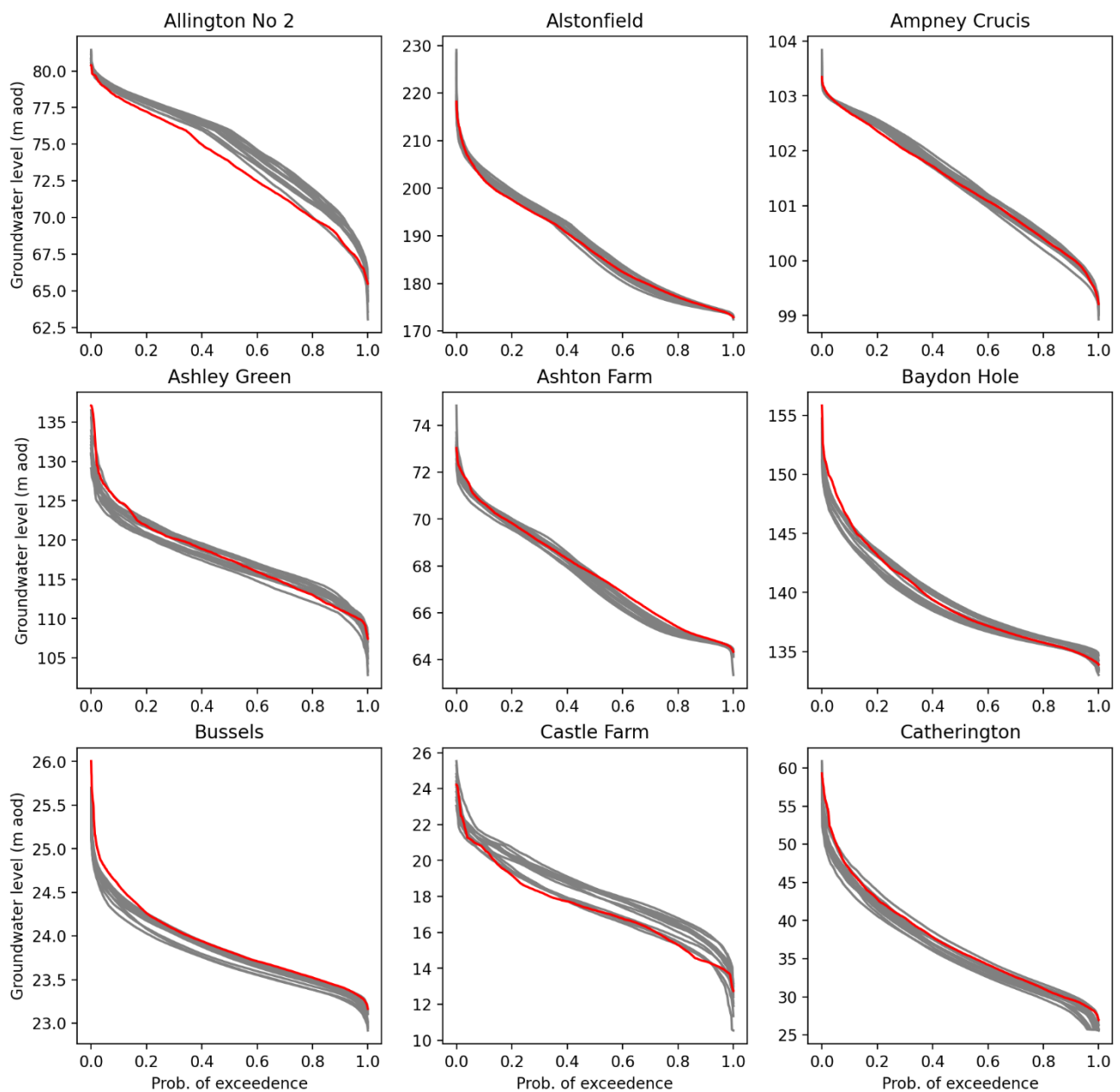


Figure S9 -- Flow duration curves (FDCs) comparing the baseline flow regime in the 12 RCM ensemble members (grey lines) to model observations (red line), 1989-2018. FDCs are featured for four hydrological models (GR4J, GR6J, PDM, G2G; rows) and eight catchments in western Scotland and Northern Ireland (79002 Nith, 83006 Ayr, 90003 Nevis, 94001 Ewe, 96002 Naver, 202002 Faughan, 205008 Lagan, 206001 Clanrye; columns). The y-axis represents river flows (cumecs) on a logarithmic scale. The absence of FDCs for G2G for 202002, 205008 and 206001 is because G2G does not cover Northern Ireland.



