

Does the weighting of climate simulations result in a more reasonable quantification of hydrological impacts?

Hui-Min Wang¹, Jie Chen^{1*}, Chong-Yu Xu^{1,2}, Hua Chen¹, Shenglian Guo¹, Ping Xie¹, Xiangquan Li¹

¹State Key Laboratory of Water Resources and Hydropower Engineering Science, Wuhan University, Wuhan, 430072, China

5 ²Department of Geosciences, University of Oslo, Oslo, Norway

Correspondence to: Jie Chen (jiechen@whu.edu.cn)

Abstract: With the increase in the number of available global climate models (GCMs), pragmatic questions come up when using them to quantify the climate change impacts on hydrology: Is it necessary to weight GCM outputs in the impact studies, and if so, how to weight them? Some weighting methods have been proposed based on the performances of GCM simulations with respect to reproducing the observed climate. However, the process from climate variables to hydrological responses is nonlinear, and thus the assigned weights based on their performances in climate simulations may not be correctly translated to hydrological responses. Assigning weights to GCM outputs based on their ability to represent hydrological simulations is more straightforward. Accordingly, the present study assigns weights to GCM simulations based on their ability to reproduce hydrological characteristics and investigates their influence on the quantification of hydrological impacts. Specifically, eight weighting schemes are used to determine the weights of GCM simulations based on streamflow series simulated by a lumped hydrological model using raw or bias-corrected GCM outputs. The impacts of weighting GCM simulations are investigated in terms of reproducing the observed hydrological regimes for the reference period (1970-1999) and quantifying the uncertainty of hydrological changes for the future period (2070-2099). The results show that when using raw GCM outputs, streamflow-based weights better represent the mean hydrograph and reduce more biases of annual streamflow than the weights calculated using climate variables. However, when applying bias correction to GCM simulations before driving the hydrological model, the streamflow-based unequal weights do not bring significant differences in the multi-model ensemble mean and uncertainty of hydrological impacts, since bias-corrected climate simulations become rather close to observations. Thus, the equal weighting method may still be a viable and conservative choice when bias correction to GCM simulations is conducted in hydrological climate change impact studies.

1 Introduction

Multi-model ensembles (MMEs) consisting of climate simulations from multiple global climate models (GCMs) have been widely used to quantify future climate change impacts and the corresponding uncertainty (Wilby and Harris, 2006; IPCC, 2013; Chiew et al., 2009; Chen et al., 2011; Tebaldi and Knutti, 2007). The number of climate models has increased rapidly, resulting in the obviously growing size of MMEs. For example, the Coupled Model Inter-comparison Project Phase 5 (CMIP5) archive contains 61 GCMs from 28 modeling institutes, with some GCMs providing multiple simulations (Taylor et al., 2012). Due to the lack of consensus on the proper way to combine simulations of a MME, the prevailing approach is the model democracy (“one model one vote”) for the sake of simplicity, where each member in an ensemble is considered to have equal ability to simulate historical and future climates. The model democracy method has been applied to many global and regional climate change impact studies (e.g., IPCC, 2014; Minville et al., 2008; Maurer, 2007). Although it has been reported that the equal average of a multi-model ensemble often outperforms any individual model in regards to the reproduction of the mean state of observed historical climate (Gleckler et al., 2008; Reichler and Kim, 2008), whether the equal weighting is a better strategy for hydrological impact studies remains to be investigated (Alder and Hostetler, 2019).

Several studies have raised concerns about the strategy of model democracy, due to the following two reasons (Lorenz et al., 2018; Knutti et al., 2017; Cheng and AghaKouchak, 2015). First, GCM simulations in an ensemble do not have identical skills at representing historical climate observations. They may perform differently in simulating future climate. GCM performances may also vary by their variables and locations (Hidalgo and Alfaro, 2015; Abramowitz et al., 2019), which further challenges the rationality of model democracy in regional impact studies. Second, equal weights imply that the individual members in an ensemble are independent of each other. However, some climate models share common modules, parts of codes, parameterizations and so on (Knutti et al., 2010; Sanderson et al., 2017). Some pairs of GCMs submitted to the CMIP5 database only differ in the spatial resolution (e.g. MPI-ESM-MR and MPI-ESM-LR; see Giorgetta et al., 2013). The replication or overlapping in these GCMs may lead to the inter-dependence of MMEs, resulting in common biases towards the replicating section and inflating confidence in the projection uncertainty (Sanderson et al., 2015; Jun et al., 2012).

With the intention of improving climate projections and reducing the uncertainty, some weighting approaches have been proposed to assign unequal weights to climate model simulations according to their performances with respect to reproducing some diagnostic metrics of historical climate observations (Murphy et al., 2004; Sanderson et al., 2017; Cheng and AghaKouchak, 2015). For example, Xu et al. (2010) apportioned weights for GCMs based on their biases to the observed data in terms of two diagnostic metrics (climatological mean and inter-annual variability) for producing probabilistic climate projections. Lorenz and Jacob (2010) used errors in the trends of temperature to evaluate climate projections and determine weights. Other criteria have also been introduced into model weighting as a complement to the performance criterion. Some examples are the convergence of climate projections for a future period (Giorgi and Mearns, 2002) and the interdependence among climate models (Sanderson et al., 2017).

Despite the different diagnostic metrics or definitions of model performances employed in these weighting methods, weights are commonly determined with respect to the ability of climate simulations at reproducing observed climate variables, such as temperature and precipitation (e.g., Chen et al., 2017; Wilby and Harris, 2006; Xu et al., 2010). However, for the impact studies, the relationship between climate variables and the impact variable is often not straightforward or explicit. In other words, the process from climate variables to their impacts may not be linear (Wang et al., 2018; Risbey and Entekhabi, 1996; Whitfield and Cannon, 2000). For example, Mpelasoka and Chiew (2009) reported that in Australia, a small change in annual precipitation can result in a several-times change in annual runoff. Thus, the weights calculated in the climate world may not be effective in the impact field.

In addition, a number of climate variables may determine the climate change impacts on a single environmental sector. For example, the runoff generation in a watershed is usually determined by precipitation, temperature, and other climate variables. Thus, it is not an easy task to determine the relative importance of each climate variable in impact studies, which is the other challenge to combining sets of weights based on different climate variables into a single set of weights for impact simulations. Previous studies have usually assumed that all variables are equally important and had an equal weight assigned to each climate variable (Xu et al., 2010; Chen et al., 2017; Zhao, 2015). However, these climate variables are usually not equally important in the impact field. For example, precipitation may be more important than temperature for a rainfall-dominated watershed, but could be different for a snowfall-dominated watershed. Thus, it may be more straightforward to calculate the weights for GCMs based on their ability to reproduce the single impact variable instead of multiple climate variables. Such a method would integrate the synthetic ability of GCMs in terms of simulating multiple climate variables to that of one impact variable. In addition, this method could also circumvent the previous problem of potential nonlinearity between climate variables and the impact variable.

Accordingly, the objectives of this study are to assign weights to GCM simulations according to their ability to represent hydrological observations, and to assess the impacts of these weighting methods on the assessment of hydrological responses to climate change. The case study was conducted over two watersheds with different climatic and hydrological characteristics. Since both bias correction and model weighting are common procedures in regional and local impact studies, this study considers two experiments (raw and bias-corrected GCM outputs) to simulate streamflows and investigate the performances of weighting methods. Seven weighting methods were used to assign unequal weights for streamflows simulated by raw or bias-corrected GCMs, respectively. The impacts of unequal weights are then assessed and compared to the equal weighting method in terms of multi-model ensemble mean and uncertainty related to the choice of a climate model.

2 Study area and data

30 2.1 Study Area

This study was conducted over two watersheds with different climate and hydrological characteristics: the rainfall-dominated Xiangjiang watershed and the snowfall-dominated Manicouagan-5 watershed (Figure 1). The Xiangjiang River is

one of the largest tributaries of the Yangtze River in central-southern China, and its drainage area is about 94 660 km² (Figure 1a). A catchment with a surface area of about 52 150 km² above the Hengyang gauged station was used in this study. The catchment is heavily influenced by the East Asian Monsoon, which causes a humid subtropical climate with hot and wet summers and mild winters. The average temperature over the catchment is about 17 °C with the coldest month averaging about 7 °C. The average annual precipitation is about 1570 mm, of which 61% falls in the wet season from April to August. The daily averaged streamflow at the Hengyang gauged station is around 1400 m³s⁻¹. The annual average of summer peak streamflow is about 4420 m³s⁻¹, mainly due to summer extreme rainfalls.

The Manicouagan-5 watershed, located in the center of the Province of Quebec, Canada, is the largest sub-basin of the Manicouagan watershed (Figure 1b). Its drainage area is about 24 610 km², most of which is covered by forest. The outlet of the Manicouagan-5 River is the Daniel-Johnson Dam. The Manicouagan-5 watershed has a continental subarctic climate characterized by long and cold winters. The average temperature over the watershed is about -3 °C, with nearly half of the year having a daily temperature below 0 °C. The average annual precipitation is about 912 mm, evenly distributed over each year. The average discharge at the outlet of the Manicouagan-5 River is about 530 m³s⁻¹. Snowmelt contributes to the peak discharge during May, whose annual average is about 2200 m³s⁻¹.

2.2 Data

This study used daily maximum and minimum temperatures and precipitation from observation and GCM simulations for both watersheds. The observed meteorological data for the Xiangjiang watershed were collected from 97 precipitation gauges and 8 temperature gauges. Streamflow series were collected from the Hengyang gauged station. For the Manicouagan-5 watershed, the observed meteorological data were extracted from the gridded dataset of Hutchinson et al. (2009), which is interpolated from daily station data using a thin-plate smoothing spline interpolation algorithm. Streamflow series were the inflows of the Daniel Johnson Dam, which were calculated using mass balance calculations. All the observation data for both watersheds cover the historical reference period (1970-1999).

For the climate simulations, maximum and minimum temperatures and precipitation of 29 GCMs were extracted from the CMIP5 archive over both watersheds (Table 1). All simulations cover both the historical reference period (1970-1999) and the future projection period (2070-2099). One Representative Concentration Pathway (RCP8.5) was used in terms of climate projections in the future period. RCP8.5 was selected because it projects the most severe increase in greenhouse gas emissions among the four RCPs, and it is often used to design conservative mitigation and adaptation strategies (IPCC, 2014).

3 Methodology

To begin the process of calculating the weight for each GCM simulation, a multi-model ensemble constructed by 29 CMIP5 GCMs was utilized to drive a calibrated hydrological model over the two watersheds. Two experiments were designed to generate the ensembles of streamflow simulations. The first experiment drives the hydrological model using raw GCM outputs

with no bias correction, while the second drives the hydrological model using bias-corrected climate simulations. Although it is not common to use raw GCM simulations for hydrological impact studies, the rationale for using them in this study is to examine the impacts of bias correction on weighting GCMs. The bias correction may adjust the relative performances between climate simulations and thus affect the determination of the relative weight for each ensemble member. Based on the ensemble of hydrological simulations from GCM outputs, eight weighting methods were employed to determine the weights of each GCM and to combine ensemble members for the assessment of hydrological climate change impacts. More detailed information is given below.

3.1 Bias correction

Since the raw outputs of GCMs are often too coarse and biased to be directly input into hydrological models for impact studies, bias correction is commonly applied to GCM outputs prior to the runoff simulation (Wilby and Harris, 2006; Chen et al., 2011; Minville et al., 2008). A distribution-based bias correction method, the daily bias correction (DBC) method of Chen et al. (2013), was used in this study. DBC is the combination of the local intensity scaling (LOCI) method (Schmidli et al., 2006) and the daily translation (DT) method (Mpelasoka and Chiew, 2009). The LOCI method was used to adjust the wet-day frequency of climate model simulated precipitation. A threshold was determined for the reference period to ensure that the simulated precipitation occurrence is identical to the observed precipitation occurrence. The same threshold was then used to correct the wet-day frequency for the future period. The DT method was used to correct biases in the frequency distribution of simulated precipitation amounts and temperature. The frequency distribution was represented by 100 percentiles ranging from the 1st to the 100th, and the correction factors were calculated for each percentile. The same correction factors were then employed to correct the distributions for the future period. The use of distribution-based biases facilitates the use of different correction factors for different levels of precipitation. Some studies have shown the advantages of distribution-based bias correction over other correction methods in the assessment of hydrological impacts (Chen et al., 2013; Teutschbein and Seibert, 2012).

