
1 Performance dependence of multi-model combination methods
2 on hydrological model calibration strategy and ensemble size

3

4 Yongjing Wan ^a, Jie Chen ^{a,b*}, Chong-Yu Xu ^c, Ping Xie ^a, Wenyan Qi ^a, Daiyuan Li ^d, Shaobo Zhang ^a

5 ^a State Key Laboratory of Water Resources and Hydropower Engineering Science, Wuhan
6 University, Wuhan 430072, China

7 ^b Hubei Key Laboratory of Water System Science for Sponge City Construction, Wuhan University,
8 Wuhan 430072, China

9 ^c Department of Geosciences, University of Oslo, P.O. Box 1047 Blindern, N-0316 Oslo, Norway

10 ^d Hydrology and Water Resources Department, Nanjing Hydraulic Research Institute, Nanjing
11 210029, China

12 * Corresponding author: Jie Chen (jiechen@whu.edu.cn)

13

14 **Abstract:**

15 The multi-model combination is a technique to improve the performances of hydrological
16 streamflow simulations. An area that has not been investigated much is the performance dependence
17 of combination techniques on the hydrological model calibration strategy and ensemble size. This
18 study aims at investigating the joint effect of the hydrological models, calibration strategies and
19 ensemble sizes on combination abilities for selecting the most appropriate multi-model combination
20 method. The ensemble members were constructed by applying four hydrological models and four
21 objective functions over 383 catchments in China. The ensemble members were combined by using
22 nine commonly used methods, which are Equal Weights (EWA), Akaike Information Criterion
23 (AICA), Bayes Information Criterion (BICA), Bates and Granger (BGA), Granger Ramanathan A,
24 B, and C (GRA, GRB, and GRC), Bayesian Model Averaging (BMA) and Multi-model Super

25 Ensemble (MMSE). The GRC is found as the best multi-model combination method for
26 hydrological simulations. Adding ensemble members by either multiple hydrological models or
27 calibration strategies could help to improve the simulation abilities. Specifically, the increase of
28 ensemble members can obviously enhance the performance of multi-model combinations when the
29 ensemble size is less than six, while only limited improvement is achieved when the ensemble size
30 is more than nine. The combination of ensemble members with various calibration strategies is hard
31 to compensate for the weakness of hydrological model structures. As well, the application of a single
32 calibration strategy in ensemble members only emphasizes single discharge periods and neglects
33 other important discharge periods. This study found that various models with different objective
34 functions are more robust and efficient. The combination performs better than any individual model
35 in terms of Nash–Sutcliffe efficiency (NSE) for approximately 70% catchments, but the multi-
36 model combination is less efficient in terms of low-flow simulations.

37 **Keywords:** multi-model combination method; ensemble modeling; calibration strategy;
38 ensemble size; joint effect

39

40 1. Introduction

41 Hydrological models are essential tools for addressing a wide spectrum of hydrological and
42 water resources problems, including water resources planning, drought and flood control, simulation
43 at ungauged locations, and impact studies for climate or land-use changes (Kotsuki et al., 2014;
44 Kudo et al., 2017; Lane et al., 2019; Wang et al., 2020; Zhang et al., 2014). Over the last decades, a
45 large number of hydrological models have been developed, ranging from lumped conceptual models
46 to physically-based distributed models (Arnold et al., 1998; Chiew et al., 2002; Edijatno et al., 1999;
47 Liang et al., 1994; Xu, 2021; Zhao et al., 1980). The performance of those models varies for diverse
48 catchments characterized by different climate, land use and topography, according to the strengths
49 and weaknesses of the modeling (Mendoza et al., 2016; Pechlivanidis et al., 2011; Vansteenkiste et
50 al., 2014a; Zhang et al., 2020). It is hard to determine a priori which model is most appropriate for
51 a given application over widely differing characteristics of catchments. A single model is not able
52 to consistently outperform the others for all catchment characteristics and heterogeneous
53 climatology (Arsenault et al., 2015; Kumar et al., 2015; Velazquez et al., 2010). Several studies
54 (Arsenault et al., 2015; Velázquez et al., 2011; Zhang et al., 2020) found that multi-model
55 combinations are more robust and efficient than their individual members with the concept of using
56 ensemble to reducing errors with an optimal bias and variance trade-off.

57 A wide range of methods can be used to generate a multi-model combination solution. The
58 simplest example is the calculation of the arithmetic mean of the input models (commonly referred
59 to as the Equal Weights Averaging (EWA)). More sophisticated techniques employ weighted
60 schemes, with differential weightings applied to each input model reflecting their relative

61 advantages or limitations. There are some popular techniques used to obtain the optimal set weights
62 for multi-model combinations like multiple linear regression (Arsenault et al., 2015; Granger and
63 Ramanathan, 1984; Kumar et al., 2015), machine learning algorithms (Jeong and Kim, 2009;
64 Shamseldin et al., 1997; Zaherpour et al., 2019), Bayesian model averaging (Neuman, 2003) and
65 Information Criterion Averaging (Akaike, 1974; Schwarz, 1978). The challenge in ensemble
66 modeling is to determine the ensemble size and to identify the best averaging method (Arsenault et
67 al., 2015; Buizza and Palmer, 1998; Kumar et al., 2015).

68 Many studies (Jeong and Kim, 2009; Shamseldin et al., 1997; Sun and Trevor, 2018; Zaherpour
69 et al., 2019) have attempted to identify the best multi-model combination method for hydrological
70 simulations. Shamseldin et al. (1997) applied three methods (simple arithmetic mean, constrained
71 ordinary least-squares weighting, and neural network) to combine four hydrological models for 11
72 catchments and found the constrained ordinary least-squares weighting and neural network are more
73 robust and efficient than the simple arithmetic mean in terms of the Nash–Sutcliffe efficiency (NSE).
74 Broderick et al. (2016) analyzed the performance of four ensemble averaging techniques using four
75 hydrological models for 37 Irish catchments. They concluded that GRA is the best ensemble
76 averaging technique, and the averaging methods performed better for the NSE as opposed to bias
77 metrics. In addition, some studies investigated the effect of ensemble size on the performance of
78 multi-model ensemble simulations. For example, Arsenault et al. (2015) compared nine multi-model
79 averaging approaches using 12 hydrographs (4 models \times 3 metrics) over 429 catchments. They
80 found that GRC performs better than other averaging methods and no catchment requiring more
81 than seven ensemble members to maximize the NSE with this method. Kumar et al. (2015)
82 compared ten different multi-model ensemble methods using eight hydrological models to select

83 the best multi-model ensemble method for the discharge estimation over a catchment of the
84 Mahanadi river basin in India. They showed that the constrained multiple linear regression is the
85 most suitable multi-model ensemble method in terms of NSE, root mean square error (RMSE) and
86 Pearson's correlation coefficient (R), and five ensembles show the best performance for the study
87 area.

88 The method of multi-model combination is usually used to extract as much information as
89 possible from a group of existing models, which may produce a better overall simulation, as each
90 simulation of the group provides specific information. In addition to the selection of the hydrological
91 models, the process of parameter identification is also a crucial step in streamflow modeling. The
92 calibration strategies (which here mean choice of objective functions) reflect the goodness of fitting
93 between hydrological model simulations and observations, which can substantially influence the
94 model parameters and the streamflow projections (Krysanova et al., 2018; Lane et al., 2019;
95 Mizukami et al., 2019; Seiller et al., 2017). The hydrological models with various structures and
96 calibration strategies have certain capacities to predict the streamflow. In general, the hydrological
97 simulations achieve maximum accuracy in terms of specific hydrological properties using a
98 particular metric, but that might limit the modeling skill in other aspects (Arsenault et al., 2015;
99 Mizukami et al., 2019; Seiller et al., 2017). For example, the most widely used calibration strategies,
100 such as Nash–Sutcliffe efficiency (NSE) and Kling–Gupta efficiency (KGE), emphasize high flow
101 events and their timing (Gupta et al., 2009; Mizukami et al., 2019; Nash and Sutcliffe, 1970). Those
102 metrics calculated on the natural logarithm of the flow values put more weight on low flows
103 (Pushpalatha et al., 2012; Seiller et al., 2017). The use of multi-model combination scheme is
104 expected to benefit from the variation of the parameter sets derived from objective functions targeted

105 at different hydrological processes to produce a better overall simulation. Although the influence of
106 calibration strategies exists in hydrological modeling, the impacts of hydrological model calibration
107 strategies on the performance of multi-model combinations and their joint effects with ensemble
108 sizes are not clear. Moreover, there is no consensus in the hydrological community in terms of the
109 selection of particular multi-model ensemble sizes to ensure good model performances.

110 This study aims to investigate the joint effect of ensemble sizes and hydrological model
111 calibration strategies on combination abilities for selecting the most appropriate multi-model
112 combination method. Specifically, nine commonly used multi-model combination techniques are
113 compared over 383 catchments in China using ensemble members derived from 4 hydrological
114 models calibrated with 4 objective functions. The rest of the paper is organized as follows. Section
115 2 presents a brief introduction of the study data and the methodology, including the hydrological
116 models, the calibration strategies, the multi-model combination techniques, and the details of the
117 evaluation technique used in this study. The results are presented in Section 3, followed by the
118 discussion and conclusion in Section 4.

119 2. Data and Methodology

120 2.1 Study region and data

121 This study used a gridded meteorological dataset ($0.5^\circ \times 0.5^\circ$) over China for the period of
122 1961–2016, which contains four climate variables, including daily precipitation, daily maximum,
123 minimum and mean air temperatures to represent observed data. This dataset was generated from
124 2,472 in-situ gauge stations across China by thin-plate spline interpolation method and GTOPO30
125 (Global 30 Arc-Second Elevation) data sampling and is considered as the latest gridded

126 meteorological data with the highest spatial resolution in China. This dataset has been commonly
127 used in many hydro-climatological studies in China (Gu et al., 2020; Li et al., 2019; Yin et al., 2020),
128 and downloaded from the China Meteorological Data Sharing Service System
129 (<http://www.cma.gov.cn>).

130 The daily streamflow series over 383 catchments in China were used (Figure 1). These
131 catchments with a wide range of climatic conditions and hydrological regimes span over all the nine
132 major river basins in China. Based on climate type and physical geography, this study region was
133 divided into four major climate regions: continental climate zone of Northwest (NW), the highland
134 climate zone of Southwest (Tibetan Plateau, SW), the temperate monsoon region of Northeast (NE),
135 and the tropical and subtropical monsoon region of Southeast (SE) (Ding, 2013; Wu et al., 2016).
136 NE is the driest region according to the average aridity index value (Figure 2). The size of the
137 catchments ranges from 612 km² to 995,343 km². The streamflow dataset covers the 1961–2016
138 period with a maximum length of 52 years and a minimum length of 22 years. The average annual
139 precipitation of the catchment varies greatly with clear gradients depending on the region. The mean
140 annual precipitation is more than 1400 mm in the southern region, while it is less than 600 mm in
141 the northern region.

142 <Figure 1>

143 <Figure 2>

144 2.2 Hydrological models

145 A wide range of hydrological models is used for different application purposes. Some studies
146 compared the performance of lumped and distributed models for outlet streamflow simulation and

147 found that two types of models may lead to similar accuracy (Kumar et al., 2015; Lobligeois et al.,
148 2014; Vansteenkiste et al., 2014b), even for quite large catchments (Merz et al., 2009). Considering
149 the scope of this study, the lumped model (the most common model type used by hydrologists for
150 water resources assessment, flood forecasting, and impact of climate change studies) was chosen,
151 and the distributed model with expensive computations was excluded. Four lumped models with
152 different complexity were used, i.e., modèle du Génie Rural à 4 paramètres Journalier (GR4J)
153 (Edijatno et al., 1999; Perrin et al., 2003), hydrological model of école de technologie supérieure
154 (HMETS) (Martel et al., 2017), simple HYDROLOG (SIMHYD) (Chiew et al., 2002), and
155 Xinanjiang (XAJ) (Zhao, 1992; Zhao et al., 1980). Those models have been widely used in
156 streamflow simulation and have been shown to be relatively efficient (Arsenault et al., 2015;
157 Broderick et al., 2016; Jones et al., 2006; Liang et al., 2013; Mathevet et al., 2020). Table 1 briefly
158 summarized the basic information for these hydrological models.

