

<<Referee 1>>

(1) The study evaluated different climate dataset source against climate stations using multiple indices and generated a synthetic dataset based on the ranks. Afterwards, the VIC model is applied as proxy validation tool to evaluate multiple datasets and generated datasets. The research is innovative and the structure of the paper is clear. Methods are valid. My only concern is about results. The performance of the VIC model in the study. It is not like what author stated “most of the climate datasets performed well”. On the contrary, in Christina and Firebag the NSE is below 0.45 for any datasets, and the worst is even below 0 which is in Pembina with NARR. The results of the model seemed unreliable. Please check the model and improve the performance of hydrological modeling.

((Reply)) We calibrated the parameters of the VIC model for the seven historical gridded climate datasets (i.e., ANUSPLIN, Alberta Township, PNWNAmet, CaPA, NARR, and two hybrid climate datasets) individually using an auto calibration method (dynamic dimensional search algorithm). Table 6 shows the Nash-Sutcliffe Efficiency (NSE) for the calibration and validation periods. Except for NARR, most of the NSE values during calibration period for Christina and Firebag are above 0.50 which is a threshold of satisfactory performance in hydrologic models as suggested by Moriasi et al. (2007). However, as indicated by the reviewer, model performance is not satisfactory for Christina and Firebag during the validation period. Accordingly, sentence has been revised in the manuscript (section 4.4). Figure 11 also shows box-whisker plots resulting from multiset-parameter hydrologic simulations that employed seven different model parameter sets (obtained through model calibration with individual climate datasets) and the same climate dataset as a forcing input data. In Figure 11, the averaged NSE values for Christina and Firebag were below 0.45 as pointed by the reviewer. However, these NSE values are different than the NSE values for calibration and validation shown in Table 6. The authors addressed more clearly how the biases in each climate dataset were estimated indirectly by the proxy validation as below.

“Under the assumption of REFRES that all of the existing climate datasets are of equal quality for hydrologic simulations, all of the calibrated parameter sets can be considered as mostly plausible parameter sets for the selected sub-basins. However, as mentioned above, intrinsic

biases exist temporally and spatially in all of the gridded climate datasets, e.g., discrepancies in the amount and spatial distribution of precipitation between the gridded climate datasets and observations. Therefore, the similarity of the gridded climate datasets in terms of magnitude, sequence, and spatial distribution of climate events relative to observations is crucial to reproduce historically observed streamflows. In addition to climate forcings, streamflows are mainly affected by geographic characteristics and physical land surface processes (e.g., infiltration and evapotranspiration), which are represented by model parametrization related to infiltration and soil properties (Demaria et al., 2007). In a hydrologic simulation, the biases in climate datasets can be compromised by model parameters that adjust hydrologic processes to observations (Harpold et al., 2017; Kirchner, 2006). That is, a calibrated parameter set may imply biases in a climate dataset. Under the assumption that the calibrated parameter sets are suitable for hydrologic simulations in each sub-basin, this study applied a multiset-parameter hydrologic simulation approach that employs all parameter sets calibrated by the seven climate datasets and the same climate dataset as a forcing input data to assess the sensitivity of the climate dataset to all feasible parameter sets. From the multiset-parameter hydrologic simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset. In other words, lower variability in the hydrologic simulations indicates higher reliability in the climate forcing dataset. The suitability of the hybrid climate dataset for improving historical hydrologic simulations was also tested by directly comparing the performances of calibration and validation for each climate dataset. Proxy validations were carried out by conducting 49 hydrologic simulations (7 climate forcing \times 7 parameter sets) for the Pembina and Christina catchment areas, whereas only 36 simulation runs were possible for Hinton, Firebag, and Clearwater sub-basins, as one of the gridded data sets (i.e., Township) did not cover the entire catchment areas of these three hydrometric stations.” (P16L11-P17L10)

(2) Section 2.1 What’s is the time duration of the climate observation data at AHCCD stations?

((Reply)) The AHCCD stations have different record lengths. For example, the longest record period is from 1840 to 2016 while the shortest period is from 1967 to 2004. As the data length are different at each AHCCD station, we selected a common period between AHCCD stations

and gridded climate datasets to estimate performance measure. The authors added this information in section 2.1.

(3) Method: Is the evaluation carried out on the whole time period and could be regarded as the average performance over the time? Is there any temporal variation of the performance for different observation dataset at different stations, and how do you consider the temporal variation of the performance?

((Reply)) The aim of REFRES is to choose a suitable climate dataset among the existing multiple historical gridded climate datasets based on the performance measures selected in this study. Each performance measure was evaluated over a whole common period at each AHCCD station. As the data lengths are different at each AHCCD station, it is not possible to consistently evaluate the temporal variation of the performance over the domain. In addition, consideration of temporal variation in performance may require a common period that covers a whole period of the hybrid climate dataset to be produced by choosing the most suitable climate dataset for a selected period. Therefore, this study only evaluated the performance averaged over a whole period to simplify the method and also to make sure that the methodology is computationally efficient.

(4) Section 3.1.3 It is not clear how the dataset is generated. Do you just choose the best one based on the evaluation over time or make a combination of several good ones?

((Reply)) Two things were considered in generating the hybrid climate data set: (i) the ranking of all datasets at each grid cell and (ii) a period of record or the availability of the gridded climate data sets. For each grid cell, the data were extracted by following the ranking (higher to lower) and data availability. For example, see the table below:

| Dataset | RANK | Period of record | Time period contributing to the hybrid climate dataset |
|----------|------|------------------|--|
| ANUSPLIN | 2 | 1950-2015 | 1950-1959 |
| Township | 1 | 1961-2016 | 1960-2016 |
| PNWNAmet | 3 | 1945-2012 | x |
| CaPA | 4 | 2002-2017 | x |

| | | | |
|------|---|-----------|---|
| NARR | 5 | 1979-2017 | x |
|------|---|-----------|---|

In the above table, the hybrid climate dataset should be a period from 1950 to 2016 which is covered by the existing climate datasets. Although Township is ranked first, Township cannot cover the period from 1950 to 1959. In this case, the data generation module in REFRES chooses the second ranked climate dataset, i.e., ANUISPLIN, to produce the hybrid climate dataset and the first ranked data for the remaining period from 1960 to 2016.

The authors addressed more clearly how the hybrid climate datasets are generated using the ranking information in DGM.

“As each climate dataset has different data periods shown in Table 1, the first ranked dataset cannot fully cover a whole target period to be extracted from a set of climate data candidates. DGM provides a systematic procedure to identify the most reliable dataset for a target region and extracts the data from the inventory of climate datasets considering the ranking and availability of each dataset for a desired period. For instance, if CaPA and ANUSPLIN ranked first and second for precipitation and the desired period is 1950 to 2016, DGM starts searching for the availability of precipitation in 1950. As CaPA is only available between 2002 to 2016, DGM reorders the rank to select ANUSPLIN as the best climate dataset available in 1950. In this way, a hybrid dataset over the period 1950 to 2016 is generated by extracting from ANUSPLIN from 1950 to 2001 and CaPA from 2002 to 2016 in this particular case.” (P14 L18-P15L2)

(5) 3.2 proxy validation “it is questionable if the hybrid climate dataset can represent a historical climate better than the individual gridded climate dataset. Utilizing a proxy validation approach (Klyszejko, 2007), this study applied streamflow records to confirm the superiority of the derived hybrid climate dataset over other existing climate datasets.” The underlying assumption is that the better input data could derive a more realistic streamflow simulation. The VIC model is calibrated against different dataset, so the calibration of parameters could offset the error from the input data. Judging the superiority through the output of a hydrological model is not straightforward and could even be misleading. How to consider the propagation of the error from the input through calibration?

((Reply)) The authors appreciate the valuable comment on the propagation of the error from the input climate data in hydrologic simulation. As the reviewer pointed, biases in climate data can be compromised or compensated by model calibration. This study indirectly estimated the impacts of the biases in climate datasets by a multiset-parameter hydrologic simulation

approach that employs all seven feasible parameter sets (obtained through calibration with the seven climate datasets separately) and seven climate dataset as a forcing data in the VIC model (i.e. 49 simulations; 7 climate forcing \times 7 parameter set). From the multiset-parameter hydrologic simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset. In other words, the lower variability in the hydrologic simulations indicates higher reliability in the climate forcing dataset as shown in Figure 11. This point has been clarified in the draft manuscript as follows:

“In a hydrologic simulation, the biases in climate datasets can be compromised by model parameters that adjust hydrologic processes to observations (Harpold et al., 2017; Kirchner, 2006). That is, a calibrated parameter set may imply biases in a climate dataset. Under the assumption that the calibrated parameter sets are suitable for hydrologic simulations in each sub-basin, this study applied a multiset-parameter hydrologic simulation approach that employs all parameter sets calibrated by the seven climate datasets and the same climate dataset as a forcing input data to assess the sensitivity of the climate dataset to all feasible parameter sets. From the multiset-parameter hydrologic simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset. In other words, lower variability in the hydrologic simulations indicates higher reliability in the climate forcing dataset.” (P16L20-P17L5)

(6) Could you specify what input you used here for the VIC model?

((Reply)) The VIC model requires several input data, i.e., climate forcing, soil, vegetation, and routing. This study used the same soil, vegetation, and routing input data as described in previous publications (Eum et al., 2014; 2017). The additional data sets used are the new climate forcing data (i.e. hybrid climate data generated in this study) comprised of daily precipitation, minimum temperature and maximum temperature climate variable.

(7) The number of Results should be 4.

((Reply)) Corrected.

(8) 3.1 Precipitation performance measures in Alberta, could you explain why ANUSPLIN and Township underestimate extreme precipitation?

((Reply)) The main reason that ANUSPLIN and Township underestimate extreme precipitation is that they employed raw station data instead of the adjusted precipitation data which are higher than the raw station data by 5 % to 20%. The authors addressed this as below,

“Interestingly, two station-based gridded climate datasets, ANUSPLIN and Township, show negative P_{bias} while PNWNAmet, CaPA, and NARR datasets have positive P_{bias} . This indicates that ANUSPLIN and Township may underestimate extreme precipitation, as they employed the raw station data instead of the adjusted precipitation data which is higher than the raw station data by 5%-20%. In contrast, other climate datasets (especially multiple sources and reanalysis data) overestimate extreme precipitation.” (P17L20-L25)

(9) Figure 10 is it a maximum, minimum or mean temperature in this figure?

((Reply)) The ranking was determined based on the performance of precipitation and temperature (minimum and maximum) individually by TOPSIS. The performance measures for both minimum and maximum temperature were employed into TOPSIS and the ranks were presented in Figure 10 (b). Figure 10 (c) showed the ranking when the performance measures for all variables were considered in TOPSIS. Please also see the following clarification text in the manuscript:

“To alleviate the erroneous output that minimum temperature is higher than maximum temperature on a certain day when producing the hybrid climate dataset by the ranking of temperature values individually, the performance measures of both minimum and maximum temperature are employed together to rank the climate datasets for temperature.” (P14L5-L8)

(10) Page 15 line 24-26 “Over the five hydrometric stations, most of the climate datasets performed well with the exception of NARR in the Pembina catchment.” Please explain why NARR in Pembina performs bad which only got -0.85 for NSE. The criterial of well or not well is quite subjective. In Hinton the model performance could be acceptable. However, in Christina and Firebag the NSE is even below 0.45 for any cases and In Pembina and Clearwater NSE below 0.7. This is not a behavioral model honestly. Is the model suitable for the river basin? If it is suitable why the NSE is low? I suggest to check the calibration of the model. Otherwise the proxy validation is not reliable.

((Reply)) In case of Pembina watershed with NAAR data set: The NSE value for calibration period (1985 to 1997) is 0.5 while it is -0.14 for the validation period (1998 -2016). There are some reasons of such a poor performance of NARR in most of the watersheds including Pembina. Since 2003, assimilation of observed precipitation data in to NARR has been discontinued and consequently, NARR overestimates precipitation (refer to section 4.1) and has

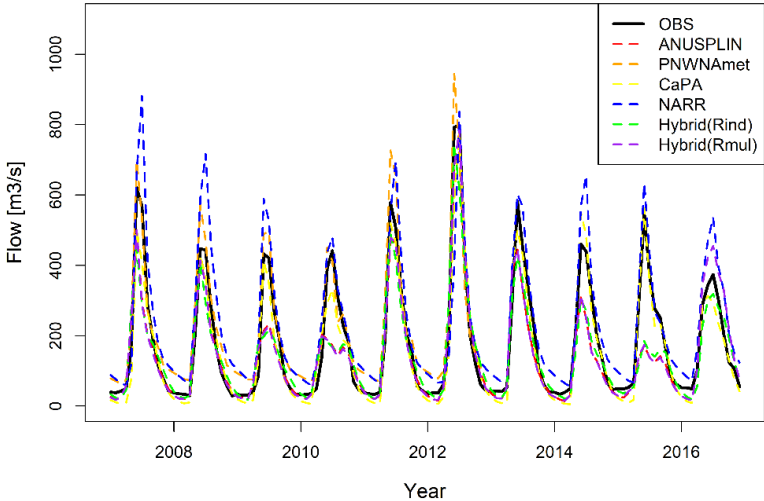
warm and cold biases in temperature (refer to section 4.2). In addition, Pembina has been recognized as a parameter-sensitive basin in Eum et al. (2014b)'s study, implying that selection of a calibration period is critical for the performance of hydrologic simulations in this watershed. These biases in NARR and the hydrologic characteristics of the basin may induce poor performance in the hydrologic simulation during the validation period in Pembina. A qualitative rating has been suggested by Moriasi et al. (2007) as shown in the table below.

| | | | |
|---------------------------|---------------------------|------------------------|-----------------|
| Very Good | Good | Satisfactory | Unsatisfactory |
| $0.75 \leq NSE \leq 1.00$ | $0.65 \leq NSE \leq 0.75$ | $0.50 < NSE \leq 0.65$ | $NSE \leq 0.50$ |

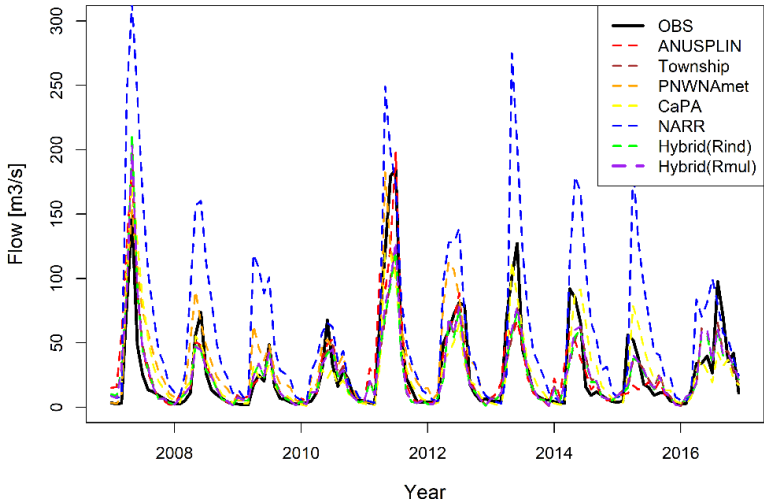
The goodness-of-fit statistics table shows modelling is satisfactory when $NSE > 0.5$. Table 6 presents Nash-Sutcliffe Efficiency (NSE) for the calibration and validation periods at the selected hydrometric stations (Hinton, Pembina, Christina, Clearwater, and Firebag) in the ARB to assess the suitability of each climate datasets as a climate forcing for hydrologic simulations. Over the five hydrometric stations, most of the climate datasets performed well with the exception of NARR in the Pembina catchment. That is, most of the NSE values in calibration for Christina and Firebag were above 0.50 which is a threshold of satisfactory performance in hydrologic models as suggested by Moriasi et al. (2007).

(11) Figure 12 is suggested to be refined it is hard to tell the difference between different experiments. Is it m³/s in the label of Y axis? There is lack of label of X axis.

((Reply)) The authors have modified Figure 12 from daily to monthly hydrograph and added another hydrographs for Pembina and x-axis has been labeled to improve visualization.



(a) Hinton



(a) Pembina

Figure 12. Monthly observed and simulated hydrographs from the gridded climate datasets at (a) Hinton and (b) Pembina

<< Referee 2 >>

<General Comments>

(1) Performance of multiple climate datasets against the ground stations.

It seems to me that the performance of the climate datasets could be affected by the interpolation method used to estimate the values at the AHCCD stations. The authors used the inverse distance squared weighting method to obtain the estimated values from all the gridded products (P8L4-5), and the Township data was shown to outperform other climate datasets for all performance measures except Pbias. I am struggling to square away in my mind that the interpolation method might favour towards the Township data because the Township data also employed inverse distance weighting and used the same (or similar) set of ECCC stations to generate the data. Thus, the Township data would most likely rank first among the climate datasets because the major deficiency of the data lies from the difference between the raw station data it used and the adjusted data in AHCCD, while the deficiencies of other climate datasets come from interpolation method, numbers of stations used, and the errors arising from the use of additional information/numerical models.

((Reply)) The authors appreciate the reviewer's valuable comments. This study investigated the performance of the five gridded climate datasets at the AHCCD stations. Among the gridded climate datasets, station-based datasets (i.e., ANUSPLN and Alberta Township) employed different numbers of observed (raw) station data depending on data availability in a given year except for PNWNAmet that set a common period from 1945 to 2012 for all stations included in the interpolation. While ANUSPLIN used the Canada-wide archive (raw) station data collected only by ECCC, the Alberta Township data has been produced on the basis of the archive (raw) station data collected by ECCC and other agencies including Alberta Environment and Parks (AEP), and Alberta Agriculture and Forestry (AF) over Alberta. Therefore, one of the possible reason for outperformance of Township dataset might be the difference in the numbers of stations (i.e. station density) employed to produce the gridded climate datasets. This point has been added to the discussion section of the manuscript, as follows:

“Among the station-based gridded climate datasets, the Township dataset outperformed other station-based gridded climate datasets. As PNWNAmet set a common period from 1945 to 2012 for all stations included in the interpolation, many stations might be left out in the data

generation processes. While ANUSPLIN used the Canada-wide archive (raw) station data collected by only ECCC, the Alberta Township data has been produced on the basis of the archive (raw) station data collected by ECCC, AEP, and AF over Alberta. Therefore, one of the possible reason for outperformance of Township dataset might be the difference in the numbers of stations (i.e. station density) employed to produce the gridded climate datasets.” (P23L2-P23L9)

(2) Superior performance of hybrid dataset over multiple existing climate datasets

I am a bit skeptical about the claim that the performance of hybrid datasets was ‘superior’ when compared to other five climate datasets (P1L30-31). By saying ‘superior’ the results should be far better than the others (e.g. a NSE value of 0.8 as compared to 0.5). In this study, I would argue that the overall performance of hybrid datasets was only marginally better than some of the existing climate datasets in most of the sub-basins. The performance of hybrid dataset, Hybrid(Rind), was even worse than ANUSPLIN at Hinton station (Figure 11). Overall, the hybrid datasets only provided comparably good NSE values as the other climate datasets.

((Reply)) The authors agree with the reviewer’s comment and agreed that ‘superior’ word may not be suitable in this context. In Table 6, the two hybrid climate datasets performed well with comparably better NSE values than other climate datasets, especially at Pembina, Clearwater, and Firebag located in the middle and lower reaches. From multiset-parameter hydrologic simulations shown in Figure 11, however, the hybrid climate datasets provided higher precision and accuracy in most of the stations except for Hinton as the reviewer pointed out. Therefore, the authors replaced the word “superior” to “utility” in the modified manuscript.

(3) Creditability of hybrid dataset in improving hydrologic simulations

(3-1) Even though the hybrid datasets provided comparably good NSE values as the other climate datasets or even higher NSE values, when examining the hydrograph in Figure 12, one can find that there are four obvious large underestimation of the peaks in 2009, 2010, 2014, and 2015 simulated by using the hybrid datasets (purple lines and potentially green lines as well). Could the authors explain what happened at Hinton station? Could the authors also show the hydrographs at other stations to see whether similar situations happened in other sub-basins?

((Reply)) The authors appreciate the reviewer’s valuable comment. The two hybrid climate datasets were produced by combining with the existing gridded climate datasets based on the performance measures. Therefore, it has an intrinsic limitation that the performance of the

hybrid dataset for a basin may resemble that of a climate dataset that is dominantly ranked first for the basin. As commented in (3-2) below, ANUSPLIN was dominantly ranked first for Hinton, consequently the hydrographs of ANUSPLIN and the hybrid datasets were similar to each other as shown in the figure below. In addition, the authors present a monthly hydrograph for Pembina where the Township data was dominantly ranked first for this basin. The hydrograph of the two hybrid climate datasets (green and purple dashed lines) are highly similar to that of Township (brown dashed line). The authors addressed the limitation in the discussion section.

“In Figure 12, the hybrid climate datasets underestimated the peak flows (in 2009, 2010, 2014, and 2015) at Hinton, and hydrograph is similar to the hydrograph produced by ANUSPLIN dataset that dominantly ranked first in this watershed. On the contrary, the hydrograph of the hybrid climate datasets at Pembina resembles that of Township that is dominantly ranked first in Pembina (refer to Table 5). These results indicate that the hybrid climate dataset has the intrinsic limitation that the performance of the hybrid dataset for a basin may closely resemble that of the climate dataset that is dominantly ranked first for the basin. However, the utility of the hybrid climate dataset can be clearly found at a whole-basin scale for a large watershed, as the added values of the hybrid climate dataset in sub-basins can be cumulated to the main stem at the downstream in the watershed.” (P23L18-P24L2)

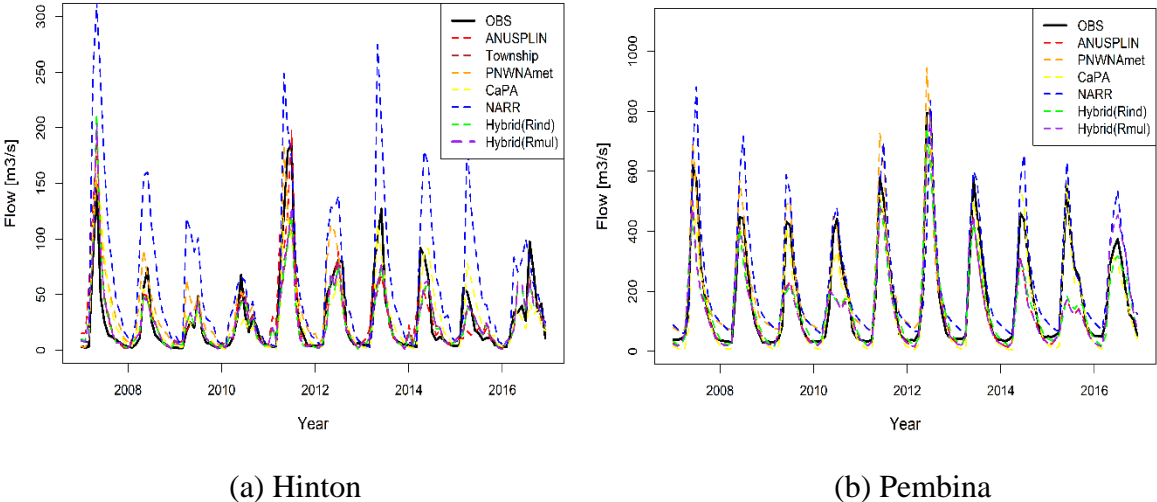


Figure 12. Monthly observed and simulated hydrographs from the gridded climate datasets at (a) Hinton and (b) Pembina

(3-2) The claim that the two hybrid datasets performed better in terms of accuracy and precision in the proxy validation (P18L28-29) could be a bit misleading. In this study, it was coincidentally that the hybrid datasets (either based on single or multiple variables) were dominantly generated from one particular climate dataset in all sub-basins (except Clearwater when using precipitation as the variable). If the authors show the breakdown of the first ranked number of grid cells for each climate dataset in each sub-basin (just like in Table 5), I would guess that over 90% of the grid cells at Hinton came from ANUSPLIN when considering the performance measures of multiple variables (Figure 9c) and almost 99% of grid cells at Pembina came from the Township data. In this regard, I would argue that the performance of the hybrid datasets shown in Figure 11 was highly resemble to the performance of the climate dataset that was dominantly generated from. I would also argue that the optimal parameter sets of the hybrid datasets would be the same (or very similar) as that of dominant climate dataset. Have the authors checked the optimal parameter sets of the hybrid datasets and the five climate datasets? Will the calibrated parameter sets of the hybrid dataset (Hybrid(Rmul)) the same as the parameter sets of Township data at Pembina, for instance? The creditability of generating a hybrid dataset might not be fully assessed at sub-basin scale, especially when the hybrid datasets were generated mainly from one particular climate dataset. I think a better assessment to reveal the usefulness of the hybrid datasets was to calibrate the model at whole-basin scale for this particular basin (e.g. calibrating at Fort McMurray using 07DA001 station). In this case, the hybrid dataset is better mixed by different climate datasets for different parts of the whole basin, thus reducing the chance of one particular climate dataset being dominant in the data generation process.