220

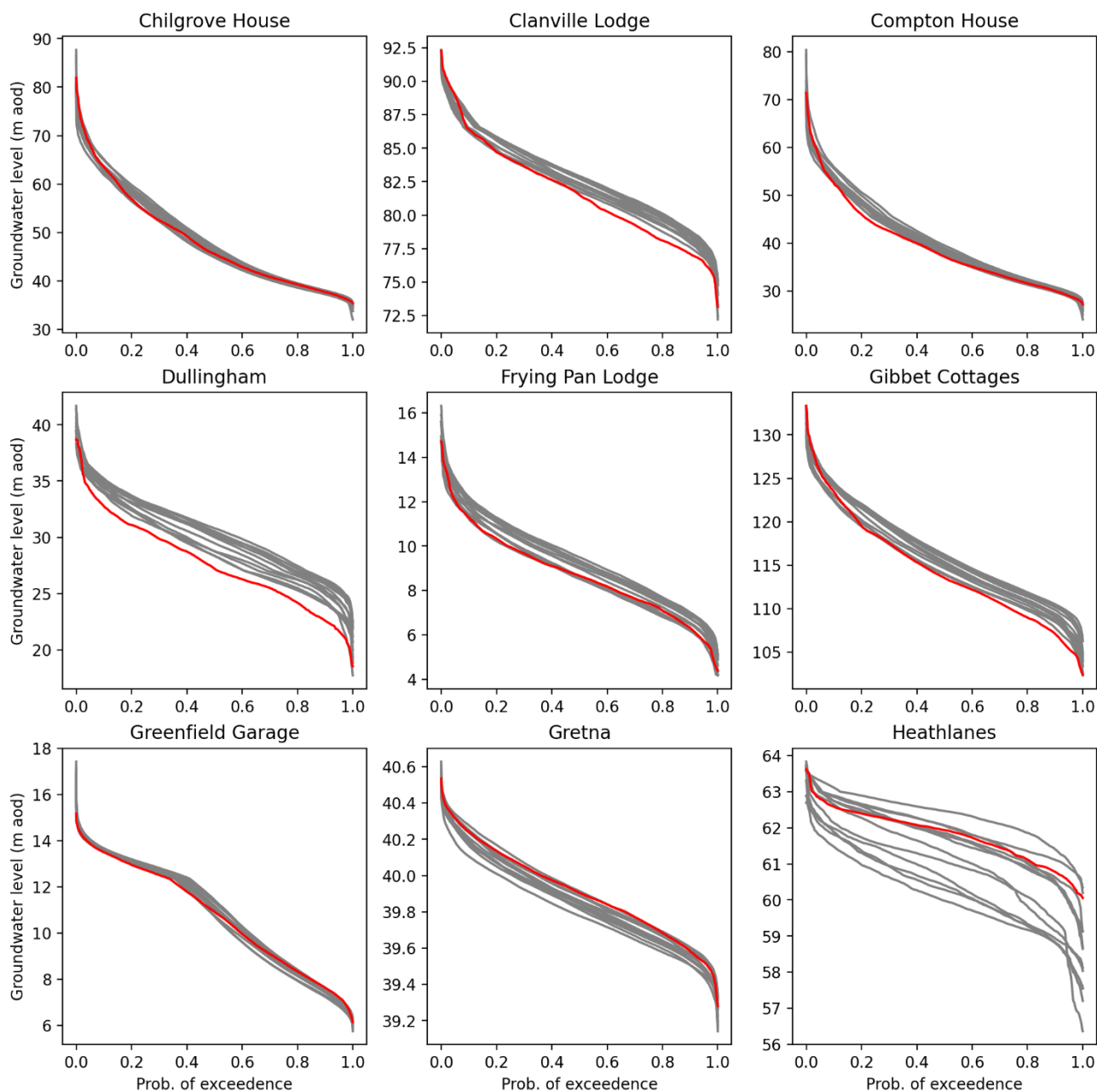
221 *Figure S10 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the*  
 222 *simrcm (grey lines) simobs (red line) simulations.*