3.2 Runoff simulation

The runoff was simulated using a lumped conceptual hydrological model, GR4J-6, which couples a snow accumulation and melt module, CemaNeige, with a rainfall-runoff model, GR4J (Arsenault et al., 2015). The CemaNeige model divides the precipitation into liquid and solid according to the daily temperature range, and generates snowmelt depending on the thermal state and water equivalent of the snowpack (Valéry et al., 2014). CemaNeige has two free parameters: the melting rate and the thermal state coefficient. The GR4J model consists of a production reservoir and a routing reservoir (Perrin et al., 2003). A portion of net rainfall (liquid precipitation with evaporation subtracted) goes into the production reservoir, whose leakage forms the effective rainfall when combined with the other proportion of net rainfall. The effective rainfall is then divided into two flow components. Ninety percent of the effective rainfall routes via a unit hydrograph and enters into the routing reservoir. The other 10% generates the direct flow through the other unit hydrograph. There is groundwater exchange with neighbouring

catchments in the direct flow and the outflow nonlinearly generated by the routing reservoir. Four free parameters in GR4J must be calibrated: the maximum capacity of the production reservoir, the groundwater exchange coefficient, the one-day-head maximum capacity of the routing reservoir and the time base of unit hydrograph.

5 The time periods of the observed data used for hydrological model calibration and validation are presented in Table 2. The shuffled complex evolution optimization algorithm (Duan et al., 1992) was employed to optimize the parameters of GR4J-6 for both watersheds. The optimized parameters were chosen to maximize the Nash-Sutcliffe Efficiency (NSE) criteria (Nash and Sutcliffe, 1970). The selected sets of parameters yield NSEs greater than 0.87 for both calibration and validation periods, indicating the reasonable performance of GR4J-6 and the high quality of the observed datasets for both watersheds.

3.3 Weighting Methods

10 Raw and bias-corrected climate simulations were input to the calibrated GR4J-6 model to generate raw and bias-corrected streamflow data series, respectively. Eight weighting methods were then employed to determine the weight of each hydrological simulation, including the equal weighting method (model democracy) and 7 unequal weighting methods. All of the unequal weighting methods are described in detail in the supplementary material so they are only briefly presented herein. Seven unequal weighting methods consist of two multiple criteria-based weighting methods and five performance-based
15 weighting methods. The two multiple criteria-based weighting methods are the reliability ensemble averaging method (REA) and the performance and interdependence skill (PI). The REA method considers both the bias of a GCM to observation in the reference period (performance criterion) and its similarity to other GCMs in the future projection (convergence criterion) (Giorgi and Mearns, 2002). The PI method weights an ensemble member according to its bias to historical observation (performance criterion) and its distance to other ensemble members in the reference period (interdependence criterion) (Knutti
20 et al., 2017; Sanderson et al., 2017). The biases and distances in the REA and PI methods were calculated based on the diagnostic metric of the climatological mean of streamflow.

The five performance-based weighting methods are the climate prediction index (CPI), upgraded reliability ensemble averaging (UREA), the skill score of the representation of the annual cycle (RAC), Bayesian model averaging (BMA), and the evaluation of the probability density function (PDF). All of these methods only consider the differences of climate simulations
25 to historical observation, but they differ in the metrics or algorithms used to determine weights. The CPI assigns weights based on the biases in the climatological mean and assumes that the simulated climatological mean follows a Gaussian distribution (Murphy et al., 2004). UREA considers biases in both the climatological mean and the inter-annual variance to determine weights (Xu et al., 2010). Both the RAC and BMA calculate weights based on monthly series. The RAC defines a skill score in simulating the annual cycle according to the relationship among the correlation coefficient, standard deviations and centered
30 root-mean-square error (Taylor, 2001). BMA combines the results of multiple models through the Bayesian theory (Duan et al., 2007; Raftery et al., 2005; Min et al., 2007). The PDF determines weights according to the overlapping area of probability density function between daily simulations and observations (Perkins et al., 2007).

Using all eight methods, the weights were calculated for each of streamflow data series simulated by raw GCM outputs and bias-corrected outputs. For a comparison, raw and bias-corrected temperature and precipitation series were also individually used to calculate climate-based weights using the above weighting methods.

3.4 Data Analysis

5 The extent of inequality of each set of weights was first investigated by the entropy of weights (Déqué and Somot, 2010). The entropy of weights reflects the extent of how a weighting method discriminates the relative reliability between GCM simulations. Next, in order to investigate the impacts of weighting GCM simulations for hydrological impact studies, unequal weights were used to combine the ensemble of hydrological simulations. The impacts of unequal weights were compared to the results obtained using the equal weighting method. The comparison focuses on three aspects: (1) the simulation of reference and future hydrological regimes; (2) the bias of the multi-model ensemble mean during the reference period; and (3) the
10 uncertainty of changes in hydrological indices between future and reference periods.

Specifically, when using the entropy of weights (Eq. (1)), the entropy reaches a maximum value when the weights are equally distributed among ensemble members. A smaller entropy indicates a larger difference among the weights of ensemble members. Thus, the entropy reflects the extent of inequality for a set of weights:

$$E = - \sum_{i=1}^N w_i \ln w_i \quad (1)$$

15 where w_i is the weight assigned to the i th ensemble member, and N is the total number of ensemble members.

Since weighting methods are usually proposed to reduce biases in the ensemble of climate simulations, the multi-model ensemble means determined by these weights are then evaluated in terms of the representation of observation during the reference period. The multi-year averages of three hydrological indices were calculated for each streamflow simulation: (1) annual streamflow; (2) peak streamflow; and (3) the center of timing of annual flow (tCMD: the occurrence day of the midpoint
20 of annual flow). The multi-model mean indices were then obtained based on the weights assigned to each simulation and compared to the indices of observation.

The influences of weighting on the uncertainty of hydrological impacts related to the choice of GCMs are investigated in terms of the changes in four hydrological indices between the reference and future periods: (1) mean annual streamflow; (2) mean streamflow during the high flow period; (3) mean streamflow during the low flow period; and (4) mean peak streamflow
25 (the periods of high and low flow are shown in Table 2). The Monte-Carlo approach was introduced to sample the uncertainty for unequally weighted ensembles (Wilby and Harris, 2006; Chen et al., 2017). The hydrological indices were randomly sampled one thousand times based on the calculated weights. For example, if a climate model simulation is assigned a weight of 0.2, the hydrological index simulated by that climate simulation has a probability of 20% to be chosen as the sample in each Monte-Carlo experiment.

4 Results

4.1 Weights of GCMs

Figure 2 presents the weights calculated based on the streamflow data series simulated by raw GCM outputs and bias-corrected outputs for 8 (one equal and 7 unequal) weighting methods over two watersheds. These results show the ability of different weighting methods to distinguish the performance or reliability of individual ensemble members. The entropy of weights was also calculated to quantify the extent of this disproportion for each set of weights (Table 3). Some weighting methods tend to aggressively discriminate the reliability of GCMs and assign differentiated weights to ensemble members, while other methods assign similar weights to each of them. Specifically, when calculating weights based on raw GCM-simulated streamflows, REA, UREA and CPI produce the weights that most radically discriminate ensemble members among all eight weighting methods for both watersheds. The RAC method generates less differentiated unequal weights, followed by the BMA and PI methods, but weights assigned by the PDF method closely resemble the equal weighting method. However, when calculating weights based on bias-corrected GCM-simulated streamflows, the inequality of weights is reduced, and all the unequal weighting methods receive a lower entropy of weights for both watersheds (Table 3). Most sets of these weights become similar to the equal weighting method, with the exception of REA and UREA for the Xiangjiang watershed, and REA for the Manicouagan-5 watershed (Fig. 2). This result was expected, as the bias correction method brings all GCM simulations to be close to the observations. The differences among GCM simulations become greatly reduced.

In addition, the weights based on the raw and bias-corrected temperature and precipitation time series of GCM simulations were also calculated and are shown in Fig. S1. For the weights based on the raw temperature and precipitation, REA, UREA and CPI still generate the most unequal weights among these weighting methods over both watersheds, as Table 3 indicates. Again, the weights become equalized when calculating weights based on bias-corrected temperature and precipitation.

4.2 Impacts on the hydrological regime

The weights determined by eight weighting methods were first utilized to combine GCM-simulated streamflow series. Figure 3 shows the weighted multi-model mean of monthly mean streamflow for the Xiangjiang watershed. The gray envelope represents the range of monthly mean streamflow simulated using 29 GCM simulations. At the reference period, streamflows simulated by raw GCMs cover a wide range (Figure 3a). However, the equal-weighted multi-model mean streamflow performs better than most of the streamflow series simulated by individual GCMs with respect to reproducing the observed streamflow; even so, the equal-weighted ensemble mean still underestimates the streamflow before the peak (January – May) and overestimates it after the peak (June – September).

For the ensemble mean combined by unequal weights, the three weighting methods that generate highly differentiated weights (REA, UREA and CPI) outperform the equal weighting method with respect to reproducing the observed monthly mean streamflow. The BMA and RAC methods improve the performance of streamflow simulations before the peak at the cost of performance after the peak, while an opposite pattern is observed when using the PI method. The PDF method generates

an ensemble mean of monthly mean streamflows almost identical to that of the equal weighting method. This is an expected result, as the PDF method assigns almost identical weights to all GCM simulations.

Weights calculated based on the raw temperature and precipitation of GCM outputs were also used to construct the ensemble mean of monthly mean streamflows (Fig. S2a,b). Particularly, the ensemble mean hydrographs combined using the REA, UREA and CPI methods largely deviate from the observation. Although REA, UREA and CPI generate highly differentiated weights when based on GCM raw temperatures, their generated ensemble mean streamflows are significantly inferior to that generated by equal weights (Fig. S2a). In addition, when using raw precipitation to calculate weights, the weighting methods perform worse than or similar to those calculated based on streamflow series (Fig. S2b). This reflects the advantage of weighting streamflow series in terms of reproducing the observed mean hydrograph.

The bias correction method can reduce the biases of precipitation and temperature in representing the mean monthly streamflow for the reference period, as indicated by the narrowed envelope (Figure 3c), although a small amount of uncertainty is still observed. The reduction of biases brings about similar weights for all GCM-simulated time series when using bias-corrected GCM-simulated streamflows. Thus, the multi-model ensemble means of monthly mean streamflow constructed by all unequal weighting method are very similar to those constructed by the equal weighting method, as shown in Figure 3c.

For the bias-corrected GCM-simulated streamflow at the future period (Figure 3d), a larger uncertainty related to the use of climate models is observed, as indicated by the wider envelope of the mean monthly streamflow. This may be because the bias of GCM outputs is non-stationary. All bias correction methods are based on a common assumption that the bias of climate model outputs is constant over time. However, this assumption may not always be true because of natural climate variability and climate sensitivity to various forcings (Hui et al., 2019; Chen et al., 2015), and most weight methods still follow the same assumption. In other words, the bias non-stationarity implies that climate models differ in their ability to simulate the climate for the future period. The weights calculated in the reference period may not be applicable in the future period. The results of this study also proved this, as all of the weighting methods project similar ensemble means of monthly mean streamflows for the future period.