((Reply)) The authors appreciate the reviewer's excellent comment. As mentioned in (3-1) above, the performance of the hybrid climate dataset is similar to that of an existing climate dataset which is dominantly ranked first for a sub-basin, and the utility of the hybrid climate dataset can be clearly demonstrated when it is applied for simulations at the whole basin scale. However, this study confirmed that the hybrid climate dataset provides a better representation of historical climatic conditions as different watersheds have different dominant gridded climate data and the proposed methodology helps to identify the appropriate dominant climate data in the derived hybrid dataset. Further, as suggested by the reviewer, we calibrated the VIC model for larger watersheds (i.e. Fort McMurray and Eymundson) to provide additional simulation results. The table below shows the NSE values calculated for ANUSPLIN and Hybrid (R_{ind}) at a few hydrometric stations in the main stream of the Athabasca River. The result shows that as

the size of watershed increases, hybrid climate dataset start performing better than the existing gridded climate dataset (in this case ANUSPLIN). This is mainly due to the fact that as the watershed area increases, the derived hybrid climate dataset is no longer dominated by a single gridded dataset. Due to the limitation of computational capacity, initially only five sub-basins were selected for proxy validation.

Nash-Sutcliffe Efficiency (NSE) of ANUSPLIN and the hybrid climate datasets at the main stream of the Athabasca River

| No | Station name/ID | Drainage area (km ²) | ANUSPLIN | | Hybrid | |
|----|-----------------------------|-------------------------------------|-------------|------------|-------------|------------|
| | | | Calibration | Validation | Calibration | Validation |
| 1 | Hinton / 07AD002 | 9,760 | 0.85 | 0.82 | 0.83 | 0.76 |
| 2 | Windfall / 07AE001 | 19,600 | 0.80 | 0.72 | 0.80 | 0.76 |
| 3 | Athabasca / 07BE001 | 74,600 | 0.78 | 0.69 | 0.77 | 0.78 |
| 4 | Fort McMurray / M07DA001 | 133,000 | 0.77 | 0.66 | 0.78 | 0.75 |
| 5 | Eymundson / S24 | 147,086 | 0.77 | 0.67 | 0.79 | 0.75 |

<Specific Comments>

(1) P8L4: How many grid points were used in the inverse distance squared weighting?

((Reply)) Four points were used for the inverse distance squared weighting method.

(2) P8L5-6: The AHCCD stations have different starting and ending points and percentage of missing values. How did the authors take care of these? Did the authors calculate the performance measures using a common period?

((Reply)) Yes, as the data lengths are different at each AHCCD station, we selected a common period between each AHCCD station and climate datasets, and neglected missing values to estimate performance measures (P6L22-24).

(3) P8L21-24: please also define i

((Reply)) Yes, we have defined i in the modified manuscript, as follows:

“ G_i and O_i represent gridded and observed climate datasets at i^{th} time step, respectively”
(P11L16-L17)

(4) P9L5: The authors mentioned 20% of all AHCCD stations were selected here but five nearest AHCCD neighbours were shown in Figure 2. Which one is correct?

((Reply)) There are two steps to select the nearest neighbors in RM. Firstly, 20% (of all AHCCD) stations are selected based on the nearest distance criteria. Then, the five nearest stations from them is finally selected by the minimum elevation difference criteria. Accordingly, Figure 2 has been modified in the revised manuscript.

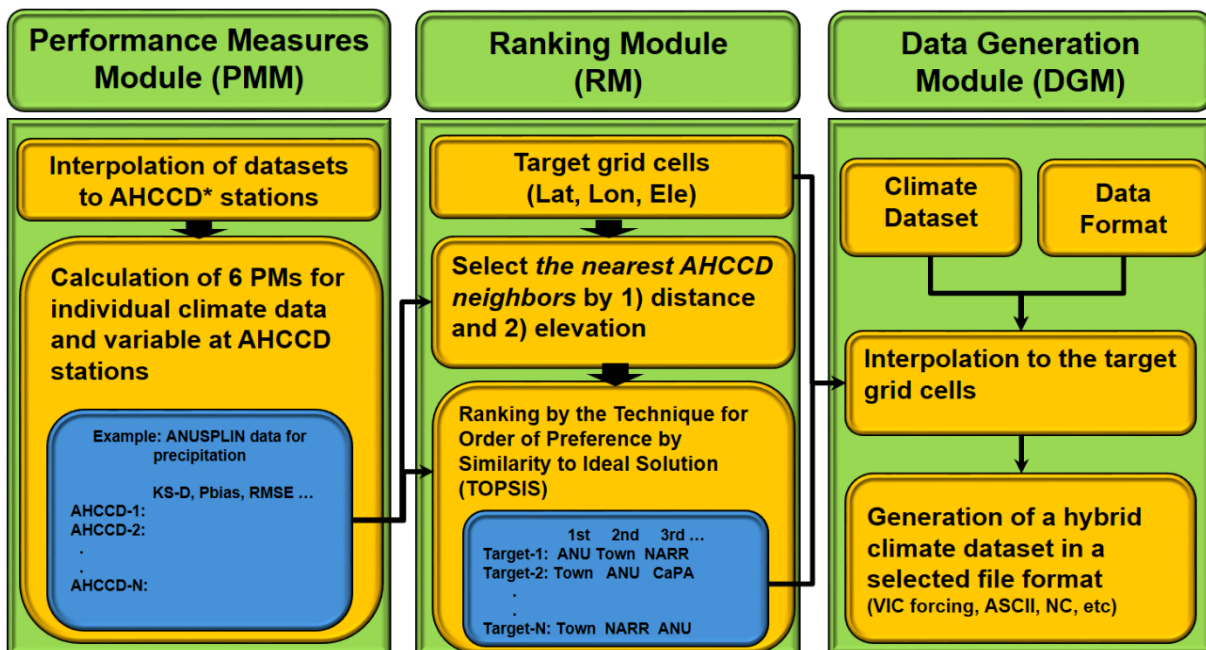


Figure 2. Structure of REFRES comprised of three modules; 1) Performance Measure Module (PMM), 2) Ranking Module (RM), and 3) Data Generation Module (DGM)

(5) P11L27-29: What did the authors mean by “the number of gridded climate datasets was optimized”? Please elaborate.

((Reply)) It has been modified as below,

“In other words, a higher number of gridded climate datasets contributing to the hybrid climate dataset within a catchment was selected to evaluate the utility of the hybrid climate data relative to the existing gridded climate datasets.” (P15L22-L24)

(6) P12L3: Why were only two hybrid datasets from the Rind and Rmul? Didn't the authors rank for precipitation and temperature separately (Rind)? (P10L12-13) I think there would be two sets of hybrid datasets based on Rind, one for precipitation only and one for temperature only, as shown in Figures 9 and 10.

((Reply)) In this study, a climate dataset consists of three variables, i.e., daily precipitation, minimum temperature, and maximum temperature. Considering the ranks from R_{ind} and R_{mul} , that is, two hybrid climate datasets was produced to be used in the proxy validation as a forcing data of the VIC model.

(7) P12L5: I assume that in this study the authors used the same version and the same VIC setup as described in Eum et al. (2017). Could the authors clarify the sources of the other meteorological variables (e.g. wind speed) required in the VIC model? Did the authors use the meteorological variables from NARR for all the climate datasets and the hybrid datasets? Did the authors use the wind speed data of the Township data itself, for instance?

((Reply)) This study used VIC version 4.2.d that has the MT-CLIM package to estimate required climate variables in VIC. Hydrologic simulations were forced by only the three daily climate variables (i.e., precipitation, minimum temperature, and maximum temperature) for the proxy validation and other climate variables including wind speed were estimated by the MT-CLIM package in VIC. Next stage of this study is to expand the number of climate variables, such as wind speed, solar radiation, etc, for further improving hydrologic simulations.

(8) P12L21: What were the calibration and validation periods in this study?

((Reply)) The calibration and validation periods were added to the modified manuscript:

“The calibration period is 1985-1997 as in Eum et al., (2017), except for CaPA that uses the period of 2003-2009 for calibration, as CaPA covers the period from 2002 to 2016. The remaining period of total record length for each climate dataset is used for validation” (P16L7-L10)

(9) P13L3-7: Table 3 shows the ‘average’ performance of each climate datasets. How did the results indicate under- or over-estimation of ‘extreme’ precipitation? Please explain.

((Reply)) The authors addressed the impacts of biases in precipitation (resulting in under or over estimation of extreme precipitation) in the discussion section of the manuscript, as follows:

“Among the station-based gridded climate datasets, the Township dataset outperformed other station-based gridded climate datasets. As PNWNAmet set a common period from 1945 to 2012 for all stations included in the interpolation, many stations might be left out in the data generation processes. While ANUSPLIN used the Canada-wide archive (raw) station data collected by only ECCC, the Alberta Township data has been produced on the basis of the archive (raw) station data collected by ECCC, AEP, and AF over Alberta. Therefore, one of the possible reason for outperformance of Township dataset might be the difference in the numbers of stations (i.e. station density) employed to produce the gridded climate datasets. In addition, PNWNAmet showed a positive P_{bias} for precipitation, especially in the mountainous areas, while ANUSPLIN, which employs similar thin plate spline interpolation, generated negative P_{bias} . PNWNAmet overestimated precipitation over the mountainous area, which considerably affects simulated low flows at Hinton in the ARB. Figure 12 shows the observed and simulated hydrographs from gridded climate datasets at (a) Hinton and (b) Pembina. It clearly shows that PNWNAmet highly overestimated the low and high, which is caused by overestimated precipitation in the drainage area of the sub-basins. As with PNWNAmet, NARR also overestimated the low and high flows, which is induced by the combined effects of overestimating precipitation and warm biases in cold temperature. The temperature bias of NARR is thus further confirmed and is consistent with the earlier finding of Eum et al., (2014) and Islam and Dery (2016).

In Figure 12, the hybrid climate datasets underestimated the peak flows (in 2009, 2010, 2014, and 2015) at Hinton, and hydrograph is similar to the hydrograph produced by ANUSPLIN dataset that dominantly ranked first in this watershed. On the contrary, the hydrograph of the hybrid climate datasets at Pembina is similar to that of Township that is dominantly ranked first in Pembina (refer to Table 5). These results indicate that the hybrid climate dataset has the intrinsic limitation that the performance of the hybrid dataset for a basin may closely resemble that of the climate dataset that is dominantly ranked first for the basin. However, the utility of the hybrid climate dataset can be clearly found at a whole-basin scale for a large watershed, as the added values of the hybrid climate dataset in sub-basins can be cumulated to the main stem at the downstream in the watershed” (P23L2-P24L2)

(10) P13L25: Should it be >800 mm/year?

((Reply)) The authors addressed this clearly as below.

“(e.g., 300 mm/year higher than the observation at the station ID 3050519)”

(11) P14L16-19: It would be better to show the breakdown of the first-ranked number of grid cells and their percentages for each sub-basin as well because the authors calibrated and validated the VIC model at sub-basin scale.

((Reply)) The authors modified Table 5 to add the information on the first ranked climate datasets for the five sub-basins and the whole Athabasca River basin.

Table 5. First ranked number of grid cells in the five sub-basins and the whole Athabasca River Basin (ARB) and their percentage for each climate dataset considering the performance measures of individual (Case A and Case B) and multi-variables (Case C, i.e., precipitation and temperature in this study). Total number of grid cells is 22,372 at 1/32° (2~3 km)

| Criteria | Basin | Climate dataset | | | | |
|--|------------|-----------------|----------------|---------------|--------------|-------------|
| | | ANUSPLIN | Township | PNWNAmet | NARR | CaPA |
| (A) Precipitation | ARB | 2985 (13%) | 17515 (78%) | 691 (3%) | 499 (2%) | 682 (3%) |
| | Hinton | 1271 (91%) | 126 (9%) | 0 (0%) | 0 (0%) | 0 (0%) |
| | Pembina | 0 (0%) | 1791 (100%) | 0 (0%) | 0 (0%) | 0 (0%) |
| | Christina | 0 (0%) | 658 (99.5%) | 3 (0.5%) | 0 (0%) | 0 (0%) |
| | Clearwater | 1474 (56%) | 252 (9.6%) | 10 (0.4%) | 682 (26%) | 215 (8%) |
| | Firebag | 129 (14%) | 750 (79%) | 9 (1%) | 0 (0%) | 64 (6%) |
| (B) Temperature (Min & Max Temp.) | ARB | 13809 (62%) | 6924 (31%) | 1639 (7%) | 0 (0%) | - |
| | Hinton | 63 (5%) | 77 (6%) | 1257 (89%) | 0 (0%) | - |
| | Pembina | 486 (27%) | 1305 (73%) | 0 (0%) | 0 (0%) | - |
| | Christina | 492 (74%) | 169 (26%) | 0 (0%) | 0 (0%) | - |
| | Clearwater | 2593 (98%) | 40 (2%) | 0 (0%) | 0 (0%) | - |

| | | | | | | |
|----------------------------|------------|---------------|----------------|-----------|-----------|---|
| (C) Multi- variables | Firebag | 924 (97%) | 28 (3%) | 0 (0%) | 0 (0%) | - |
| | ARB | 8049 (36%) | 14323 (64%) | 0 (0%) | 0 (0%) | - |
| | Hinton | 1271 (91%) | 126 (9%) | 0 (0%) | 0 (0%) | - |
| | Pembina | 0 (0%) | 1791 (100%) | 0 (0%) | 0 (0%) | - |
| | Christina | 109 (16%) | 552 (84%) | 0 (0%) | 0 (0%) | - |
| | Clearwater | 2574 (98%) | 59 (2%) | 0 (0%) | 0 (0%) | - |
| | Firebag | 536 (56%) | 416 (44%) | 0 (0%) | 0 (0%) | - |

(12) P15L12: Again, I think there should be three different hybrid datasets.

((Reply)) Based on the response mentioned in (6), I believe the reviewer fully understands the definition of a climate dataset.

(13) P15L19: Same as the above comment. If only two hybrid datasets were implemented, could the authors clarify which Rind was used?

((Reply)) Please refer to the response provided for (6) and (12).

(14) P15L20-22: It was shown that NARR did not perform well in temperature (Section 3.2). Why did the authors still combine CaPA precipitation with NARR temperature for the proxy validation? Would such combination be unfair to CaPA performance? The performance of CaPA should be assessed by combining with the temperature data of all other climate datasets.

((Reply)) As both CaPA and NARR datasets are produced from climate model-based outputs, authors thought that it will be more logical to supplement the CaPA precipitation data with temperature data from another similar type of dataset (i.e., NAAR). The performance evaluation of CaPA data when supplemented with different temperature data is beyond the scope of this study.

(15) P16L4-9: What was the validation period for other climate datasets? For better comparison with CaPA, I think the authors could show the NSE results calculated from 2010 to 2016 for all the climate datasets.

((Reply)) Please refer to the reply of (8) and P21L25-P22L5.

“The validation period of CaPA is only six years from 2010 to 2016, as CaPA data are only available between 2002 to 2016. This might be a reason why CaPA produced the highest NSE (accuracy) among the climate datasets used in this study. Therefore, the results of CaPA need to be considered carefully otherwise they might be misleading. In this context, the CaPA dataset was excluded from further assessment of the precision and accuracy even though all of the results of CaPA were included in Figure 11 for reference only.” (P22L6-L11)

(16) P16L12: The VIC performance using NARR did not get positive NSE even after calibration. This means that no optimal parameter sets could be identified using NARR and the parameter sets could be anywhere in the parameter space. I wonder how such unidentified parameter sets could still produce fair NSE values when it was used with other climate datasets (Figure 11). I would expect a long lower whisker (just like the case in CaPA). Otherwise, I would think that the errors from the climate dataset were greatly compensated by the parameter uncertainties during the calibration. Could the authors explain what happened at Pembina?

((Reply))

The reviewer 1 has raised the same issue on the results in the performance of NARR in Pembina. In the case of Pembina watershed with NAAR dataset, the NSE value for the calibration period (1985 to 1997) is 0.5 while it is -0.14 for the validation period (1998 -2016). There are some reasons for such a poor performance of NARR in most of the watersheds including Pembina. Since 2003, assimilation of observed precipitation data in to NARR has been discontinued and consequently, NARR overestimates precipitation (refer to section 4.1) and has warm and cold biases in temperature (refer to section 4.2), resulting in highly overestimating flows (refer to Figure 12). In addition, Pembina has been recognized as a parameter-sensitive basin in Eum et al. (2014b)'s study, implying that selection of a calibration period is critical for the performance of hydrologic simulations in this watershed. These biases in NARR and the hydrologic characteristics of the basin may induce poor performance in the hydrologic simulation during the validation period in Pembina. As the reviewer commented, the NARR parameter set produced fair NSEs in simulations forced by the other climate datasets except for CaPA and PNWNAmet. Such result indicates that 1) all of parameter sets used in this study were calibrated reasonably and 2) climate forcing input data plays a more crucial role in hydrolog simulations as any parameter sets did not produce a fair NSE value from NARR in Pembina. The authors

addressed the impacts of NARR on hydrologic simulations in the discussion section of the manuscript, as follows:

“Literature has demonstrated that NARR, a reanalysis-based climate dataset, can be an alternative as a climate forcing dataset for hydrologic simulations in data sparse regions (Choi et al., 2009; Praskievicz and Bartlein, 2014; Islam and Dery, 2016). In this study, the NARR dataset performed quite well in high-elevation regions (Hinton in this study) while it did not perform so well in the middle and lower reaches, i.e., lower-elevation watersheds. NARR performed especially poorly in the Pembina sub-basin, a region where hydrologic simulations are highly sensitive to model parameters (Eum et al., 2014b). In Figure 11 (b), however, the NARR parameter set produced fair NSE values in hydrologic simulations forced by the other climate datasets except for CaPA and PNWNAmet. Such result indicates that 1) all of parameter sets used in this study were calibrated reasonably and 2) climate forcing input data plays a more crucial role in hydrologic simulations as any parameter sets did not produce a fair NSE value from NARR in Pembina.” (P24L19-P25L3)

<Remarks>

(1) P2L20: should be “may not produce” not “may not produces”

((Reply)) Corrected

(2) P4L4: should be “the aims of this study are” not “the aims of this study is”

((Reply)) Corrected

(3) P4L32: should be “Peace River” not “Peasce River”

((Reply)) Corrected

(4) P9L5: should be “criteria” not “citeria”

((Reply)) Corrected

(5) P19L19-21: please update the reference. Christensen and Lettenmaier (2007) has been published in HESS already, not HESSD.

((Reply)) Corrected

(6) P20L16-18: missing the name of journal

((Reply)) Corrected

(7) P20L19: should be “Dibike, Y.” not “Yonas, D.”

((Reply)) Corrected

(8) Table 6: should there be two hybrid datasets of Rind?

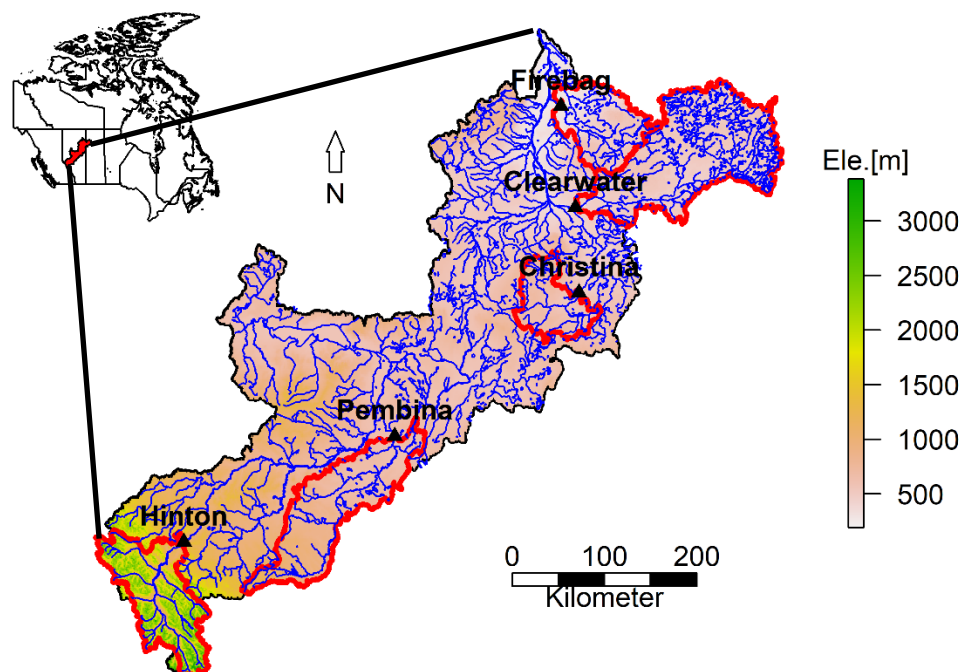
((Reply)) Based on the reply above, I believe the reviewer fully understands how the hydrologic simulations were conducted with two hybrid climate datasets (i.e., R_{ind} and R_{mul}).