223

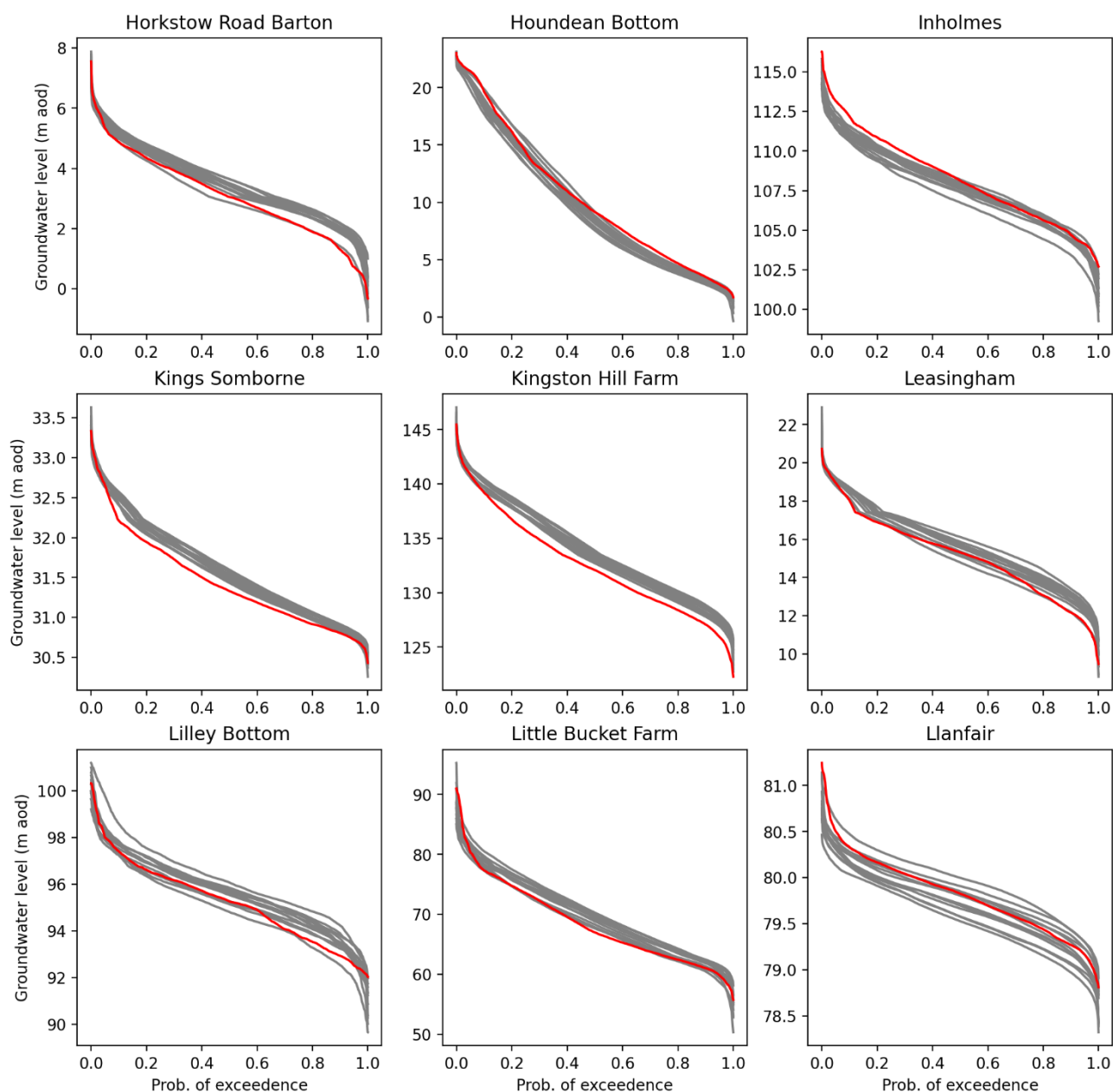
224 *Figure S11 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the*  
 225 *simrcm (grey lines) simobs (red line) simulations.*