159 The four models have different numbers of parameters and are different in model structures
160 and underlying mechanisms. For example, the physical process is described in more detailed and
161 complex mechanisms in HMETS, SIMHYD and XAJ than in the most parsimonious structure GR4J
162 with only four free parameters. The main feature of the runoff generation of HMETS and XAJ is
163 using the saturation excess flow mechanism based on the soil moisture content of the aeration zone
164 reaching its field capacity. While SIMHYD considers both infiltration excess runoff and saturation
165 excess runoff in streamflow production calculated by an interception store, a nonlinear soil moisture
166 store. For the simulation of evaporation, XAJ uses a three-layer evaporation model, while HMETS
167 and SIMHYD use a one-layer model. Additionally, GR4J and HMETS consider the incorporation
168 of groundwater exchange by surface water–groundwater interaction functions, but XAJ and

169 SIMHYD do not have this consideration.

170 Since GR4J, XAJ and SIMHYD do not simulate snow accumulation or melt processes, a snow
171 module (CEMANEIGE) with 2 free parameters (Valéry et al., 2014) was combined with the original
172 model to make it applicable in seasonally snow-covered catchments in northern China. The basic
173 inputs of these four models are catchment-averaged precipitation and temperature/potential
174 evapotranspiration over the entire basin for the computation of discharge. The potential
175 evapotranspiration was estimated using a temperature-based method proposed by Oudin et al.
176 (2006a).

177 <Table 1>

178 2.3 Calibration and evaluation metrics

179 This study used four objective functions to calibrate the four hydrological models over 383
180 catchments. The four calibration strategies are the widely-used Nash–Sutcliffe efficiency (NSE)
181 (Nash and Sutcliffe, 1970), the NSE computed on natural logarithm and square root of the flow
182 values (NSE(ln) and NSE(sqrt)), and percent bias (PBias). NSE, NSE(ln) and NSE(sqrt) all range
183 from negative infinity to 1, with a value of 1 indicating perfect fitting. The PBias' value being closer
184 to zero indicates the better simulating performances.

185 Different calibration strategies were included in the combination since they emphasize
186 different aspects of hydrological streamflow properties. The original NSE without discharge
187 transformation puts great emphasis on high flows (Li et al., 2019; Mizukami et al., 2019). The
188 natural logarithm of discharge transformation (NSE(ln)) is to optimize the performance for low flow
189 segments (Seiller et al., 2017). The analysis of NSE(sqrt) well balances simulated streamflow

190 without too much emphasis on low or high flow (Oudin et al., 2006b), and the PBias emphasizes
191 the total water balance. Regardless of the objective function, the hydrological models optimized the
192 parameter using the Shuffled Complex Evolution-University of Arizona (SCE-UA) algorithm
193 (Duan et al., 1992). A cross-validation method that divided the complete record of each catchment
194 (Arsenault et al., 2015; Yang et al., 2020) into odd and even years for model calibration and
195 validation was used, to reduce the influence from the non-stationarity of hydro-climatological
196 conditions. In this process, the first year in the calibration period was used for model warm-up.

197 In addition to using NSE, NSE(ln), NSE(sqrt) and PBias as evaluation metrics to represent the
198 overall performance of simulations, the high and low flows were also analyzed based on the
199 discharge segments of flow duration curves (FDC) to identify the influence of ensemble method on
200 the performance of various flow components. Following previous studies (Laaha and Blöschl, 2006;
201 Pfannerstill et al., 2014; Yilmaz et al., 2008), the flow exceedance probability of 70% was used to
202 represent the low flow, and the mid-flow segment was shifted from 20% to 70%. The very high-
203 flow range was defined between 0% and 5%, and the high-flow range between 5% and 20%. The
204 performance of the model simulations within these FDC segments was analyzed using PBias
205 (PBiasFSV), noted by PBiasFSV-5, PBiasFSV-20, PBiasFSV-mid, and PBiasFSV-low. The basic
206 information of those metrics which were used for evaluating the performance of the different
207 hydrograph phases was shown in Table 2.

208 <Table 2>

209 2.4 Multi-model averaging methods

210 There are several methods available in the literature for developing the multi-model

211 combination. Here, we compared the performance of nine commonly used deterministic ensemble
212 techniques for creating multi-model ensembles. The selected methods include Equal Weights
213 Averaging (EWA), Akaike Information Criterion Averaging (AICA), Bayes Information Criterion
214 Averaging (BICA), Bates and Granger Averaging (BGA), Granger Ramanathan A Averaging (GRA),
215 Granger Ramanathan B Averaging (GRB), Granger Ramanathan C Averaging (GRC), Bayesian
216 Model Averaging (BMA) and Multi-model Super Ensemble (MMSE). The general model for
217 averaging methods can be expressed as:

$$218 \quad Q_{ens} = \sum_{i=1}^n W_i \cdot Q_{simi}$$

219 where Q_{ens} is the ensemble simulation, Q_{simi} is the individual simulation, and W_i is the weight
220 of the ensemble member.

221 A description of the selected averaging methods was given in the Appendix. The basic
222 characteristics of the nine multi-model averaging techniques were summarized in Table 3.

223 <Table 3>

224 3. Results

225 3.1 Performance of the individual members

226 The performances of the 16 ensemble members over 383 catchments in the calibration and
227 validation were analyzed using the evaluation of NSE, NSE(ln), NSE(sqrt), and PBias representing
228 the average simulated abilities and using the PBias of four different FDC segments (PBiasFSV-5,
229 PBiasFSV-20, PBiasFSV-mid and PBiasFSV-low) representing the high and low flows, as presented
230 in Figure 3. The result shows that the evaluation metrics in the validation period are consistent with

231 those in the calibration period, which indicates the hydrological simulations are robust and
232 transferable.

233 The ensemble members with calibration strategies closely related to the evaluation metrics
234 work best, as expected. Models calibrated with NSE(sqrt) maintain good performances in terms of
235 NSE and NSE(ln) evaluation with median values around 0.7 and 0.8. All ensemble members show
236 the absolute value of PBias being less than 10%, but the models calibrated with PBias yield inferior
237 performance in terms of NSE-based evaluation metrics. According to the FDC segment evaluation
238 values (Figures 3e-3h), the models calibrated with NSE show good performance for high flow
239 segments that obtain the best score for 38% and 42% of catchments in terms of the PBiasFSV-5 and
240 PBiasFSV-20, respectively (Table 4). The models calibrated with NSE(ln) put high weight on low
241 flows and obtain the best score at a frequency of 42% catchments. Models calibrated with NSE(sqrt)
242 emphasize middle flows, which obtain the best score at a frequency of 38% catchments. Those
243 results confirmed the specialization of the objective functions for specific parts of the hydrographs.

244 This figure also demonstrates that hydrologic models perform differently with respect to
245 different calibration strategies and evaluation metrics. For example, calibrated with NSE, both
246 HMETS and GR4J with median NSE(ln) value of 0.6 perform better than SIMHYD and XAJ with
247 median NSE(ln) value of 0.3 and 0.5. In addition, HMETS and GR4J tend to overestimate the middle
248 and low flow segments, while the other two models (i.e., SIMHYD and XAJ) tend to underestimate
249 those flows. Comparing the hydrological models, XAJ generally performs best for the evaluation
250 metrics, followed by HMETS, GR4J and SIMHYD.

251 When looking at four sub-regions based on the median value of evaluation metrics (Figure 4),

252 different models show similar spatial patterns in terms of NSE-based metrics with better
253 performances in the wetter southern catchments (SW and SE) than in drier northern catchments
254 (NW and NE). However, the performance of the four hydrological models is not the same. GR4J
255 and SIMHYD perform better than XAJ and HMETS in terms of NSE-based metrics over SW when
256 using NSE(ln) as the objective function.

257 <Figure 3>

258 <Figure 4>

259 <Table 4>

260 3.2 Performance of the multi-model averaging methods

261 The nine multi-model combination methods were applied to calculate the optimal weights
262 based on the observed and the 16 simulated streamflow series. The performance of those methods
263 over 383 catchments in the validation period is presented in Figure 5. The result shows that the
264 differences of the ensemble simulations are evident not just in how well the averaging methods but
265 also in the evaluation metrics used.

266 For NSE metric, almost all multi-model combination methods show similar performances with
267 median values around 0.8 except for EWA, BGA, and BMA with median values less than 0.75. For
268 NSE(ln) metric, EWA, BGA, and BMA perform better than others with median values around 0.8,
269 followed by GRA, GRB and GRC with median values around 0.75. Almost all averaging methods
270 show similar performances in terms of NSE(sqrt) metric with median values exceeding 0.81. For
271 PBias and PBiasFSV-5, GRA, GRB, GRC and MMSE show similar performances and outperform

272 others, followed by AICA and BICA. For PBiasFSV-20, PBiasFSV-mid and PBiasFSV-low, GRC
273 performs the best, followed by GRA, GRB and MMSE. In general, all simulations tend to
274 underestimate the high flow but overestimating the mid and low flows.

275 Among various averaging methods, GRA, GRB, GRC, and MMSE show similar performance,
276 since they derived from the same optimal weighting group. The AICA and BICA perform poorly in
277 terms of NSE(ln) and PBiasFSV-low metrics compared with other averaging methods, as they put
278 almost all weights on the individual member with minimum RMSE. It can also be seen that the
279 BMA, EWA, and BGA without bias-correction show worse performances in terms of PBias-based
280 measures, especially for PBias, PBiasFSV-5 and PBiasFSV-low.

281 Overall, the GRA, GRB, GRC, and MMSE consistently outperformed other combination
282 methods in terms of the eight evaluation metrics, and GRC provided the best performances,
283 especially in terms of PBias and PBiasFSV. Therefore, GRC is considered as the best multi-model
284 averaging method for hydrological simulations in this study.

285 <Figure 5>

286 3.3 Impact of the hydrological model, calibration strategy, and ensemble size on multi- 287 model combinations

288 This section investigated the joint influence of the hydrological models, calibration strategies
289 and ensemble size on the performance of GRC. The selection of ensemble members is based on
290 multiple hydrological models calibrated with a single objective function, a single hydrological
291 model calibrated with multiple objective functions, and multiple hydrological models calibrated
292 with multiple objective functions. Figure 6 shows the median performance value over 383

293 catchments in the validation period to represent the overall ability of those combinations.

294 The performance of multiple hydrological models calibrated with a single objective function
295 varies depending on the selection of objective function and is improved along with the increase of
296 hydrological model numbers (areas numbered 2 in Figure 6). Areas numbered 3 in Figure 6 show
297 the performances of a hydrological model with multiple objective functions. Though using multiple
298 objective functions could improve the combination abilities, the performances still depend on the
299 selection of hydrological models, as each hydrological model has its own advantages and limitations.
300 For example, the combinations based on SIMHYD generally perform worse than others in terms of
301 NSE and NSE(sqrt) but perform better in terms of NSE(ln), PBias and PBiasFSV-low. In addition,
302 a single model is not able to consistently outperform others for all catchments with various hydro-
303 climatic regimes. The combinations of multiple hydrological models calibrated with various
304 objective functions taking advantage of all ensemble members generally perform more robust and
305 efficient than combinations only using a single calibration or a single hydrological model in terms
306 of all NSE-based metrics (areas numbered 4 in Figure 6).

307 Overall, either using multiple hydrological models or multiple calibration strategies could
308 improve the multi-model combination performances in terms of NSE-based metrics but is less
309 efficient in PBias-based evaluation metrics. The above results also indicate that the influence of
310 increasing hydrological models on combination is larger than increasing model calibration strategies,
311 and the effect of the individual simulation is decreasing along with the increase of ensemble
312 members. The qualitative comparison of different combinations and individual simulations would
313 be analyzed later (Section 3.4).

314

<Figure 6>

315 The combinations with ensemble sizes ranging from 2 to 16 ensemble members are generated
316 by re-sampling the ensemble members 100 times from all combination members to investigate the
317 effect of ensemble size on simulating abilities. Figure 7 shows the relationship between the
318 ensemble size and the combination performances. The 0.05 and 0.25 quantiles of NSE, NSE(ln) and
319 NSE(sqrt) values correspond to the poor efficiency modeling, while 0.75 and 0.95 quantiles
320 correspond to simulations with good performance. The median performance is represented by
321 quantile 0.50 (red lines). The absolute values of PBias and four PBiasFSV being infinite to zero
322 indicating better performances. This result indicates that the performance is improved along with
323 the increase in ensemble sizes, and the effect of ensemble size on the low quantile values is larger
324 than that on the high quantile. In general, the performance of the multi-model combinations is
325 sensitive to the ensemble size and the selection of ensemble members when the member is less than
326 six. The influence of enlarging ensemble numbers on simulation performances could be ignored
327 when the ensemble member is more than nine. In other words, when the ensemble member is more
328 than nine, neither enriching the hydrological model nor increasing calibration strategies show
329 limited improvement in simulating abilities, as larger ensemble sizes mean each individual member
330 with smaller weights and limited influence on the combination.