Figure 4 presents the same information as Figure 3 but for the Manicouagan-5 watershed. Nearly half of the monthly mean streamflow time series simulated by raw GCM outputs have delayed peak (June) compared to the observed one (May) at the reference period, which leads to the delayed peak streamflow of the weighted multi-model mean streamflows for all weighting methods (Figure 4a). Nonetheless, when using raw GCM-simulated streamflow series to calculate weights, the multi-model mean streamflows perform better than or similar to those simulated using GCM raw temperature and precipitation data (Fig. S2c). However, for the bias-corrected streamflow series, the uncertainty of monthly streamflows simulated by individual bias-corrected GCMs is largely reduced and the problem of delayed peak streamflow is corrected (Figure 4c). Similar to the case in the Xiangjiang watershed, all unequally weighted multi-model mean streamflows are identical to that of the equal weighting method. For the future period, although the uncertainty of single bias-corrected GCM-simulated streamflows increases (Figure 4d), there are still very little differences among the future multi-model mean streamflows combined by different weighting methods.

4.3 Bias in multi-model mean

In order to quantify the performance of weighting methods with respect to reproducing the multi-model ensemble mean, biases of the multi-model ensemble mean relative to corresponding observation were calculated for the reference period in terms of three hydrological indices (mean annual streamflow, mean peak streamflow and mean center of timing of annual flow; tCMD). A smaller bias represents a better performance. Figure 5 presents the biases of weighted multi-model mean indices over the Xiangjiang watershed. For the streamflows simulated using raw GCM outputs, the weighting methods show varied performance in terms of reproducing observed indices (Figure 5a-c). Except for the PI method, the unequal-weighted multi-model means more or less outperform the equal weighting method in terms of reducing biases in mean annual streamflow and mean center timing, while an opposite result is observed in mean peak streamflow. This may be because only the mean value (climatological mean or monthly mean series) was used as the evaluation metric when determining weights, while peak or extreme values were not considered. Additionally, weights calculated based on the raw temperature and precipitation of GCM outputs were used to calculate multi-model mean indices for comparison (Fig. S3a-c). When using raw temperature series of GCMs to determine weights, they often bring about more biases in mean annual streamflow and tCMD. The weights based on raw precipitation show some superiority in reducing bias in mean peak streamflow. However, when using bias-corrected GCM-simulated streamflows to calculate weights (Figure 5d-f), the biases in multi-model mean indices are much less varied among different weighting methods. This is similar to the previous results of hydrological regimes.

For the case in the Manicouagan-5 watershed, twenty-five of the 29 streamflow series simulated by raw GCMs have larger mean annual streamflows and mean peak streamflows than those of the observations, and 26 series generate delayed tCMD. This leads to the overestimation of multi-model mean indices for all weighting methods (Figure 6a-c). Compared to the equal weighting method, all unequal weighting methods overcome this overestimation more or less. The three weighting methods that generate highly differentiated weights (REA, UREA and CPI) notably reduce biases for all three hydrological indices. For most weights calculated based on raw temperature and precipitation of GCM outputs (Fig. S3d-f), a certain improvement on mean indices was also observed (the only exception is raw precipitation-based PDF weights). Compared to weights calculated using streamflow series, nearly all weights based on GCM-simulated streamflows reduce more biases than those based on temperature and precipitation. However, when using bias-corrected GCM-simulated streamflows (Fig. 6d-f), again, all weighting methods generate very similar mean indices to the equal weighting method, since the biases among different GCM-simulated streamflows have been largely reduced by the bias correction method.

4.4 Impacts on uncertainty

In addition to the multi-model ensemble mean, the impacts of weighting GCM simulations on uncertainty of hydrological responses also need to be assessed. Thus, this study also evaluated how unequal weighting methods affect the uncertainty of hydrological impacts related to the choice of GCMs. Figures 7 and 8 present the box plots of changes in 4 hydrological indices (mean annual streamflow, mean streamflow during the high/low flow periods and mean peak streamflow) between the

reference and future periods. The box plots of the equal weighting method are depicted using 29 values simulated by each climate simulation, while the box plots of 7 unequal weighting methods are constructed using 1,000 values sampled by the Monte-Carlo approach based on assigned weights. For example, a simulation with 2-times the weight as another one will occur 2-times as often as that one in the 1,000 samples of Monte-Carlo experiments. While the 1,000 samples still only consist of 5 the 29 values, the occurrence of each value reflects its possibility to be chosen and presents the uncertainty related to the choice of GCMs determined by assigned weights.

Figure 7 presents the uncertainty of hydrological changes for the Xiangjiang watershed. When using raw GCM-simulated streamflows (Figure 7a-d), depending on the weighting methods, unequal weights show the varying effects on the uncertainty. Both the PDF and PI methods suggest similar uncertainties to those of the equal weighting method for all four hydrological 10 indices. The BMA and RAC methods generate slightly larger uncertainty for the change in mean annual streamflow and slightly smaller uncertainty of the change in low streamflow. The two weighting methods that generate the most differentiated weights (REA and UREA) largely reduce the uncertainty and increase the changes of the upper and lower probabilities for all four hydrological variables. The impacts of weights calculated based on raw GCM temperature and precipitation series were also analyzed (Fig. S4a-d). When calculating weights based on raw temperature, REA, UREA and CPI tend to aggressively reduce 15 the uncertainty in mean high streamflow and peak streamflow. Precipitation-based weights show similar influences on uncertainty as weights based on streamflows. However, for the bias-corrected GCM-simulated streamflows (Figure 7e-h), the uncertainty of changes in the four hydrological indices is similar among all weighting methods.

Figure 8 presents the uncertainty of hydrological impacts in terms of four hydrological indices over the Manicouagan-5 watershed. For weights calculated using raw GCM-simulated streamflows (Figure 8a-d), only UREA clearly reduces the 20 uncertainty for mean annual streamflow. The REA, UREA and CPI methods reduce the uncertainty for mean low streamflow and decrease its value of upper probability. There are few differences in the uncertainty of mean high streamflow and peak streamflow among all weighting methods. However, when using bias-corrected GCM-simulated streamflows (Figure 8e-h), again, the uncertainty of changes in all four hydrological indices is very similar among most of the weighting methods. Only CPI suggests slight increases in changes of the lower probability.

25 4.5 Out-of-sample Testing

In the above assessments for weighting methods except their impacts on uncertainty, the weighting methods are mostly evaluated in terms of their performances to simulate observations in the reference period. This kind of assessments has been referred to as “in-sample” testing (Herger et al., 2018). But the performances of weighting methods in the future period (“out-of-sample”) may also need to be investigated. However, there is no observations to be compared with in the future period. 30 Thus, an out-of-sample testing was then performed by conducting model-as-truth experiments (Herger et al., 2018; Abramowitz et al., 2019). In model-as-truth experiments, the output of each climate model was regarded as the “truth” in turn and the outputs of the remaining 28 climate models were used as simulations to this “truth” model. Then, the weights were re-

calculated for these remaining models. Since there is a “truth” at the future period in this case, the performances of weighting methods can be evaluated in terms of reproducing the future “truth”.

Figure 9 shows the results of out-of-sample testing over the Xiangjiang watershed for biases of weighted multi-model mean hydrological indices, which are the same as those in Fig. 5. The left and right sides of each stick respectively represent the biases at the reference and future periods when one climate model is regarded as the truth. Similar to Fig. 5, the bias of weighted mean being closer to 0 means that the corresponding weighting method performs better. In general, the results of out-of-sample testing are similar to those where historical observations are used. For the experiment of streamflows simulated by raw GCM outputs, Fig. 9a-c shows that unequally weighted means more or less become closer to the truth simulation than those of equal weighting for both reference and future periods. The unequal streamflow-based weights can help to reduce the biases. In particular, the three methods with the most differentiated weights (REA, UREA and CPI) reduce more biases of annual streamflow when compared with other methods, in that the ranges of the biases calculated by these three methods are narrower and closer to 0 when different simulations are used as the truth. In addition, although the biases in the future period tend to be larger than those in the reference period, the weighted means still have a slight improvement in most cases. However, for the experiment of using bias-corrected GCM outputs to simulate streamflows, as shown by the similar patterns among equal and unequal weighting methods (Fig. 9d-f), the unequally weighted multi-model means have similar biases to those of using equal weighting method at both reference and future periods. In addition, the results of out-of-sample testing over the Manicouagan-5 watershed are shown in Fig. 10, and generally, they are also similar to the results of using observations (Fig. 6).

5 Discussion

In addition to the equal weighting method, which is a normal strategy for handling multi-model ensembles, many studies have proposed various unequal weighting methods for impact studies (e.g., Giorgi and Mearns, 2002; Sanderson et al., 2017; Xu et al., 2010; Min et al., 2007; Murphy et al., 2004). Most of these methods calculate weights based on the reliability of GCM simulations relative to observed climates, or at least adopt their reliability as one of their weighting criteria. In other words, the performances of GCM simulations are usually evaluated by comparing them to observed climate using certain metrics. However, this method may have two problems. First, the trade-off between multiple climate variables related to the impact variable remains uncertain, which leads to difficulty in obtaining a single set of weights for impact studies. Second, the relationship between climate variables and the impact variable is often non-linear and not explicit, which may jeopardize the validity and reasonableness of climate-based weights in the impact studies. Some examples are the weights based on temperature in the experiment of raw GCM-simulated streamflows in the Xiangjiang watershed, which lead to obviously biased multi-model mean hydrographs at the reference period. But using the weights calculated based on raw GCM precipitation does not lead to such biases. This may be because the runoff generation in the Xiangjiang watershed is dominated more by rainfall than temperature. Therefore, weights calculated using temperature may not reflect a GCMs’ reliability that is relevant to

hydrological responses. On the contrary, for the snow-dominated Manicouagan-5 watershed, the snowmelt-driven spring flood is an important characteristic of its hydrological regime, and both temperature and precipitation conditions have large influences on this process. Thus, weights based on temperature and precipitation do not lead to obviously biased multi-model mean hydrographs in this case. Furthermore, over both watersheds, most weights calculated using raw GCM-simulated streamflows reduce more biases of the mean annual streamflow than those based on raw temperature and precipitation. This is as expected, because weights based on streamflows directly reflect how GCM simulations conform to the observed streamflow and are not affected by the non-linear relationship between climate variables and impact variables. Generally, in the experiment of simulating streamflows using raw GCMs, weights calculated based on streamflows not only circumvent the above two problems, they also bring about fewer biases in mean annual streamflow for the multi-model means.

Since bias correction methods are routinely applied to GCM outputs for hydrological impact assessments, this study considered two experiments where raw and bias-corrected GCM-simulated streamflows were separately used to determine weights. The performances of weighting methods are separately examined for the two experiments. Although the equal weighting is often used by default to combine bias-corrected ensembles in hydrological impact studies, whether unequal weighting is necessary still remains to be investigated (Alder and Hostetler, 2019). As shown in Figures 3 and 4, biases in the simulated mean monthly streamflows are greatly reduced for the reference period after bias correction. This change in biases affects the ability of most unequal weighting methods to discriminate the performances of climate simulations. In this experiment, all of the weighting methods assign similar weights to all simulations (as indicated by the decline of entropy of weights calculated by each weighting method). This is because climate simulations become rather close to each other in the reference period, and all weighting methods except REA in this study only rely on reference performances (which means that they lose the ability to discriminate the performances of climate simulations). As to the REA method, even though it considers future projections in its convergence criterion when calculating weights and its weights are still the most differentiated for the bias-corrected ensemble (as shown in Fig. 2), they bring little impacts on the final results of the multi-model mean. In addition, the PI method considers independency among simulations, but it only relies on reference values which have been tuned by the bias-correction method. The ability of independent criterion may be affected because of the bias correction. In general, in this experiment, compared to the equal weighing method, unequal weighting methods do not bring about much disparateness to the results of hydrological impacts. The out-of-sample testing also manifested the same phenomena. Therefore, it is still viable to attend to the bias-corrected ensembles with the equal weighting method.