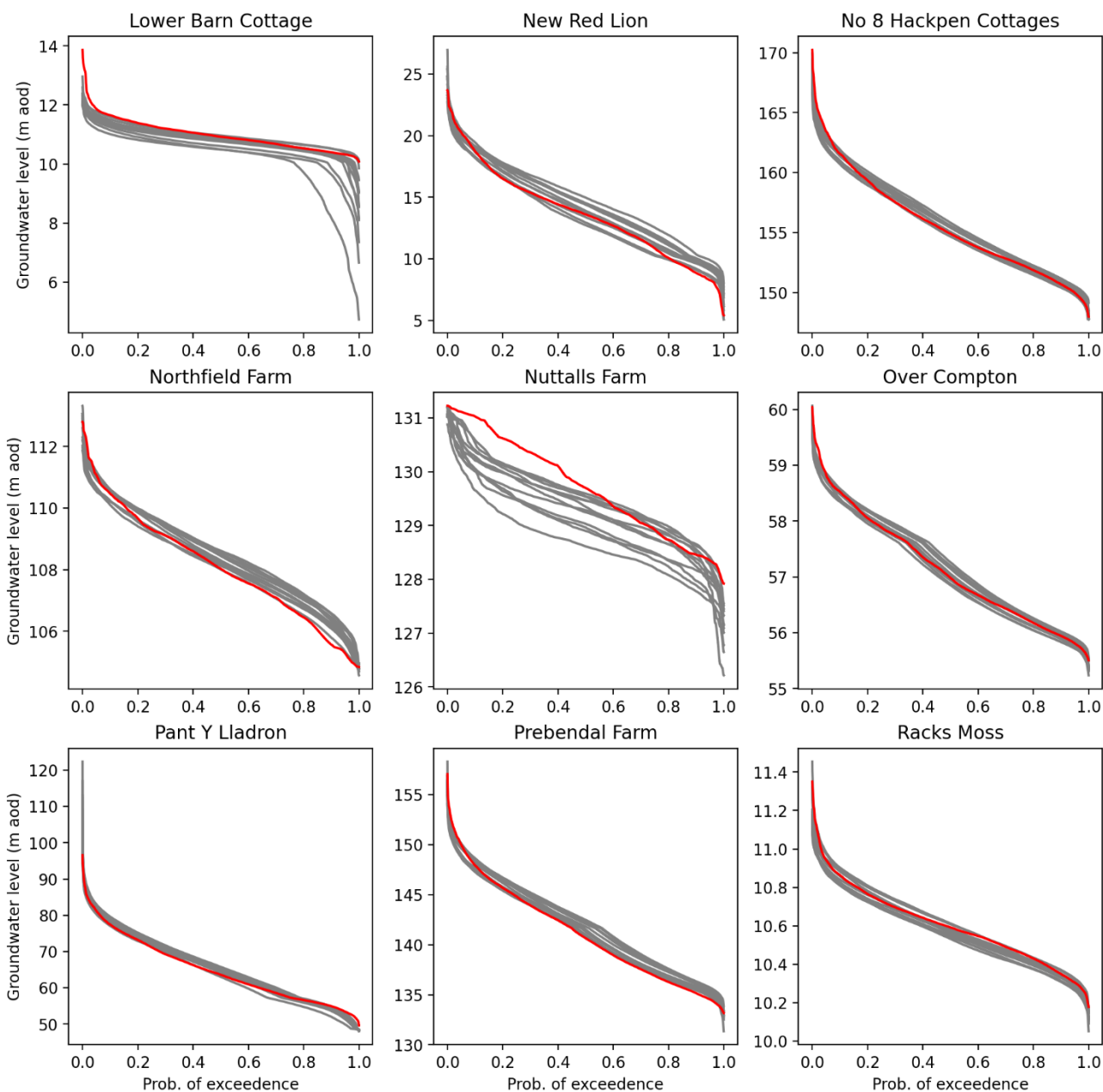
226

227 *Figure S12 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the*  
 228 *simrcm (grey lines) simobs (red line) simulations.*



229

230 *Figure S13 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the*  
 231 *simrcm (grey lines) simobs (red line) simulations.*



232

233 *Figure S14 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the*  
 234 *simrcm (grey lines) simobs (red line) simulations.*