331

<Figure 7>

332 In order to investigate whether all ensemble members contribute to the performance of the
333 combination, the frequency of the individual member being selected in the best-performed
334 combinations of all catchments is calculated and plotted in Figure 8. Here, the ensemble sizes of the

335 combinations are nine since the results from Figure 7 concluded that the effect of increasing
336 ensemble numbers more than nine could be ignored. Figure 8 shows that all individual members
337 have similar chances of being selected in the best-performed combinations in terms of different
338 evaluation metrics, indicating all ensemble members contribute to enhancing the combination
339 simulation.

340 <Figure 8>

341 3.4 Comparison of the multi-model combination and ensemble members

342 The performances of GRC and the ensemble members were compared using the cumulative
343 distributions of NSE, NSE(ln), NSE(sqrt), absolute values of PBias and PBiasFSV over the 383
344 catchments in the validation period (Figure 9). Here, five combinations are represented, including
345 four hydrological models calibrated with a specified objective function (i.e., H4-NSE, H4-NSE(ln),
346 H4-NSE(sqrt) and H4-PBias) and four hydrological models calibrated with four calibration
347 strategies (i.e., H4C4).

348 It is apparent that the multi-model combinations, either using multiple hydrological models or
349 using multiple calibration strategies, outperform the individual member in terms of those
350 performance metrics, except for NSE(ln) and PBiasFSV-low. The combinations of four hydrological
351 models calibrated with four calibration strategies perform the best and show 85% catchments have
352 NSE values exceeding 0.6, and 75% catchments have NSE values exceeding 0.7 (Figure 9a). In
353 addition, 91% catchments have NSE(sqrt) values exceeding 0.6, and 73% catchments have NSE(sqrt)
354 values exceeding 0.7 (Figure 9b). The absolute value of PBias is less than 10 for more than 95%
355 catchments, and that is less than 5 for 78% catchments. In addition, the absolute values of PBiasFSV-

356 5, PBiasFSV-20, and PBiasFSV-mid are less than 10 for 55%, 80% and 80% of catchments,
357 respectively. However, the combinations are less efficient in terms of NSE(ln) and PBiasFSV-low
358 evaluation than the individual members using NSE(ln) calibration, since the averaging methods
359 based on minimum RMSE are more sensitive to flood peaks than low-flow period.

360 <Figure 9>

361 Table 5 shows the rate at which the averaging simulation surpasses the best individual member
362 in terms of four average performance metrics and four FDC segment metrics over the 383
363 catchments in the validation period. Here, the “best individual member” value is selected from the
364 16 members independently for each catchment. The results show that four hydrological models
365 calibrated with four calibrations (H4C4) perform more robustly than other combinations in terms of
366 various performance metrics. The combination of four hydrological models calibrated with a
367 specified objective only performs efficiently for the specific parts of the hydrographs. The
368 combinations perform better than the best individual member for about 70% catchments in terms of
369 NSE and for about 55% catchments in terms of NSE(sqrt). On the contrary, combinations show less
370 efficiency for performance values in terms of NSE(ln) and PBias.

371 Figures 10(a) and 10(b) show the geographic distribution of the NSE value in the validation
372 period from the best individual member and combinations. The spatial patterns are consistent with
373 better performances in the wetter southern catchments than in the drier northern catchments. For
374 catchments with annual precipitation larger than 600 mm, 82% catchments obtain NSE value
375 exceeding 0.7, and the combination outperforms the best individual in 82% cases. However, for
376 catchments with annual precipitation less than 600 mm, only 51% catchments obtain NSE value

377 exceeding 0.7, and the combinations outperform the best individual only in 45% cases (Figures 10c
378 and 10d). Since hydrological simulations are big challenges for the arid catchments, especially when
379 using lumped models, which cannot represent the rainfall and loss variability that tends to be higher.
380 The multi-model combinations have less benefit when the ensemble members cannot accurately
381 simulate the streamflow.

382 <Table 5>

383 <Figure 10>

384 4. Discussion and conclusion

385 The multi-model combination is a technique widely used to improve the performance of
386 hydrological streamflow simulations, as it extracts potentially useful information from a group of
387 existing models. The commonly used multi-model combination usually takes information from
388 different hydrological models while neglects the additional information from calibration strategies.
389 This study investigated the joint effect of hydrological models, calibration strategies and ensemble
390 sizes on combination abilities for selecting the appropriate multi-model combination method. This
391 study compared different multi-model ensemble methods by combining the simulated discharge
392 with four hydrological models and four different calibration strategies.

393 Generally, the hydrological model performances were similar in terms of geographic
394 distribution, with better performances in wetter catchments than drier catchments. Nevertheless,
395 there were some obvious differences between those models. GR4J and HMETS with surface water–
396 groundwater interaction functions performed better than XAJ and SIMHYD in terms of NSE(ln)
397 value. This is consistent with Pushpalatha et al. (2011) and Fleckenstein et al. (2006), who found

398 that the lumped model with the incorporation of groundwater exchange functions improves the
399 predictions of low flows.

400 One of the other goals was to compare the nine multi-model averaging methods. Overall, GRA,
401 GRB, GRC, and MMSE with similar methods for optimal weight showed consistent performances
402 and performed better than other combination methods. GRC with bias-correction showed better
403 performances in terms of PBias and PBiasFSV metrics than other combination methods and was
404 considered as the best multi-model averaging method for hydrological simulations in this study. As
405 AICA and BICA put all weights on the individual model with a minim root-mean-square error, they
406 performed similarly to the best-performing individual model. BGA, without the consideration of the
407 biases from ensemble members, only performed slightly better than EWA, which is consistent with
408 the previous study (Arsenault et al., 2015). This study also found that BMA, which was widely used
409 in previous ensemble studies (Li et al., 2019; Zhang et al., 2020), performs similarly to BGA in
410 terms of PBias-based evaluation metrics. BMA is recommended to apply on ensemble simulations,
411 whereas it is sometimes difficult to effectively remove bias from the predictions of complex
412 streamflow simulation (Madadgar and Moradkhani, 2014), which is related to the poor performance
413 in terms of PBias evaluations.

414 In addition, this study found that adding ensemble members by either increasing hydrological
415 models or increasing calibration strategies could improve the simulation abilities, but the
416 combinations of taking advantage of different hydrological models and objective functions were
417 more robust and efficient in terms of different hydrological properties and hydro-climatic regimes.
418 In addition, this study found that the increase of ensemble members could obviously improve the

419 multi-model combination performance when the ensemble size is less than six, but has limited
420 effects when the ensemble size is more than nine. This is consistent with Kumar et al. (2015) and
421 Arsenault et al. (2015), who found that four and seven ensemble members are efficient for multi-
422 model combinations.

423 Comparing the combinations and individual simulations, the combinations of different
424 hydrological models with various objective functions performed better than any individual model
425 for around 70% of 383 catchments in terms of NSE scores. In contrast, the combinations were less
426 efficient than the best individual members in terms of low-flow simulations. This study also found
427 that the multi-model combinations perform better for wetter catchments than for drier catchments.
428 The frequencies of ensemble simulation outperforming the best individual simulation in terms of
429 NSE were 82% catchments in wetter regions with precipitation more than 600 mm/year and 45%
430 catchments in drier regions with precipitation less than 600 mm/year.

431 Acknowledgments

432 This work was partially supported by the National Key Research and Development Program of
433 China (No. 2017YFA0603704), the National Natural Science Foundation of China (Grant No.
434 52079093), the Hubei Provincial Natural Science Foundation of China (Grant No. 2020CFA100),
435 and the Overseas Expertise Introduction Project for Discipline Innovation (111 Project) funded by
436 Ministry of Education and State Administration of Foreign Experts Affairs P.R. China (Grant No.
437 B18037). The authors wish to thank the China Meteorological Data Sharing Service System for
438 providing gauged precipitation for China.

439 Appendix: Multi-model ensemble methods

440 A description of the multi-model averaging methods is presented here.

441 1. Equal Weights Averaging (EWA)

442 In the EWA (Shamseldin et al., 1997), the equal weight is simply assigned to each of the
443 ensemble members. This is expressed mathematically as:

$$444 \quad w = \frac{1}{N} \quad (A1)$$

445 where N is the number of ensemble models.

446 2. Akaike and Bayes information criteria averaging (AICA and BICA)

447 AICA (Akaike, 1974) and BICA (Schwarz, 1978) methods combine ensemble members based
448 on both performance and model complexity. Weight represents a trade-off between reducing the
449 simulated error while tending toward less complex, which is calculated as:

$$450 \quad w = \frac{\exp\left(-\frac{1}{2}I\right)}{\sum_{i=1}^N \exp\left(-\frac{1}{2}I_i\right)} \quad (A2)$$

451 where I is the information criterion estimated based on the mean of the logarithm of the ensemble
452 member variances and the number of calibration parameters, and is calculated as:

$$453 \quad I = -2 \log(L) + q(p) \quad (A3)$$

454 where L and $q(p)$ are the maximum likelihood of ensemble member and the penalty term,
455 respectively.

456 The difference between the AICA and BICA methods lies in the penalty term calculation. The
457 penalty terms of AICA and BICA are estimated by equations (A4) and (A5):

$$458 \quad q = 2p \quad (A4)$$

459 $q = p \log(k)$ (A5)

460 where p donates the number of calibrated parameters in the members, and k donates the sample size
461 (here is the number of time steps).

462 3. Bates and Granger Averaging (BGA)

463 The BGA (Bates and Granger, 1969) method aims to produce a combined ensemble by
464 minimizing the Root Mean Square Error (RMSE) between the observations and simulations. The
465 model weighting vector is estimated according to:

466
$$W = \frac{1/RMSE^2}{\sum_{i=1}^N 1/RMSE_i^2}$$
 (A6)

467 4. Granger Ramanathan A, B and C (GRA, GRB and GRC)

468 The GRA (Granger and Ramanathan, 1984) approach minimizes the RMSE setting weights
469 based on the ordinary least squares (OLS) algorithm. The weights are estimated by:

470
$$W = (Q_{sim}^T Q_{sim})^{-1} Q_{sim}^T Q_{obs}$$
 (A7)

471 where Q_{obs} is the observation. The GRB variant is similar to the GRA method, but the OLS
472 algorithm is constrained such that the sum of the weights to unity. The GRC variant is unconstrained,
473 but the averaged streamflow is bias-corrected through the use of a constant term.

474 5. Bayesian Model Averaging (BMA)

475 BMA (Neuman, 2003) determines the weights of each member through the use of the ensemble
476 members' probability distribution functions (PDFs). The combined distribution is bias-corrected,
477 and the difference between the distributions is minimized. According to BMA, the posterior

478 probability of the predictand (y) is described as:

479
$$p(y | Q_{obs}) = \sum_{i=1}^N P(Q_{simi} | Q_{obs})p(y | Q_{simi}, Q_{obs}) \quad (A8)$$

480 where $p(y | Q_{simi}, Q_{obs})$ is the posterior predictive distribution of y on the condition of the given
481 sample Q_{obs} and each individual model Q_{simi} ; $P(Q_{simi} | Q_{obs})$ is the optimal model on the
482 condition of the given sample Q_{obs} denoted the weight (w_i) of each ensemble member. The mean
483 and variance of y are given:

484
$$E[y | Q_{obs}] = \sum_{i=1}^N P(Q_{simi} | Q_{obs}) \int_{-\infty}^{+\infty} yp(y | Q_{simi}, Q_{obs})dy = \sum_{i=1}^N w_i\eta_i \quad (A9)$$

485
$$\text{Var}[y | Q_{obs}] = \sum_{i=1}^N w_i(\eta_i - \sum_{i=1}^N w_i\eta_i)^2 + \sum_{i=1}^N w_i\sigma_i^2 \quad (A10)$$

486
$$w_i = P(Q_{simi} | Q_{obs}) \quad (A11)$$

487 where w_i is the weight of the ensemble member and sum to unite; η and σ are the expectation
488 and variance of y , respectively, on the condition of the given sample Q_{obs} .