(9) Figure 1: should be “precipitation” not” precipitation”

((Reply)) Corrected

(10) Figure 3: this figure could be combined with Figure 8 to reduce the numbers of figures (or the other way round). Otherwise, the authors should provide the geographical information about the basin on the map to facilitate the understanding of the international readers (e.g. elevation, latitude and longitude, a mini map showing the geographical location of the basin in Canada). Also, it would be better to show the river network of the basin.

((Reply)) The authors modified Figure 3 to provide the geographical information of the ARB as the reviewer suggested.



(11) Figure 9: there are too much unnecessary white space between the labels, the figures, and the legend. Consider squeezing the white space to make the figure more compact.

((Reply)) Corrected

(12) Figure 11: should there be two hybrid datasets of Hybrid(Rind)?

((Reply)) Again, I believe the reviewer fully understands how the hydrologic simulations were conducted with two hybrid climate datasets.

<<Short comment from David Thompson >>

<General Comments>

The study evaluates five climate datasets; ANUSPLIN, Alberta Township, PNWNAmet, CaPA, and NARR. The method can be divided in three major parts: (a) comparing climate datasets (identified in the method section of the manuscript as “Performance Measure Module”), (b) ranking the gridded datasets based on their performance measures (identified in the method section of the manuscript as “Ranking Module”), and (c) further evaluating climate datasets and their ranking using the VIC hydrological model (identified in the method section of the manuscript as “Proxy validation”).

Each part of the method section raises concerns as follows:

Part 1:

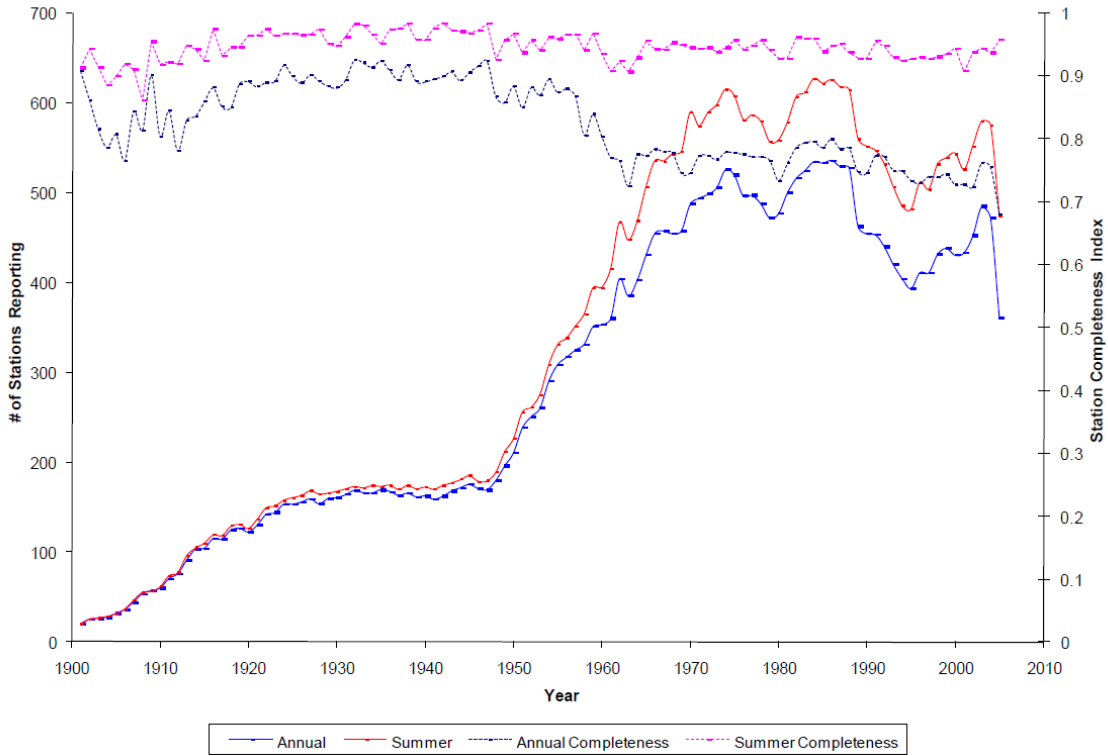
In the first part of the methodology, five climate datasets were compared. Three of them (ANUSPLIN, Alberta Township, and PNWNAmet) are climate datasets which are originally generated based on interpolation, and the other two (CaPA and NARR) are generated based on models and satellite technologies. The accuracy of all the datasets is compared to the (observed stations) Adjusted and Homogenized Canadian Climate Data (AHCCD). The main concern is how the authors did this comparison? The study states that “the inverse distance squared weighting method was applied to obtain the values at the AHCCD stations from all the gridded climate datasets. Then, performance measures were calculated by comparing the interpolated values with the data collected at AHCCD stations.” This raises major concerns about the method used as follows:

1-1) First and foremost, the ANUSPLIN, Alberta Township, and PNWNAmet climate datasets were originally generated/interpolated based on “the same source of observed data (AHCCD).” If they are slightly different in the interpolated values, this is simply due to:

a. different generation (updated version) of AHCCD were used to interpolate the data (Vincent et al., 2002, 2012; Mekis and Vincent, 2011). This implies that if one dataset illustrates slightly poor performance compared to the others, it doesn’t mean it is still the poor choice as they are continuously being updated.

((Reply)) The station-based gridded climate datasets included in this study were not generated based on the adjusted values at the AHCCD stations only, but they employed the raw archive of station data available in a given year to produce gridded climate datasets. While there are only 45 AHCCD stations for precipitation within Alberta, for example, the Alberta Township

dataset employed more than 500 stations (refer to the figure below) to produce the gridded climate data in the 1970s within Alberta. Similarly, for ANUSPLIN, quality-controlled, but unadjusted, station data from the National Climate Data Archive of Environment and Climate Change Canada data (Hutchinson et al., 2009) were interpolated onto the high-resolution grid using thin plate splines. Station density varies over time with changes in station availability, peaking in the 1970s with a general decrease towards the present day (Hutchinson et al., 2009). Thus, the number of stations active across Canada between 1950 and 2011 ranged from 2000 to 3000 for precipitation and 1500 to 3000 for air temperature (Hopkinson et al., 2011). In other words, the station-based gridded climate datasets have been produced based on different station densities which varied spatially and temporally and by applying different set of rules for inclusion of stations in interpolation. Thus, the number of stations included in each dataset is significantly different apart from differences in the interpolation techniques. Therefore, differences are expected in the interpolated values at a location using different gridded climate datasets.



Number of stations used in the interpolation scheme of the Alberta Township dataset

(Source: https://agriculture.alberta.ca/acis/docs/Methodology-and-data-sources-for-interpolated-data-y2019_m03_d27.pdf)

b. the three climate datasets have been generated based on different interpolation techniques. Therefore, the errors/uncertainties might be associated with the interpolation techniques. In this regard, even if one assumes the three climate datasets were generated using the same version of AHCCD at the time of comparison (which is not the case here), the interpolation method of each individual dataset should have been used to estimate unknown points based on known points. That is the way to evaluate the performance of each climate dataset generated by different interpolation techniques. Instead, authors used their own interpolation method (inverse distance squared weighting method) “to obtain the values at the AHCCD stations”, “Then, performance measures were calculated by comparing the interpolated values with the data collected at AHCCD stations.” This means the error found in one dataset could be associated with the interpolation techniques used, - not the original datasets. This could be one of the reasons the Alberta Township climate datasets illustrate better accuracy compared to others. Because the Alberta Township climate datasets have been generated based on different versions of the Inverse Distance Weighting method including “the inverse distance squared weighting method” which was used by the authors to do the evaluation.

((Reply)) If all of the climate datasets were generated from the same set of stations data (e.g., only AHCCD), the skill of interpolation techniques can be evaluated as the reviewer commented. However, the three station-based climate datasets have not used the same source of AHCCD stations as commented above in 1-1 a). Due to the limitation of data availability in a given year, each station-based climate dataset investigated in this study employed different numbers of raw station data. For example, ANUSPLIN used the number of stations ranging from 2000 to 3000 for precipitation and from 1500 to 3000 for temperature.

Although ANUSPLIN and PNWNAmet used the same interpolation approach, i.e., thin-plate smoothing spline, it was found in this study that the performance of ANUSPLIN was much better than that of PNWNAmet. The reason of this difference in performance is that ANUSPLIN used all of Canada-wide archive (raw) station data collected by ECCC in a given year while PNWNAmet employed only stations which cover a common period from 1945 to 2002. Therefore, the different number of stations employed in these two climate datasets may induce the different performances in this study. In addition, the Alberta Township dataset has been produced by the archive (raw) station data collected by ECCC, Alberta Environment and Parks (AEP), and Alberta Agriculture and Forestry (AF) over Alberta. It means, additional stations were used in the Alberta Township data for interpolation, so that the accuracy of the dataset was

also improved. In other words, the performance of the station-based climate datasets included in this study considerably depends on the station density employed in interpolation rather than only on the interpolation techniques used.

(1-2) The authors should avoid comparing apples with oranges when the two CaPA and NARR datasets obtained from models and satellites were compared to the ANUSPLIN, Alberta Township, and PNWNAmets datasets obtained from interpolation techniques. This comparison was done based on the observed detests (AHCCD) which was originally used to generate the ANUSPLIN, Alberta Township, and PNWNAmets datasets. Each point of comparison has been initially used as a centre point to generate the ANUSPLIN, Alberta Township, and PNWNAmets datasets, which can result in high correlations between three as well as the AHCCD datasets due the “existing spatial dependency.” The point values should have been used for evaluations which are “spatially independent.” Otherwise, there is no point in comparing the three interpolated climate datasets with CaPA and NARR which were originally generated to address a poor monitoring network density.

((Reply)) As commented in 1-1, the station-based climate datasets used the archive of raw station data and not the only adjusted values at AHCCD stations. It means the raw but quality controlled observations were used at stations (number of stations used are much greater than only AHCCD stations). Unfortunately, it cannot be guaranteed that the station-based climate datasets are spatially independent with the AHCCD stations as the raw station values at the same location of AHCCD stations might be included in interpolation schemes of each climate dataset. However, it is sure that each station-based climate dataset has been produced using their own spatial structures i.e., different station densities in data generation processes and thus they are unique. On the other hand, CaPA is an amalgamation of rain gauge data, radar data and output from a numerical weather prediction model whereas the NARR data is an amalgamation of NCEP Eta Model (32km/45 layer) output with the Regional Data Assimilation System (RDAS). The archive of raw station data were employed in developing both of these products. As shown in Table 1 (manuscript), the climate datasets used in this study have several inconsistencies with respect to spatial domain, data length, number of climate variables, and spatial resolution. In past, large-scale modelling studies have combined multiple climate datasets to cover the entire study domain or period of record for all the required climate variables, usually without evaluating the performance of different climate datasets for the modelled regions. Thus, the ultimate aim of this study is to suggest a framework that systematically combines multiple climate datasets. In this context, it is meaningful to rank all of the climate datasets and to produce a performance-based hybrid climate dataset to enhance

the performance of numerical models. This study also proved that the hybrid climate dataset provides better representation of historical climatic conditions and thus, enhance the reliability of hydrologic simulations.

Part 2: The gridded datasets have been ranked based on their performance measures.

However:

(2-1) We can not necessarily assign a high performance rank to a grid cell just because of being highly correlated with a nearby station - neither due to its distance nor elevation.

((Reply)) The idea is to rank all of available climate datasets (i.e., five datasets included in this study) based on their varying performance spatially. The performance is determined by comparing the interpolated values against the observed values (at several locations within the study area, AHCCD stations). Various performance measures have been used for ranking instead of just using correlation coefficient. Furthermore, as mentioned in (1-1) above, all of five considered climate datasets are different and unique as they have employed different numbers of climate stations (also varied over time) and generation techniques. Two of the datasets (CaPA and NARR) are very different from three station-based datasets as they also employed the output of weather prediction numerical models in addition to assimilation of station based data. When several datasets are available but there are considerable differences between them, it is reasonable to compare them against the observations to determine their accuracy and thus preference for use.

(2-2) The ranking concept may not be still valid considering some of the comments mentioned in part 1.

((Reply)) Based on the response provided above (part 1, 1-1, 1-2; and part 2, 2-1), authors strongly believe that the reviewer now has a better understanding of the five existing gridded climate datasets and how the suggested methodology could help the researchers in identifying the best data for their study area. The suggested methodology will also help in constructing a reliable, gap filled data for larger regional areas which are otherwise affected by the available data domains.

Part 3: Further evaluation of climate datasets and their ranking have been done using the VIC hydrological model.

(3-1) Five VIC models have been calibrated corresponding to each individual climate dataset. How can you justify associating the errors to the climate data rather than to “the calibration parameters and/or the calibration process, and/or the model structure”?

arbitrary adjustment of parameters might have been done to compensate for the errors in the input climate data - which has been done for each VIC model separately.

((Reply)) Authors agrees that in hydrologic simulations, the biases in climate datasets can be compensated or compromised by model parameters that adjust hydrologic processes to observations. That is, a calibrated parameter set may imply biases in a climate dataset. Therefore, this study applied a multiset-parameter hydrologic simulation approach that employs all parameter sets calibrated by the seven climate datasets and the same climate dataset as a forcing input data to assess the sensitivity of the climate dataset to all feasible parameter sets. From the multiset-parameter hydrologic simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset. The results showed the hybrid climate dataset provides a better representation of historical hydrologic simulations compared to the results of individual climate datasets. The authors also clarified this in the revised manuscript as below,

“As mentioned above, however, intrinsic biases exist temporally and spatially in all of the gridded climate datasets, e.g., discrepancies in the amount and spatial distribution of precipitation between the gridded climate datasets and observations. Therefore, the similarity of the gridded climate datasets in terms of magnitude, sequence and spatial distribution of climate events relative to observations is crucial to reproduce historical observed streamflows. In addition to climate forcings, streamflows are mainly affected by geographic characteristics and physical land surface processes (e.g., infiltration and evapotranspiration), which are represented by model parametrization related to infiltration and soil properties (Demaria et al., 2007). In a hydrologic simulation, the biases in climate datasets can be compromised by model parameters that adjust hydrologic processes to observations (Harpold et al., 2017; Kirchner, 2006). That is, a calibrated parameter set may imply biases in a climate dataset. Under the assumption that the calibrated parameter sets are suitable for hydrologic simulations in each sub-basin, this study applied a multiset-parameter hydrologic simulation approach that employs all parameter sets calibrated by the seven climate datasets and the same climate dataset as a forcing input data to assess the sensitivity of the climate dataset to all feasible parameter sets. From the multiset-parameter hydrologic simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset. In other words, the lower variability in the hydrologic simulations indicates higher reliability in the climate forcing dataset.”(P16L12-P17L4)

(3-2) It has been mentioned in the manuscript that “The proxy validation also confirmed the superior performance of hybrid climate datasets compared with the other five individual climate datasets investigated in this study.” However, the results of the proxy validation (in Table 6) confirm otherwise. Maybe even going one step further, and ask this question whether the two climate datasets; ANUSPLIN, Alberta Township can confirm that there is no need to generate another dataset called “hybrid climate dataset”. Overall, I agree the use of various available data sources in hydrological modeling and qualifying them through alternative simulation scenarios prior to calibration of the model parameters (e.g., Faramarzi et al., 2015), but we need way more rigorous method and justification than what are used in this study to introduce ‘a reference climate dataset’ for a province.

((Reply)) The accuracy of the historical gridded climate datasets considerably depends on employed station density which varies with time and region. As commented in (1-1), all of the climate datasets have employed different station densities, methods, and techniques in the processes of data generation, thus they are quite different. Therefore, there is a need to evaluate their performance before application so that an informed decision could be made before their application. Having several products of varying quality may pose serious concerns especially when these are applied without understanding the differences, reliability and accuracies. Further, the performance of a dataset may vary with region and hence requires such assessment for each study area as the results presented here cannot be generalized for the entire data domains. As commented in (1-2), numerical modelers have suffered from the inconsistency of available climate datasets in spatial domain and resolution, data length, and climate variables. In this context, the Athabasca River basin is a good test-bed because the whole domain cannot be covered by the Alberta Township data which was dominantly ranked first. Combining Alberta Township with ANUSPLIN simply for the Athabasca River basin instead of generating the hybrid climate dataset, as commented by the reviewer, we may neglect added values of ANUSPLIN within the domain of Alberta Township. Further, we may neglect added values of other climate datasets available within the basin. Therefore, we suggested the REFERENCE Reliability Evaluation System (REFRES) that systematically produces a performance-based hybrid climate dataset. For the Clearwater sub-basin, all of five climate datasets contribute to generating the hybrid climate data for precipitation (refer to Table 5), resulting in relatively a larger improvement in hydrologic simulations as shown in Table 6. In addition, the other reviewer also suggested that the usefulness of the hybrid climate dataset can be clearly found

at the whole basin scale instead of a sub-basin scale as the added values may be accumulated at the main stream over the entire ARB. The authors conducted additional analysis to simulate the entire basin and computed NSE values at a few hydrometric stations in the main stream of the Athabasca River (refer to the table below). The results showed that the hybrid climate dataset performs better than the existing gridded climate dataset (in this case ANUSPLIN) as the drainage areas are larger. This is mainly due to the fact that as the watershed area increases, the derived hybrid climate dataset is no longer dominated by a single gridded dataset. We also addressed these results in the reply of (3-2) in AC2.

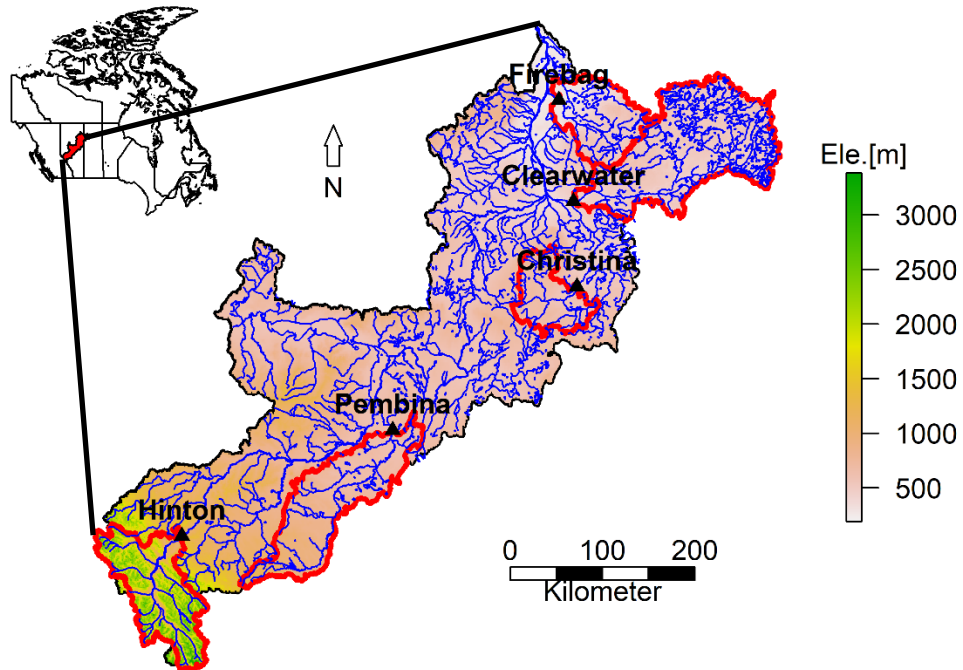
Nash-Sutcliffe Efficiency (NSE) of ANUSPLIN and the hybrid climate datasets at the main stream of the Athabasca River

| No | Station name/ID | Drainage area (km ²) | ANUSPLIN | | Hybrid | |
|----|-----------------------------|-------------------------------------|-------------|------------|-------------|------------|
| | | | Calibration | Validation | Calibration | Validation |
| 1 | Hinton / 07AD002 | 9,760 | 0.85 | 0.82 | 0.83 | 0.76 |
| 2 | Windfall / 07AE001 | 19,600 | 0.80 | 0.72 | 0.80 | 0.76 |
| 3 | Athabasca / 07BE001 | 74,600 | 0.78 | 0.69 | 0.77 | 0.78 |
| 4 | Fort McMurray / M07DA001 | 133,000 | 0.77 | 0.66 | 0.78 | 0.75 |
| 5 | Eymundson / S24 | 147,086 | 0.77 | 0.67 | 0.79 | 0.75 |

<Specific Comments>

Authors may consider using coordinate systems for figures, especially Fig. 3 and 8 that can help readers to locate the study area and better investigate its climate.

((Reply)) The authors modified Figure 3 to provide the geographical information of the ARB as the reviewer suggested.



<<Short Comments from Fuad Yassin>>

<General Comments>

This study addresses a relevant topic, particularly in Canada, where there is a huge limitation of reliable high-density observed climate data. Although I find the study very interesting, I have two important general comments that need better clarification.

(1) The first comment is that why other important data sources ignored in this study? If you look at the study of Wong et al. (2017), they demonstrated that GPCC and CRU data are good candidates in Canada compared to NARR. In their study, NARR was found to be the worst data set, and it is not clear why it is accounted in this study, while GPCC and CRU data present unique data globally with long-term and high-temporal resolution data. I believe a better explanation about this is needed, and accounting GPCC and CRU data would provide greater insight for the audience.

((Reply)) Wong et al. (2017) intercompared multiple climate datasets for only precipitation at monthly time step while this study did for precipitation and temperature at daily scale. Both GPCC and CaPA provide daily precipitation at the global and North America domains. However, GPCC has a coarser resolution ($1.0^\circ = \sim 100\text{km}$) while CaPA provides a higher resolution, at 10km, with a better monitoring network in Canada. Therefore, CaPA has been selected in this study. In addition, CRU has been also excluded as it provides monthly climate datasets. REFRES has a flexible structure to include a new climate dataset when available. For example, the Climate Forecast System Reanalysis (CFSR) dataset will be included in the next version of REFRES.

(2) My second observation is that why only few streamflow stations are used for proxy validation? My understanding is that there are many streamflow stations in the study area, especially around headwaters where huge variability and magnitude of precipitation expected.

((Reply)) Yes, there are other hydrometric stations in the upper reach in the ARB. The five sub-basins were selected for the proxy validation based on three criteria: a) hydrometric record length, b) location defined by upper, middle and lower reaches (Northern River Basin Study, 2002), and c) the number of gridded climate datasets used to generate a hybrid climate dataset for the catchment area of the selected hydrometric station. Based on first criteria, hydrometric stations with a short-period of record and/or severe data gaps were excluded for the proxy validation. There are several stations in the lower watersheds with shorter record length (as they

were installed between 2000 and 2010). In addition several stations are only being operated during open water season (i.e., summer) and do not have any observations during winter seasons and they have also been excluded. For example, the Windfall station (ID: 07AE001) has been excluded because it has hydrometric record only during open-water period. The other reviewer also suggested that the usefulness of the hybrid climate dataset can be clearly found at the whole basin scale instead of a sub-basin scale. Accordingly, we calibrated the VIC model for the whole basin to provide additional results at the main stream of the Athabasca River as shown in the table below. The results showed that the hybrid climate dataset performs better than the existing gridded climate dataset as the drainage areas are larger. We also addressed these results in the reply of (3-2) in AC2 and SC1.