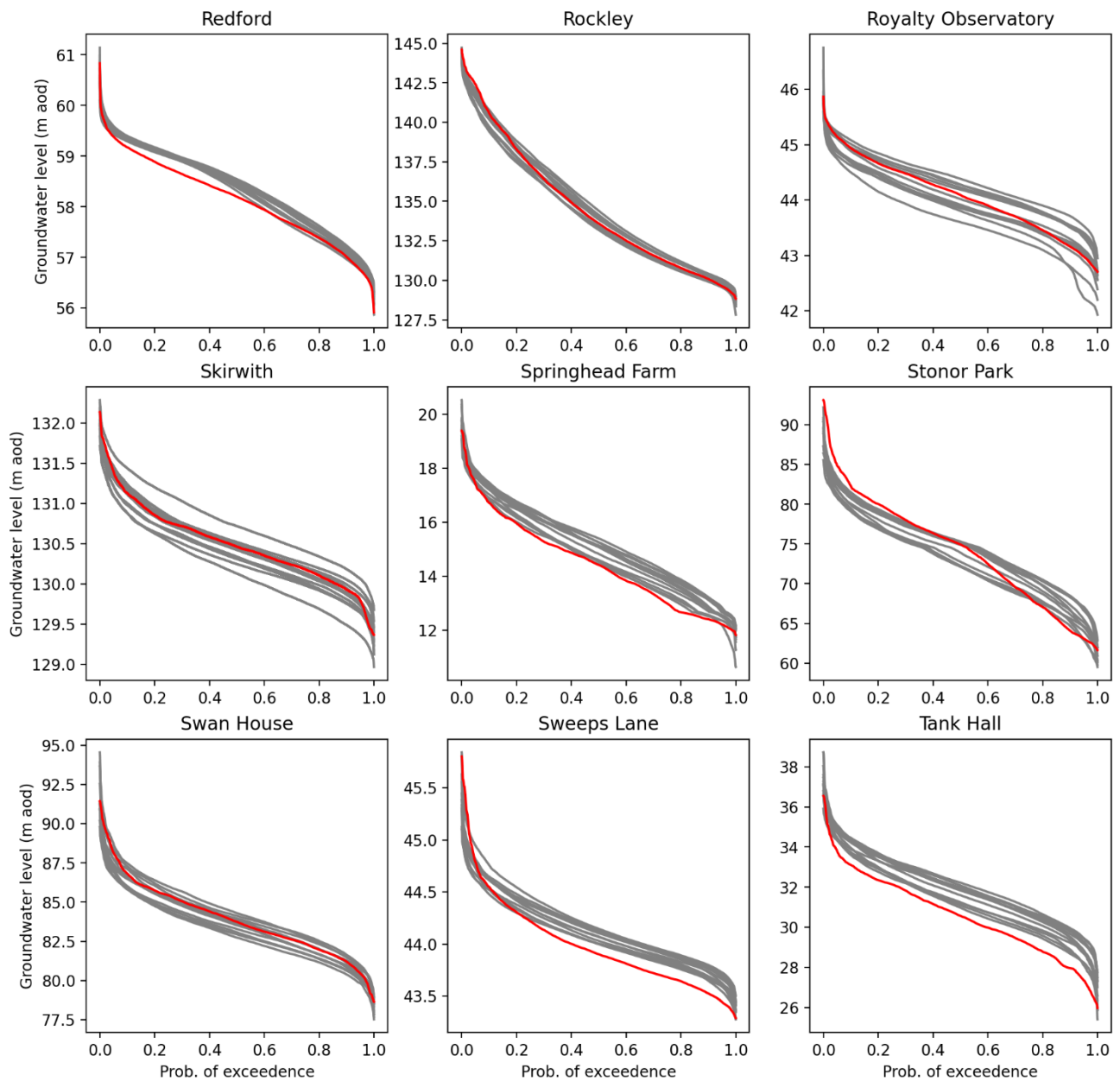
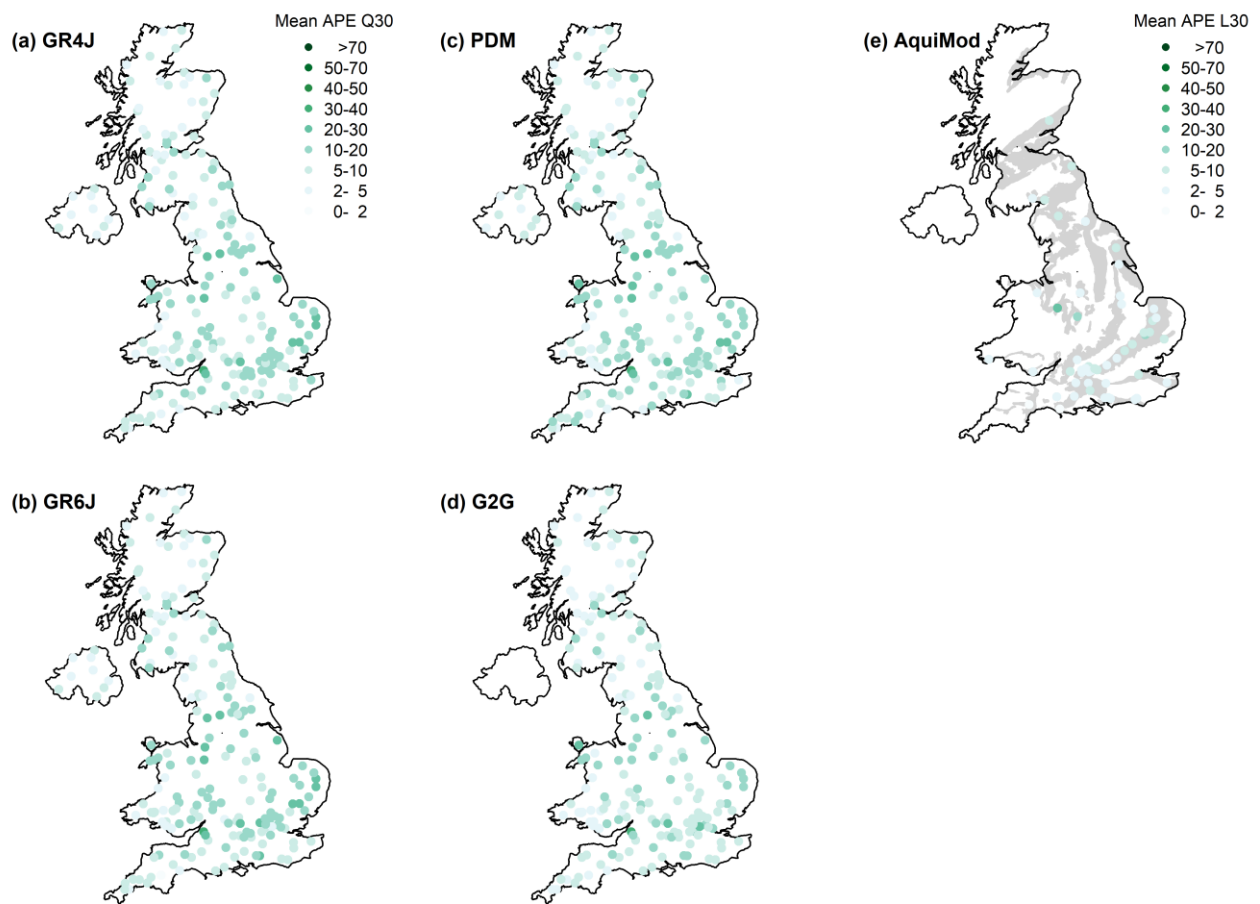