489 6. Multi-model Super Ensemble (MMSE)

490 MMSE (Krishnamurti et al., 2000) uses the logic of bias reduction along with variance
491 reduction through using the mean of observational value and the combination of ensemble member,
492 respectively. According to this method, the ensemble discharge is estimated as:

493
$$Q_{MMSE,j} = \bar{Q}_{obs} + \sum_{i=1}^N w_i[Q_{sim,i,j} - \bar{Q}_{sim,i}] \quad (A12)$$

494 where w_i is the weight of i th model, which is estimated by the unconstrained least square technique
495 (Eq. (A1)).

496 References

497 Akaike, H., 1974. A new look at the statistical model identification. IEEE Transactions on Automatic
498 Control, 19(6): 716-723. <https://doi.org/10.1109/TAC.1974.1100705>

499 Arnold, J.G., Srinivasan, R., Muttiah, R.S., Williams, J.R., 1998. LARGE AREA HYDROLOGIC
500 MODELING AND ASSESSMENT PART I: MODEL DEVELOPMENT1. JAWRA Journal of the
501 American Water Resources Association, 34(1): 73-89. [https://doi.org/10.1111/j.1752-
502 1688.1998.tb05961.x](https://doi.org/10.1111/j.1752-1688.1998.tb05961.x)

503 Arsenault, R., Gatién, P., Renaud, B., Brissette, F., Martel, J., 2015. A comparative analysis of 9 multi-
504 model averaging approaches in hydrological continuous streamflow simulation. Journal of
505 Hydrology, 529: 754-767. <http://dx.doi.org/10.1016/j.jhydrol.2015.09.001>

506 Bates, J.M., Granger, C.W.J., 1969. The Combination of Forecasts. Journal of the Operational Research
507 Society, 20(4): 451-468. <https://doi.org/10.1057/jors.1969.103>

508 Broderick, C., Matthews, T.K.R., Wilby, R.L., Bastola, S., Murphy, C., 2016. Transferability of
509 hydrological models and ensemble averaging methods between contrasting climatic periods. Water
510 Resources Research, 52(10): 8343-8373. <https://doi.org/10.1002/2016WR018850>.

511 Buizza, R., Palmer, T.N., 1998. Impact of Ensemble Size on Ensemble Prediction. Monthly Weather
512 Review, 126(9): 2503-2518. [https://doi.org/10.1175/1520-
513 0493\(1998\)126<2503:IOESOE>2.0.CO;2](https://doi.org/10.1175/1520-0493(1998)126<2503:IOESOE>2.0.CO;2)

514 Chiew, F.H.S., Peel, M.C., Western, A.W., Singh, V.P., Frevert, D., 2002. Application and testing of the
515 simple rainfall-runoff model SIMHYD. Water Resources Publications, Colorado, USA.

516 Ding, Y., 2013. China Climate. Science Press, Beijing, China.

517 Duan, Q., Sorooshian, S., Gupta, V., 1992. Effective and efficient global optimization for conceptual
518 rainfall-runoff models. Water Resources Research, 28(4): 1015-1031.
519 <https://doi.org/10.1029/91WR02985>

520 Edijatno, De Oliveira Nascimento, N., Yang, X., Makhlof, Z., Michel, C., 1999. GR3J: a daily
521 watershed model with three free parameters. Hydrological Sciences Journal, 44(2): 263-277.
522 <https://doi.org/10.1080/02626669909492221>

523 Fleckenstein, J.H., Niswonger, R.G., Fogg, G.E., 2006. River-aquifer interactions, geologic heterogeneity,
524 and low-flow management. Ground Water, 44(6): 837-852. [https://doi.org/10.1111/j.1745-
525 6584.2006.00190.x](https://doi.org/10.1111/j.1745-6584.2006.00190.x)

526 Granger, C.W.J., Ramanathan, R., 1984. Improved methods of combining forecasts. Journal of
527 Forecasting, 3(2): 197-204. <https://doi.org/10.1002/for.3980030207>

528 Gu, L. et al., 2020. On future flood magnitudes and estimation uncertainty across 151 catchments in
529 mainland China. International Journal of Climatology, n/a(n/a). <https://doi.org/10.1002/joc.6725>

530 Gupta, H.V., Kling, H., Yilmaz, K.K., Martinez, G.F., 2009. Decomposition of the mean squared error
531 and NSE performance criteria: Implications for improving hydrological modelling. Journal of
532 Hydrology, 377(1): 80-91. <https://doi.org/10.1016/j.jhydrol.2009.08.003>

533 Jeong, D.I., Kim, Y.-O., 2009. Combining single-value streamflow forecasts – A review and guidelines
534 for selecting techniques. Journal of Hydrology, 377(3): 284-299.
535 <https://doi.org/10.1016/j.jhydrol.2009.08.028>

536 Jones, R.N., Chiew, F., Boughton, W.C., Lu, Z.J.A.i.W.R., 2006. Estimating the sensitivity of mean
537 annual runoff to climate change using selected hydrological models. Advances in Water Resources,
538 29(10): 1419-1429. <https://doi.org/10.1016/j.advwatres.2005.11.001>

539 Kotsuki, S., Tanaka, K., Watanabe, S., 2014. Projected hydrological changes and their consistency under
540 future climate in the Chao Phraya River Basin using multi-model and multi-scenario of CMIP5

541 dataset. *Hydrological Research Letters*, 8(1): 27-32. <https://doi.org/10.3178/hr1.8.27>

542 Krishnamurti, T.N. et al., 2000. Multimodel Ensemble Forecasts for Weather and Seasonal Climate.
543 *Journal of Climate*, 13(23): 4196-4216. [https://doi.org/10.1175/1520-0442\(2000\)013<4196:MEFFWA>2.0.CO;2](https://doi.org/10.1175/1520-0442(2000)013<4196:MEFFWA>2.0.CO;2)

544

545 Krysanova, V. et al., 2018. How the performance of hydrological models relates to credibility of
546 projections under climate change. *Hydrological Sciences Journal*, 63(5): 696-720.
547 <https://doi.org/10.1080/02626667.2018.1446214>

548 Kudo, R., Yoshida, T., Masumoto, T., 2017. Nationwide assessment of the impact of climate change on
549 agricultural water resources in Japan using multiple emission scenarios in CMIP5. *Hydrological*
550 *Research Letters*, 11(1): 31-36. <https://doi.org/10.3178/hr1.11.31>

551 Kumar, A., Singh, R., Jena, P.P., Chatterjee, C., Mishra, A., 2015. Identification of the best multi-model
552 combination for simulating river discharge. *Journal of Hydrology*, 525(525): 313-325.
553 <http://dx.doi.org/10.1016/j.jhydrol.2015.03.060>

554 Laaha, G., Blöschl, G., 2006. Seasonality indices for regionalizing low flows. *Hydrological Processes*,
555 20(18): 3851-3878. <https://doi.org/10.1002/hyp.6161>

556 Lane, R. et al., 2019. Benchmarking the predictive capability of hydrological models for river flow and
557 flood peak predictions across over 1000 catchments in Great Britain. *Hydrology and Earth System*
558 *Sciences*, 23(10): 4011-4032. <https://doi.org/10.5194/hess-23-4011-2019>

559 Li, X., Chen, J., Xu, C., Li, L., Chen, H., 2019. Performance of Post-Processed Methods in Hydrological
560 Predictions Evaluated by Deterministic and Probabilistic Criteria. *Water Resources Management*,
561 33(9): 3289-3302. <https://doi.org/10.1007/s11269-019-02302-y>

562 Liang, X., Lettenmaier, D.P., Wood, E.F., Burges, S.J., 1994. A simple hydrologically based model of
563 land surface water and energy fluxes for general circulation models. *Journal of Geophysical*
564 *Research: Atmospheres*, 99(D7): 14415-14428. <https://doi.org/10.1029/94JD00483>

565 Liang, Z., Wang, D., Guo, Y., Zhang, Y., Dai, R., 2013. Application of Bayesian Model Averaging
566 Approach to Multimodel Ensemble Hydrologic Forecasting. *Journal of Hydrologic Engineering*,
567 18(11): 1426-1436. [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000493](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000493)

568 Lobligeois, F., Andréassian, V., Perrin, C., Tabary, P., Loumagne, C., 2014. When does higher spatial
569 resolution rainfall information improve streamflow simulation? An evaluation using 3620 flood
570 events. *Hydrol. Earth Syst. Sci.*, 18(2): 575-594. <https://hess.copernicus.org/articles/18/575/2014/>

571 Madadgar, S., Moradkhani, H., 2014. Improved Bayesian multimodeling: Integration of copulas and
572 Bayesian model averaging. *Water Resources Research*, 50(12): 9586-9603.
573 <https://doi.org/10.1002/2014WR015965>

574 Martel, J.-L., Demeester, K., Brissette, F.P., Arsenaul, R., Poulin, A., 2017. HMET: a simple and efficient
575 hydrology model for teaching hydrological modelling, flow forecasting and climate change impacts.
576 *The International journal of engineering education*, 33(4): 1307-1316.
577 <https://dialnet.unirioja.es/servlet/articulo?codigo=6897050>

578 Mathevet, T., Gupta, H., Perrin, C., Andréassian, V., Le Moine, N., 2020. Assessing the performance and
579 robustness of two conceptual rainfall-runoff models on a worldwide sample of watersheds. *Journal*
580 *of Hydrology*, 585: 124698. <https://doi.org/10.1016/j.jhydrol.2020.124698>

581 Mendoza, P.A. et al., 2016. How do hydrologic modeling decisions affect the portrayal of climate change
582 impacts? *Hydrological Processes*, 30(7): 1071-1095. <https://doi.org/10.1002/hyp.10684>

583 Merz, R., Parajka, J., Blöschl, G., 2009. Scale effects in conceptual hydrological modeling. *Water*
584 *Resources Research*, 45(9): W09405. <https://doi.org/10.1029/2009WR007872>

585 Mizukami, N. et al., 2019. On the choice of calibration metrics for “high-flow” estimation using
586 hydrologic models. *Hydrology and Earth System Sciences*, 23(6): 2601-2614.
587 <https://doi.org/10.5194/hess-23-2601-2019>

588 Nash, J.E., Sutcliffe, J.V., 1970. River flow forecasting through conceptual models part I — A discussion
589 of principles. *Journal of Hydrology*, 10(3): 282-290. [https://doi.org/10.1016/0022-1694\(70\)90255-](https://doi.org/10.1016/0022-1694(70)90255-6)
590 [6](https://doi.org/10.1016/0022-1694(70)90255-6)

591 Neuman, S.P., 2003. Maximum likelihood Bayesian averaging of uncertain model predictions. *Stochastic*
592 *Environmental Research and Risk Assessment*, 17(5): 291-305. [https://doi.org/10.1007/s00477-](https://doi.org/10.1007/s00477-003-0151-7)
593 [003-0151-7](https://doi.org/10.1007/s00477-003-0151-7)

594 Oudin, L., Andréassian, V., Mathevet, T., Perrin, C., Michel, C., 2006a. Dynamic averaging of rainfall -
595 runoff model simulations from complementary model parameterizations. *Water Resources Research*,
596 42(7): W07410. <https://doi.org/10.1029/2005WR004636>

597 Oudin, L., Perrin, C., Mathevet, T., Andréassian, V., Michel, C., 2006b. Impact of biased and randomly
598 corrupted inputs on the efficiency and the parameters of watershed models. *Journal of Hydrology*,
599 320(1): 62-83. <https://doi.org/10.1016/j.jhydrol.2005.07.016>

600 Pechlivanidis, I.G., Jackson, B., McIntyre, N., Wheater, H.S., 2011. Catchment scale hydrological
601 modelling: a review of model types, calibration approaches and uncertainty analysis methods in the
602 context of recent developments in technology and applications. *Global Nest Journal*, 13(3): 193-
603 214.