Nash-Sutcliffe Efficiency (NSE) of ANUSPLIN and the hybrid climate datasets at the main stream of the Athabasca River

| No | Station name/ID | Drainage area (km ²) | ANUSPLIN | | Hybrid | |
|----|-----------------------------|-------------------------------------|-------------|------------|-------------|------------|
| | | | Calibration | Validation | Calibration | Validation |
| 1 | Hinton / 07AD002 | 9,760 | 0.85 | 0.82 | 0.83 | 0.76 |
| 2 | Windfall / 07AE001 | 19,600 | 0.80 | 0.72 | 0.80 | 0.76 |
| 3 | Athabasca / 07BE001 | 74,600 | 0.78 | 0.69 | 0.77 | 0.78 |
| 4 | Fort McMurray / M07DA001 | 133,000 | 0.77 | 0.66 | 0.78 | 0.75 |
| 5 | Eymundson / S24 | 147,086 | 0.77 | 0.67 | 0.79 | 0.75 |

1 **Hybrid climate datasets from a climate data evaluation system and their impacts on**
2 **hydrologic simulations for the Athabasca River basin in Canada**

3
4
5
6
7
8

Hyung-II Eum¹, and Anil Gupta^{1,2}

For submission to Hydrology and Earth System Sciences (HESS)

H.-I. Eum (Corresponding author, email: hyung.eum@gov.ab.ca)

¹ Alberta Environment and Parks, Environment Monitoring and Science Division, 3535 Research Road NW, Calgary, Canada, T2L 2K8

² Department of Geomatics Engineering, University of Calgary, 2500 University Drive NW, Calgary, Alberta, Canada

1 **Abstract**

2 A reliable climate dataset is a backbone for modeling the essential processes of the water cycle and
3 predicting future conditions. Although a number of gridded climate datasets are available for North
4 American continent, which provides reasonable estimates of climatic conditions in the region, there are
5 inherent inconsistencies in these available climate datasets (e.g., spatial_ and temporal_varying data
6 accuracies, meteorological parameters, lengths of records, spatial coverage, temporal resolution, etc). These
7 inconsistencies raise questions as to which datasets are the most suitable for the study area and how to
8 systematically combine these datasets to produce a reliable climate dataset for climate studies and
9 hydrological modeling. This study suggested a framework, called the reference reliability evaluation system
10 (REFRES), that systematically ranks multiple climate datasets to generate a hybrid climate dataset for a
11 region. To demonstrate the usefulness of the proposed framework, REFRES was applied to produce a
12 historical hybrid climate dataset for the Athabasca River basin in Alberta, Canada. A proxy validation was
13 also conducted to prove the applicability of the generated hybrid climate datasets to hydrologic simulations.
14 This study evaluated five climate datasets, including station-based gridded climate datasets (ANUSPLIN,
15 Alberta Township, and PNWNAmet), a multi-source gridded dataset (Canadian Precipitation Analysis -
16 CaPA), and a reanalysis-based dataset (NARR). The results showed that the gridded climate interpolated
17 from station data performed better than multi-source and reanalysis based climate datasets. For the
18 Athabasca River basin, Township and ANUSPLIN were ranked first for precipitation and temperature,
19 respectively. The proxy validation also confirmed the utility of hybrid climate datasets in hydrologic
20 simulations, compared with the other five individual climate datasets investigated in this study. These
21 results indicate that the hybrid climate dataset provides the best representation of historical climatic
22 conditions and thus, enhances the reliability of hydrologic simulations.

23
24 **Key words:** Historical gridded climate data, reference reliability evaluation system, hydrological
25 simulation, Athabasca River basin, proxy validation

Deleted: the

Deleted:

Deleted: length

Deleted:

Deleted: a valid question

Deleted: determines a ranking of

Deleted: mostly

Deleted: superior performance

Deleted: a better

Deleted: enhancing

Deleted: Reference

Deleted: Proxy

1 **1. Introduction**

2 A reliable historical climate dataset is essential in understanding the climatic and hydrological
3 characteristics of a watershed, as it is a crucial forcing input data for simulating key processes of the water
4 and energy cycles in impact models (Deacu et al., 2012; Essou et al., 2016; Wong et al., 2017). Although
5 climate monitoring networks have advanced over the last decades, poor network density still exists,
6 especially in western mountainous and northern parts of Canada. Moreover, climate observations are often
7 spatially interpolated to cover ungauged regions, which may cause unexpected erroneous model predictions
8 as a consequence of the sparse measurements network, especially for mountainous areas affected by
9 orographic effects (Rinke et al., 2004; Wang and Lin, 2015).

Deleted: to understand

Deleted: cycle

Deleted: a

Deleted: the

10 As advances in numerical hydrologic and hydrodynamic modeling have increased the capability and
11 reliability in simulating complex natural processes to detect anthropogenic and natural climate changes, a
12 need for temporally_ and spatially_ reliable climate data has also been grown to accommodate the
13 requirements of input data for numerical models (Shen et al., 2010; Shrestha et al., 2012; Islam and Dery,
14 2017). For instance, process-based distributed hydrologic models have a grid-based structure that requires
15 input data for each grid cell. However, a simple spatial interpolation of observational station data to all
16 model grid cells may not produce a reliable input forcing dataset for hydrologic models, particularly in a
17 region with a sparse gauging network. A reliable historical climate dataset is also crucial in climate change
18 studies when used for statistical downscaling techniques that employ the relationships between observations
19 and outputs of global (or regional) climate models to produce climate forcing at regional or local scales.
20 Since the resolution of products from a statistical downscaling technique usually corresponds to that of the
21 historical climate dataset (Werner and Cannon, 2016; Eum and Cannon, 2017), the availability of
22 temporally_ and spatially_ reliable, historical climate data is essential for climate-related impact studies
23 (Christensen and Lettenmaier, 2007; Kay et al., 2009; Gutmann et al., 2014; Eum et al., 2016).

Deleted: produces

Deleted:

24 A number of high-resolution gridded climate datasets have been developed for various applications
25 such as inter-comparison studies (Eum et al., 2014a; Wong et al., 2017) and hydrologic modeling (Choi et

1 al., 2009; Eum et al., 2016). There are various types of gridded climate datasets available for the North
2 American region; 1) station-based interpolated, 2) station-based multiple-source, and 3) reanalysis-based
3 multiple-source (Wong et al., 2017). By interpolation of observational station data, long-term gridded
4 climate datasets have been produced over various domains defined by stations incorporated such as Canada-
5 wide Australia National University's spline (ANUSPLIN, Hutchison et al., 2009), the Alberta Township
6 data (Shen et al., 2001), and the PCIC NorthWest North America meteorological (PNWNAmet) dataset
7 (Werner et al., 2019). The Canadian Precipitation Analysis (CaPA) system, a multiple source-based climate
8 dataset, has been developed to produce near real-time precipitation analyses (6-hr accumulated precipitation)
9 over North America at 15 km resolution which has been further improved to 10km resolution (Lespinas et
10 al., 2015). North American Regional Reanalysis (NARR), one of the reanalysis-based datasets derived from
11 a regional climate model (~32km), has been tested as an alternative climate dataset (Choi et al., 2009;
12 Praskievicz and Bartlein, 2014; Essou et al., 2016; Islam and Dery, 2017).

Deleted: It is an example of a multiple source-based climate dataset.

13 In most of the large-scale modelling studies, multiple climate data sets were combined to cover the
14 entire modelling domain for all the required climate variables, usually without evaluating the performance
15 of different climate datasets for the modelled regions (Faramarzi et al., 2015; Shrestha et al., 2017; Wong
16 et al., 2017). The lack of performance indicators for available climate datasets may cause inappropriate
17 application of these datasets for various large scale studies, resulting in unreliable outputs, e.g., considerable
18 bias in statistical downscaling studies. Therefore, selecting reliable gridded climate data for a study area is
19 crucial for any hydrological or climate-related studies (Werner and Cannon, 2016; Eum et al., 2014a; 2017).
20 Eum et al. (2014a) intercompared three gridded climate datasets (ANUSPLIN, NARR, and CaPA) for the
21 Athabasca River Basin (ARB) and found that data accuracy varies spatially and temporally over the basin
22 mainly due to the heterogeneity of spatial density of the observational climate network in the basin and
23 limited data assimilation. Wong et al. (2017) also intercompared gridded precipitation datasets derived from
24 different data sources over Canada. Few studies have attempted to incorporate spatially-varied performance
25 measures of various climate datasets to produce a complete long-term historical climate dataset for a study

Deleted:

Deleted: , however

Deleted: an

Deleted: varying

Deleted:

Deleted: measure

Deleted: in producing

1 region (Faramarzi et al., 2015; Shrestha et al., 2017). In addition, no systematic framework has been
2 developed yet that could be employed by climatic and hydrologic studies.

Deleted:), however

3 Therefore, this study provides a framework, called REFERENCE Reliability Evaluation System
4 (REFRES), to systematically determine the ranking of multiple climate datasets based on their performance
5 and generate a hybrid climate dataset for a study region by extracting the best candidate (based on the
6 ranking) from multiple climate datasets available in a repository. Several performance measures were
7 identified and calculated by comparing to the Adjusted and Homogenized Canadian Climate Data (AHCCD)
8 over western Canada. Based on the performance measures, the climate datasets were ranked to generate a
9 hybrid climate dataset for the area of interest (target area). A hybrid dataset for two climate variables -
10 precipitation and temperature, key forcing for hydrological modeling, was produced for a period of record

Deleted: to

11 that is fully covered by the multiple climate datasets. To validate the applicability of the hybrid climate
12 dataset, a proxy validation approach was employed by comparing simulated streamflows derived from the
13 generated hybrid climate data and other available climate datasets to recorded streamflows at various
14 hydrometric stations in the Athabasca River basin (ARB). Streamflows were simulated using a hydrologic
15 model (Variable Infiltration Capacity, VIC) calibrated and forced by individual climate datasets and the
16 generated hybrid climate dataset. Therefore, the aims of this study are 1) to develop a methodology (i.e.,
17 reference reliability evaluation system, REFRES) to compare and rank multiple gridded climate datasets
18 based on the proposed performance measures and to generates the hybrid climate dataset, and 2) to validate
19 the hybrid climate dataset using the proxy validation approach for the Athabasca River basin as a case study
20 to confirm the applicability of hybrid climate dataset to hydrologic simulations.

Deleted: full

Deleted: overlapped or

Deleted: using

Deleted: is

Deleted: superiority

22 2. Climate data

23 2.1 Adjusted and Homogenized Canadian Climate Data (AHCCD)

24 Climate station observations in Canada are available from the national climate data and information
25 archive of Environment and Climate Change Canada (ECCC, <http://climate.weather.gc.ca/>). Besides the

1 variable number of observations due to frequent changes in operations including discontinuation of stations,
2 the observations are also subject to various errors from undercatch of solid precipitation, orographic effects,
3 and malfunction of measurements (Mekis and Hogg, 1999; Rinke et al., 2004).

4 Mekis and Vincent (2011) adjusted daily rainfall and snowfall data, considering wind undercatch,
5 evaporation, and wetting losses corresponding to the types of gauges for 450 stations over Canada. The
6 most recent version released in 2016 provides the adjusted precipitation observations, expanded to 464
7 precipitation stations. Vincent et al. (2012) produced the 2nd generation of homogenized daily temperature
8 by adjusting the time series at 120 synoptic stations to account for a nation-wide change in observing time
9 and homogenizing discontinuities over 338 temperature (daily minimum and maximum) stations in Canada.
10 The adjusted and homogenized Canadian Climate Data (AHCCD) are available through Environment and
11 Climate Change Canada (<http://ec.gc.ca/dccha-ahccd/default.asp?lang=En&n=B1F8423>).

Deleted: expanding

12 Considering that archived raw station data were used to produce the historical gridded climate datasets
13 used in our study, the evaluation of performance at the AHCCD stations is more meaningful because the
14 AHCCD data were adjusted to account for the known measurement issues in the raw station data. For
15 example, the adjusted precipitation data are higher by 5 % to 20 %, varying with topographic characteristics
16 (Mekis and Vincent, 2011). Therefore, the AHCCD dataset is recognized as the best estimate of actual
17 climate variables in Canada, and consequently used in a number of climate-related studies (Asong et al.,
18 2015; Eum et al., 2014a; Shook and Pomeroy, 2012; Wong et al., 2017). As large-scale watersheds in Alberta
19 are crossing the province, e.g., the Peace River and Athabasca River basins, this study evaluated the
20 performance of the historical gridded climate datasets at the AHCCD stations within British Columbia (BC),
21 Alberta (AB), and Saskatchewan (SK) (190 and 129 stations for precipitation and temperature, respectively,
22 in Figure 1). The AHCCD stations have different record lengths. For example, the longest record period is
23 from 1840 to 2016 while the shortest period is from 1967 to 2004. As the data lengths are different at each
24 AHCCD station, we selected a common period between each AHCCD station and climate dataset to
25 estimate performance measures.

Deleted: %

Deleted: Peacece

1 Figure 1. AHCCD stations within the British Columbia (BC), Alberta (AB), and Saskatchewan (SK)
2 provinces

4 2.2 Historical gridded climate datasets

5 In general, the available historical gridded climate dataset can be divided into three categories; 1)
6 station-based, 2) multiple source-based, and 3) reanalysis-based. In this study, five high-resolution gridded
7 climate datasets available for Alberta were selected (Table 1) to evaluate their performance and include in
8 the generation of a hybrid climate dataset for Alberta.

Deleted: were included

9 Table 1. High-resolution gridded historical climate datasets used in this study

11 2.2.1 Station-based datasets

12 Hutchinson et al. (2009) produced a Canada-wide daily climate dataset at 10 km resolution from 1961
13 to 2003 by the Australia National University's trivariate thin-plate smoothing spline (ANUSPLIN)
14 technique to model the complex spatial patterns (e.g., large variations in ground elevation and station
15 density over Canada) of daily weather data. Hopkinson et al. (2011) updated the existing ANUSPLIN
16 dataset by reducing residuals and extended the daily weather data from 1950 to 2011. Recently,
17 ANUSPLIN data were extended until 2015 for three climate variables, i.e., daily precipitation, minimum
18 and maximum air temperature, which were interpolated with 7,514 surface-based observations (archive
19 data) of Environment Canada. However, the numbers of stations included in interpolation varied year to
20 year, ranging from 2,000 to 3,000 for precipitation and from 1,500 to 3,000 for air temperature. The
21 ANUSPLIN data generated by Natural Resource Canada (NRCan) have been used as the source data to
22 compare climate products (Eum et al., 2014a; Wong et al., 2017), evaluate the accuracy of regional climate
23 models (Eum et al., 2012), and to model hydrologic regimes (Islam and Dery, 2017; Eum et al., 2017;
24 Dibike et al., 2018).

Deleted: and include 3-

Deleted: .g.

Deleted: and

Deleted: ranged

Deleted: inter-

Deleted: to

Deleted: modeling

1 Similar to the ANUSPLIN dataset, Pacific Climate Impacts Consortium (PCIC) also generated daily
2 precipitation, minimum and maximum air temperature, and wind speed from 1945 to 2012 at 1/16 degree
3 (6~7km) resolution using a thin-plate smoothing spline technique over Northwest North America, called
4 the PCIC North West North America meteorological (PNWNAmet, Werner et al., 2019) dataset
5 (https://data.pacificclimate.org/portal/gridded_observations/map/). While ANUSPLIN utilized a varying
6 number of gauge stations depending on availability of observations in a given year, PNWNAmet set a
7 common period from 1945 to 2012 for all stations included in the interpolation over regularly spaced grid
8 cells within the domain. The PNWNAmet dataset was developed to produce forcing data for an updated
9 version of the Variable Infiltration Capacity model with glaciers (VIC-GL). In addition to precipitation, and
10 minimum and maximum temperature, PNWNAmet includes wind speed, which considerably affects vital
11 hydrologic processes, especially evapotranspiration, sublimation, and snow transport (i.e., snow blowing).
12 Because the AHCCD dataset provides only daily precipitation and temperature, wind speed was excluded
13 in this study.

14 Alberta Agriculture and Forestry (AF) produced the Alberta Township data
15 (<http://agriculture.alberta.ca/acis/township-data-viewer.jsp>) from 1961 to 2016 at approximately 10km
16 (Alberta Township grid) resolution using a hybrid inverse distance weighting (IDW) process (Shen et al.,
17 2001) for daily precipitation, minimum and maximum temperature, relative humidity, wind speed, and solar
18 radiation. The archive (raw) station data collected by ECCC, Alberta Environment and Parks (AEP), and
19 AF over Alberta were used in producing the Township dataset. The Township data used various effective
20 radiuses (60 km to 200 km) to ensure a sufficient number of gauge stations in IDW. When there is no station
21 within 200 km, it is assumed that the nearest station represents the climate conditions of the Township
22 center. The domain of Township data covers most of Alberta except the mountainous regions while both
23 ANUSPLIN and PNWNAmet cover all of western Canada (refer to Table 1). Therefore, one of the
24 limitations of the Township dataset is its application to a large watershed spanning Alberta and other
25 neighboring provinces.

1

2 **2.2.2 Multiple source-based dataset**

3 As an operational system, the Meteorological Service of Canada initiated the Canadian Precipitation
4 Analysis (CaPA) in 2003 to produce superior gridded precipitation data over North America at 10 km
5 resolution (Lespinas et al., 2015), especially for regions with poor observational networks (Mahfouf et al.,
6 2007). CaPA employs an optimum interpolation technique that requires properties of error statistics among
7 observations and a first guess, i.e., background field (Garand and Grassotti, 1995). A short-term forecast of
8 6-hr accumulated precipitation from the Canadian Meteorological Centre (CMC) regional Global
9 Environmental Multiscale (GEM) model (Côté et al., 1998a; 1998b) is used in CaPA as the background
10 field. The assimilated precipitation from the Canadian weather radar network and 33 US radars near the
11 border are used as additional observations to generate analysis error among multiple sources of observations
12 and the background precipitation. Zhao (2013) tested the applicability of CaPA for hydrologic modelling in
13 the Canadian Prairies and proved its usefulness in data-sparse regions and the winter season. In addition,
14 CaPA has been widely used in agricultural and hydrologic applications (Deacu et al., 2012; NIDIS, 2015).
15 Eum et al. (2014a) further addressed some of the limitations of CaPA, i.e., lack of air temperature which is
16 one of the primary drivers in hydrologic modeling and shorter data length (only from 2002 to 2017), for
17 model calibration and validation. Using 6-hr accumulated precipitation CaPA products, in this study, daily
18 accumulated precipitation was generated over western Canada by adjusting the time zone from Universal
19 Time Coordinated (UTC) to Mountain Time (MT).

Deleted:

Deleted:)

20

21 **2.2.3 Reanalysis-based dataset**

22 Reanalysis products are another common type of gridded dataset used in climate and hydrologic
23 studies. The North American Regional Reanalysis (NARR) was developed to create a long-term set of
24 dynamically consistent 3-hourly climate data from 1979 to 2003 at a regional scale ($0.3^\circ = \sim 32\text{km}$) for the
25 North America domain (Mesinger et al., 2006). By utilizing advanced land-surface modeling and data

1 assimilation through the Eta Data Assimilation System (EDAS), NARR improved the National Centers for
2 Environmental Prediction/National Center for Atmospheric Research (NCEP/NCAR) global reanalysis data.
3 NARR cycled every 3 hours to produce a climate dataset from 1979 to the current year. Choi et al. (2009)
4 tested the applicability of NARR for hydrologic modeling in Manitoba for a region with a poor monitoring
5 network density. However, the NARR dataset after 2004 is not consistent with that of prior years (i.e., 1979
6 to 2003) because assimilation of precipitation observations was discontinued in 2003 (Eum et al., 2014a).
7 Using the 3-hr NARR climate data, daily precipitation, and minimum and maximum temperature were
8 calculated by adjusting the time zone to MT from the original NARR dataset (UTC zone).

Deleted: is

Deleted: ,

3. Methodology

3.1 Reference Reliability Evaluation System (REFRES)

12 This study suggests a REFERENCE Reliability Evaluation System (REFRES) that consists of three
13 main modules (refer to Figure 2): 1) a performance measure module (PMM) to evaluate various
14 performance measures for each climate dataset, 2) a ranking module (RM) to identify the most reliable
15 climate data for a target grid cell using a multi-criteria decision-making technique based on the performance
16 measures provided by PMM, and 3) a data generation module (DGM) to produce a hybrid climate dataset
17 by selecting the most reliable climate dataset based on the ranking provided by the RM (ranking model).
18 These three modules are seamlessly integrated and exchange the required data and information to generate
19 a hybrid climate dataset. The next section provides further details on each module.

Deleted:);

Deleted: RM.

Deleted: more

20 Figure 2. Structure of REFRES comprise of three modules; 1) Performance Measure Module (PMM), 2)
21 Ranking Module (RM), and 3) Data Generation Module (DGM)

Moved (insertion) [1]

Formatted: English (Canada)

3.1.1 Performance Measure Module (PMM)

24 AHCCD is a point (station) dataset while the other climate datasets used in this study (refer to Table
25 1) are regularly spaced gridded datasets with varying time period, spatial resolution, and coverage (i.e.,

Moved down [2]: Figure 2. Structure of REFRES comprised of three modules; 1) Performance Measure Module (PMM), 2) Ranking Module (RM), and 3) Data Generation Module (DGM)¶

Formatted: English (United States)

Formatted: Font: Not Bold

1 domain). Therefore, the inverse distance squared weighting method was applied to obtain the values at the
 2 AHCCD stations from all the gridded climate datasets. Then, performance measures were calculated by
 3 comparing the interpolated values with the data collected at AHCCD stations. The choice of the
 4 performance measures is vital in REFRES₂ as the ranking of climate datasets entirely depends on included
 5 performance measures. In this study, performance measures were selected based on three criteria: 1)
 6 distribution, 2) sequencing, and 3) spatial pattern. Distribution-related performance is assessed by [the](#)
 7 Kolmogorov-Smirnov D statistic (D_{KS}) and standard deviation ratio (σ_{ratio}). Sequence-related performance
 8 is assessed by [the](#) percentage of bias (P_{bias}), root mean square error (RMSE), and temporal correlation
 9 coefficient (TCC). Spatial pattern-related performance is evaluated by [the](#) pattern correlation coefficient
 10 (PCC) as shown in Eq. (1) to Eq. (5). The equations of TCC and PCC are identical but TCC is calculated
 11 with [the](#) daily time series of climate variables and PCC is obtained by the mean annual precipitation and
 12 temperature of the AHCCD stations over a target domain. Therefore, PCC varies with the user specified
 13 target domain.

$$14 \quad D_{KS} = \sup |F_G(x) - F_O(x)| \quad (1)$$

$$15 \quad \sigma_{ratio} = \{(\sigma_G / \sigma_O) - 1\} \quad (2)$$

$$16 \quad P_{bias} = \frac{\sum_{i=1}^N (G_i - O_i)}{\sum_{i=1}^N O_i} \times 100 \quad (3)$$

$$17 \quad RMSE = \sqrt{\frac{\sum_{i=1}^N (G_i - O_i)^2}{N}} \quad (4)$$

$$18 \quad TCC, PCC = \frac{\sum_{i=1}^N (G_i - \bar{G})(O_i - \bar{O})}{\sqrt{\sum_{i=1}^N (G_i - \bar{G})^2} \sqrt{\sum_{i=1}^N (O_i - \bar{O})^2}} \quad (5)$$

19 where σ_G and σ_O are the standard deviation of gridded and observed climate datasets, G_i and O_i represent
 20 gridded and observed climate datasets [at \$i\$ th time step](#), respectively; F is the empirical distribution function
 21 of a climate dataset; σ is standard deviation; \bar{G} and \bar{O} represent the mean of gridded and observed
 22 climate datasets, respectively and N is a total number of data points. These six performance measures were

Deleted: G

Deleted: O

1 calculated for all the selected climate datasets and variables at each AHCCD station. Figure 2 (blue box in
2 PMM) shows an example of 6 PMs calculated for the precipitation variable using the ANUSPLIN gridded
3 data. Thus, 15 tables (5 climate datasets \times 3 variables) were generated by PMM and transferred to the RM.