Despite the choices of variables used to calculate weights, the establishment of any weighting method involves subjective choices of diagnostic metrics, its translation to performance measurement, and normalization to weights (Knutti et al., 2017; Santer et al., 2009). For example, in the RAC method, the correlation coefficient and standard deviation are used as diagnostic metrics, and GCM skills are measured through the translation of a fourth-order formulation. The skill scores are then divided by their sum to be normalized. Any of these steps can ultimately affect the property of a weighting method. For example, the REA, UREA and CPI methods are inclined to generate more differentiated weights, while other methods assign more similar weights to ensemble members. All of these aspects in weighting methods are often predefined without detailed examination

or based on expert experience and, thus, can actually introduce several layers of subjective uncertainty. An improper weighting method may even cause a risk of reducing projection accuracy (Weigel et al., 2010), and extremely aggressive weighting may conceal the uncertainty rather than reduce it (Chen et al., 2017). Thus, notwithstanding the equal weighting is not a perfect solution, model weighting methods should be used with cautions and the results of equal weighting should be presented along with those of unequal weighting methods.

Moreover, some risks may exist in the usage of weighting methods in impact studies. Firstly, weights are generally assigned to climate simulations in a static way (i.e. weights in the future period are the same as those in the reference period). This usage shares the same assumption with bias-correction methods that the performances of GCM simulations are stable and stationary. However, some studies have shown that model skills are nonstationary in a changing climate (Weigel et al., 2010; Miao et al., 2016), and models with better performance in the reference period do not necessarily provide more realistic signals of climate change (Reifen and Toumi, 2009; Knutti et al., 2010). The way to deal with the dynamic reliability of climate models deserves further studies. Secondly, many researchers and end-users in hydrological impacts only consider one diagnostic metric to determine weights, such as the climatological mean (e.g., Wilby and Harris, 2006; Chen et al., 2017). It is not clear whether reducing the bias of one specific metric can transfer to other metrics. The weights calculated using the raw GCM-simulated streamflows in the Xiangjiang watershed are one negative example, where the bias in mean annual streamflow is reduced while the bias in the mean peak streamflow is enlarged. Some studies have also shown similar problems (Jun et al., 2012; Santer et al., 2009). For example, Jun et al. (2012) demonstrated that there is little relationship between a GCMs' ability to reproduce mean temperature state and trend of temperature. Actually, a set of metrics can be introduced to determine weights (e.g., Sanderson et al., 2017). Some studies suggested using calibrated multiple metrics because it can improve the rationality of weighted multi-model mean (Knutti et al., 2017; Lorenz et al., 2018), while some argued that multiple metrics form another level of uncertainty within weighting methods (Christensen et al., 2010). Thus, the best way to choose proper metrics and synthesize performances in multiple metrics still remains in doubt and deserves further research.

There is a limitation in the hydrological modeling in this study. Only large watersheds were considered, as well as a lumped hydrological model. When using a lumped model, the nonlinear relationship between the climate variables and the impact variable (streamflow) may not be sufficiently revealed. Spatial differences between different climate simulations only affect the basin-averaged inputs to the hydrological model but not directly affect the process of runoff generation and streamflow routing (Lebel et al., 1987). Temporal variations of climate simulations may be partially reduced by the lumped hydrological model as well. With the help of other more sophisticated hydrological models (such as distributed models), the differences between climate-based weights and streamflow-based weights may become more obvious. For the experiment of raw GCM-simulated streamflows, the weights based on streamflow perform better than those based on climate variables. This may be related to large differences among climate simulations. But in the experiment of streamflows simulated using bias-corrected GCM outputs, that no much discrepancy is seen in the performances between unequal and equal weighting may be partly because only a simple hydrological model is used. In other words, the remaining differences among corrected climate

simulations may not be well presented in streamflow simulations when a lumped hydrological model is used in such large watersheds.

6 Conclusion

In order to weight climate models based on the impact variable and to quantify its influences on the impact assessment, this study assigns weights to an ensemble of 29 CMIP5 GCMs over two watersheds through a group of weighting methods based on GCM-simulated streamflow time series. Streamflow series are simulated by separately inputting the raw and bias-corrected GCM simulations to hydrological models. Using streamflows to determine weights is straightforward and can avoid the difficulty of combining weights based on multiple climate variables for impact studies. The influences of these unequal weights on the assessment of hydrological impacts were then investigated and compared to the common strategy of model democracy.

This study concludes that for the streamflows simulated using raw GCM outputs without bias correction, using unequal weights has some advantages over the equal weighting method in simulating observed hydrographs and reducing the biases of multi-model means in mean annual streamflow. In particular, the weights calculated based on streamflows can reduce more biases of multi-model mean annual streamflow and better reproduce observed hydrographs, compared with the weights calculated based on climate variables. However, when using bias-corrected GCM outputs to simulate streamflow, GCM simulations were brought close to the observations by the bias correction method. Consequently, the weights assigned to climate simulations become similar to each other, resulting in similar multi-model means and uncertainty of hydrological impacts for all unequal weighting methods. Therefore, the equal weighting method is still a conservative and viable option for combining the bias-corrected multi-model ensembles. Or, if an unequal weighting method is applied, it is better to present it with a detailed explanation of the weighting procedure, as well as the results of using equal weighting method to end-users.

Data availability

The climate simulation data can be accessed from the CMIP5 archive (<https://esgf-node.llnl.gov/projects/esgf-llnl/>, last access: 3 June 2019). The observation data in the Xiangjiang and Manicouagan-5 are not publicly available due to the restrictions of data providers, but can be requested by contacting the corresponding author.

25 Author contributions

JC conceived the original idea, and HMW and JC designed the methodology. JC and HC collected the data. HMW developed the model code and performed the simulations, with some contributions from XL. HMW, JC, CYX, SG and PX contributed to the interpretation of results. HMW wrote the paper, and JC, CYX, SG and PX revised the paper.

Competing interests

The authors declare that they have no conflict of interest.

Acknowledgements

This work was partially supported by the National Natural Science Foundation of China (Grant No. 51779176, 51539009, 91547205), the Overseas Expertise Introduction Project for Discipline Innovation (111 Project) funded by Ministry of Education and State Administration of Foreign Experts Affairs P.R. China (Grant No. B18037), the Thousand Youth Talents Plan from the Organization Department of CCP Central Committee (Wuhan University, China) and the Research Council of Norway (FRINATEK Project 274310). The authors would like to acknowledge the World Climate Research Program Working Group on Coupled Modelling, and all climate modeling institutions listed in Table 1 for making GCM outputs available. We also thank Hydro-Québec and the Changjiang Water Resources Commission for providing observation data in the Manicouagan-5 and Xiangjiang watersheds, respectively.

References

- Abramowitz, G., Herger, N., Gutmann, E., Hammerling, D., Knutti, R., Leduc, M., Lorenz, R., Pincus, R., and Schmidt, G. A.: ESD Reviews: Model dependence in multi-model climate ensembles: weighting, sub-selection and out-of-sample testing, *Earth System Dynamics*, 10, 91-105, <https://doi.org/10.5194/esd-10-91-2019>, 2019.
- Alder, J. R., and Hostetler, S. W.: The Dependence of Hydroclimate Projections in Snow-Dominated Regions of the Western United States on the Choice of Statistically Downscaled Climate Data, *Water Resources Research*, 55, 2279-2300, <https://doi.org/10.1029/2018wr023458>, 2019.
- Arsenault, R., Gatién, P., Renaud, B., Brissette, F., and Martel, J.-L.: A comparative analysis of 9 multi-model averaging approaches in hydrological continuous streamflow simulation, *Journal of Hydrology*, 529, 754-767, <https://doi.org/10.1016/j.jhydrol.2015.09.001>, 2015.
- Chen, J., Brissette, F. P., Poulin, A., and Leconte, R.: Overall uncertainty study of the hydrological impacts of climate change for a Canadian watershed, *Water Resources Research*, 47, W12509, <https://doi.org/10.1029/2011wr010602>, 2011.
- Chen, J., Brissette, F. P., Chaumont, D., and Braun, M.: Performance and uncertainty evaluation of empirical downscaling methods in quantifying the climate change impacts on hydrology over two North American river basins, *Journal of Hydrology*, 479, 200-214, <https://doi.org/10.1016/j.jhydrol.2012.11.062>, 2013.
- Chen, J., Brissette, F. P., and Lucas-Picher, P.: Assessing the limits of bias-correcting climate model outputs for climate change impact studies, *Journal of Geophysical Research: Atmospheres*, 120, 1123-1136, <https://doi.org/10.1002/2014jd022635>, 2015.

- Chen, J., Brissette, F. P., Lucas-Picher, P., and Caya, D.: Impacts of weighting climate models for hydro-meteorological climate change studies, *Journal of Hydrology*, 549, 534-546, <https://doi.org/10.1016/j.jhydrol.2017.04.025>, 2017.
- Cheng, L., and AghaKouchak, A.: A methodology for deriving ensemble response from multimodel simulations, *Journal of Hydrology*, 522, 49-57, <https://doi.org/10.1016/j.jhydrol.2014.12.025>, 2015.
- 5 Chiew, F. H. S., Teng, J., Vaze, J., Post, D. A., Perraud, J. M., Kirono, D. G. C., and Viney, N. R.: Estimating climate change impact on runoff across southeast Australia: Method, results, and implications of the modeling method, *Water Resources Research*, 45, <https://doi.org/10.1029/2008wr007338>, 2009.
- Christensen, J. H., Kjellström, E., Giorgi, F., Lenderink, G., and Rummukainen, M.: Weight assignment in regional climate models, *Climate Research*, 44, 179-194, <https://doi.org/10.3354/cr00916>, 2010.
- 10 Déqué, M., and Somot, S.: Weighted frequency distributions express modelling uncertainties in the ENSEMBLES regional climate experiments, *Climate Research*, 44, 195-209, <https://doi.org/10.3354/cr00866>, 2010.
- Duan, Q., Sorooshian, S., and Gupta, V.: Effective and efficient global optimization for conceptual rainfall-runoff models, *Water Resources Research*, 28, 1015-1031, <https://doi.org/10.1029/91WR02985>, 1992.
- Duan, Q., Ajami, N. K., Gao, X., and Sorooshian, S.: Multi-model ensemble hydrologic prediction using Bayesian model averaging, *Advances in Water Resources*, 30, 1371-1386, <https://doi.org/10.1016/j.advwatres.2006.11.014>, 2007.
- 15 Giorgetta, M. A., Jungclaus, J., Reick, C. H., Legutke, S., Bader, J., Böttinger, M., Brovkin, V., Crueger, T., Esch, M., Fieg, K., Glushak, K., Gayler, V., Haak, H., Hollweg, H.-D., Ilyina, T., Kinne, S., Kornblueh, L., Matei, D., Mauritsen, T., Mikolajewicz, U., Mueller, W., Notz, D., Pithan, F., Raddatz, T., Rast, S., Redler, R., Roeckner, E., Schmidt, H., Schnur, R., Segschneider, J., Six, K. D., Stockhause, M., Timmreck, C., Wegner, J., Widmann, H., Wieners, K.-H., Claussen, M.,
- 20 Marotzke, J., and Stevens, B.: Climate and carbon cycle changes from 1850 to 2100 in MPI-ESM simulations for the Coupled Model Intercomparison Project phase 5, *Journal of Advances in Modeling Earth Systems*, 5, 572-597, <https://doi.org/10.1002/jame.20038>, 2013.
- Giorgi, F., and Mearns, L. O.: Calculation of Average, Uncertainty Range, and Reliability of Regional Climate Changes from AOGCM Simulations via the “Reliability Ensemble Averaging” (REA) Method, *Journal of Climate*, 15, 1141-1158, [https://doi.org/10.1175/1520-0442\(2002\)015<1141:coaura>2.0.co;2](https://doi.org/10.1175/1520-0442(2002)015<1141:coaura>2.0.co;2), 2002.
- 25 Gleckler, P. J., Taylor, K. E., and Doutriaux, C.: Performance metrics for climate models, *Journal of Geophysical Research*, 113, D06104, <https://doi.org/10.1029/2007jd008972>, 2008.
- Herger, N., Abramowitz, G., Knutti, R., Angéilil, O., Lehmann, K., and Sanderson, B. M.: Selecting a climate model subset to optimise key ensemble properties, *Earth System Dynamics*, 9, 135-151, <https://doi.org/10.5194/esd-9-135-2018>, 2018.
- 30 Hidalgo, H. G., and Alfaro, E. J.: Skill of CMIP5 climate models in reproducing 20th century basic climate features in Central America, *International Journal of Climatology*, 35, 3397-3421, <https://doi.org/10.1002/joc.4216>, 2015.
- Hui, Y., Chen, J., Xu, C. Y., Xiong, L., and Chen, H.: Bias nonstationarity of global climate model outputs: The role of internal climate variability and climate model sensitivity, *International Journal of Climatology*, 39, 2278-2294, <https://doi.org/10.1002/joc.5950>, 2019.