604 Perrin, C., Michel, C., Andréassian, V., 2003. Improvement of a parsimonious model for streamflow
605 simulation. *Journal of Hydrology*, 279(1): 275-289. [https://doi.org/10.1016/S0022-1694\(03\)00225-](https://doi.org/10.1016/S0022-1694(03)00225-7)
606 [7](https://doi.org/10.1016/S0022-1694(03)00225-7)

607 Pfannerstill, M., Guse, B., Fohrer, N., 2014. Smart low flow signature metrics for an improved overall
608 performance evaluation of hydrological models. *Journal of Hydrology*, 510: 447-458.
609 <http://dx.doi.org/10.1016/j.jhydrol.2013.12.044>

610 Pushpalatha, R., Perrin, C., Le Moine, N., Mathevet, T., Andréassian, V., 2011. A downward structural
611 sensitivity analysis of hydrological models to improve low-flow simulation. *Journal of Hydrology*,
612 411(1): 66-76. <https://doi.org/10.1016/j.jhydrol.2011.09.034>

613 Pushpalatha, R., Perrin, C., Moine, N.L., Andréassian, V., 2012. A review of efficiency criteria suitable
614 for evaluating low-flow simulations. *Journal of Hydrology*, 420-421: 171-182.
615 <https://doi.org/10.1016/j.jhydrol.2011.11.055>

616 Schwarz, G., 1978. Estimating the dimension of a model. *Ann. Statist.*, 6(2): 461-464.

617 Seiller, G., Roy, R., Anctil, F., 2017. Influence of three common calibration metrics on the diagnosis of
618 climate change impacts on water resources. *Journal of Hydrology*, 547: 280-295.
619 <https://doi.org/10.1016/j.jhydrol.2017.02.004>

620 Shamseldin, A.Y., O'Connor, K.M., Liang, G.C., 1997. Methods for combining the outputs of different
621 rainfall-runoff models. *Journal of Hydrology*, 197(1): 203-229. [https://doi.org/10.1016/S0022-](https://doi.org/10.1016/S0022-1694(96)03259-3)
622 [1694\(96\)03259-3](https://doi.org/10.1016/S0022-1694(96)03259-3)

623 Sun, W., Trevor, B., 2018. Multiple model combination methods for annual maximum water level
624 prediction during river ice breakup. *Hydrological Processes*, 32(3): 421-435.
625 <https://doi.org/10.1002/hyp.11429>

626 Valéry, A., Andréassian, V., Perrin, C., 2014. ‘As simple as possible but not simpler’: What is useful in a
627 temperature-based snow-accounting routine? Part 2 – Sensitivity analysis of the Cemaneige snow
628 accounting routine on 380 catchments. *Journal of Hydrology*, 517: 1176-1187.

629 <https://doi.org/10.1016/j.jhydrol.2014.04.058>

630 Vansteenkiste, T. et al., 2014a. Intercomparison of hydrological model structures and calibration
631 approaches in climate scenario impact projections. *Journal of Hydrology*, 519: 743-755.
632 <https://doi.org/10.1016/j.jhydrol.2014.07.062>

633 Vansteenkiste, T. et al., 2014b. Intercomparison of five lumped and distributed models for catchment
634 runoff and extreme flow simulation. *Journal of Hydrology*, 511: 335-349.
635 <https://doi.org/10.1016/j.jhydrol.2014.01.050>

636 Velazquez, J.A., Anctil, F., Perrin, C., 2010. Performance and reliability of multimodel hydrological
637 ensemble simulations based on seventeen lumped models and a thousand catchments. *Hydrology
638 and Earth System Sciences*, 14(11): 2303-2317. <https://doi.org/10.5194/hess-14-2303-2010>

639 Velázquez, J.A., Anctil, F., Ramos, M.H., Perrin, C., 2011. Can a multi-model approach improve
640 hydrological ensemble forecasting? A study on 29 French catchments using 16 hydrological model
641 structures. *Advances in Geosciences*, 29: 33-42. <https://doi.org/10.5194/adgeo-29-33-2011>

642 Wang, H.-M., Chen, J., Xu, C.-Y., Zhang, J., Chen, H., 2020. A Framework to Quantify the Uncertainty
643 Contribution of GCMs Over Multiple Sources in Hydrological Impacts of Climate Change. *Earth's
644 Future*, 8(8): e2020EF001602. <https://doi.org/10.1029/2020EF001602>

645 Wu, Y., Wu, S.-Y., Wen, J., Xu, M., Tan, J., 2016. Changing characteristics of precipitation in China
646 during 1960 - 2012. *International Journal of Climatology*, 36(3): 1387-1402.
647 <https://doi.org/10.1002/joc.4432>

648 Xu, C.-y., 2021. Issues influencing accuracy of hydrological modeling in a changing environment. *Water
649 Science and Engineering*. <https://doi.org/10.1016/j.wse.2021.06.005>

650 Yang, W. et al., 2020. Temporal and spatial transferabilities of hydrological models under different
651 climates and underlying surface conditions. *Journal of Hydrology*, 591: 125276.
652 <https://doi.org/10.1016/j.jhydrol.2020.125276>

653 Yilmaz, K.K., Gupta, H.V., Wagener, T., 2008. A process-based diagnostic approach to model evaluation:
654 Application to the NWS distributed hydrologic model. *Water Resources Research*, 44(9): W09417.
655 <https://doi.org/10.1029/2007WR006716>

656 Yin, J. et al., 2020. Projected changes of bivariate flood quantiles and estimation uncertainty based on
657 multi-model ensembles over China. *Journal of Hydrology*, 585: 124760.
658 <https://doi.org/10.1016/j.jhydrol.2020.124760>

659 Zaherpour, J. et al., 2019. Exploring the value of machine learning for weighted multi-model combination
660 of an ensemble of global hydrological models. *Environmental Modelling and Software*, 114: 112-
661 128. <https://doi.org/10.1016/j.envsoft.2019.01.003>

662 Zhang, J. et al., 2020. Combining Postprocessed Ensemble Weather Forecasts and Multiple Hydrological
663 Models for Ensemble Streamflow Predictions. *Journal of Hydrologic Engineering*, 25(1): 04019060.
664 [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0001871](https://doi.org/10.1061/(ASCE)HE.1943-5584.0001871)

665 Zhang, Y., Vaze, J., Chiew, F.H.S., Teng, J., Li, M., 2014. Predicting hydrological signatures in ungauged
666 catchments using spatial interpolation, index model, and rainfall-runoff modelling. *Journal of
667 Hydrology*, 517: 936-948. <https://doi.org/10.1016/j.jhydrol.2014.06.032>

668 Zhao, R.-J., 1992. The Xinanjiang model applied in China. *Journal of Hydrology*, 135(1): 371-381.
669 [https://doi.org/10.1016/0022-1694\(92\)90096-E](https://doi.org/10.1016/0022-1694(92)90096-E)

670 Zhao, R.J., Zuang, Y., Fang, L., Liu, X., Zhang, Q., 1980. The Xinanjiang Model.

671

Replies to Referee #1

Performance dependence of multi-model combination methods on hydrological model calibration strategy and ensemble size

Yongjing Wan, Jie Chen, Chong-Yu Xu, Ping Xie, Wenyan Qi, Daiyuan Li, Shaobo Zhang

We would like to thank the referee for constructive comments and suggestions. All comments are very helpful for improving this paper and beneficial to our research in general. We have provided detailed point-by-point responses to each comment below and have revised the manuscript accordingly. For clarity, comments are given in black, and our responses are given in blue. Please note that the page and line numbers mentioned in reviewers' comments refer to the original version, while in our reply, they refer to the revised version.

The authors have paid careful attention to the reviewers' comments and significantly improved the paper. I have a few minor suggestions in annotations in the attached file. The paper, especially the latest revisions, would benefit from a proof reading for English quality.

Reply: We sincerely thank the reviewer for the positive evaluation of our manuscript and for providing insightful comments. All comments in annotations in the attachment have been addressed as follows. The English has been carefully revised.

Line 103. Perhaps this is expected, if the metric on which the multi-model average is assessed is some generic metric like NSE. If the metric is more specific like high flows, then we do not expect the combination derived from using different objective functions to work better than simply using a high-flow objective. This is a critical point that needs to be made clear.

Reply: We agree with the reviewer. The use of multi-model combination scheme is expected to benefit from the variation of the parameter sets derived from objective functions targeted at different hydrological processes for producing a better overall simulation.

We have made a clear declaration in line 105 of the revised manuscript.

Line 195. This is poorly written. Needs a new sentence; and needs written more clearly and checked by an English speaker.

Reply: We have made a clear declaration in line 195 of the revised manuscript.

Line 277. Unclear what this means.

Reply: Sorry for the confusion. This sentence has been modified as: Among various averaging methods, GRA, GRB, GRC, and MMSE show similar performance, since they derived from the same optimal weighting group.

We have clarified this in line 276 of the revised manuscript.

Line 319. Random selection from all combinations?

Reply: Yes, and we have clarified this in line 316 of the revised manuscript.

Line 331. not the right word here? "significant", or "much"?

Reply: Thanks, and we have verified this sentence.

Line 335. does not make sense to me. Do you mean "whether all ensemble members contribute to the performance of the combination"

Reply: Yes, and we have made a clear declaration in line 332 of the revised manuscript.

Line 381. Especially when using lumped models, which cannot represent the rainfall and loss variability that tends to be higher in arid areas.

Reply: Thanks for the supplement explanation. We have added this in the revised manuscript in line 378.

Line 428. This gives limited insight into where and why the performances were better. At least, the effect of catchment area should be explored too.

Reply: Thanks for the suggestion. We have analyzed the relationship between the combination performance and catchment area, as shown in figure R1. The result indicates that there is no obvious relationship between two of them. Thus, this result was not shown in the revised manuscript.

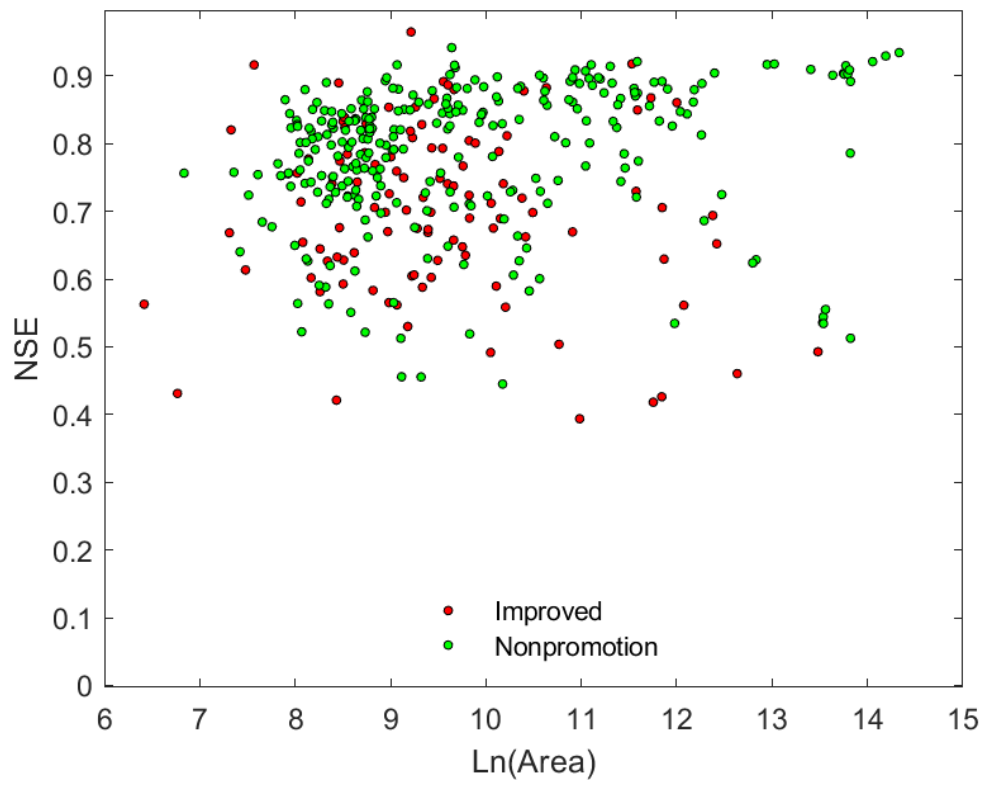


Figure R1. The relationship between the NSE value and catchment area. The green/red markers represent the combination performed better/worse than the best individual member.