Deleted: ranking module (RM).

5 3.1.2 Ranking Module (RM)

6 The function of the ranking module is to select the appropriate AHCCD stations for a given target grid
7 cell and to rank all the gridded data sets based on the six performance measures calculated in the previous
8 module. For a given target cell, AHCCD stations are selected based on two criteria: distance and elevation.
9 Firstly, 20% (of all AHCCD) stations are selected based on the nearest distance criteria, which were then
10 again reduced by the five nearest stations based on the minimum elevation difference criteria. Then the
11 performance measures are averaged over the selected AHCCD stations to represent the skill of each climate
12 dataset for the given target grid cell.

Deleted: ,

Deleted: First

Deleted: criteria

Deleted: to half

Deleted:

13 As multiple performance measures are employed in this study, there are situations when a climate
14 dataset may perform well for some measures but not for others. Therefore, a multi-criteria decision-making
15 (MCDM) technique is required to systematically rank all of the climate datasets while considering multiple
16 performance measures. This study applied a multi-criteria decision-making technique called the Technique
17 for Order of Preference by Similarity to Ideal Solution (TOPSIS, Hwang and Yoon 1981) to systematically
18 determine the order of preference for all climate datasets at each target grid cell. TOPSIS calculates the
19 geometric distance between alternatives and an ideal solution defined by the best performance on each
20 criterion from the alternatives, and then determines the best and worst alternatives based on the distance.
21 TOPSIS has been successfully applied to watershed management for multi-criteria problems (Jun et al.,
22 2013; Lee et al., 2013). TOPSIS starts with the averaged performance measures, $(x_{ij})_{m \times n}$ for the i^{th} alternative
23 (climate dataset in this study) and j^{th} criterion (i.e., a performance measure). A weighted normalized decision
24 matrix, $(t_{ij})_{m \times n}$ is given by

Deleted: in

Deleted: ,

Deleted:),

Deleted: determine

Deleted:), and a

$$(t_{ij})_{m \times n} = (w_j n_{ij})_{m \times n} \quad i = 1, 2, \dots, m; \quad j = 1, 2, \dots, n \quad (6)$$

$$n_{ij} = \frac{x_{ij}}{\sum_{i=1}^m x_{ij}^2} \quad (7)$$

where, m and n are the total number of alternatives and criteria, respectively, n_{ij} is normalized matrix by Eq. (7), and w_j represents weighting on the j^{th} criterion. Under the assumption that all performance measures are important, this study used an equal weighting. Then, Euclidean distances (d_{iw} and d_{ib}) of climate datasets from the best (A_b) and worst (A_w) conditions were calculated respectively by Eq. (8) to Eq. (11)

$$A_w = \{ \{ \max(t_{ij} | i = 1, 2, \dots, m) | j \in J_-, \langle \min(t_{ij} | i = 1, 2, \dots, m) | j \in J_+ \rangle \} \equiv \{ t_{wj} | j = 1, 2, \dots, n \} \quad (8)$$

$$A_b = \{ \{ \min(t_{ij} | i = 1, 2, \dots, m) | j \in J_-, \langle \max(t_{ij} | i = 1, 2, \dots, m) | j \in J_+ \rangle \} \equiv \{ t_{bj} | j = 1, 2, \dots, n \} \quad (9)$$

$$d_{iw} = \sqrt{\sum_{j=1}^n (t_{ij} - t_{wj})^2} \quad i = 1, 2, \dots, m \quad (10)$$

$$d_{ib} = \sqrt{\sum_{j=1}^n (t_{ij} - t_{bj})^2} \quad i = 1, 2, \dots, m \quad (11)$$

Where, t_{bj} and t_{wj} are the best and worst decision matrices determined by Eq. (8) and (9), respectively, and J_+ and J_- represent criteria that have a positive and a negative impact on performance. For example, TCC and PCC are in J_+ while D_{KS} , σ_{ratio} , P_{bias} , and RMSE are in J_- . Using the Euclidean distances, the order of preference for all climate datasets was determined by the similarity (S_{iw}) to the worst condition in Eq. (15).

$$s_{iw} = \frac{d_{ib}}{d_{iw} + d_{ib}}, \quad 0 \leq s_{iw} \leq 1, \quad i = 1, 2, \dots, m \quad (15)$$

$s_{iw} = 1$, when the alternative is equal to the best condition (A_b) and $s_{iw} = 0$ if the alternative is equal to the worst condition (A_w). In other words, a higher s_{iw} represents higher preference among alternatives. As we evaluate the performance measures (criteria) for individual climate variables, TOPSIS can be applied to decide the preference of climate datasets considering the performance measures for either individual or multiple variables. In this study, TOPSIS provides two types of ranking information by using performance measures from i) individual climate variable and ii) all climate variables. That is, one is the ranking for

Deleted: fairly

Deleted: matrix

Deleted: When

Deleted: ,

Deleted:),

Deleted: evaluated

1 precipitation and temperature separately (R_{ind}) and the other is the ranking for multiple variables (R_{mul}). For
2 example, in this study, R_{ind} was determined by a 5×6 decision matrix (5 climate datasets and 6 performance
3 measures) for precipitation and temperature individually, while R_{mul} was determined by a 4×18 decision
4 matrix (4 climate datasets excluding CaPA that provides only precipitation by 18 performance measures
5 from three variables). To alleviate the erroneous output that minimum temperature is higher than maximum
6 temperature on a certain day when producing the hybrid climate dataset by the ranking of temperature
7 values individually, the performance measures of both minimum and maximum temperature are employed
8 together to rank the climate datasets for temperature.

Deleted: an

10 3.1.3 Data Generation Module (DGM)

11 DGM extracts the most reliable climate data for a user-specified target region based on the ranking
12 information obtained from the RM. The tool is flexible enough to provide output in various common
13 formats, i.e., NetCDF, ASCII (text) or in the specific format of a numerical model. As all of the historical
14 gridded climate datasets have been tested and employed in numerous climatic and hydrologic studies, an
15 assumption was made in generating the hybrid climate dataset that all of the climate datasets are equally
16 qualified for inclusion but the final selection can be determined by the proven superiority evaluated through
17 the performance measures. Under this assumption, the available datasets can be combined systematically
18 based on the rank (performance) of each dataset at target grid cells. As each climate dataset has different
19 data periods shown in Table 1, the first ranked dataset cannot fully cover a whole target period to be
20 extracted from a set of climate data candidates. DGM provides a systematic procedure to identify the most
21 reliable dataset for a target region and extracts the data from the inventory of climate datasets considering
22 the ranking and availability of each dataset for a desired period. For instance, if CaPA and ANUSPLIN
23 ranked first and second for precipitation and the desired period is 1950 to 2016, DGM starts searching for
24 the availability of precipitation in 1950. As CaPA is only available between 2002 to 2016, DGM reorders
25 the rank to select ANUSPLIN as the best climate dataset available in 1950. In this way, a hybrid dataset

Deleted:

Deleted:),

Deleted: has

Deleted: Thus,

Deleted: when

1 over the period 1950 to 2016 is generated by extracting from ANUSPLIN from 1950 to 2001 and CaPA
2 from 2002 to 2016 in this particular case. Once the best climate datasets are extracted over all the target
3 grid cells (study domain), the hybrid climate dataset is produced in a user-defined format. This study
4 generated the hybrid climate datasets in the form of the VIC forcing input format to be directly employed
5 into the hydrologic model.

Deleted:

7 **3.2 Proxy validation**

Deleted: ¶

8 Although the AHCCD dataset has been adjusted to provide better estimates of actual precipitation and
9 temperature, it contains statistical artifacts that include inevitable errors from sequential data processes that
10 can be propagated in the derived hybrid climate dataset. Given that the AHCCD stations, the reference
11 dataset for the performance measures, are not regularly distributed and have especially poor density in the
12 northern parts of the study area (refer to Figure 1), it is questionable if the hybrid climate dataset can
13 represent a historical climate better than the individual gridded climate dataset. Utilizing a proxy validation

Deleted: , with

14 approach (Klyszejko, 2007), this study applied streamflow records to validate the utility of the derived
15 hybrid climate dataset over other existing climate datasets in hydrologic simulations. In this study, the proxy
16 validation was conducted using an existing hydrologic model (Eum et al., 2017), Variable Infiltration
17 Capacity (VIC, Liang et al., 1994), for the Athabasca River basin (ARB). The VIC model was further
18 refined at 1/32° (2~3 km) for a finer spatial resolution and to better simulate the complex river network in

Deleted: confirm

Deleted: superiority

19 the Lower Athabasca River basin. Five of the catchment areas listed in Table 2 were selected for the proxy
20 validation based on three criteria: i) hydrometric record length, ii) location defined by upper, middle and
21 lower reaches (Northern River Basin Study, 2002), and iii) the number of gridded climate datasets used to
22 generate a hybrid climate dataset for the catchment area of the selected hydrometric station. In other words,

Deleted: at 1/32° (2~3 km)

Deleted: a

Deleted: b

Deleted: c

Deleted: The

23 a higher number of gridded climate datasets contributing to the hybrid climate dataset within a catchment
24 was selected to evaluate the utility of the hybrid climate data relative to the existing gridded climate datasets,

Deleted: was optimized for each catchment area

Deleted: maximize the evaluation of the hydrologic simulations using

Deleted: multiple

25 Hinton is located near the headwaters of ARB, which are characterized by mountainous topography and

Deleted:

1 snow- and glacier-ice melt dominated hydrologic regimes. Pembina is one of the major rivers in the middle
2 reach. The other three stations (Christina, Clearwater above Christina and Firebag) are located in the lower
3 reach, which is a water-limited (dry) region due to a higher amount of evapotranspiration (Eum et al.,
4 2014b). The sub-basins of Hinton, Firebag, and Clearwater include a partial area outside of the Township
5 data domain, thus inducing a higher or lower number of climate datasets in the derived hybrid dataset.
6 A total of seven climate datasets (five individual and two hybrid climate datasets from the R_{ind} and R_{mul}) are
7 available to calibrate the VIC hydrologic model parameter set related to soil properties and routing. The
8 calibration period is 1985-1997 as in Eum et al., (2017), except for CaPA that uses the period of 2003-2009
9 for calibration, as CaPA covers the period from 2002 to 2016. The remaining period of total record length
10 for each climate dataset is used for validation. More details on calibration can be found in Eum et al. (2017).
11 Under the assumption of REFRES that all of the existing climate datasets are of equal quality for hydrologic
12 simulations, all of the calibrated parameter sets can be considered as mostly plausible parameter sets for
13 the selected sub-basins. However, as mentioned above, intrinsic biases exist temporally and spatially in all
14 of the gridded climate datasets, e.g., discrepancies in the amount and spatial distribution of precipitation
15 between the gridded climate datasets and observations. Therefore, the similarity of the gridded climate
16 datasets in terms of magnitude, sequence, and spatial distribution of climate events relative to observations
17 is crucial to reproduce historically observed streamflows. In addition to climate forcings, streamflows are
18 mainly affected by geographic characteristics and physical land surface processes (e.g., infiltration and
19 evapotranspiration), which are represented by model parametrization related to infiltration and soil
20 properties (Demaria et al., 2007). In a hydrologic simulation, the biases in climate datasets can be
21 compromised by model parameters that adjust hydrologic processes to observations (Harpold et al., 2017;
22 Kirchner, 2006). That is, a calibrated parameter set may imply biases in a climate dataset. Under the
23 assumption that the calibrated parameter sets are suitable for hydrologic simulations in each sub-basin, this
24 study applied a multiset-parameter hydrologic simulation approach that employs all parameter sets
25 calibrated by the seven climate datasets and the same climate dataset as a forcing input data to assess the

Deleted: , and

Deleted: number

Deleted: have an

Deleted: As

Deleted: , however

Deleted: characteristics

Deleted: the similarity of the gridded

Deleted: in terms of magnitude, sequence

Deleted: spatial distribution of

Deleted: events relative

Deleted: observations is crucial to reproduce historical observed streamflows. Therefore, the variability of hydrologic simulations may signify the reliability

1 sensitivity of the climate dataset to all feasible parameter sets. From the multiset-parameter hydrologic
 2 simulations, the bias in a climate dataset can be estimated indirectly by quantifying the variability in
 3 hydrologic simulations derived from the feasible calibrated parameter sets under a climate forcing dataset.
 4 In other words, lower variability in the hydrologic simulations indicates higher reliability in the climate
 5 forcing dataset. The suitability of the hybrid climate dataset for improving historical hydrologic simulations
 6 was also tested by directly comparing the performances of calibration and validation for each climate
 7 dataset. Proxy validations were carried out by conducting 49 hydrologic simulations (7 climate forcing × 7
 8 parameter sets) for the Pembina and Christina catchment areas, whereas only 36 simulation runs were
 9 possible for Hinton, Firebag, and Clearwater sub-basins, as one of the gridded data sets (i.e., Township) did
 10 not cover the entire catchment areas of these three hydrometric stations.

Deleted: used to produce

Deleted: multiple calibrated parameter sets, i.e., the lower variability indicates higher reliability. For instance, all of hydrologic simulations from the calibrated parameter sets perform well for a

Deleted: and show low

Deleted: , it indicates that the

Deleted: may provide reliable historical

Deleted: does

12 4. Results

Deleted: ¶
3

13 4.1 Precipitation performance measures in Alberta

Deleted: 3

14 Although the performance measures were calculated for 190 AHCCD stations in western Canada, the
 15 target area of this study is in Alberta, where only 45 stations are located. Therefore, the results for the 45
 16 AHCCD stations are given in this study. Table 3 shows spatially-averaged performance measures for
 17 precipitation. The Township data outperformed other climate datasets for all performance measures except
 18 P_{bias} . ANUSPLIN is the second best climate dataset for Alberta. All climate datasets underestimate the
 19 standard deviation of observed daily precipitation (i.e., negative σ_{ratio}), especially PNWNAmet and CaPA
 20 which underestimated by 34 % and 39 %, respectively. Interestingly, two station-based gridded climate
 21 datasets, ANUSPLIN and Township, show negative P_{bias} while PNWNAmet, CaPA, and NARR datasets
 22 have positive P_{bias} . This indicates that ANUSPLIN and Township may underestimate extreme precipitation,
 23 as they employed the raw station data instead of the adjusted precipitation data which is higher than the raw
 24 station data by 5%-20%. In contrast, other climate datasets (especially multiple sources and reanalysis data)
 25 overestimate extreme precipitation. These results are consistent with findings in Eum et al. (2014a) that

Deleted:

Deleted: underestimate

Deleted: and

1 CaPA and NARR overestimate extreme precipitation events by overly reflecting the orographic effects on
2 precipitation in western Alberta.

Deleted: western

3 Figure 4 shows the temporal correlation coefficient (TCC) data averaged over the AHCCD stations in
4 Alberta to investigate the similarity between historical precipitation datasets employed in this study. As
5 expected, station-based climate datasets (i.e., ANUSPLIN, PNWNAmet, and Township) showed better
6 TCCs than CaPA and NARR. The TCC between ANUSPLIN and Township was the highest among climate
7 datasets except for the observations (i.e., OBS), even though they incorporated different interpolation
8 techniques. PNWNAmet showed the highest TCC with ANUSPLIN because they both are based on thin
9 plate spline interpolation. TCCs between CaPA and other climate datasets are similar, as CaPA is produced
10 from multiple sources such as GEM's outputs and weather radar networks of Canada and US. NARR, the
11 reanalysis-based climate dataset, showed higher TCC with CaPA than with other datasets, as it is assimilated
12 with multiple sources of observations.

Deleted:)

13 Maps of each performance measure are shown in Figure 5. It is evident from the spatial variability that
14 the ANUSPLIN and Township datasets outperformed the other datasets in D_{KS} throughout Alberta. In the
15 mountainous region of southwest Alberta, most of the climate datasets performed poorly in P_{bias} , σ_{ratio} ,
16 RMSE, and PCC, resulting mainly from the sparse observation network and inconsistent observations near
17 the Canada-US border. PNWNAmet highly overestimates the mean annual precipitation in the mountainous
18 area (e.g., 300 mm/year higher than that observed at station ID 3050519), which may considerably affect
19 simulated streamflows originating in mountainous headwaters and further downstream.

Deleted: mountains

Deleted: >

Deleted: for

21 **4.2 Air temperature performance measures in Alberta**

Deleted: 3

22 The performance measures for air temperature averaged over 37 AHCCD stations in Alberta are
23 presented in Table 4. As CaPA provides only precipitation, it was excluded in the assessment for temperature.

24 All of the performance measures for temperature are better than those for precipitation except P_{bias} . NARR
25 is highly biased as it underestimates minimum and maximum temperatures, which might be an attribute of

Deleted: Performance

1 discontinuation of observation assimilation since 2003 (Eum et al., 2014a). ANUSPLIN and Township
2 showed an almost perfect linear relationship (TCC) with the observations (i.e., > 0.97 for all of the climate
3 datasets). The performance measures for maximum temperature are better than those for minimum
4 temperature as maximum temperature is dominated by mainly large-scale heat waves while minimum
5 temperature is affected by local physical processes, e.g., topography and surface conditions (Eum et al.,
6 2012). NARR showed less skill in capturing these local effects due to the coarse spatial resolution (~32km)
7 compared to other station-based climate datasets. As with precipitation, the maps of performance measures
8 for minimum and maximum temperature presented in Figure 6 and Figure 7 showed that data from the
9 mountainous areas performed poorly in most of the performance measures. NARR showed positive and
10 negative P_{bias} for minimum and maximum temperature, respectively, in the mountainous region, indicating
11 that NARR has a warm bias in extreme cold temperatures and a cold bias in extreme warm temperatures.

Deleted: temperature while

Deleted: temperature

13 **4.3 Ranking of climate datasets in the ARB**

Deleted: 3

14 The geospatial information (i.e., latitude, longitude, and elevation) of 22,372 grid cells within the ARB
15 was extracted from the Canadian digital elevation data provided by Natural Resources Canada (refer to
16 <https://open.canada.ca/data/dataset/7f245e4d-76c2-4caa-951a-45d1d2051333>). Using this information, the
17 RM in REFRES ranked the five climate datasets by TOPSIS for each grid cell. Table 5 presents the first-
18 ranked number of grid cells and their percentage for each climate dataset according to the performance
19 measures of individual variables (Case A and Case B) and multi-variables (Case C), i.e., precipitation and
20 (minimum and maximum) temperature in this study.

Deleted: the

21 For precipitation, the Alberta township dataset was ranked first in most of the grid cells within the
22 basin (78%) for the whole ARB, followed by ANUSPLIN (13%), PNWNAmet (3%), CaPA (3%), and
23 NARR (2%). However, the Township data domain covers only 83% of the ARB within Alberta; the
24 remaining 17% of the watershed area that lies on the outside the province is not covered (Figure 8). The
25 Township dataset was ranked first for almost 95% of grid cells within its domain, indicating that the

Deleted: %),

1 Township dataset overwhelmingly outperformed other climate datasets for precipitation. Township was
2 dominantly ranked first for the subbasins (Pembina and Christina) within the Township domain.

3 For temperature, ANUSPLIN was ranked first (in 62% grid cells) for the whole ARB, followed by
4 Township (31%) and PNWNAmets (7%). In the upper and middle reaches, i.e., Hinton and Pembina,
5 PNWNAmets and Township were mostly ranked first, respectively, while ANUSPLIN outperformed other
6 climate datasets for the subbasins in the lower reach. When considering the performance measures for
7 multiple variables simultaneously, the Township dataset was ranked first, followed by ANUSPLIN for 64%
8 and 36% of the grid cells for the whole ARB. Figure 9 shows maps of the first-ranked climate datasets for
9 each case in Table 5, i.e., individual variable (Case A and B) and multi-variables (Case C). Due to the
10 limited spatial coverage of the Township dataset, other climate datasets were ranked first in the headwaters
11 of the ARB and the area of the river basin in Saskatchewan. For instance, ANUSPLIN and PNWNAmets
12 were ranked first in the headwaters, while no specific climate dataset dominated in Saskatchewan for
13 precipitation (refer to Figure 9A). For temperature, ANUSPLIN outperformed in the northern part (middle
14 and lower reaches of the ARB) due to outstanding performance of the P_{bias} performance measure for
15 minimum temperature as shown in Table 4 and Figure 6(b). For multi-variables, Township was mostly
16 ranked first within its domain and ANUSPLIN was ranked first outside the Township dataset domain and
17 also for a small part of lower reach area in the ARB.

18 Figure 10 shows the percentage of each climate dataset at each rank for the three cases (e.g. A, B, and
19 C in Table 5). For precipitation (Case A), Township overwhelmed other climate datasets. The second
20 alternative was ANUSPLIN in the majority of grid cells in the ARB. PNWNAmets, NARR and CaPA were
21 mostly ranked 3rd, 4th and 5th, respectively. For temperature (Case B), ANUSPLIN was ranked mostly first
22 and Township was a distinct second choice in the majority of grid cells, followed by PNWNAmets and
23 NARR. For multi-variables (Case C), Township and ANUSPLIN were the first and second choices in the
24 majority of grid cells in the ARB, respectively.

Deleted:),

Deleted: and the

Deleted: on the contrary,

Deleted: Multi

Deleted: choice

1 As two different hybrid climate datasets were generated using the ranking information from single-
2 and multi-variable approaches, i.e., Hybrid (R_{ind}) and Hybrid (R_{mul}), further investigation is required to
3 identify which hybrid climate dataset may provide better performance and consequently will be
4 recommended for future climate-related studies. A proxy validation approach was applied using both
5 generated hybrid climate datasets to validate the utility of one dataset over the other.