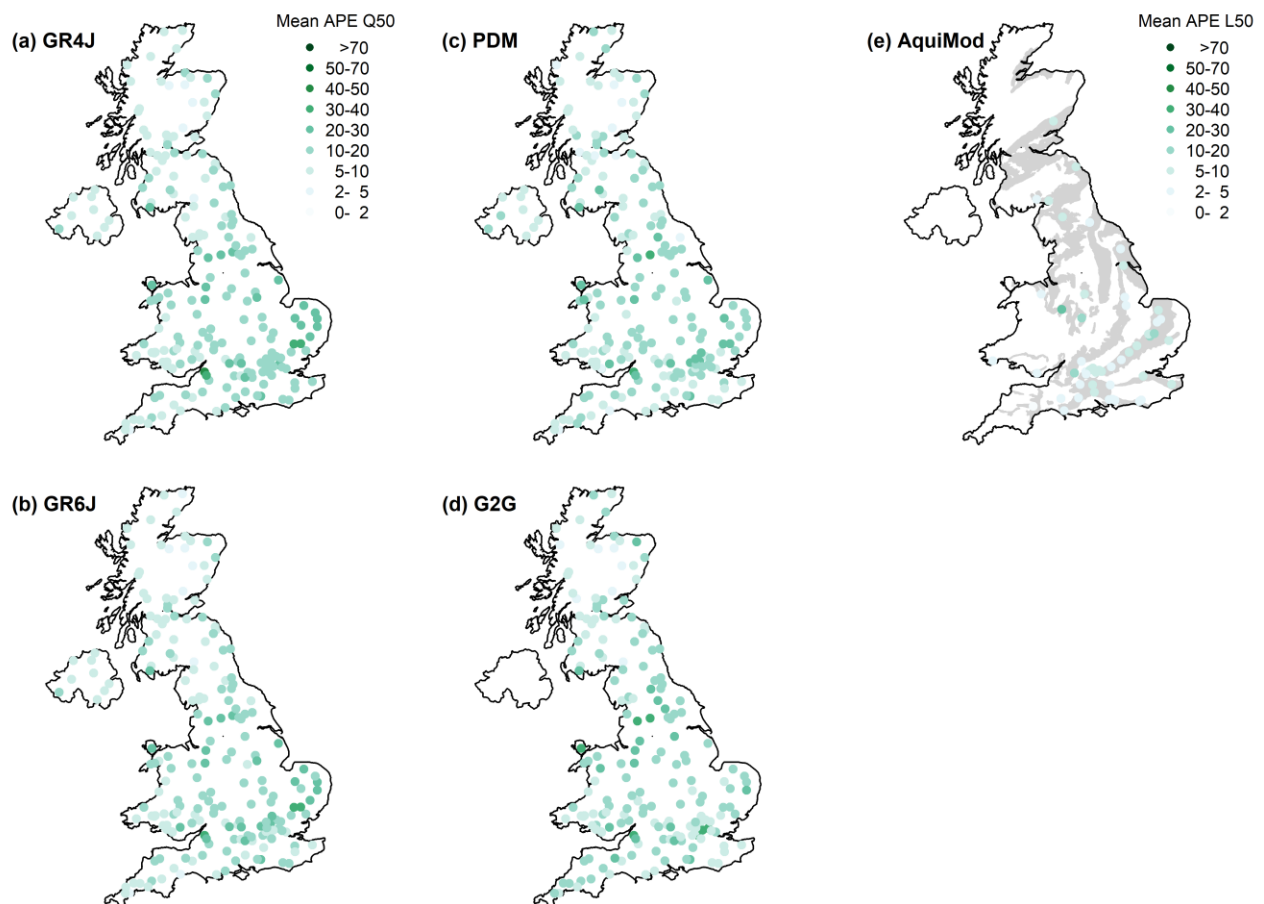