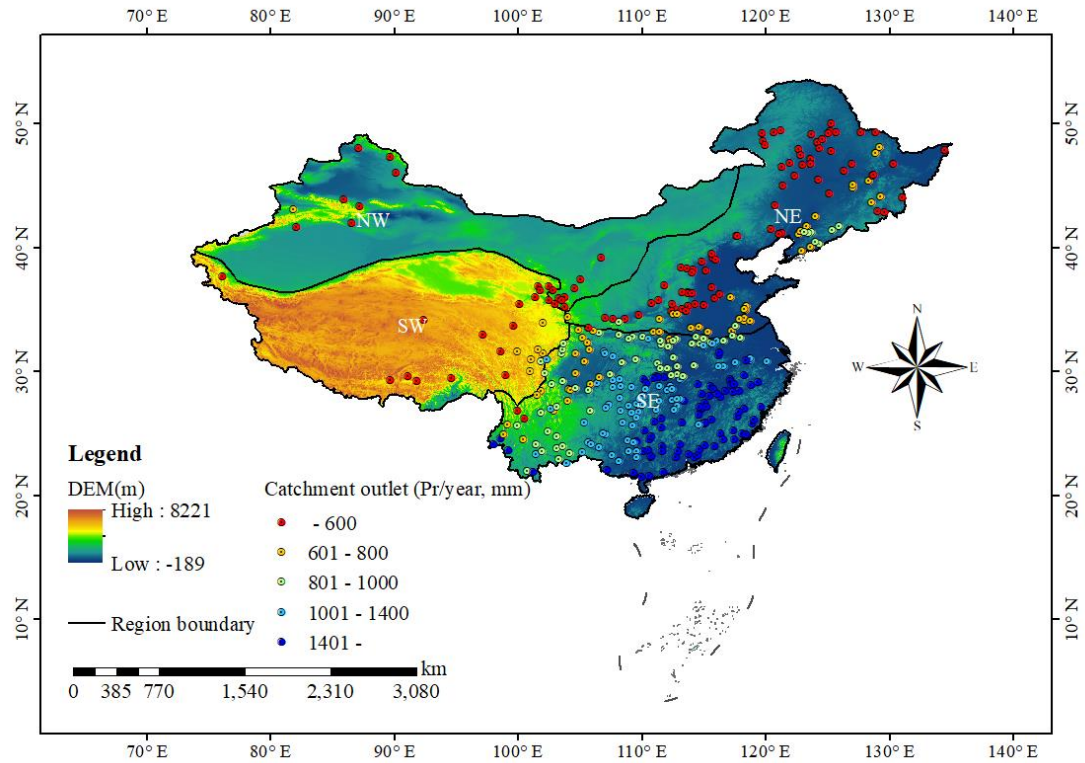


Figure 1. Spatial distribution of the outlets and total annual precipitation (Pr, mm) for 383 catchments in China.

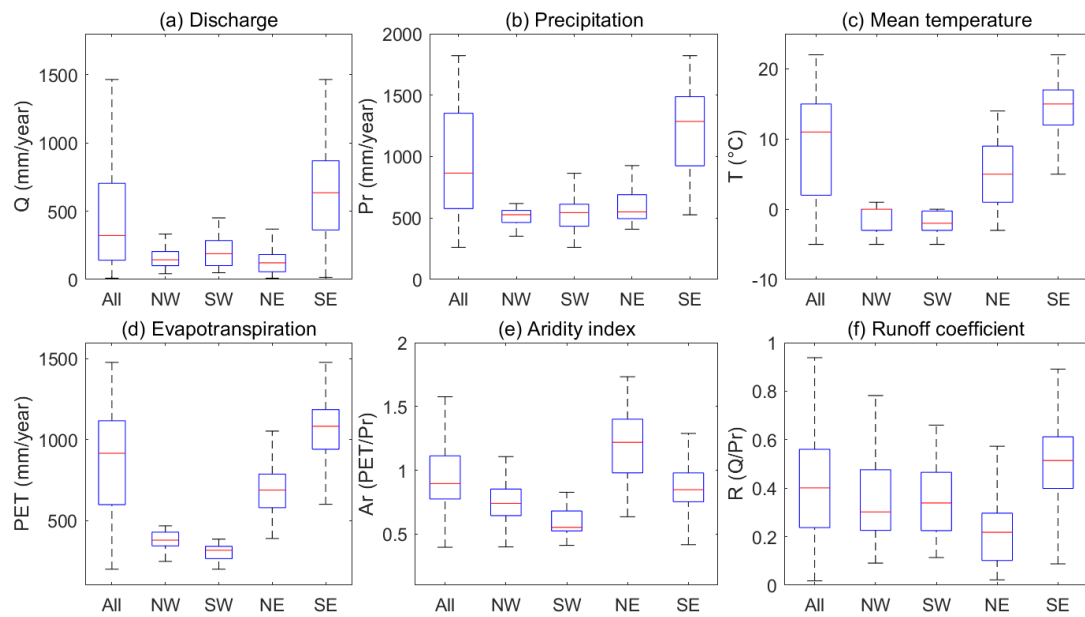


Figure 2. Characteristics of the catchments over different regions. The red line in the boxplots represents the median

value, the ends of the boxes represent the 25th and 75th percentiles, the whiskers represent the values at the 5th and

95th percentiles, and outliers are not shown.

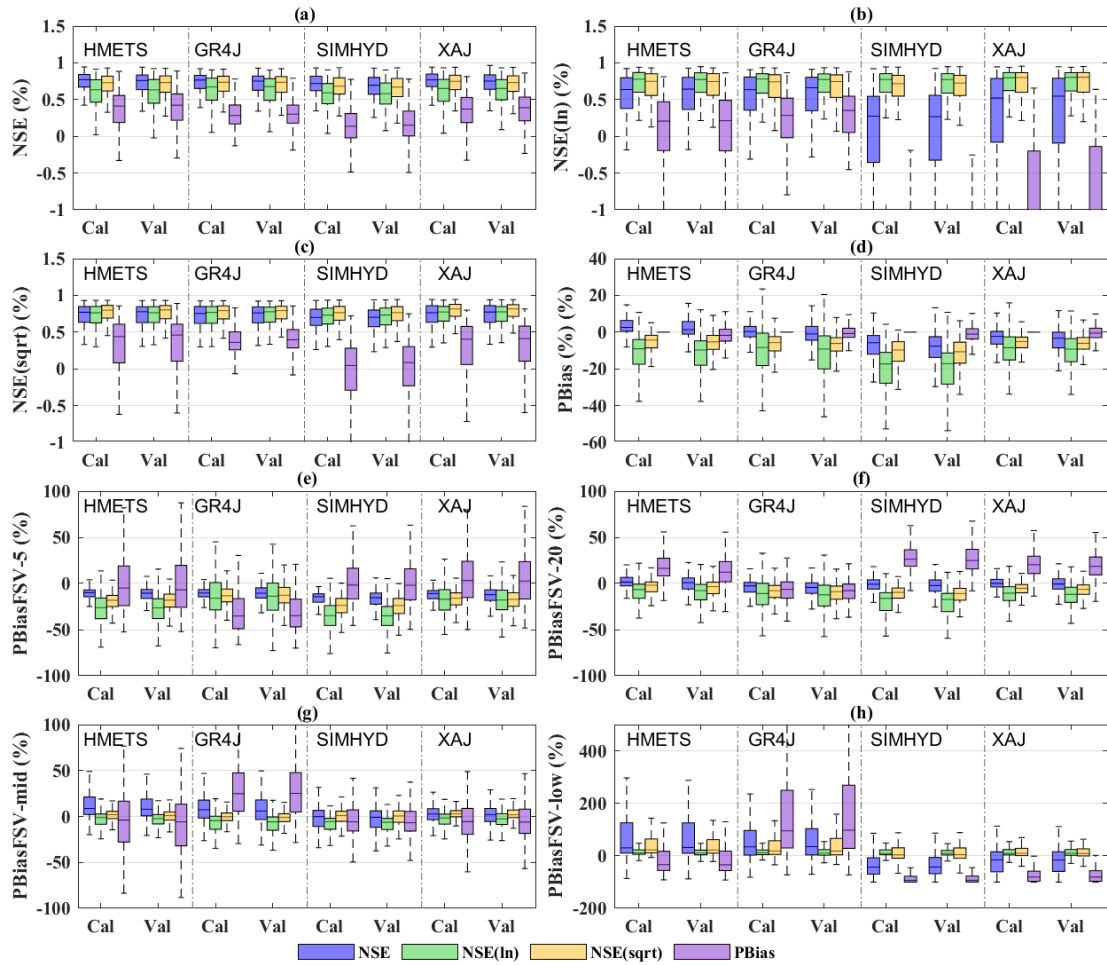


Figure 3. NSE, NSE(ln), NSE(sqrt), PBias, and four PBiasFSV values for 16 ensemble members (model/objective function pairs). The boxplots indicate the spread of performance values of the evaluation metrics over 383 catchments in the calibration and validation periods.

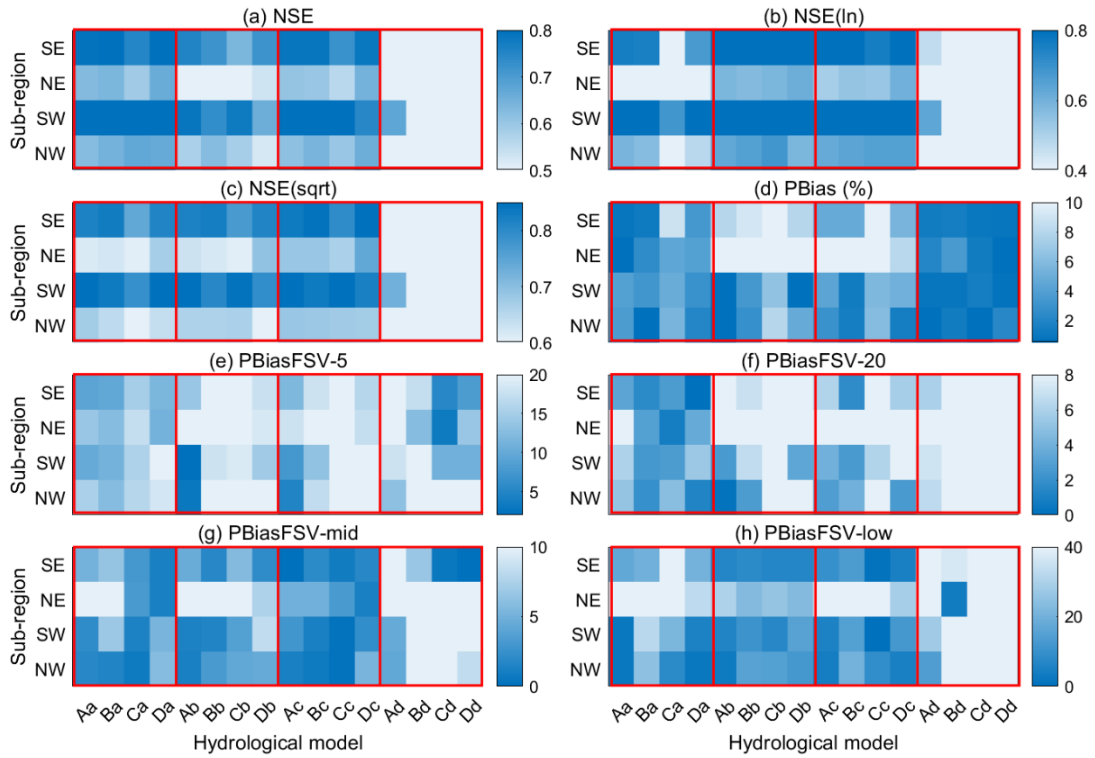


Figure 4. The median value of NSE, NSE(ln), NSE(sqrt), absolute of PBias and four PBiasFSV values in the validation period for the catchments over the four sub-regions. The four hydrological models GR4J, HMETS, SIMHYD and XAJ are denoted by A, B, C and D, respectively. The four objective functions NSE, NSE(ln), NSE(sqrt), and PBias are denoted by a, b, c and d, respectively.