Deleted: superior

Deleted: the

Deleted: determine

Deleted: superiority

7 4.4 Proxy validation of generated hybrid climate datasets

Deleted: 3

8 In addition to the five gridded climate datasets, the two hybrid climate datasets were implemented for
9 proxy validation using the VIC model. In contrast to the station-based climate datasets, both CaPA and
10 NARR were produced from climate models and multiple sources of observations, consequently showing a
11 higher correlation with each other as shown in Figure 4. Since CaPA also provides only precipitation, this
12 study combined precipitation of CaPA with the NARR temperature to prepare the CaPA climate forcing
13 dataset for the proxy validation. Table 6 presents the Nash-Sutcliffe Efficiency (NSE) for the calibration
14 and validation periods at the selected hydrometric stations (Hinton, Pembina, Christina, Clearwater, and
15 Firebag) in the ARB to assess the suitability of each climate dataset as a climate forcing input data for

Deleted: Since CaPA

Deleted: datasets

16 hydrologic simulations. Over the five hydrometric stations, most of the climate datasets performed well
17 with the exception of NARR in the Pembina catchment. Most of NSE values in calibration for Christina
18 and Firebag were above 0.50, which is the threshold of satisfactory performance in hydrologic models as
19 suggested by Moriasi et al. (2007). However, model performance is not satisfactory but acceptable for
20 Christina and Firebag during the validation period. The two hybrid climate datasets performed well, with

Deleted: slightly

21 comparably good and better NSE values than other climate datasets, especially at Pembina, Clearwater, and
22 Firebag, located in the middle and lower reaches. Figure 11 presents the boxplots of NSEs obtained through
23 the multiset-parameter VIC simulations. The NSE ranges were obtained from multiple VIC simulations,
24 with each climate dataset used as climate forcing for all the plausible model parameter sets, which were
25 calibrated with seven climate datasets, individually. The values above each boxplot represent the averaged

Deleted: (

Deleted:)

1 value of the NSEs over the multiset-parameter hydrologic simulations. A narrower range of NSE values
2 represents a higher precision for a climate dataset and a higher averaged NSE value means higher accuracy.
3 Therefore, a climate dataset showing both a higher averaged NSE and a narrow range of NSEs indicates
4 that it is a relatively more appropriate and reliable climate forcing dataset for hydrologic simulations.

Deleted: multiple

5 At Hinton, all of the climate datasets showed satisfactory NSE values for accuracy, while ANUSPLIN,
6 Hybrid(R_{ind}), and Hybrid(R_{mul}) showed better precision. The validation period of CaPA is only six years
7 from 2010 to 2016, as CaPA data are only available between 2002 to 2016. This might be a reason why
8 CaPA produced the highest NSE (accuracy) among the climate datasets used in this study. Therefore, the
9 results of CaPA need to be considered carefully otherwise they might be misleading. In this context, the
10 CaPA dataset was excluded from further assessment of the precision and accuracy even though all of the

Deleted: a

Deleted: , which

Deleted: that

11 results of CaPA were included in Figure 11 for reference only. Hybrid(R_{mul}) and ANUSPLIN showed the
12 highest accuracy as forcing data, followed by Hybrid(R_{ind}), PNWNAmet, and NARR. In the Pembina and
13 Christina catchments, the Hybrid(R_{ind}), Hybrid(R_{mul}), and Township datasets had the highest precision and
14 accuracy. NARR produced negative NSEs at Pembina, indicating it is not reliable or suitable as a forcing

Deleted: in

Deleted: and

15 dataset. For Clearwater, Hybrid(R_{ind}) is the top performer, followed by Hybrid(R_{mul}), ANUSPLIN,
16 PNWNAmet, and NARR. Clearwater had the highest number of climate datasets combined in the hybrid
17 climate dataset within the basin for precipitation as shown in Figure 9. Interestingly, the precision of NARR
18 is similar to that of CaPA because they shared the temperature data from NARR. For Firebag, Hybrid(R_{ind})

Deleted: has

Deleted: used

Deleted: same

Deleted: the

19 also showed top performance in both precision and accuracy, followed by Hybrid(R_{mul}), ANUSPLIN,
20 PNWNAmet, and NARR. Overall, Hybrid(R_{ind}) showed the best accuracy and precision at all hydrometric
21 stations, indicating that it has the potential not only to improve historical hydrologic simulations but also
22 to be used as reference data for statistical downscaling of climate change projections in the province.

Deleted: a

23
24
25

Deleted: 4

1 **5. Discussion**

2 Among the station-based gridded climate datasets, the Township dataset outperformed other station-
3 based gridded climate datasets. As PNWNAmet set a common period from 1945 to 2012 for all stations
4 included in the interpolation, many stations might be left out in the data generation processes. While
5 ANUSPLIN used the Canada-wide archive (raw) station data collected by only ECCC, the Alberta
6 Township data has been produced on the basis of the archive (raw) station data collected by ECCC, AEP,
7 and AF over Alberta. Therefore, one of the possible reason for outperformance of Township dataset might
8 be the difference in the numbers of stations (i.e. station density) employed to produce the gridded climate
9 datasets. In addition, PNWNAmet showed a positive P_{bias} for precipitation, especially in the mountainous
10 areas, while ANUSPLIN, which employs similar thin plate spline interpolation, generated negative P_{bias} .
11 PNWNAmet overestimated precipitation over the mountainous area, which considerably affects simulated
12 low flows at Hinton in the ARB. Figure 12 shows the observed and simulated hydrographs from gridded
13 climate datasets at (a) Hinton and (b) Pembina. It clearly shows that PNWNAmet highly overestimated the
14 low and high, which is caused by overestimated precipitation in the drainage area of the sub-basins. As with
15 PNWNAmet, NARR also overestimated the low and high flows, which is induced by the combined effects
16 of overestimating precipitation and warm biases in cold temperature. The temperature bias of NARR is thus
17 further confirmed and is consistent with the earlier finding of Eum et al., (2014) and Islam and Dery (2016).

18 In Figure 12, the hybrid climate datasets underestimated the peak flows (in 2009, 2010, 2014, and
19 2015) at Hinton, and hydrograph is similar to the hydrograph produced by ANUSPLIN data set that
20 dominantly ranked first in this watershed. On the contrary, the hydrograph of the hybrid climate datasets at
21 Pembina is similar to that of Township that is dominantly ranked first in Pembina (refer to Table 5). These
22 results indicate that the hybrid climate dataset has the intrinsic limitation that the performance of the hybrid
23 dataset for a basin may closely resemble that of the climate dataset that is dominantly ranked first for the
24 basin. However, the utility of the hybrid climate dataset can be clearly found at a whole-basin scale for a

Deleted: Among the station-based gridded climate datasets

Deleted: enormously

Deleted: six

Deleted: (excluding Township)

Deleted: , the headwater region in the ARB.

Deleted: flows during winter

Deleted: Hinton.

Deleted: overestimates

Deleted: during winter

1 large watershed, as the added values of the hybrid climate dataset in sub-basins can be cumulated to the
2 main stem at the downstream in the watershed.

3 Among the station-based gridded climate datasets, ANUSPLIN and Township employed a different
4 number of stations depending on their periods of record. Therefore, there is an inconsistency in these climate
5 datasets over time. For example, the Township dataset employed only 300~400 stations in the 1960s, but
6 has increased to 400~500 since 1970. A change-point analysis of these datasets may provide some useful
7 information to end-users with respect to when and where changes occurred, which will help in establishing
8 spatial and temporal accuracies of these datasets (Eum et al. 2014a). Further, PNWNAmet employed the
9 same number of stations over time to avoid the above mentioned inconsistency, but this study found that it
10 induced overestimation of precipitation in data-poor regions such as mountainous regions in Alberta. As
11 the hybrid climate datasets are generated from the multiple historical gridded datasets, they may also have
12 the same inconsistencies identified in other datasets. The proxy validation, however, demonstrated that the
13 generated hybrid climate datasets can improve the performance of hydrologic simulations.

14 This study identified the preference order of all gridded climate datasets based on the performance
15 measures evaluated at the AHCCD stations, therefore the ranking somewhat relies on the spatial distribution
16 of the AHCCD stations. As shown in Figure 1, the density of AHCCD stations varies across western Canada,
17 and is low in the cold climates of mountainous and northern areas. Therefore, the ranking could further be
18 improved with a more uniform density of AHCCD stations over western Canada.

19 Literature has demonstrated that NARR, a reanalysis-based climate dataset, can be an alternative as a
20 climate forcing dataset for hydrologic simulations in data sparse regions (Choi et al., 2009; Praskievicz and
21 Bartlein, 2014; Islam and Dery, 2016). In this study, the NARR dataset performed quite well in high-
22 elevation regions (Hinton in this study) while it did not perform so well in the middle and lower reaches,
23 i.e., lower-elevation watersheds. NARR performed especially poorly in the Pembina sub-basin, a region
24 where hydrologic simulations are highly sensitive to model parameters (Eum et al., 2014b). In Figure 11
25 (b), however, the NARR parameter set produced fair NSE values in hydrologic simulations forced by the

Deleted: while

Deleted: it

Deleted: issue

Deleted: and

Deleted: climate

Deleted: Much literature

Deleted: Especially,

Deleted: other than climate forcings (Eum et al., 2014a).

1 other climate datasets except for CaPA and PNWNAmet. Such result indicates that 1) all of parameter sets
2 used in this study were calibrated reasonably and 2) climate forcing input data plays a more crucial role in
3 hydrologic simulations as any parameter sets did not produce a fair NSE value from NARR in Pembina.
4 CaPA was more suitable than NARR for the selected sub-basins in this study, which indicates that CaPA
5 might be a better alternative in low station-density regions such as the ARB. However, since the validation
6 period in this study is only 7 years from 2010 to 2016, a longer data period is necessary to validate the
7 suitability of CaPA as indicated in Eum et al. (2014a) and Wong et al. (2017).

Formatted: Font: 12 pt, Font color: Auto

Deleted: Since

Deleted: however,

Deleted: .(

Deleted: .(

8 In the proxy validation, Hybrid(R_{ind}) performed well in the Clearwater sub-basin where the highest
9 number of climate datasets were combined in the generated hybrid climate datasets. The Township dataset,
10 which mostly ranked first within its spatial domain, partially covers the drainage area of Clearwater, so that
11 the generated hybrid climate dataset, Hybrid(R_{ind}), is composed of many climate datasets in this sub-basin.
12 In a traditional approach to hydrological modelling for Clearwater, either the Township dataset might be
13 completely excluded (as it does not cover the entire Clearwater watershed), or potentially, combined with
14 other gridded climate datasets to cover the entire watershed. However, combining different climate datasets
15 to construct the climate forcing for a larger region requires an evaluation of the datasets to identify the order
16 of preference for such aggregation when multiple choices are available. Therefore, this study suggested the
17 REFRES methodology to systematically compare all-available climate datasets for a region to produce a
18 hybrid climate dataset that covers a desired period of record and spatial domain by considering the order of
19 preference for combining various climate datasets at each grid cell. The proxy validation approach also
20 confirmed the utility of a generated hybrid climate dataset over other data sets, especially in hydrologic
21 simulations.

Deleted:)

Deleted: could be

Deleted: superiority

23 6. Summary and concluding remarks

Deleted: 5

24 This study suggested a framework called reference reliability evaluation system (REFRES) to
25 systematically generate a performance-based hybrid climate dataset from multiple climate datasets for a

1 region. The hybrid dataset was found to more reliable for hydrological modelling. The REFRES is
2 composed of three modules; 1) performance measures, 2) ranking, and 3) data generation. The suggested
3 framework was applied to the ARB as a test-bed and generated two hybrid climate datasets from single-
4 (R_{ind}) and multi-variable (R_{mul}) approaches by evaluating the performance of five available gridded climate
5 datasets: station-based gridded climate datasets (i.e. ANUSPLIN, Alberta Township, and PNWNAmets), a
6 multi-source dataset (CaPA), and a reanalysis-based dataset (NARR). A hydrologic modelling-based proxy
7 validation approach was applied to demonstrate the applicability of the hybrid climate dataset generated for
8 the five sub-basins in the ARB. The results showed that

- 9 - Among the five climate datasets, the station-based climate datasets performed better than multi-
10 source- and reanalysis-based datasets. The Township dataset, in particular, outperformed other
11 climate datasets in the selected performance measures over northern Alberta.
- 12 - Most of the climate datasets performed poorly in the mountainous areas of southwest Alberta, due
13 to a sparse observation network, orographic effects, topographic complexity, and inconsistencies in
14 observation between Canada and the US.
- 15 - As a result of REFRES' application for the ARB, the Township and ANUSPLIN datasets are mostly
16 ranked the highest among the five climate datasets for precipitation and temperature, respectively.
- 17 - In the proxy validation, two hybrid climate datasets, Hybrid(R_{ind}) and Hybrid(R_{mul}), performed
18 better in terms of precision and accuracy as forcing data for hydrologic simulations.
- 19 - Hybrid(R_{ind}) especially outperformed other climate datasets in the Clearwater sub-basin where the
20 highest number of climate datasets were combined in generating Hybrid(R_{ind}) for precipitation. This
21 indicates that the hybrid climate dataset generated by REFRES may lead to more reliable
22 hydrologic simulations, resulting in improved hydrologic predictions.

23 This study provided the preference order of climate datasets available in Alberta, which may be useful
24 for modelers and decision-makers as to which climate dataset is the most suitable for their studies and
25 projects. Furthermore, this study demonstrated that the hybrid climate dataset produced by REFRES is more

Deleted:

Deleted:

Deleted: including mountainous areas,

Deleted: inconsistency

Deleted: clearly

Deleted: and thus

1 representative of historical climatic conditions. Therefore, the hybrid climate dataset is recommended to be
2 used as a reference dataset for statistical downscaling and hydrologic model forcing, resulting in more
3 reliable high-resolution climatic and hydrologic projections.

Deleted: the

Deleted: should

5 Acknowledgements

6 The authors would like to thank the Natural Resources Canada, Alberta Agriculture and Forest, the
7 Pacific Climate Impacts Consortium (PCIC), Environment and Climate Change Canada, and
8 NOAA/OAR/ESRL PSD for providing the historical gridded climate datasets.

Deleted: National

10 References

- 11 Asong, Z. E., Khaliq, M. N. and Wheeler, H. S.: Regionalization of precipitation characteristics in the
12 Canadian Prairie Provinces using large-scale atmospheric covariates and geophysical attributes,
13 Stochastic Environmental Research and Risk Assessment, 29(3), 875–892, 2015.
- 14 Choi, W., Kim, S. J., Rasmussen, P. F. and Moore, A. R.: Use of the North American Regional Reanalysis
15 for hydrological modeling in Manitoba, Can. Water Resour. J., 34, 13–36, 2009.
- 16 Christensen, N. S. and Lettenmaier, D. P.: A multimodel ensemble approach to assessment of climate
17 change impacts on the hydrology and water resources of the Colorado River Basin, Hydrology and
18 Earth System Sciences, 11(4), 1417–1434, 2007.
- 19 Côté, J., Desmarais, J.-G., Gravel, S., Méthot, A., Patoine, A., Roch, M. and Staniforth, A.: The
20 operational CMC–MRB global environmental multiscale (GEM) model. Part II: Results, Monthly
21 Weather Review, 126(6), 1397–1418, 1998a.
- 22 Côté, J., Gravel, S., Méthot, A., Patoine, A., Roch, M. and Staniforth, A.: The operational CMC–MRB
23 global environmental multiscale (GEM) model. Part I: Design considerations and formulation,
24 Monthly Weather Review, 126(6), 1373–1395, 1998b.

Deleted: Discussions

- 1 Deacu, D., Fortin, V., Klyszejko, E., Spence, C. and Blanken, P. D.: Predicting the Net Basin Supply to
2 the Great Lakes with a Hydrometeorological Model, *Journal of Hydrometeorology*, 13(6), 1739–
3 1759, doi:10.1175/JHM-D-11-0151.1, 2012.
- 4 Demaria, E.M, Nijssen, B., Wagener, T.: Monte Carlo sensitivity analysis of land surface parameters
5 using the variable infiltration capacity model, *Journal of Geophysical Research*, 112, D11113, 2007
- 6 Dibike, Y., Eum, H.-I. and Prowse, T.: Modelling the Athabasca watershed snow response to a changing
7 climate, *Journal of Hydrology: Regional Studies*, 15, 134–148, doi:10.1016/j.ejrh.2018.01.003,
8 2018.
- 9 Essou, G. R. C., Sabarly, F., Lucas-Picher, P., Brissette, F. and Poulin, A.: Can Precipitation and
10 Temperature from Meteorological Reanalyses Be Used for Hydrological Modeling?, *Journal of*
11 *Hydrometeorology*, 17(7), 1929–1950, doi:10.1175/JHM-D-15-0138.1, 2016.
- 12 Eum, H.-I. and Cannon, A. J.: Intercomparison of projected changes in climate extremes for South Korea:
13 application of trend preserving statistical downscaling methods to the CMIP5 ensemble,
14 *International Journal of Climatology*, 37(8), 3381–3397, doi:10.1002/joc.4924, 2017.
- 15 Eum, H.-I., Dibike, Y. and Prowse, T.: Climate-induced alteration of hydrologic indicators in the
16 Athabasca River Basin, Alberta, Canada, *Journal of Hydrology*, 544, 327–342,
17 doi:10.1016/j.jhydrol.2016.11.034, 2017.
- 18 Eum, H.-I., Dibike, Y. and Prowse, T.: Comparative evaluation of the effects of climate and land-cover
19 changes on hydrologic responses of the Muskeg River, Alberta, Canada, *Journal of Hydrology:*
20 *Regional Studies*, 8, 198–221, doi:10.1016/j.ejrh.2016.10.003, 2016.
- 21 Eum, H.-I., Dibike, Y., Prowse, T. and Bonsal, B.: Inter-comparison of high-resolution gridded climate
22 data sets and their implication on hydrological model simulation over the Athabasca Watershed,
23 Canada, [Hydrological Processes](#), 28(14), 4250–4271, doi:10.1002/hyp.10236, 2014a.

- 1 Eum, H.-I., [Dibike, Y.](#) and Prowse, T.: Uncertainty in modelling the hydrologic responses of a large
2 watershed: a case study of the Athabasca River basin, Canada, *Hydrological Processes*, 28(14),
3 4272–4293, doi:10.1002/hyp.10230, 2014b.
- 4 Eum, H.-I., Gachon, P., Laprise, R. and Ouarda, T.: Evaluation of regional climate model simulations
5 versus gridded observed and regional reanalysis products using a combined weighting scheme,
6 *Climate Dynamics*, 38(7-8), 1433–1457, doi:10.1007/s00382-011-1149-3, 2012.
- 7 Faramarzi, M., Srinivasan, R., Irvani, M., Bladon, K. D., Abbaspour, K. C., Zehnder, A. J. B. and Goss,
8 G. G.: Setting up a hydrological model of Alberta: Data discrimination analyses prior to calibration,
9 *Environmental Modelling & Software*, 74, 48–65, doi:10.1016/j.envsoft.2015.09.006, 2015.
- 10 Garand, L. and Grassotti, C.: Toward an objective analysis of rainfall rate combining observations and
11 short-term forecast model estimates, *Journal of Applied Meteorology*, 34, 1962–1977, 1995.
- 12 Gutmann, E., Pruitt, T., Clark, M. P., Brekke, L., Arnold, J. R., Raff, D. A. and Rasmussen, R. M.: An
13 intercomparison of statistical downscaling methods used for water resource assessments in the
14 United States, *Water Resources Research*, 50(9), 7167–7186, doi:10.1002/2014WR015559, 2014.
- 15 [Harpold, A. A., Kaplan, M. L., Klos, P. Z., Link, T., McNamara, J. P., Rajagopal, S., Schumer, R. and](#)
16 [Steele, C. M.: Rain or snow: hydrologic processes, observations, prediction, and research needs,](#)
17 [Hydrology and Earth System Sciences](#), 21(1), 1–22, doi:10.5194/hess-21-1-2017, 2017.
- 18 Hopkinson, R. F., McKenney, D. W., Milewska, E. J., Hutchinson, M. F., Papadopol, P. and Vincent, L.
19 A.: Impact of Aligning Climatological Day on Gridding Daily Maximum–Minimum Temperature
20 and Precipitation over Canada, *Journal of Applied Meteorology and Climatology*, 50(8), 1654–
21 1665, doi:10.1175/2011JAMC2684.1, 2011.
- 22 Hutchinson, M. F., McKenney, D. W., Lawrence, K., Pedlar, J. H., Hopkinson, R. F., Milewska, E. and
23 Papadopol, P.: Development and Testing of Canada-Wide Interpolated Spatial Models of Daily
24 Minimum–Maximum Temperature and Precipitation for 1961–2003, *Journal of Applied*
25 *Meteorology and Climatology*, 48(4), 725–741, doi:10.1175/2008JAMC1979.1, 2009.

- 1 Hwang, C.L. and Yoon, K.: Multiple attribute decision making: methods and applications. Springer, New
2 York, 1981.
- 3 Islam, S. U. and Déry, S. J.: Evaluating uncertainties in modelling the snow hydrology of the Fraser River
4 Basin, British Columbia, Canada, *Hydrology and Earth System Sciences*, 21(3), 1827–1847,
5 doi:10.5194/hess-21-1827-2017, 2017.
- 6 Jun, K. S., Chung, E.-S., Kim, Y.-G. and Kim, Y.: A fuzzy multi-criteria approach to flood risk
7 vulnerability in South Korea by considering climate change impacts, *Expert Systems with*
8 *Applications*, 40(4), 1003–1013, 2013.
- 9 Kay, A. L., Davies, H. N., Bell, V. A. and Jones, R. G.: Comparison of uncertainty sources for climate
10 change impacts: flood frequency in England, *Climatic Change*, 92(1-2), 41–63,
11 doi:10.1007/s10584-008-9471-4, 2009.
- 12 [Kirchner, J. W.: Getting the right answers for the right reasons: Linking measurements, analyses, and](#)
13 [models to advance the science of hydrology: GETTING THE RIGHT ANSWERS FOR THE](#)
14 [RIGHT REASONS, *Water Resources Research*, 42\(3\), doi:10.1029/2005WR004362, 2006.](#)
- 15 Klyszejko, E. S.: Hydrologic Validation of Real-Time Weather Radar VPR Correction Methods,
16 University of Waterloo., 2007.
- 17 Lee, G., Jun, K.-S. and Chung, E.-S.: Integrated multi-criteria flood vulnerability approach using fuzzy
18 TOPSIS and Delphi technique, *Natural Hazards and Earth System Science*, 13(5), 1293–1312,
19 doi:10.5194/nhess-13-1293-2013, 2013.
- 20 Lespinas, F., Fortin, V., Roy, G., Rasmussen, P. and Stadnyk, T.: Performance Evaluation of the Canadian
21 Precipitation Analysis (CaPA), *Journal of Hydrometeorology*, 16(5), 2045–2064, doi:10.1175/JHM-
22 D-14-0191.1, 2015.
- 23 Liang, X., Lettenmaier, D. P., Wood, E. F. and Burges, S. J.: A simple hydrologically based model of
24 land surface water and energy fluxes for general circulation model, *Journal of Geophysical Research*,
25 99(D7), 14,415–14,428, doi:10.1029/94JD00483, 1994.