Figure S15 – Groundwater level duration curves (GLDCs) for the period 1989-2018 using the *simrcm* (grey lines) *simobs* (red line) simulations.



247

248      *Figure S16 -- Comparison of river flows and groundwater levels exceeded 30% of the time (Q30) in*  
249      *model observations and RCM ensemble baseline, 1989-2018. Colour scale indicates the mean of 12*  
250      *absolute percent errors (APEs) between Q30 in model observations and Q30 in each of 12*  
251      *ensemble members. Results are presented for each of the four hydrological models and one borehole*  
252      *model: (a) GR4J; (b) GR6J; (c) PDM; (d) G2G; (e) AquiMod. Note: AquiMod levels expressed relative*  
253      *to the minimum level prior to calculating APEs, to remove influence of arbitrarily high datums.*



254

255 *Figure S17 -- Comparison of river flows and groundwater levels exceeded 50% of the time (Q50) in*  
 256 *model observations and RCM ensemble baseline, 1989-2018. Colour scale indicates the mean of 12*  
 257 *absolute percent errors (APEs) between Q50 in model observations and Q50 in each of 12 RCM*  
 258 *ensemble members. Results are presented for each of the four hydrological models and one borehole*  
 259 *model: (a) GR4J; (b) GR6J; (c) PDM; (d) G2G; (e) AquiMod. Note: AquiMod levels expressed relative*  
 260 *to the minimum level prior to calculating APEs, to remove influence of arbitrarily high datums.*