- Hutchinson, M. F., McKenney, D. W., Lawrence, K., Pedlar, J. H., Hopkinson, R. F., Milewska, E., and Papadopol, P.: Development and Testing of Canada-Wide Interpolated Spatial Models of Daily Minimum–Maximum Temperature and Precipitation for 1961–2003, *Journal of Applied Meteorology and Climatology*, 48, 725-741, <https://doi.org/10.1175/2008jamc1979.1>, 2009.
- 5 IPCC: Evaluation of Climate Models, in: *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change*, edited by: Stocker, T. F., Qin, D., Plattner, G.-K., Tignor, M. M. B., Allen, S. K., Boschung, J., Nauels, A., Xia, Y., Bex, V., and Midgley, P. M., Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 741-866, 2013.
- IPCC: Summary for Policymakers, in: *Climate Change 2014 – Impacts, Adaptation and Vulnerability: Part A: Global and Sectoral Aspects: Working Group II Contribution to the IPCC Fifth Assessment Report*, edited by: Barros, V. R., Field, C. B., Dokken, D. J., Mastrandrea, M. D., Mach, K. J., Bilir, T. E., Chatterjee, M., Ebi, K. L., Estrada, Y. O., Genova, R. C., Girma, B., Kissel, E. S., Levy, A. N., MacCracken, S., Mastrandrea, P. R., and White, L. L., Cambridge University Press, Cambridge, 1-32, 2014.
- 10 Jun, M., Knutti, R., and Nychka, D. W.: Spatial Analysis to Quantify Numerical Model Bias and Dependence, *Journal of the American Statistical Association*, 103, 934-947, <https://doi.org/10.1198/016214507000001265>, 2012.
- Knutti, R., Furrer, R., Tebaldi, C., Cermak, J., and Meehl, G. A.: Challenges in Combining Projections from Multiple Climate Models, *Journal of Climate*, 23, 2739-2758, <https://doi.org/10.1175/2009jcli3361.1>, 2010.
- Knutti, R., Sedláček, J., Sanderson, B. M., Lorenz, R., Fischer, E. M., and Eyring, V.: A climate model projection weighting scheme accounting for performance and interdependence, *Geophysical Research Letters*,
 20 <https://doi.org/10.1002/2016gl072012>, 2017.
- Lebel, T., Bastin, G., Obled, C., and Creutin, J. D.: On the accuracy of areal rainfall estimation: A case study, *Water Resources Research*, 23, 2123-2134, <https://doi.org/10.1029/WR023i01p02123>, 1987.
- Lorenz, P., and Jacob, D.: Validation of temperature trends in the ENSEMBLES regional climate model runs driven by ERA40, *Climate Research*, 44, 167-177, <https://doi.org/10.3354/cr00973>, 2010.
- 25 Lorenz, R., Herger, N., Sedláček, J., Eyring, V., Fischer, E. M., and Knutti, R.: Prospects and Caveats of Weighting Climate Models for Summer Maximum Temperature Projections Over North America, *Journal of Geophysical Research: Atmospheres*, 123, 4509-4526, <https://doi.org/10.1029/2017jd027992>, 2018.
- Maurer, E. P.: Uncertainty in hydrologic impacts of climate change in the Sierra Nevada, California, under two emissions scenarios, *Climatic Change*, 82, 309-325, <https://doi.org/10.1007/s10584-006-9180-9>, 2007.
- 30 Miao, C., Su, L., Sun, Q., and Duan, Q.: A nonstationary bias-correction technique to remove bias in GCM simulations, *Journal of Geophysical Research: Atmospheres*, 121, 5718-5735, <https://doi.org/10.1002/2015jd024159>, 2016.
- Min, S. K., Simonis, D., and Hense, A.: Probabilistic climate change predictions applying Bayesian model averaging, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 365, 2103-2116, <https://doi.org/10.1098/rsta.2007.2070>, 2007.

- Minville, M., Brissette, F., and Leconte, R.: Uncertainty of the impact of climate change on the hydrology of a nordic watershed, *Journal of Hydrology*, 358, 70-83, <https://doi.org/10.1016/j.jhydrol.2008.05.033>, 2008.
- Mpelasoka, F. S., and Chiew, F. H. S.: Influence of Rainfall Scenario Construction Methods on Runoff Projections, *Journal of Hydrometeorology*, 10, 1168-1183, <https://doi.org/10.1175/2009jhm1045.1>, 2009.
- 5 Murphy, J. M., Sexton, D. M., Barnett, D. N., Jones, G. S., Webb, M. J., Collins, M., and Stainforth, D. A.: Quantification of modelling uncertainties in a large ensemble of climate change simulations, *Nature*, 430, 768-772, <https://doi.org/10.1038/nature02771>, 2004.
- Nash, J. E., and Sutcliffe, J. V.: River flow forecasting through conceptual models part I — A discussion of principles, *Journal of Hydrology*, 10, 282-290, [https://doi.org/10.1016/0022-1694\(70\)90255-6](https://doi.org/10.1016/0022-1694(70)90255-6), 1970.
- 10 Perkins, S. E., Pitman, A. J., Holbrook, N. J., and McAneney, J.: Evaluation of the AR4 Climate Models' Simulated Daily Maximum Temperature, Minimum Temperature, and Precipitation over Australia Using Probability Density Functions, *Journal of Climate*, 20, 4356-4376, <https://doi.org/10.1175/jcli4253.1>, 2007.
- Perrin, C., Michel, C., and Andréassian, V.: Improvement of a parsimonious model for streamflow simulation, *Journal of Hydrology*, 279, 275-289, [https://doi.org/10.1016/s0022-1694\(03\)00225-7](https://doi.org/10.1016/s0022-1694(03)00225-7), 2003.
- 15 Raftery, A. E., Gneiting, T., Balabdaoui, F., and Polakowski, M.: Using Bayesian Model Averaging to Calibrate Forecast Ensembles, *Monthly Weather Review*, 133, 1155-1174, <https://doi.org/10.1175/mwr2906.1>, 2005.
- Reichler, T., and Kim, J.: How Well Do Coupled Models Simulate Today's Climate?, *Bulletin of the American Meteorological Society*, 89, 303-312, <https://doi.org/10.1175/bams-89-3-303>, 2008.
- Reifen, C., and Toumi, R.: Climate projections: Past performance no guarantee of future skill?, *Geophysical Research Letters*, 20 36, <https://doi.org/10.1029/2009gl038082>, 2009.
- Risbey, J. S., and Entekhabi, D.: Observed Sacramento Basin streamflow response to precipitation and temperature changes and its relevance to climate impact studies, *Journal of Hydrology*, 184, 209-223, [https://doi.org/10.1016/0022-1694\(95\)02984-2](https://doi.org/10.1016/0022-1694(95)02984-2), 1996.
- Sanderson, B. M., Knutti, R., and Caldwell, P.: A Representative Democracy to Reduce Interdependency in a Multimodel Ensemble, *Journal of Climate*, 28, 5171-5194, <https://doi.org/10.1175/jcli-d-14-00362.1>, 2015.
- 25 Sanderson, B. M., Wehner, M., and Knutti, R.: Skill and independence weighting for multi-model assessments, *Geoscientific Model Development*, 10, 2379-2395, <https://doi.org/10.5194/gmd-10-2379-2017>, 2017.
- Santer, B. D., Taylor, K. E., Gleckler, P. J., Bonfils, C., Barnett, T. P., Pierce, D. W., Wigley, T. M., Mears, C., Wentz, F. J., Bruggemann, W., Gillett, N. P., Klein, S. A., Solomon, S., Stott, P. A., and Wehner, M. F.: Incorporating model quality information in climate change detection and attribution studies, *Proceedings of The National Academy of Sciences of the United States of America*, 106, 14778-14783, <https://doi.org/10.1073/pnas.0901736106>, 2009.
- 30 Schmidli, J., Frei, C., and Vidale, P. L.: Downscaling from GCM precipitation: a benchmark for dynamical and statistical downscaling methods, *International Journal of Climatology*, 26, 679-689, <https://doi.org/10.1002/joc.1287>, 2006.

- Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *Journal of Geophysical Research: Atmospheres*, 106, 7183-7192, <https://doi.org/10.1029/2000jd900719>, 2001.
- Taylor, K. E., Stouffer, R. J., and Meehl, G. A.: An Overview of CMIP5 and the Experiment Design, *Bulletin of the American Meteorological Society*, 93, 485-498, <https://doi.org/10.1175/bams-d-11-00094.1>, 2012.
- 5 Tebaldi, C., and Knutti, R.: The use of the multi-model ensemble in probabilistic climate projections, *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 365, 2053-2075, <https://doi.org/10.1098/rsta.2007.2076>, 2007.
- Teutschbein, C., and Seibert, J.: Bias correction of regional climate model simulations for hydrological climate-change impact studies: Review and evaluation of different methods, *Journal of Hydrology*, 456-457, 12-29, <https://doi.org/10.1016/j.jhydrol.2012.05.052>, 2012.
- 10 Valéry, A., Andréassian, V., and Perrin, C.: ‘As simple as possible but not simpler’: What is useful in a temperature-based snow-accounting routine? Part 2 – Sensitivity analysis of the Cemaneige snow accounting routine on 380 catchments, *Journal of Hydrology*, 517, 1176-1187, <https://doi.org/10.1016/j.jhydrol.2014.04.058>, 2014.
- Wang, H.-M., Chen, J., Cannon, A. J., Xu, C.-Y., and Chen, H.: Transferability of climate simulation uncertainty to hydrological impacts, *Hydrology and Earth System Sciences*, 22, 3739-3759, <https://doi.org/10.5194/hess-22-3739-2018>, 2018.
- 15 Weigel, A. P., Knutti, R., Liniger, M. A., and Appenzeller, C.: Risks of Model Weighting in Multimodel Climate Projections, *Journal of Climate*, 23, 4175-4191, <https://doi.org/10.1175/2010jcli3594.1>, 2010.
- Whitfield, P. H., and Cannon, A. J.: Recent Variations in Climate and Hydrology in Canada, *Canadian Water Resources Journal*, 25, 19-65, <https://doi.org/10.4296/cwrj2501019>, 2000.
- 20 Wilby, R. L., and Harris, I.: A framework for assessing uncertainties in climate change impacts: Low-flow scenarios for the River Thames, UK, *Water Resources Research*, 42, W02419, <https://doi.org/10.1029/2005wr004065>, 2006.
- Xu, Y., Gao, X., and Giorgi, F.: Upgrades to the reliability ensemble averaging method for producing probabilistic climate-change projections, *Climate Research*, 41, 61-81, <https://doi.org/10.3354/cr00835>, 2010.
- 25 Zhao, Y.: Investigation of uncertainties in assessing climate change impacts on the hydrology of a Canadian river watershed, Thèse de doctorat électronique, École de technologie supérieure, Montréal, 2015.