1 Mahfouf, J.-F., Brasnett, B. and Gagnon, S.: A Canadian Precipitation Analysis (CaPA) Project:
2 Description and Preliminary Results, *ATMOSPHERE-OCEAN*, 45(1), 1–17,
3 doi:10.3137/ao.v450101, 2007.

4 Mekis, E. and Hogg, W. D.: Rehabilitation and analysis of Canadian daily precipitation time series,
5 *Atmosphere-Ocean*, 37(1), 53–85, doi:10.1080/07055900.1999.9649621, 1999.

6 Mekis, É. and Vincent, L. A.: An Overview of the Second Generation Adjusted Daily Precipitation
7 Dataset for Trend Analysis in Canada, *Atmosphere-Ocean*, 49(2), 163–177,
8 doi:10.1080/07055900.2011.583910, 2011.

9 Mesinger, F., DiMego, G., Kalnay, E., Mitchell, K., Shafran, P. C., Ebisuzaki, W., Jović, D., Woollen, J.,
10 Rogers, E., Berbery, E. H., Ek, M. B., Fan, Y., Grumbine, R., Higgins, W., Li, H., Lin, Y., Manikin,
11 G., Parrish, D. and Shi, W.: North American Regional Reanalysis, *Bulletin of the American
12 Meteorological Society*, 87(3), 343–360, doi:10.1175/BAMS-87-3-343, 2006.

13 [Moriasi, D. N., Arnold, J. G., Van Liew, M. W., Bingner, R. L., Harmel, R. D. and Veith, T. L.: Model
14 evaluation guidelines for systematic quantification of accuracy in watershed simulations,
15 Transactions of the ASABE, 50\(3\), 885–900, 2007.](#)

16 NIDIS, U.S. Drought Portal. NOAA, 2015 [Available online at <http://www.drought.gov>.]

17 Praskievicz, S. and Bartlein, P.: Hydrologic modeling using elevationally adjusted NARR and
18 NARCCAP regional climate-model simulations: Tucannon River, Washington, *Journal of
19 Hydrology*, 517, 803–814, doi:10.1016/j.jhydrol.2014.06.017, 2014.

20 Rinke, A., Marbaix, P. and Dethloff, K.: Internal variability in Arctic regional climate simulations: case
21 study for the SHEBA year, *Climate research*, 27(3), 197–209, doi:doi:10.3354/cr027197, 2004.

22 Shen, S. S., Dzikowski, P., Li, G. and Griffith, D.: Interpolation of 1961–97 daily temperature and
23 precipitation data onto Alberta polygons of ecodistrict and soil landscapes of Canada, *Journal of
24 applied meteorology*, 40(12), 2162–2177, 2001.

Formatted

1 Shen, Y., Xiong, A., Wang, Y. and Xie, P.: Performance of high-resolution satellite precipitation products
2 over China, *Journal of Geophysical Research*, 115(D2), doi:10.1029/2009JD012097, 2010.

3 Shrestha, R. R., Cannon, A. J., Schnorbus, M. A. and Zwiers, F. W.: Projecting future nonstationary
4 extreme streamflow for the Fraser River, Canada, *Climatic Change*, 145(3-4), 289–303,
5 doi:10.1007/s10584-017-2098-6, 2017.

6 Shrestha, R. R., Schnorbus, M. A., Werner, A. T. and Berland, A. J.: Modelling spatial and temporal
7 variability of hydrologic impacts of climate change in the Fraser River basin, British Columbia,
8 Canada, *Hydrological Processes*, 26(12), 1840–1860, 2012.

9 Shook, K. and Pomeroy, J.: Changes in the hydrological character of rainfall on the Canadian prairies,
10 *Hydrological Processes*, 26(12), 1752–1766, 2012.

11 Vincent, L. A., Wang, X. L., Milewska, E. J., Wan, H., Yang, F. and Swail, V.: A second generation of
12 homogenized Canadian monthly surface air temperature for climate trend analysis:
13 HOMOGENIZED CANADIAN TEMPERATURE, *Journal of Geophysical Research:*
14 *Atmospheres*, 117(D18), n/a–n/a, doi:10.1029/2012JD017859, 2012.

15 Wang, X. L. and Lin, A.: An algorithm for integrating satellite precipitation estimates with in situ
16 precipitation data on a pentad time scale: BLENDED PENTAD PRECIPITATION DATA, *Journal*
17 *of Geophysical Research: Atmospheres*, 120(9), 3728–3744, doi:10.1002/2014JD022788, 2015.

18 Werner, A. T. and Cannon, A. J.: Hydrologic extremes – an intercomparison of multiple gridded
19 statistical downscaling methods, *Hydrology and Earth System Sciences*, 20(4), 1483–1508,
20 doi:10.5194/hess-20-1483-2016, 2016.

21 Werner, A., Schnorbus, M., Shrestha, R., Cannon, A., Zwiers, F., Dayon, G., and Anslow, F.: A long-
22 term, temporally consistent, gridded daily meteorological dataset for northwestern North America,
23 *Scientific Data*, 6, 180299, 2019.

1 Wong, J. S., Razavi, S., Bonsal, B. R., Wheeler, H. S. and Asong, Z. E.: Inter-comparison of daily
2 precipitation products for large-scale hydro-climatic applications over Canada, *Hydrology and*
3 *Earth System Sciences*, 21(4), 2163–2185, doi:10.5194/hess-21-2163-2017, 2017.

4 Zhao, K.: Validation of the Canadian Precipitation Analysis (CaPA) for hydrological modelling in the
5 Canadian Prairies, University of Manitoba (Canada)., 2013.

6
7
8
9
10

1 **Table 1. High-resolution gridded historical climate datasets used in this study**

| Dataset | Full name | Variable | Type | Period | Resolution | Domain | Institution |
|----------|---|----------------------------------|-----------------------|-----------|-----------------------------|--|--|
| ANUSPLIN | Australia National University Spline | PRCP, TMX, TMN | Station-based | 1950-2015 | 10 km, Daily | Canada | Natural Resource Canada (NRCan) |
| Township | Alberta Township | PRCP, TMX, TMN, Tave, WS, RH, SR | Station-based | 1961-2016 | 10km, Daily | Alberta | Alberta Agriculture and Forestry |
| PNWNAmet | PCIC NorthWest North America meteorological dataset | PRCP, TMX, TMN, WS | Station-based | 1945-2012 | 1/16 degree (6~7 km), Daily | Western Canada (BC, AB, SK) and Alaska | Pacific Climate Impacts Consortium |
| CaPA | Canadian Precipitation Analysis | PRCP | Multiple source-based | 2002-2017 | 10 km, 6-hr | North America | Canadian Meteorological Centre |
| NARR | North American Regional Reanalysis | PRCP, Tair, WS, RH, SR, GH, etc* | Reanalysis-based | 1979-2017 | 32km, 3-hr | North America | National Oceanic and Atmospheric Administration (NOAA) |

2 PRCP: precipitation, TMX: maximum temperature, TMN: minimum temperature, Tave: average
 3 temperature, Tair: air temperature, WS: wind speed, RH: relative humidity, SR: solar radiation, GH:
 4 Geopotential Height

5 *: Refer to <https://www.esrl.noaa.gov/psd/data/gridded/data.narr.monolevel.html> for details

6
7

1

Table 2. Characteristics of hydrometric stations selected in this study

| Station name | Station ID | Record length | Drainage (km ²) | Reach |
|-------------------------------|------------------|---------------|-----------------------------|--------|
| Hinton | 07AD002 | 1961-2016 | 9,760 | Upper |
| Pembina | 07BC002 | 1957-2016 | 13,100 | Middle |
| Christina | S29 (07CE002) | 1982-2016 | 4,836 | Lower |
| Clearwater above Christina | S42 (07CD005) | 1966-2016 | 18,061 | Lower |
| Firebag | S27 (07DC001) | 1971-2016 | 5,980 | Lower |

2

3

Table 3. Performance measures averaged over AHCCD stations in Alberta for precipitation

| Performance measure | Climate Dataset | | | | |
|-------------------------|-----------------|----------|-------|-------|----------|
| | ANUSPLIN | PNWNAmet | CaPA | NARR | Township |
| D _{KS} | 0.09 | 0.62 | 0.60 | 0.42 | 0.09 |
| σ_{ratio} | -0.17 | -0.34 | -0.39 | -0.28 | -0.03 |
| P _{bias} | -7.05 | 5.80 | 3.02 | 2.43 | -6.73 |
| RMSE | 2.02 | 2.50 | 2.59 | 3.53 | 1.07 |
| TCC | 0.87 | 0.81 | 0.77 | 0.53 | 0.95 |
| PCC | 0.87 | 0.80 | 0.73 | 0.74 | 0.93 |

4

5

1 Table 4. Performance measures averaged over the AHCCD stations in Alberta for minimum and
 2 maximum temperature

| Performance measure | Climate Dataset | | | | | | | |
|---------------------|-----------------|-------|----------|-------|---------|--------|----------|-------|
| | ANUSPLIN | | PNWNAmet | | NARR | | Township | |
| | Tmin | Tmax | Tmin | Tmax | Tmin | Tmax | Tmin | Tmax |
| D_{KS} | 0.03 | 0.02 | 0.05 | 0.04 | 0.12 | 0.08 | 0.03 | 0.02 |
| σ_{ratio} | -0.01 | -0.01 | -0.03 | -0.03 | -0.03 | -0.03 | -0.01 | -0.02 |
| P_{bias} | -0.43 | -0.28 | 22.90 | -3.89 | -306.52 | -14.09 | 7.33 | -0.86 |
| RMSE | 1.48 | 1.25 | 1.97 | 1.82 | 4.40 | 3.47 | 1.31 | 0.97 |
| TCC | 0.99 | 0.99 | 0.98 | 0.99 | 0.96 | 0.97 | 0.99 | 0.99 |
| PCC | 0.91 | 0.98 | 0.87 | 0.95 | 0.71 | 0.78 | 0.93 | 0.98 |

3

4

1 Table 5. First ranked number of grid cells in the five sub-basins and the whole Athabasca River Basin
 2 (ARB) and their percentages for each climate dataset, considering the performance measures of individual
 3 (Case A and Case B) and multi-variables (Case C, i.e., precipitation and temperature in this study). Total
 4 number of grid cells is 22,372 at 1/32° (2~3 km)

Deleted: percentage

| Criteria | Basin | Climate dataset | | | | |
|------------------------|--|-----------------|----------------|---------------|--------------|-------------|
| | | ANUSPLIN | Township | PNWNAmet | NARR | CaPA |
| (A) Precipitation | ARB | 2985 (13%) | 17515 (78%) | 691 (3%) | 499 (2%) | 682 (3%) |
| | Hinton | 1271 (91%) | 126 (9%) | 0 (0%) | 0 (0%) | 0 (0%) |
| | Pembina | 0 (0%) | 1791 (100%) | 0 (0%) | 0 (0%) | 0 (0%) |
| | Christina | 0 (0%) | 658 (99.5%) | 3 (0.5%) | 0 (0%) | 0 (0%) |
| | Clearwater | 1474 (56%) | 252 (9.6%) | 10 (0.4%) | 682 (26%) | 215 (8%) |
| | Firebag | 129 (14%) | 750 (79%) | 9 (1%) | 0 (0%) | 64 (6%) |
| | (B) Temperature (Min & Max Temp.) | ARB | 13809 (62%) | 6924 (31%) | 1639 (7%) | 0 (0%) |
| | Hinton | 63 (5%) | 77 (6%) | 1257 (89%) | 0 (0%) | = |
| | Pembina | 486 (27%) | 1305 (73%) | 0 (0%) | 0 (0%) | |
| | Christina | 492 (74%) | 169 (26%) | 0 (0%) | 0 (0%) | = |
| | Clearwater | 2593 (98%) | 40 (2%) | 0 (0%) | 0 (0%) | = |
| | Firebag | 924 (97%) | 28 (3%) | 0 (0%) | 0 (0%) | = |
| (C) Multi-variables | ARB | 8049 (36%) | 14323 (64%) | 0 (0%) | 0 (0%) | - |

Formatted: Line spacing: single

Inserted Cells

Formatted: Line spacing: single

Formatted Table

Formatted: Line spacing: single

Merged Cells

Formatted: Line spacing: single

Formatted: Line spacing: single

Formatted: Line spacing: single

Inserted Cells

Formatted: Line spacing: single

Inserted Cells

| | | | | | |
|-------------------|-----------------------------|------------------------------|-------------------------|-------------------------|---|
| <u>Hinton</u> | <u>1271</u> <u>(91%)</u> | <u>126</u> <u>(9%)</u> | <u>0</u> <u>(0%)</u> | <u>0</u> <u>(0%)</u> | = |
| <u>Pembina</u> | <u>0</u> <u>(0%)</u> | <u>1791</u> <u>(100%)</u> | <u>0</u> <u>(0%)</u> | <u>0</u> <u>(0%)</u> | = |
| <u>Christina</u> | <u>109</u> <u>(16%)</u> | <u>552</u> <u>(84%)</u> | <u>0</u> <u>(0%)</u> | <u>0</u> <u>(0%)</u> | = |
| <u>Clearwater</u> | <u>2574</u> <u>(98%)</u> | <u>59</u> <u>(2%)</u> | <u>0</u> <u>(0%)</u> | <u>0</u> <u>(0%)</u> | = |
| <u>Firebag</u> | <u>536</u> <u>(56%)</u> | <u>416</u> <u>(44%)</u> | <u>0</u> <u>(0%)</u> | <u>0</u> <u>(0%)</u> | = |

1

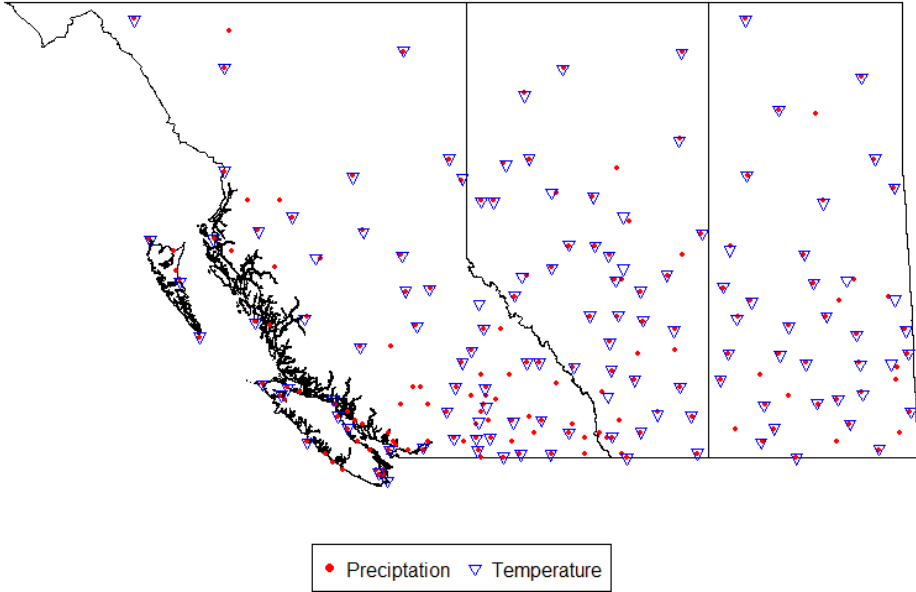
2

1 Table 6. Nash-Sutcliffe Efficiency (NSE) for the calibration and validation periods at five sub-basins in
 2 ARB for the climate datasets investigated in this study

| Climate forcing | Hinton | | Pembina | | Christina | | Clearwater | | Firebag | |
|----------------------|--------|------|---------|-------|-----------|------|------------|------|---------|------|
| | Cal. | Val. | Cal. | Val. | Cal. | Val. | Cal. | Val. | Cal. | Val. |
| ANU | 0.88 | 0.83 | 0.61 | 0.64 | 0.52 | 0.46 | 0.76 | 0.54 | 0.61 | 0.49 |
| SPLIN | | | | | | | | | | |
| Township | - | - | 0.62 | 0.66 | 0.54 | 0.49 | - | - | - | - |
| PNWNA met | 0.82 | 0.81 | 0.53 | 0.54 | 0.40 | 0.35 | 0.73 | 0.59 | 0.65 | 0.48 |
| CaPA | 0.89 | 0.90 | 0.53 | 0.61 | 0.55 | 0.44 | 0.74 | 0.74 | 0.51 | 0.53 |
| NARR | 0.84 | 0.79 | 0.50 | -0.14 | 0.39 | 0.34 | 0.75 | 0.42 | 0.44 | 0.32 |
| Hybrid (R_{ind}) | 0.82 | 0.78 | 0.61 | 0.66 | 0.55 | 0.49 | 0.78 | 0.67 | 0.60 | 0.52 |
| Hybrid (R_{mul}) | 0.89 | 0.83 | 0.61 | 0.65 | 0.54 | 0.48 | 0.77 | 0.53 | 0.59 | 0.47 |

3

4



1
2
3
4

Figure 1. AHCCD stations within the BC, AB, and SK provinces

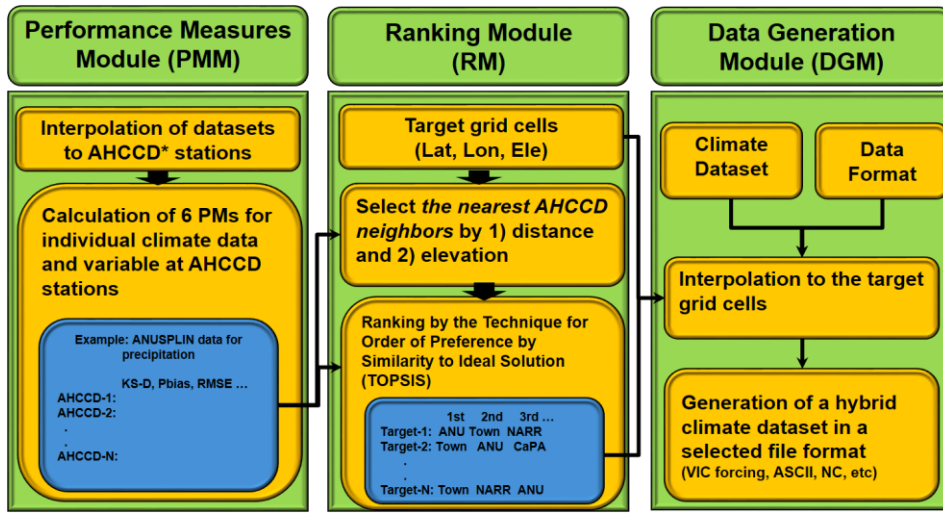


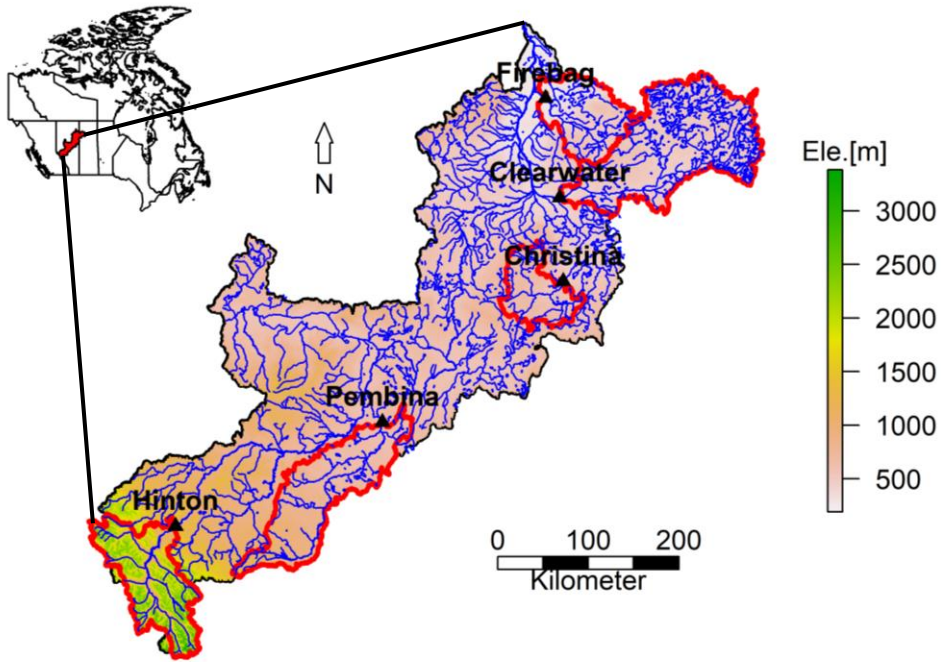
Figure 2. Structure of REFRES comprised of three modules: 1) Performance Measure Module (PMM), 2) Ranking Module (RM), and 3) Data Generation Module (DGM)

Moved (insertion) [2]

Formatted: English (United States)

Formatted: Font: Not Bold

Formatted: Centered



1

2 Figure 3. Geographical information on the five sub-basins (red line) selected in the Athabasca River basin

3

for the proxy validation

4

Moved up [1]: Figure 2.

Formatted: English (Canada)

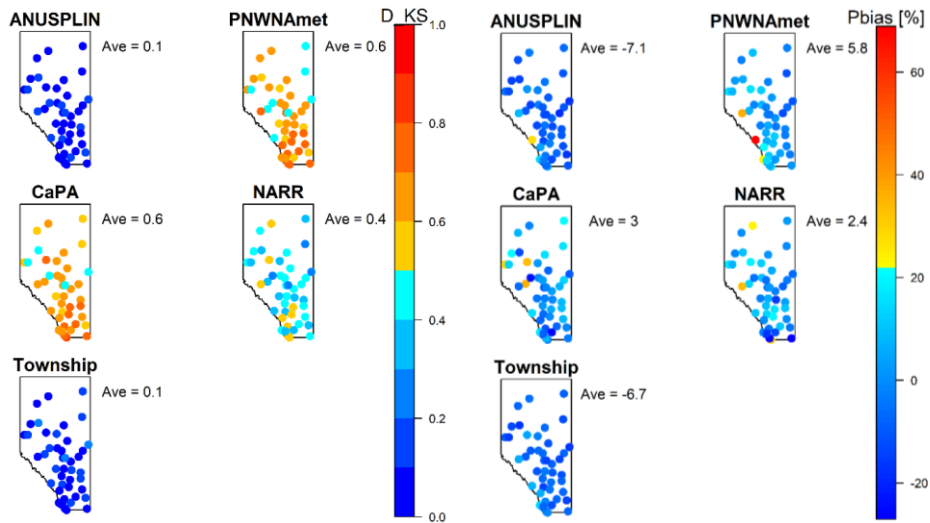
Deleted: Five

| | OBS* | ANUSPLIN | PNWNAme | CaPA | NARR | Township |
|----------|------|----------|---------|------|------|----------|
| OBS* | 1 | 0.87 | 0.81 | 0.77 | 0.53 | 0.95 |
| ANUSPLIN | 0.87 | 1 | 0.84 | 0.81 | 0.61 | 0.86 |
| PNWNAmet | 0.81 | 0.84 | 1 | 0.81 | 0.65 | 0.78 |
| CaPA | 0.77 | 0.81 | 0.81 | 1 | 0.76 | 0.81 |
| NARR | 0.53 | 0.61 | 0.65 | 0.76 | 1 | 0.55 |
| Township | 0.95 | 0.86 | 0.78 | 0.81 | 0.55 | 1 |

1
2
3
4

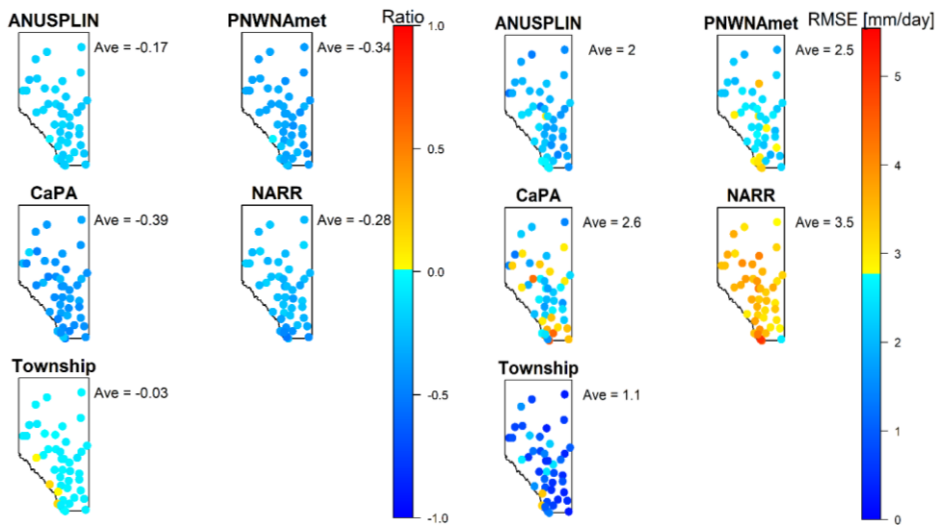
Figure 4. Temporal Correlation Coefficient (TCC) between historical precipitation data.