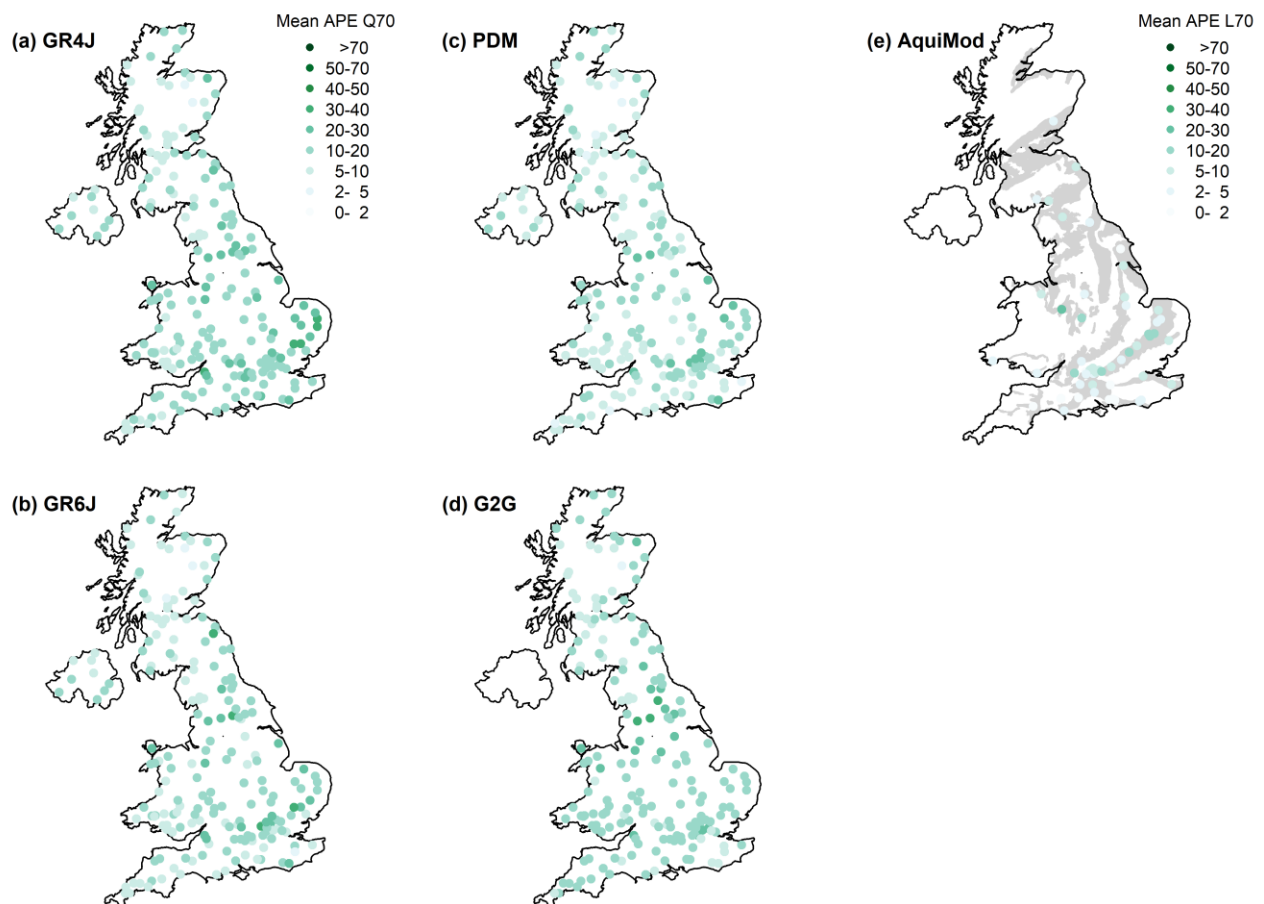


Figure S18 -- Comparison of river flows and groundwater levels exceeded 70% of the time (Q70) in model observations and RCM ensemble baseline, 1989-2018. Colour scale indicates the mean of 12 absolute percent errors (APEs) between Q70 in model observations and Q70 in each of 12 RCM ensemble members. Results are presented for each of the four hydrological models and one borehole model: (a) GR4J; (b) GR6J; (c) PDM; (d) G2G; (e) AquiMod. Note: AquiMod levels expressed relative to the minimum level prior to calculating APEs, to remove influence of arbitrarily high datums.

## Supplementary References

Santos, L., Thirel, G., and Perri, C.,: Technical note: Pitfalls in using log-transformed flows within the KGE criterion. *Hydrology and Earth System Sciences*, 22, 4583–4591, 2018.

<https://doi.org/10.5194/hess-22-4583-2018>