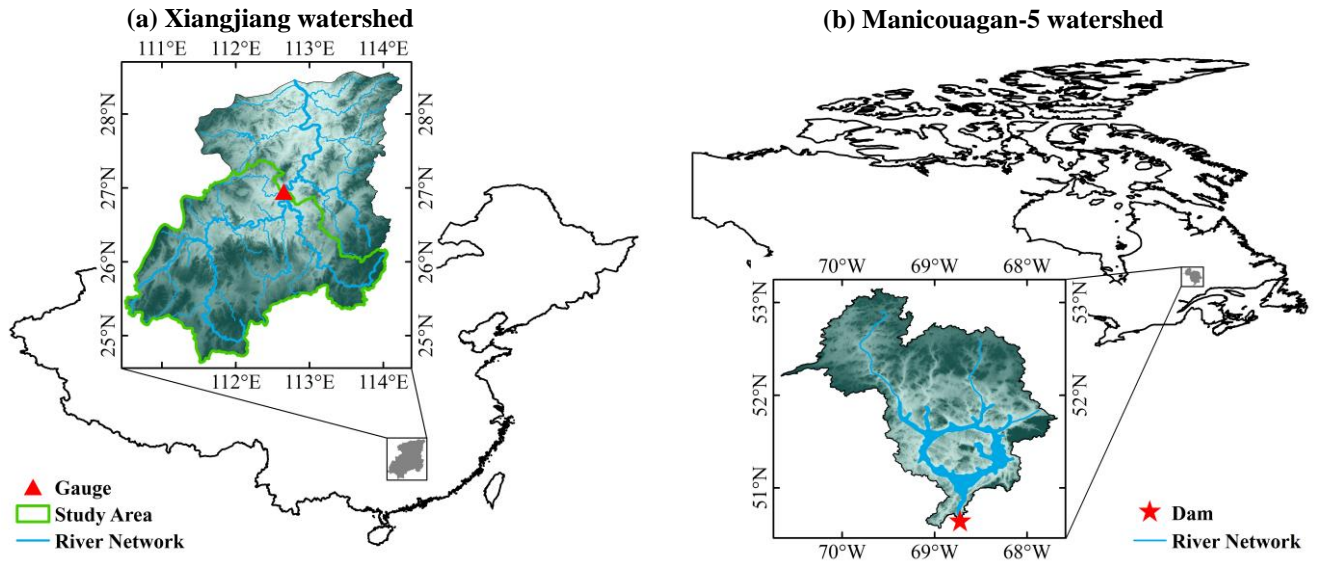


Figure 1. Locations of the (a) Xiangjiang and (b) Manicouagan-5 watersheds. (The study area in the Xiangjiang watershed is one of its sub-basins as the green boundary.)

Table 1. Information about the 29 GCMs used.

No.	Model name	Resolution (Lon. × Lat.)	Institution
1	ACCESS1.0	1.875 × 1.25	Commonwealth Scientific and Industrial Research Organization (CSIRO) and Bureau of Meteorology (BOM), Australia
2	ACCESS1.3	1.875 × 1.25	
3	BCC-CSM1.1	2.8 × 2.8	Beijing Climate Center, China Meteorological Administration
4	BCC-CSM1.1(m)	1.125 × 1.125	
5	BNU-ESM	2.8 × 2.8	College of Global Change and Earth System Science, Beijing Normal University
6	CanESM2	2.8 × 2.8	Canadian Centre for Climate Modelling and Analysis
7	CCSM4	1.25 × 0.94	US National Centre for Atmospheric Research
8	CESM1(CAM5)	1.25 × 0.94	National Science Foundation, Department of Energy, NCAR, USA
9	CMCC-CMS	1.875 × 1.875	Centro Euro-Mediterraneo per I Cambiamenti Climatici
10	CMCC-CM	0.75 × 0.75	
11	CMCC-CESM	3.75 × 3.7	
12	CNRM-CM5	1.4 × 1.4	Centre National de Recherches Météorologiques and Centre Européen de Recherche et Formation Avancée en Calcul Scientifique
13	CSIRO-Mk3.6.0	1.8 × 1.8	Commonwealth Scientific and Industrial Research Organization and Queensland Climate Change Centre of Excellence
14	FGOALS-g2	1.875 × 1.25	LASG, Institute of Atmospheric Physics, Chinese Academy of Sciences, and CESS, Tsinghua University
15	GFDL-CM3	2.5 × 2.0	NOAA Geophysical Fluid Dynamics Laboratory
16	GFDL-ESM2G	2.5 × 2.0	
17	GFDL-ESM2M	2.5 × 2.0	
18	INM-CM4	2.0 × 1.5	Russian Institute for Numerical Mathematics
19	IPSL-CM5A-LR	3.75 × 1.9	Institut Pierre Simon Laplace
20	IPSL-CM5A-MR	2.5 × 1.25	
21	IPSL-CM5B-LR	3.75 × 1.9	
22	MIROC-ESM-CHEM	2.8 × 2.8	Japan Agency for Marine-Earth Science and Technology, Atmosphere and Ocean Research Institute (The University of Tokyo), and National Institute for Environmental Studies
23	MIROC-ESM	2.8 × 2.8	
24	MIROC5	1.4 × 1.4	Atmosphere and Ocean Research Institute (The University of Tokyo), National Institute for Environmental Studies, and Japan Agency for Marine-Earth Science and Technology
25	MPI-ESM-LR	1.875 × 1.875	Max Planck Institute for Meteorology
26	MPI-ESM-MR	1.875 × 1.875	
27	MRI-ESM1	1.125 × 1.125	Meteorological Research Institute
28	MRI-CGCM3	1.1 × 1.1	
29	NorESM1-M	2.5 × 1.875	Norwegian Climate Centre

Table 2. Nash-Sutcliffe Efficiency (NSE) of hydrological models in the calibration and validation periods.

Country	Watershed name	Area (km ²)	High flow	Low flow	Calibration period	NSE calibration	Validation period	NSE validation
China	Xiangjiang	52150	Apr-Jun	Jul-Nov	1975-1987	0.912	1988-2000	0.871
Canada	Manicouagan-5	24610	Mar-Jul	Aug-Feb	1970-1979	0.926	1980-1989	0.881

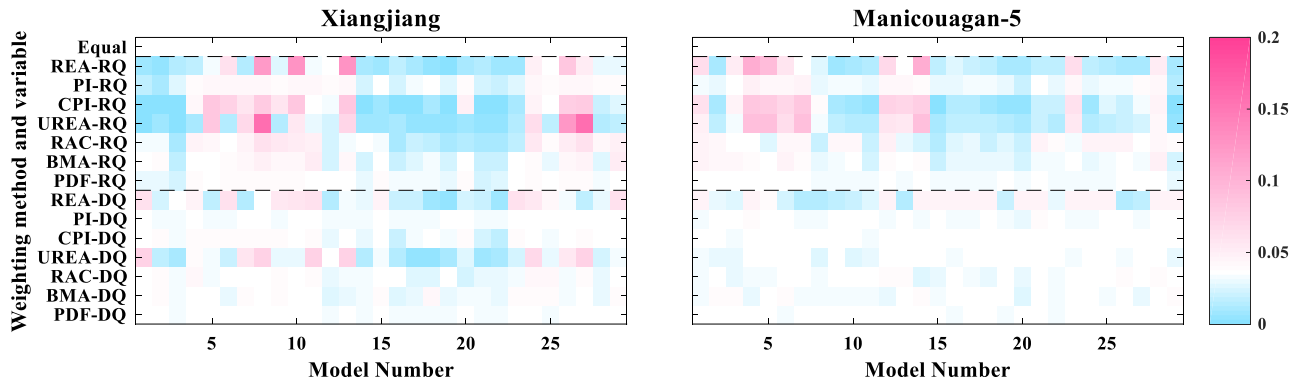


Figure 2. Weights assigned by equal weighting and 7 unequal weighting methods based on raw climate model data-simulated streamflow (RQ) and bias corrected data-simulated streamflow (DQ) for two watersheds. (Equal weight is presented in white, weights greater than equal are presented in red, and weights less than equal in blue.)

5

Table 3. The entropy of weights calculated by equal weighting and 7 unequal weighting methods based on raw climate model data-simulated streamflow (RQ) and bias corrected data-simulated streamflow (DQ) for two watersheds. The entropy of weights calculated based on raw and bias-corrected temperature (RT and DT) and precipitation (RP and DP) are also presented for comparison.

	Xiangjiang watershed						Manicouagan-5 watershed					
	RT	RP	RQ	DT	DP	DQ	RT	RP	RQ	DT	DP	DQ
REA	2.45	3.04	2.93	3.05	3.18	3.22	2.87	3.11	3.06	3.12	3.30	3.29
PI	3.34	3.35	3.33	3.37	3.37	3.37	3.34	3.34	3.34	3.36	3.36	3.37
CPI	2.46	2.92	2.86	3.37	3.36	3.35	2.99	3.12	3.00	3.37	3.37	3.37
UREA	2.72	3.00	2.73	3.33	3.22	3.15	3.02	3.15	3.10	3.33	3.35	3.36
RAC	3.37	3.35	3.25	3.37	3.36	3.36	3.37	3.36	3.32	3.37	3.36	3.36
BMA	3.34	3.36	3.33	3.36	3.36	3.36	3.35	3.36	3.35	3.37	3.36	3.36
PDF	3.36	3.37	3.36	3.37	3.37	3.37	3.37	3.37	3.37	3.37	3.37	3.37
Equal	3.37						3.37					

5

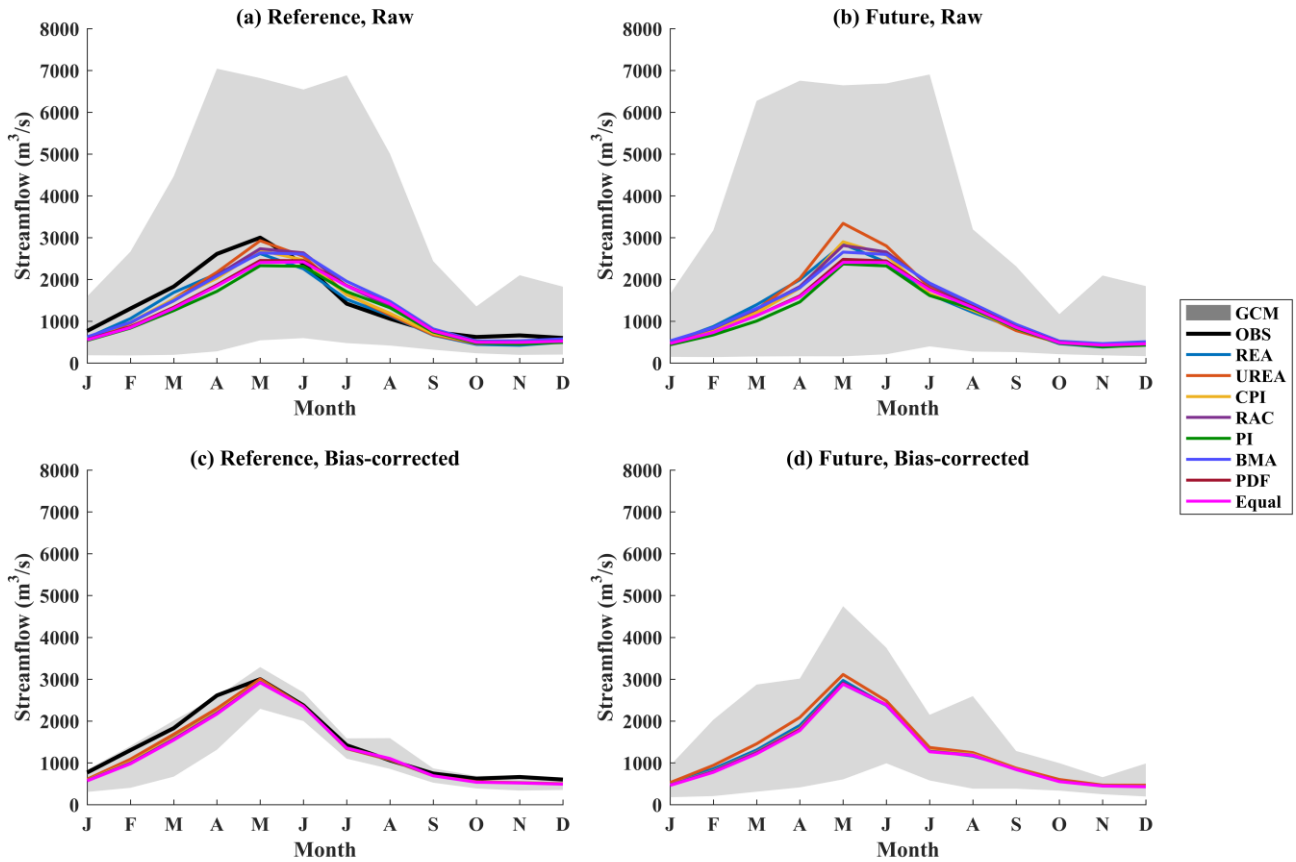


Figure 3. The envelope of monthly mean streamflows simulated by 29 raw and bias-corrected GCM outputs and the multi-model ensemble means of monthly mean streamflows weighted by 8 weighting methods based on GCM-simulated streamflows over the Xiangjiang watershed for the reference and future periods (OBS = the hydrograph simulated from meteorological observation).