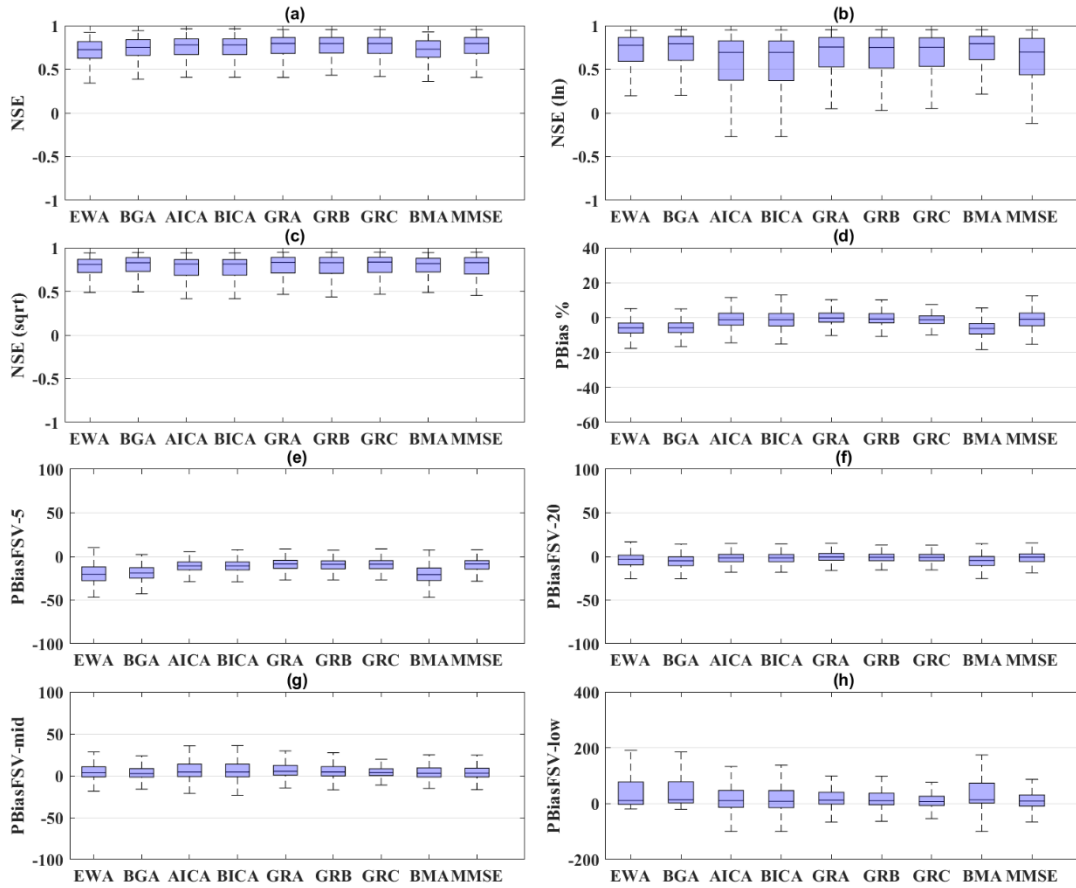


Figure 5. NSE, NSE(ln), NSE(sqrt), PBias, and four PBiasFSV values over 383 catchments for nine multi-model ensemble techniques with combinations of the 16 ensemble members in the validation period.

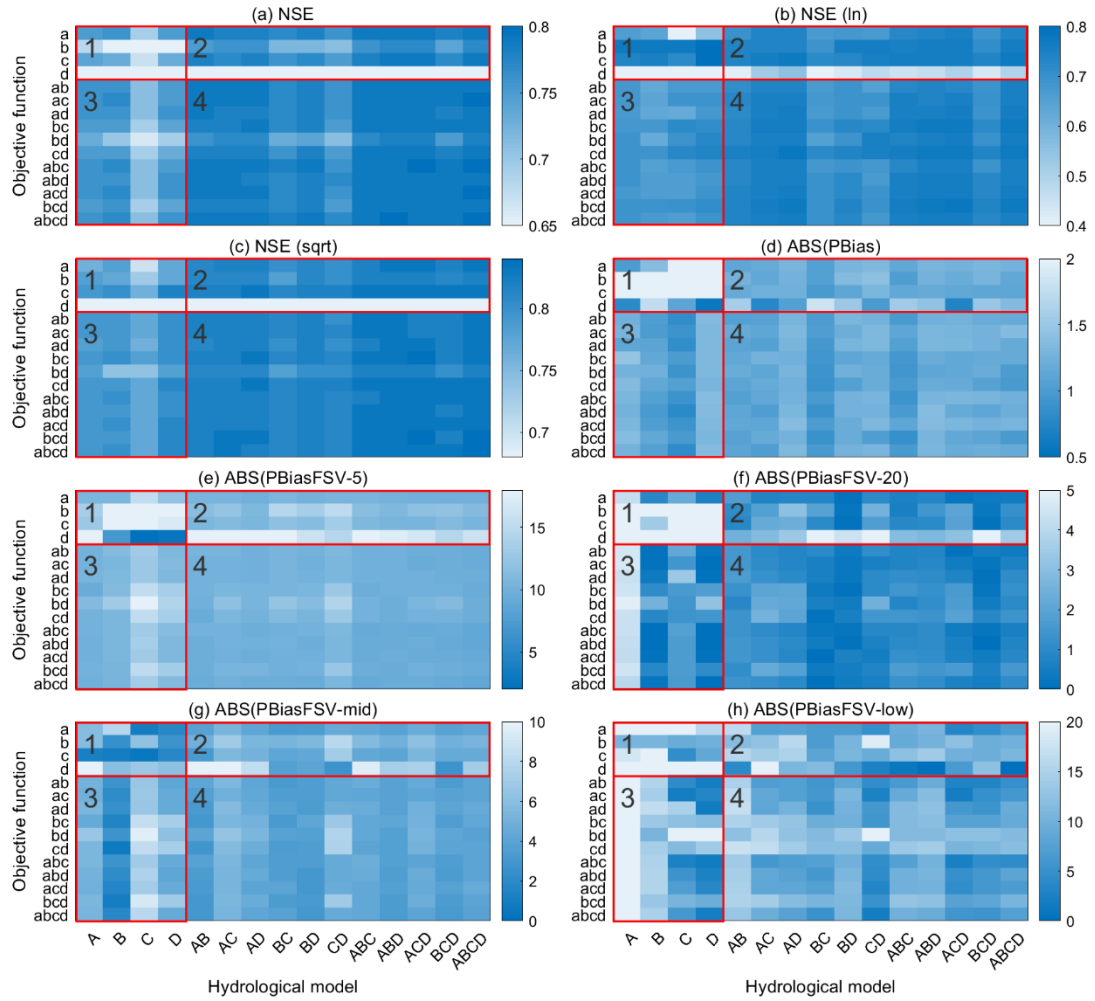


Figure 6. The median value of NSE, NSE(ln), NSE(sqrt), and absolute values of PBias four PBiasFSV over 383 catchments in the validation period based on different combinations in terms of ensemble members. The four hydrological models GR4J, HMETs, SIMHYD and XAJ are denoted by A, B, C and D, respectively. The four objective functions NSE, NSE(ln), NSE(sqrt), and PBias are denoted by a, b, c and d, respectively.

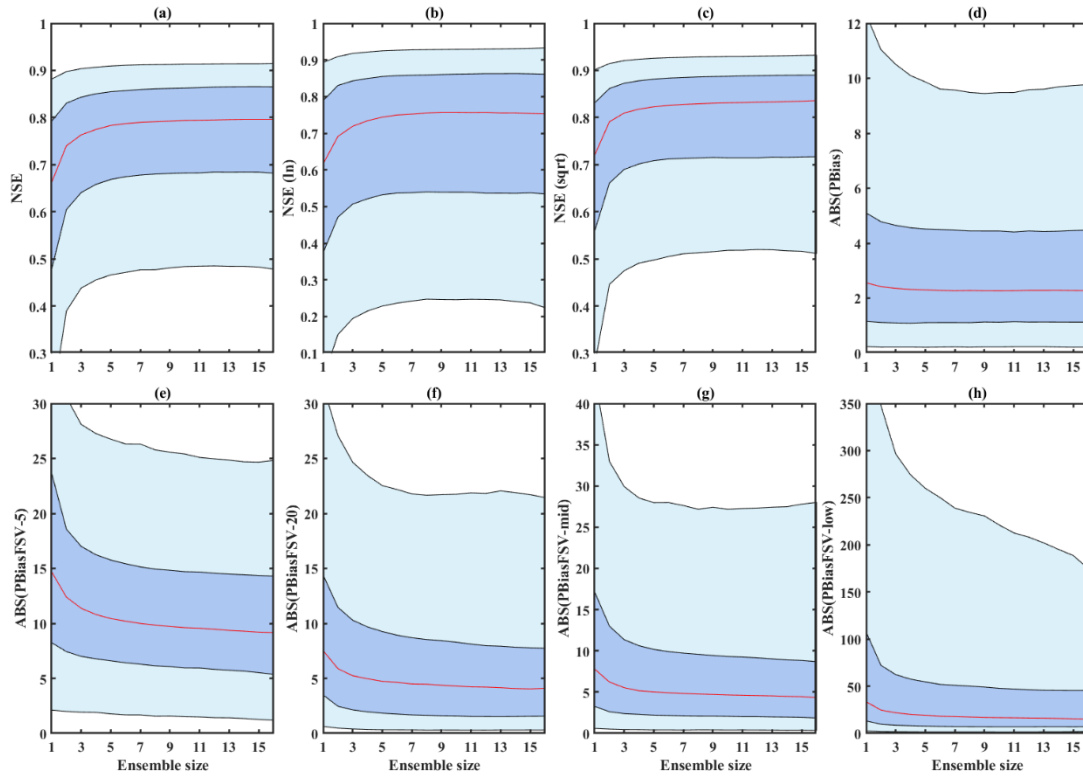


Figure 7. The relationship between the ensemble size and the performance metrics (NSE, NSE(ln), NSE(sqrt), absolute values of PBias and four PBiasFSV) over 383 catchments in the validation period. The light-colored envelopes: the 0.05 and 0.95 interquartile range values. The dark-colored envelopes: the 0.25 and 0.75 interquartile range values. The red lines are the median values.

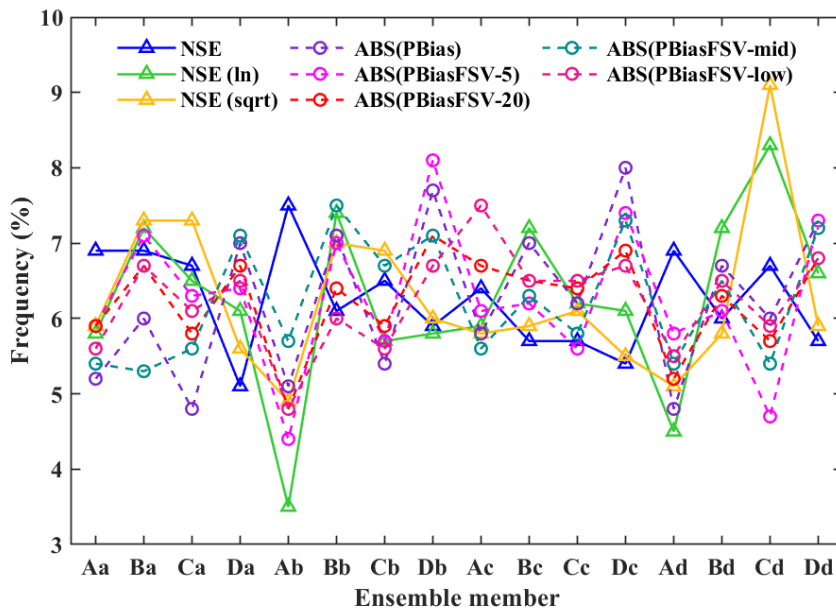


Figure 8. Frequency of individual members is selected in the best combination scheme in terms of different

evaluation metrics (NSE, NSE(ln), NSE(sqrt), absolute values of PBias and four PBiasFSV). The four hydrological models GR4J, HMETS, SIMHYD and XAJ are denoted by A, B, C and D, respectively. The four objective functions NSE, NSE(ln), NSE(sqrt), and PBias are denoted by a, b, c and d, respectively.