*: AHCCD data



(a) D_{KS}

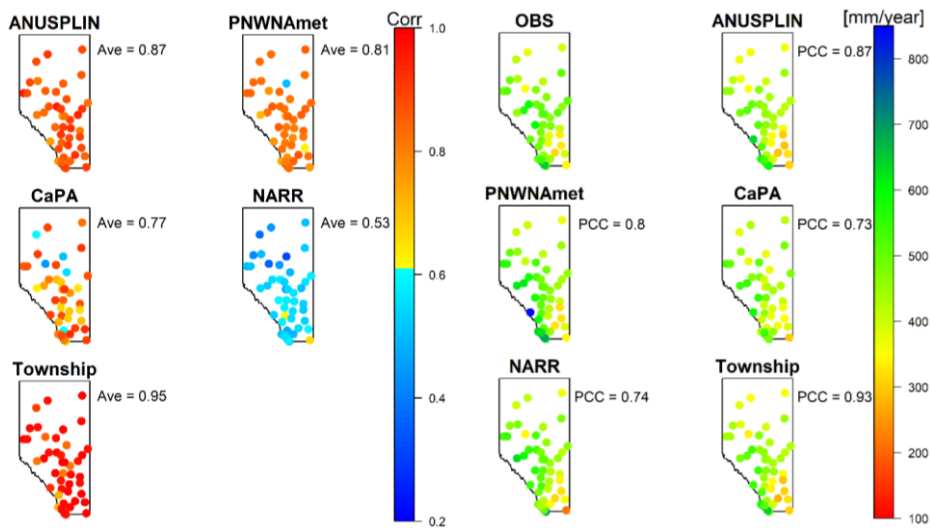
(b) P_{bias}



(c) σ_{ratio}

(d) RMSE

Figure 5. Maps of performance measures for AHCCD precipitation stations in Alberta



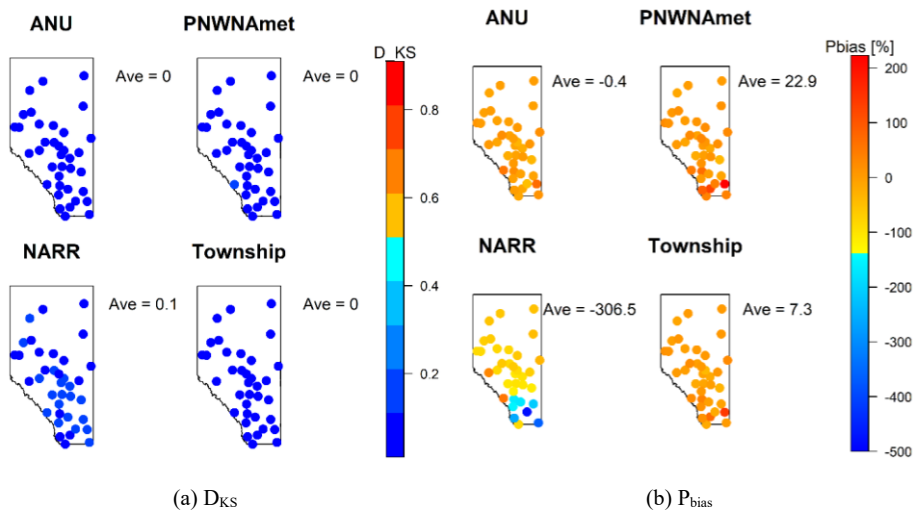
(e) TCC

(f) Mean annual precipitation

Figure 5. Continued

1
2
3
4

1
2
3



4
5
6

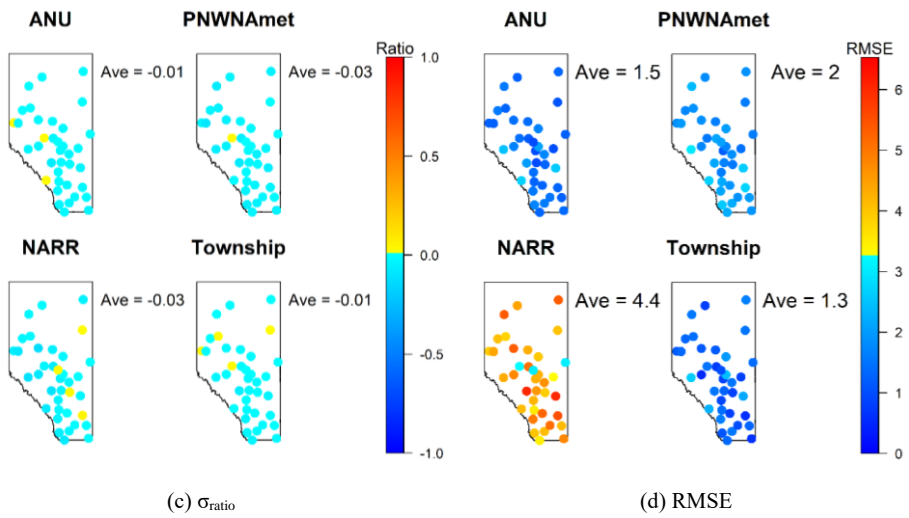
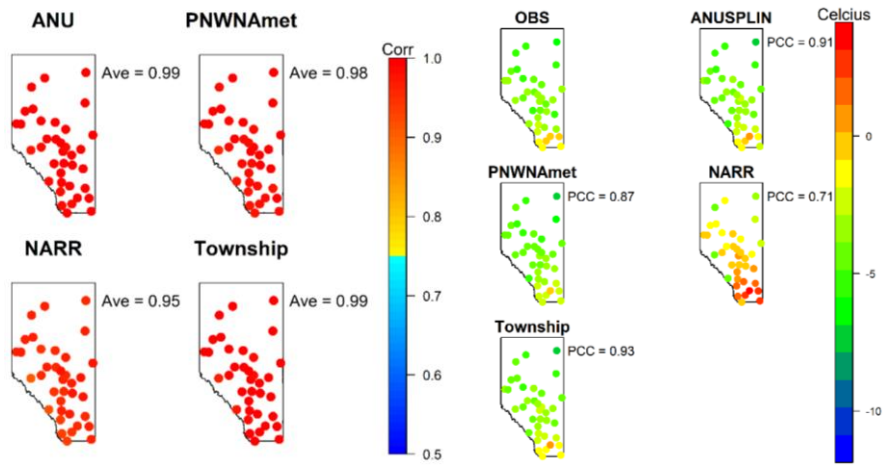


Figure 6. Maps of performance measures for minimum temperature over the AHCCD stations in Alberta

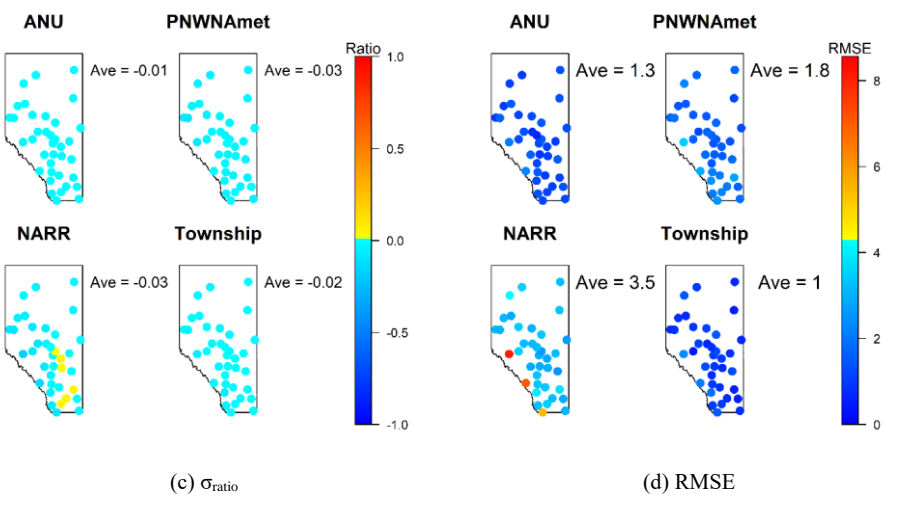
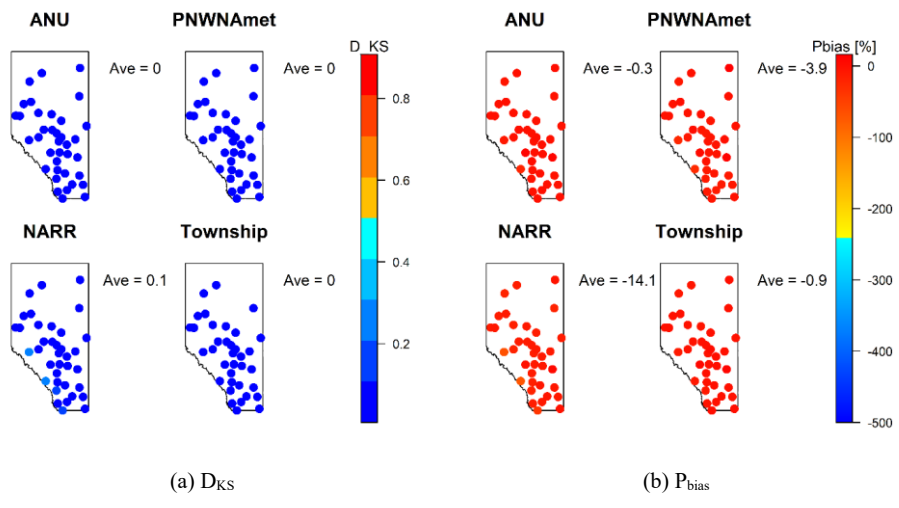


(e) TCC

(f) Mean annual minimum temperature

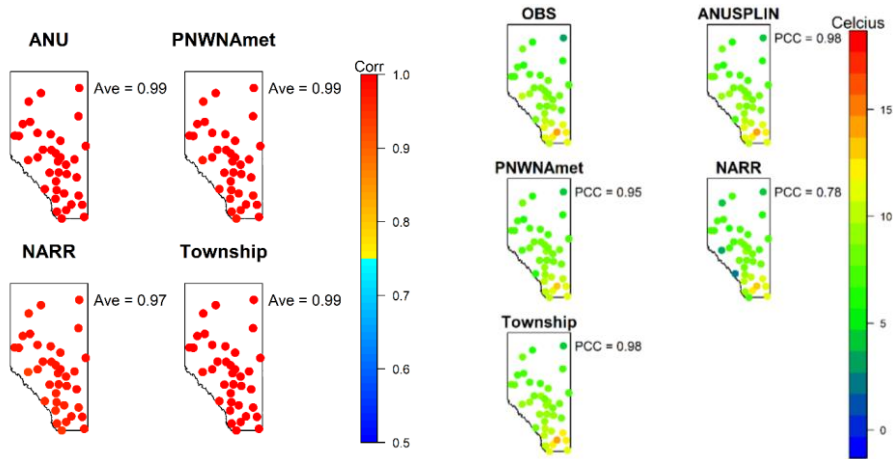
Figure 6. Continued

1
2
3
4



5 Figure 7. Maps of performance measures for maximum temperature over the AHCCD stations in Alberta

6

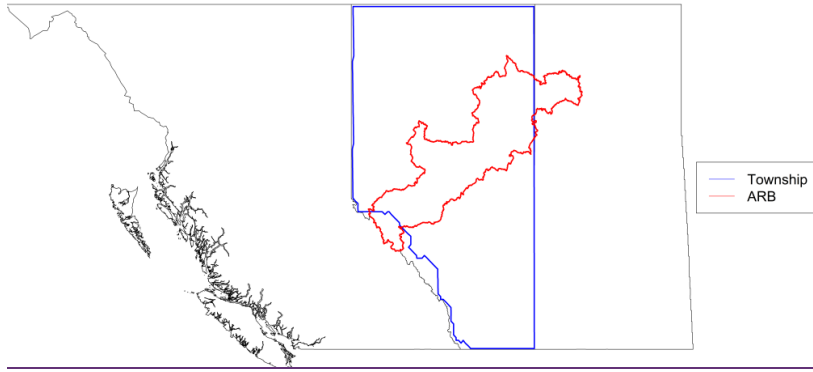


(e) TCC

(f) Mean annual maximum temperature

Figure 7. Continued

1
2
3
4
5

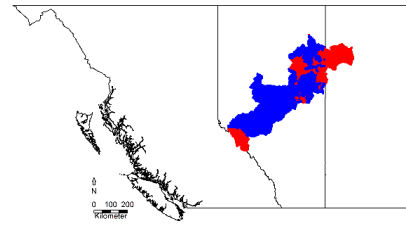
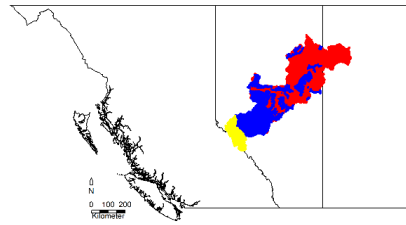
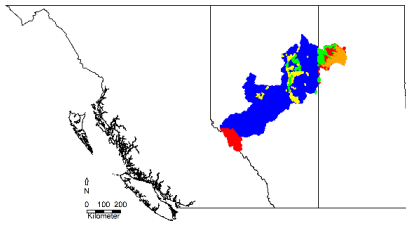


1
2 Figure 8. Domain of the Township dataset (blue line) and the boundary of the Athabasca River basin (red
3 line)
4

(A)Precipitation

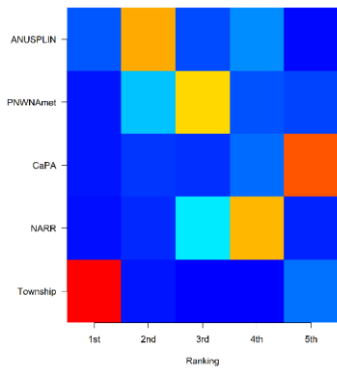
(B)Temperature

(C)Multi-variables

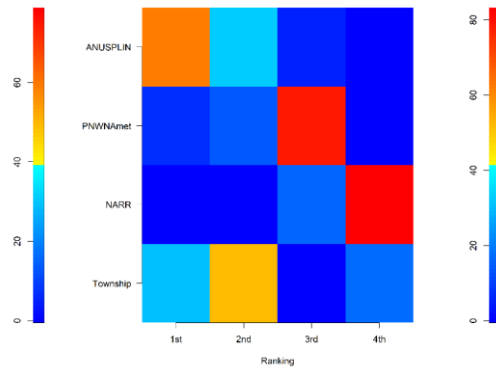


● ANUSPLIN ● Township ● PNWNAmet ● NARR ● CaPA

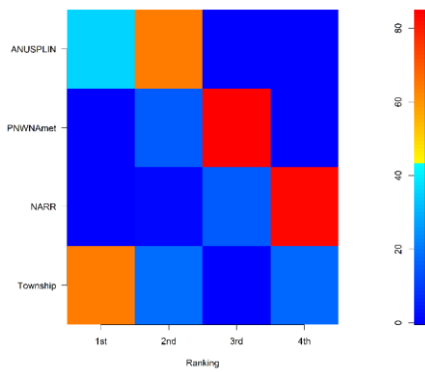
Figure 9. Maps of the first-ranked climate datasets in ARB for the individual variable (A and B) and multi-variables (C)



(a) Precipitation



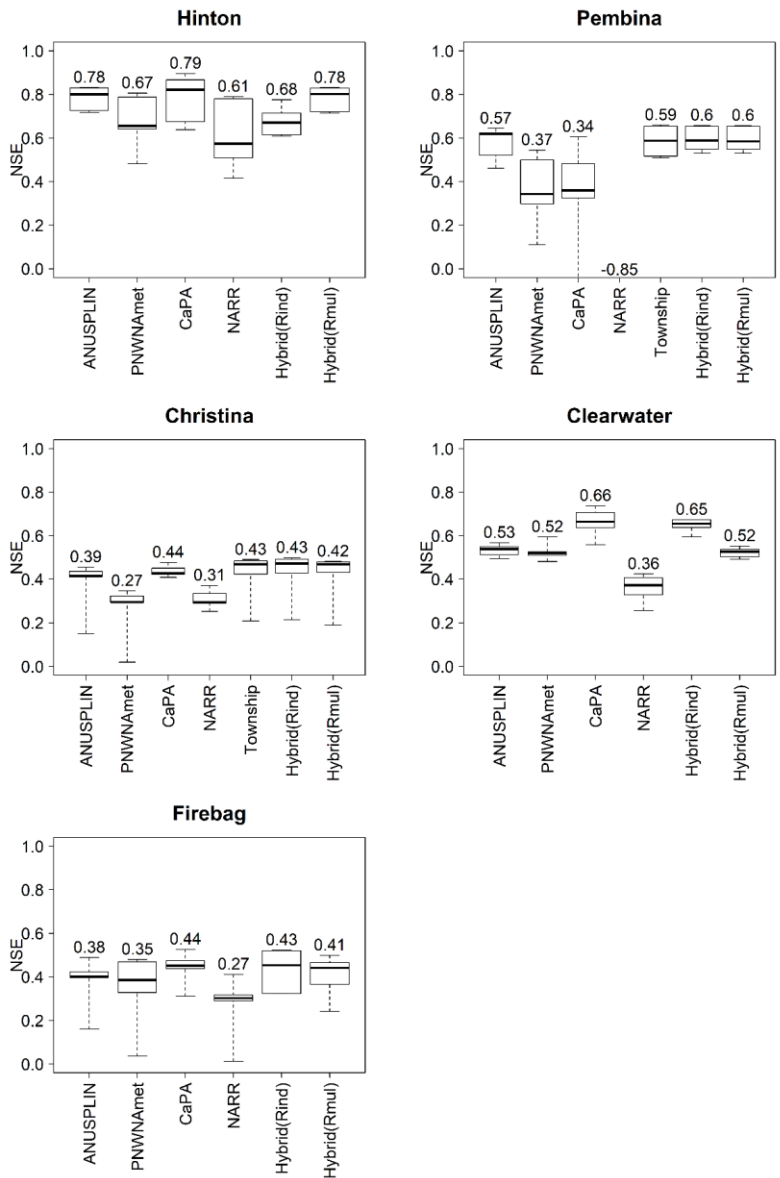
(b) Temperature



(c) Multi-variables

Figure 10. Percentage of climate datasets on each rank for R_{ind} and R_{mul}

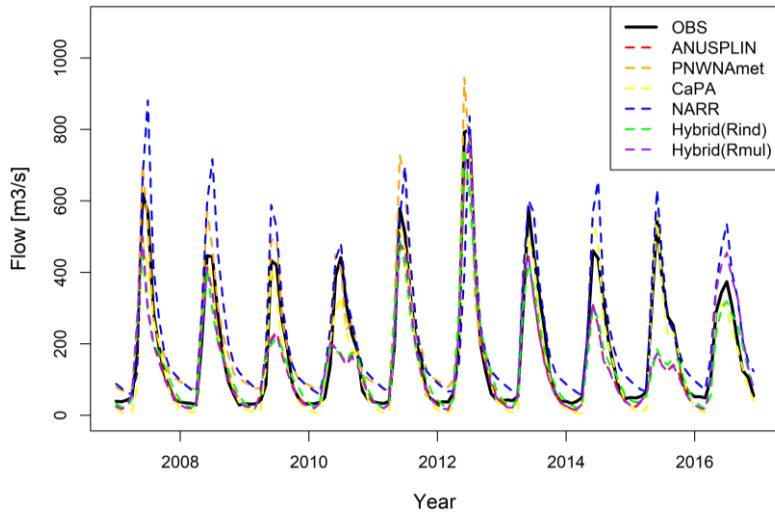
1
2
3
4
5
6
7



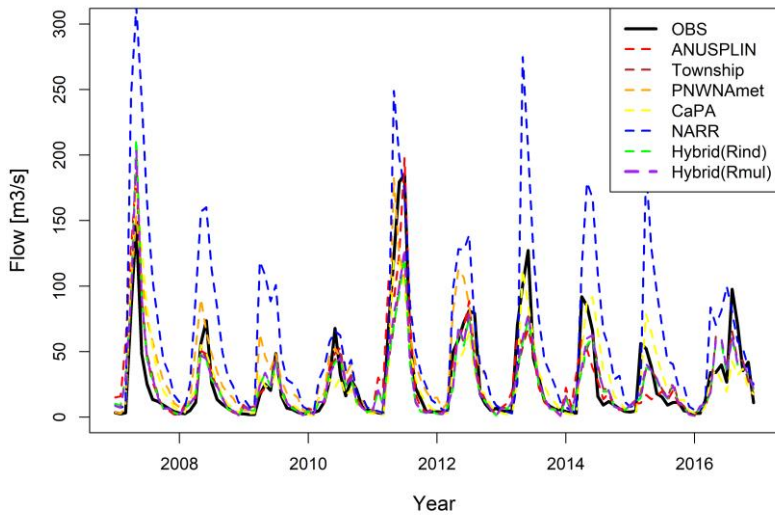
1

2 Figure 11. Boxplots of the NSEs of the proxy validation at the five sub-basins in ARB. The values
 3 above each boxplot represent the average over NSEs of the proxy validation.

Deleted: NSE



(a) Hinton



(a) Pembina

Figure 12. Monthly observed and simulated hydrographs from the gridded climate datasets at (a) Hinton and (b) Pembina

Deleted: Daily

Deleted: six