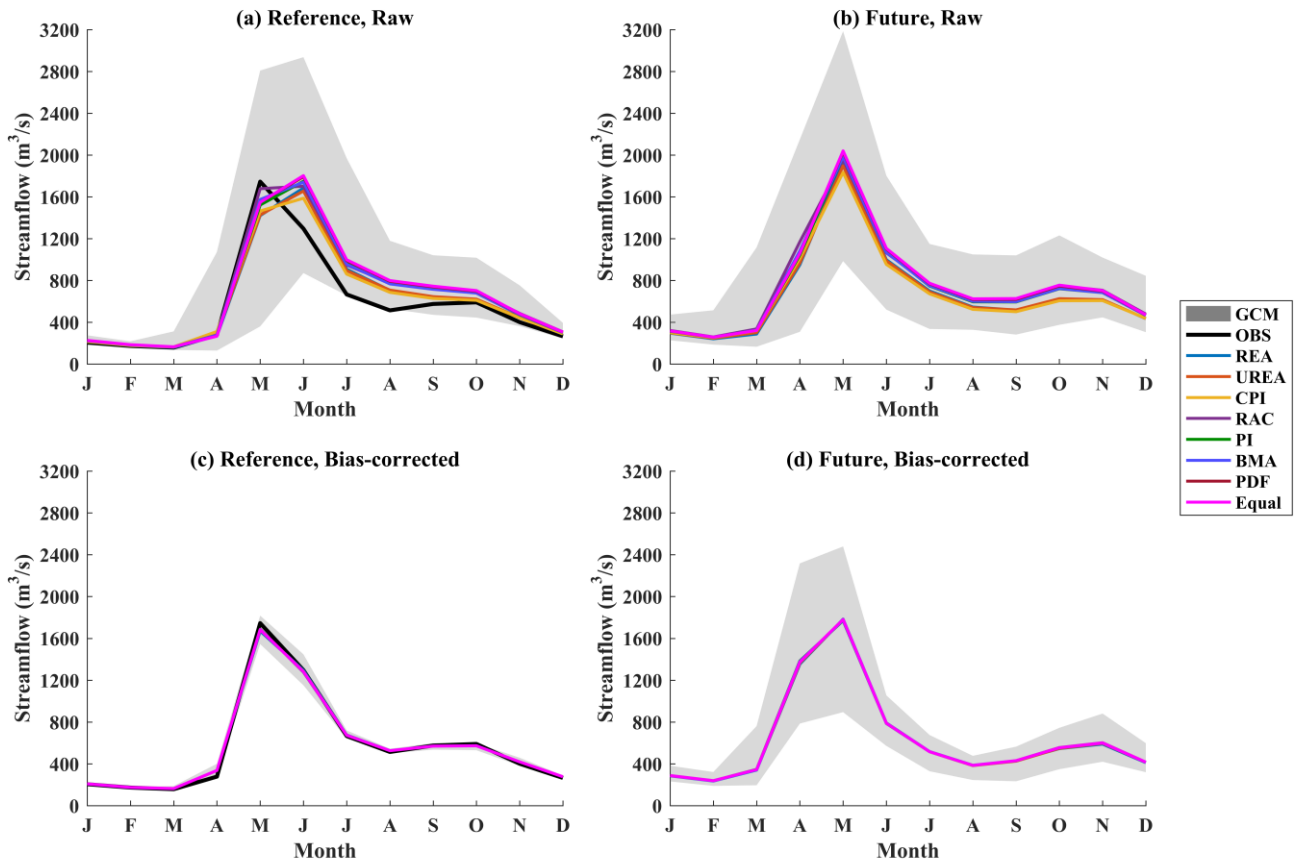


Figure 4. The same as Fig. 3 but for the Manicouagan-5 watershed.

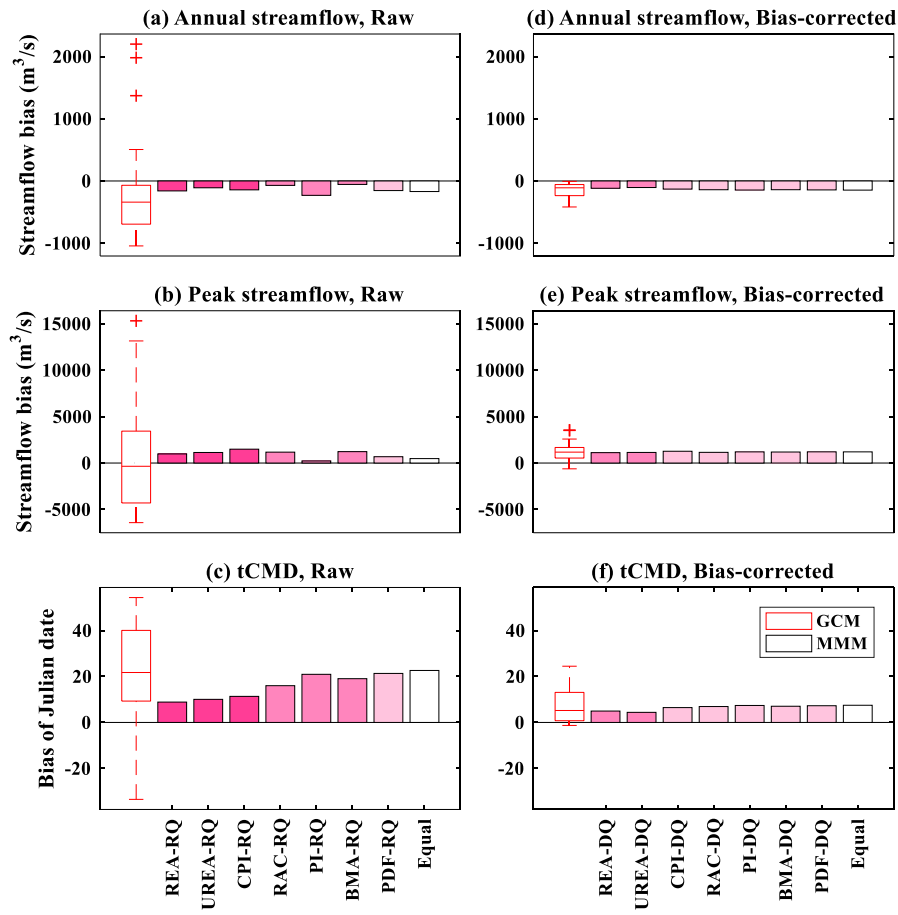


Figure 5. Bias in mean annual streamflow, mean peak streamflow and mean center of timing of annual flow (tCMD) simulated using 29 raw or bias-corrected GCM outputs and the multi-model means (MMM) combined by weights based on raw (RQ) and bias-corrected (DQ) GCM-simulated streamflows in the Xiangjiang watershed in the reference period. (The depth of pink in the MMM bars represents the level of inequality of weights, as indicated in Table 3.)

5

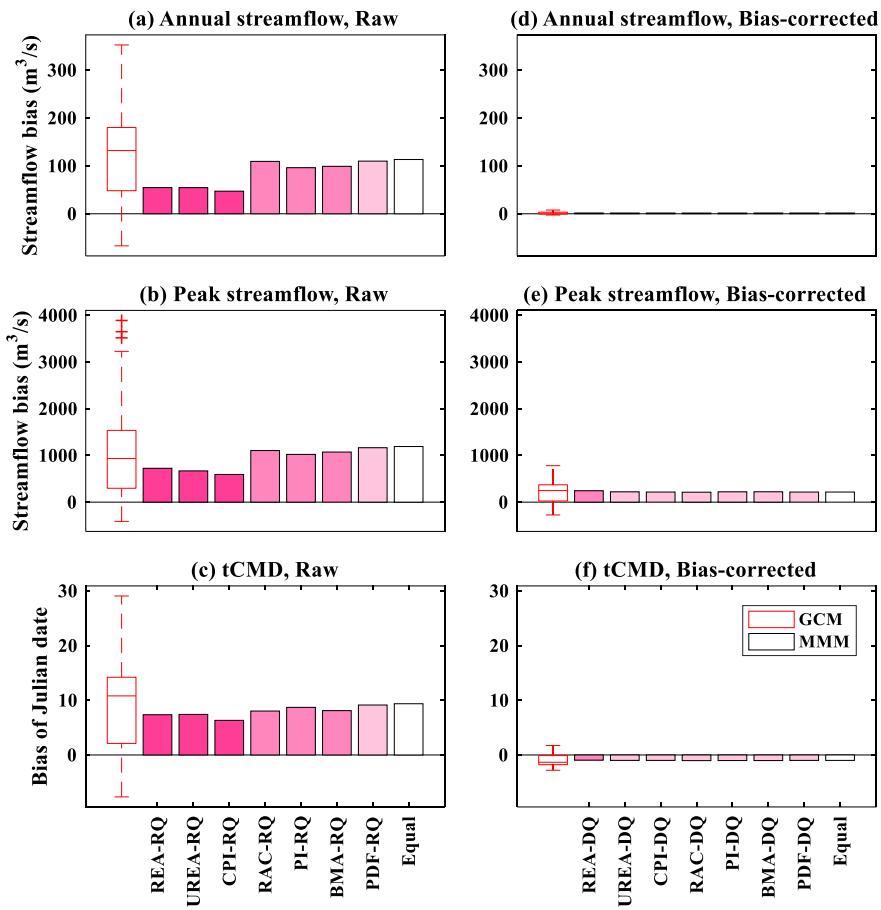


Figure 6. The same as Fig. 5 but for the Manicouagan-5 watershed.