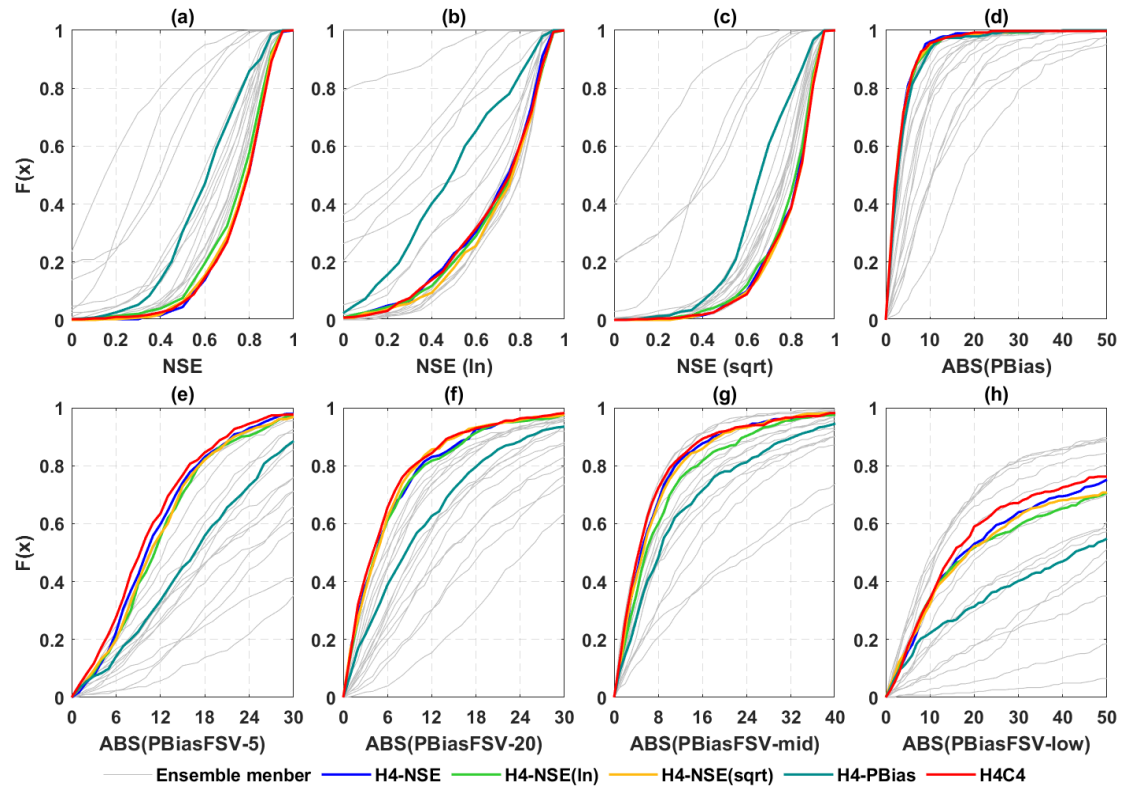


Figure 9. Cumulative distributions of NSE, NSE(ln), NSE(sqrt), absolute values of PBias and four PBiasFSV over the 383 catchments in the validation period.

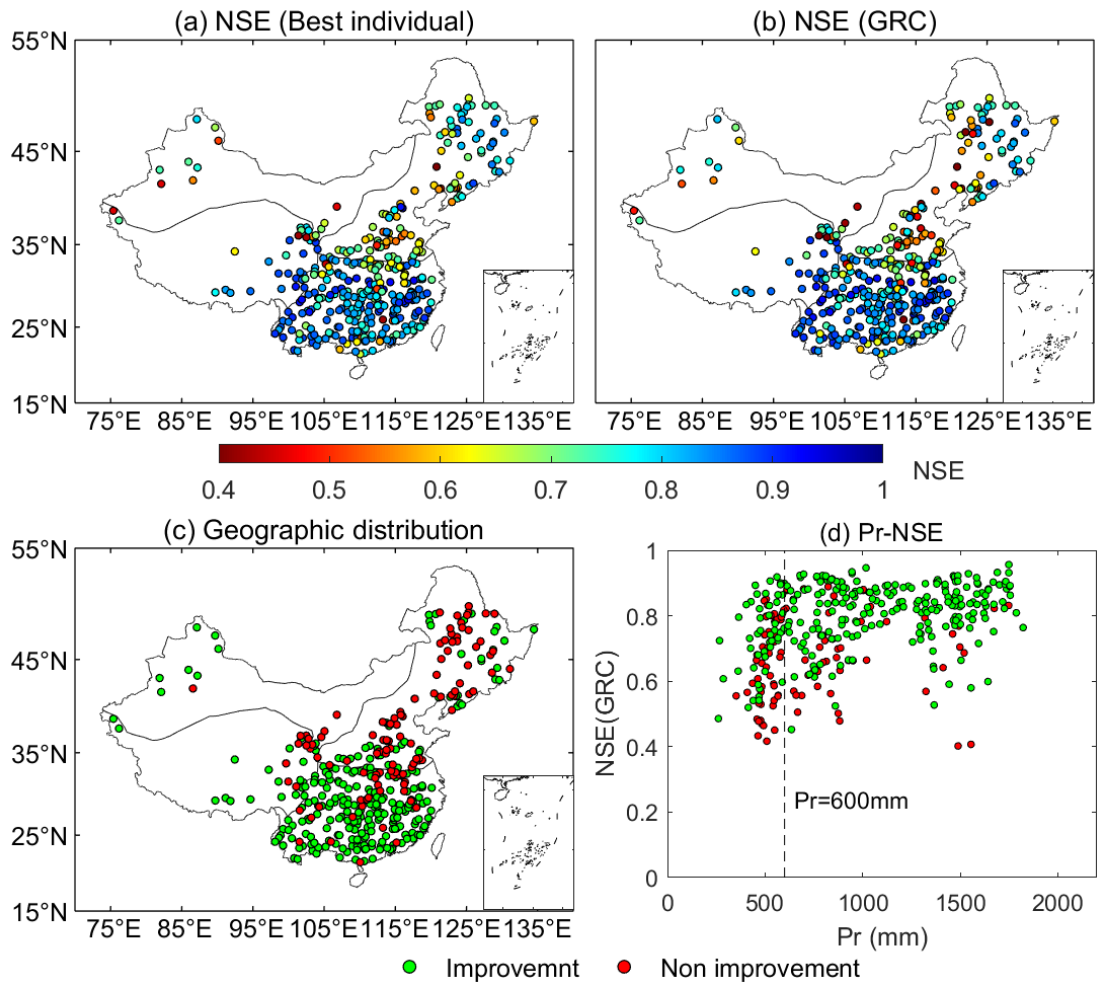


Figure 10. Geographic distribution of the NSE value of (a) the best individual member and (b) the combination over the 383 catchments in the validation period. Comparison of the best individual and GRC averaging method (c) on geographic distribution and (d) related to annual precipitation. The green/red markers present GRC produced results better/worse than the best individual member.

Table 1. Structures of the four lumped conceptual rainfall-runoff models

ID	Model	Number of parameters	Characteristics of the model	References
A	GR4J	6	The effective rainfall is partitioned as a 10:90 split representing direct runoff and delayed runoff, a nonlinear production reservoir with two-unit hydrographs, a routing reservoir	Edijatno et al. (1999); Perrin et al. (2003)
B	HMETS	21	Generation of hypodermic flow and groundwater flow with two linear reservoirs, a routing module with two unit hydrographs, a snowmelt module, an evapotranspiration calculation module	Martel et al. (2017)
C	SIMHYD	11	Two linear reservoirs for the calculation of interflow and base-flow, a nonlinear routing reservoir, an evapotranspiration calculation module	Chiew et al. (2002)
D	XAJ	17	Linear reservoirs for surface flow routing, two recession coefficients for interflow and groundwater flow routing, three-layer evapotranspiration system	Zhao (1992); Zhao et al. (1980)

Table 2. Performance metrics for the evaluation of different phases of the hydrograph.

Performance metric	Sensitive hydrograph phase	Range
NSE	Peak and discharge dynamic	$-\infty \sim 1$
NSE(ln)	Low-flow and discharge dynamic	$-\infty \sim 1$
NSE(sqrt)	Discharge dynamic	$0 \sim 1$
PBias	Overall water balance	$-\infty \sim \infty$
PBiasFSV-5	Tendencies of overestimation and underestimation for FDC very high-segment volume	$-\infty \sim \infty$
PBiasFSV-20	Tendencies of overestimation and underestimation for FDC high-	$-\infty \sim \infty$

	segment volume	
PBiasFSV-mid	Tendencies of overestimation and underestimation for FDC mid-segment volume	$-\infty \sim \infty$
PBiasFSV-low	Tendencies of overestimation and underestimation for FDC very low-segment volume	$-\infty \sim \infty$

Table 3. Basic characteristics of the nine multi-model averaging techniques used in this study

Combination method	Acronym	Method description	References
Equal Weights	EWA	Unweighted average	Shamseldin et al. (1997)
Akaike's Information Criterion	AICA	Mean of the logarithm of the member variances added a penalty equalling to double the number of calibrated parameters	Akaike (1974)
Bayes Information Criterion	BICA	Mean of the logarithm of the member variances added a penalty equalling to the number of calibrated parameters times the logarithm of the number of time steps	Schwarz (1978)
Bates and Granger	BGA	Minimizing the Root Mean Square Error	Bates and Granger (1969)
Granger	GRA	Based on ordinary least squares (OLS) algorithm	Granger and

Ramanathan A				Ramanathan (1984)
Granger	GRB	Weights based on ordinary least squares (OLS) algorithm and constrained its sum to unity		Granger and Ramanathan (1984)
Ramanathan-B				
Granger	GRC	Based on ordinary least squares (OLS) algorithm and bias-corrected the results		Granger and Ramanathan (1984)
Ramanathan-C				
Bayesian Model Averaging	BMA	Determining weights by probability distribution functions (PDFs)		Neuman (2003)
Multi-model Super Ensemble	MMSE	Based on ordinary least squares (OLS) algorithm and using the logic of bias reduction with respect to individual member models along with variance reduction in simulation		Krishnamurti et al. (2000)

Table 4. The frequency (%) of individual members obtaining the best score for the performance metrics in the validation period.

Model member	NSE	NSE(ln)	NSE(sqrt)	PBias	PBiasFSV			
					5	20	mid	low
GR4J-NSE	17.2	0	0.8	11.2	13.8	7.8	2.9	3.9
GR4J-NSE(ln)	0.8	9.7	0.8	2.1	6.8	3.7	7.6	12.8
GR4J-NSE(sqrt)	5.5	1.3	16.4	3.7	5.7	3.7	10.4	4.4
GR4J-PBias	0.3	0	0	15.4	1.6	6	3.1	2.1
HMETS-NSE	20.1	0	1.6	11.5	12.5	11	4.4	4.2
HMETS-NSE(ln)	0.5	18.8	1.8	2.9	2.1	4.2	9.1	7.6
HMETS-NSE(sqrt)	4.4	3.1	17.8	4.7	4.4	11.7	8.4	3.1

HMETS-PBias	0	0	0	7	6.5	8.4	3.1	5
SIMHYD-NSE	7	0	1.8	4.4	1.6	9.7	4.7	5.2
SIMHYD-NSE(ln)	0.5	17.5	2.1	0.3	0.5	1.6	5.7	13.6
SIMHYD-NSE(sqrt)	2.3	2.3	7.8	1.8	0.8	1.6	7.3	8.9
SIMHYD-PBias	0	0	0	7	12.8	1.8	4.2	1
XAJ-NSE	26.6	2.1	5.2	7.3	9.7	13.3	5.7	3.4
XAJ-NSE(ln)	0.3	23.5	1.3	2.1	7.3	2.6	7.6	13.8
XAJ-NSE(sqrt)	13.6	21.7	42.6	2.6	2.3	8.1	12	7.8
XAJ-PBias	0.8	0	0	15.9	11.5	5	3.7	3.1

Table 5. Comparison of the multi-model combination and the best individual member for each of the 383 catchments in the validation period. Here, H4-NSE, H4-NSE(ln), H4-NSE(sqrt) and H4-PBias donate four hydrological models calibrated with a specified objective function, and H4C4 donates four hydrological models calibrated with four calibration strategies.

Acronym	NSE	NSE(ln)	NSE(sqrt)	PBias	PBiasFSV			
					5	20	mid	low
H4-NSE	70	11.2	45.2	10.7	14.1	13.8	4.4	8.4
H4-NSE(ln)	32.4	18.3	35.2	11.2	10.4	13.1	4.2	5.7
H4-NSE(sqrt)	60.1	17.8	56.4	11.7	13.3	12.5	7	6
H4-PBias	1.8	1	1.6	7.3	3.1	7.8	5	7.6
H4C4	71	17	55.4	12.3	20.4	13.6	8.4	8.4

Highlights:

- Four hydrological models calibrated with four objective functions are compared.
- The Granger Ramanathan average variant C (GRC) method performs the best.
- Using more than nine ensemble members does not further improve performance.
- Combinations of models and objective functions are better than the single model and objective.
- Averaging outperforms the ensemble members except in low-flow simulations.

Yongjing Wan: Conceptualization, Software, Formal analysis, Writing-Original Draft.

Jie Chen: Conceptualization, Writing-Review & Editing, Project administration,
Funding acquisition.

Chong-Yu Xu: Conceptualization, Writing-Review & Editing.

Ping Xie: Writing-Review & Editing.

Wenyan Qi: Writing-Review & Editing.

Daiyuan Li: Writing-Review & Editing.

Shaobo Zhang: Data curation.