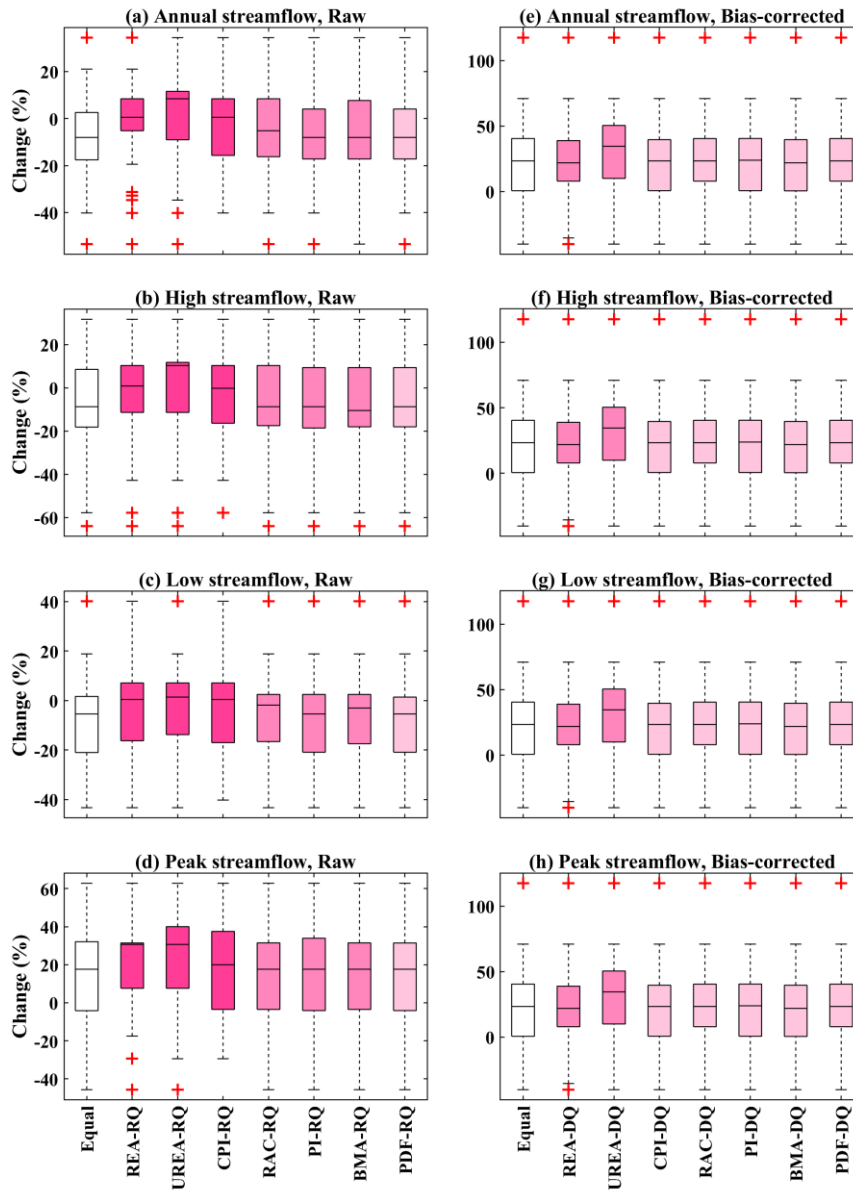


Figure 7. Box plot of changes in four hydrological indices calculated by raw or bias-corrected GCM-simulated streamflows in the Xiangjiang watershed. The changes of hydrological variables were sampled through the Monte-Carlo approach based on the weights calculated using raw (RQ) or bias-corrected (DQ) GCM-simulated streamflows. (The depth of pink represents the level of inequality of the weights.)

5

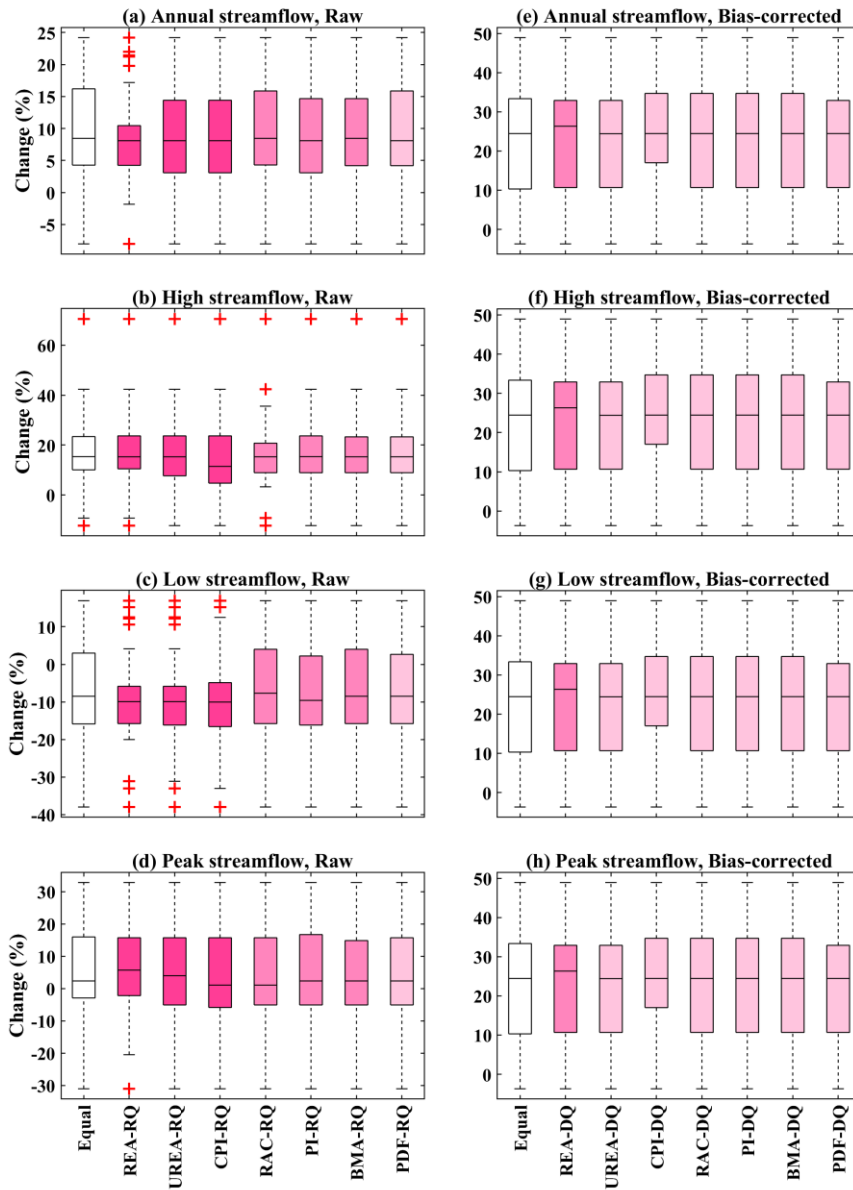


Figure 8. The same as Fig. 7 but for the Manicouagan-5 watershed.

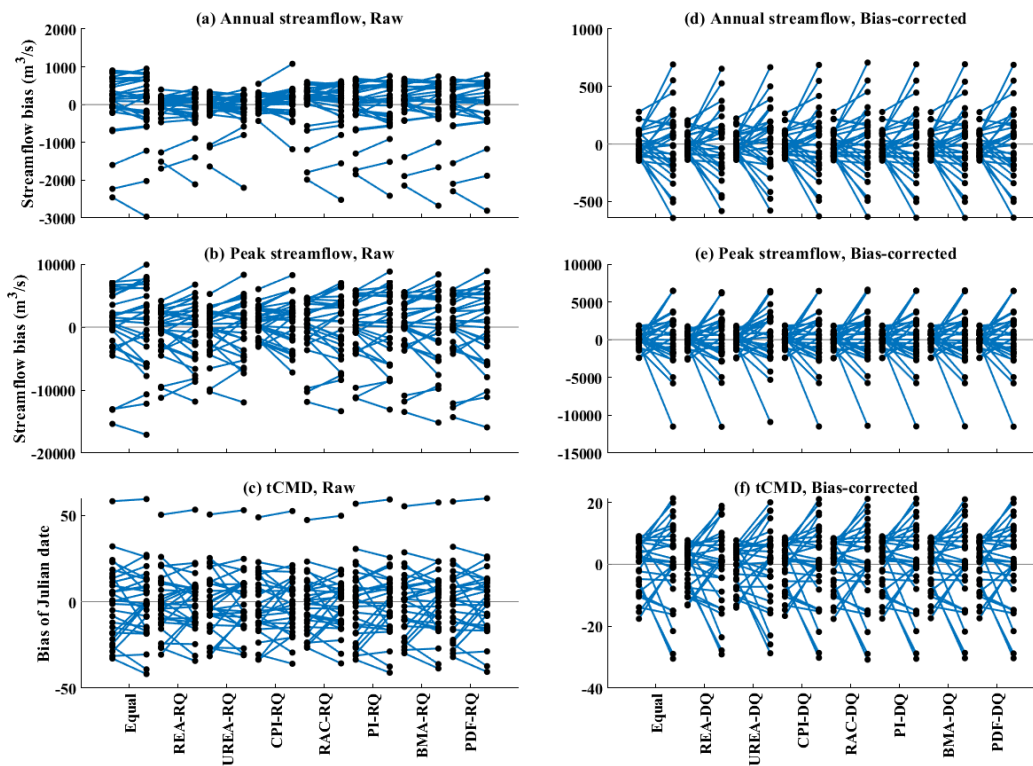


Figure 9. Bias in mean annual streamflow, mean peak streamflow and mean center of timing of annual flow (tCMD) of weighted multi-model ensemble mean in the out-of-sample testing over the Xiangjiang watershed. Twenty-nine sticks of each weighting method represent the results when each of 29 climate models was regarded as the “truth” in turn, and the left and right points in each stick represent the bias for the reference and future periods, respectively.

5

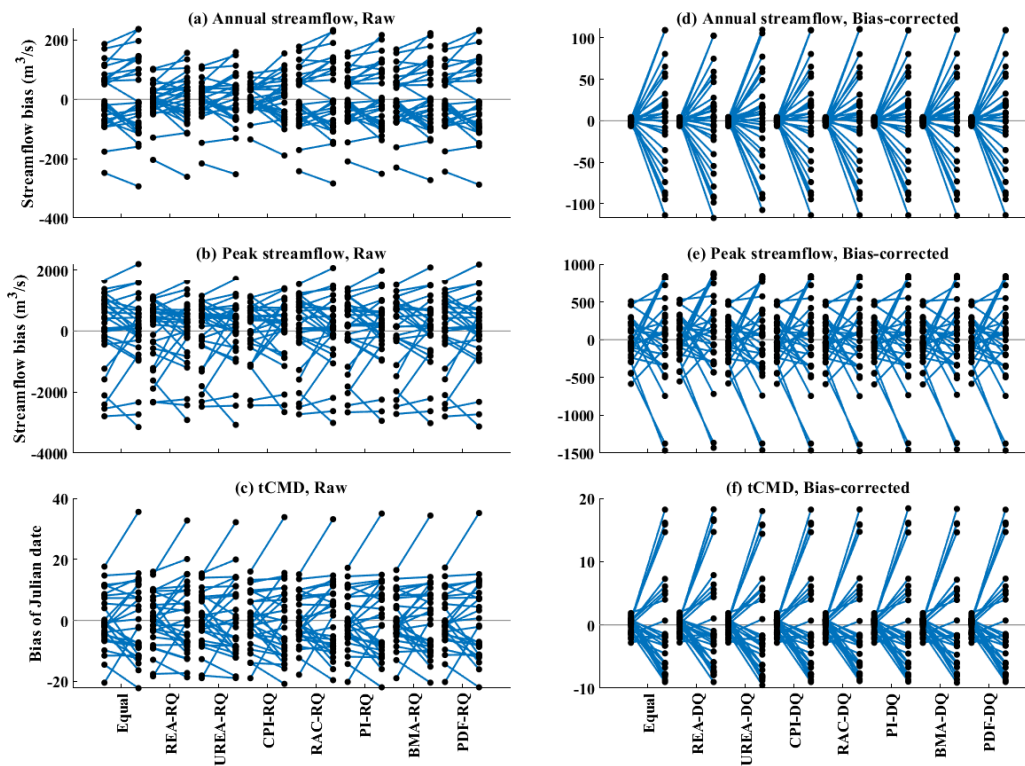


Figure 10. The same as Fig. 9 but for the Manicouagan-5 watershed.