1 **<<Reviewer 1>>**

2

3 **(1) I think the authors have modified the manuscript accordingly and respond to the**
4 **comments reasonably. Major concern is still the performance of the models. The authors**
5 **have pointed out the NSE above 0.5 could be regarded as satisfactory, However, they used**
6 **different models. There are numerous studies showing that the NSE above 0.8 with VIC.**
7 **Since the performance for Christina and Firebag is not satisfactory as admitted by the**
8 **authors, what are the possible reasons for this? Could it be the possible that the VIC model**
9 **is not suitable for the study area? The authors failed to give answers to those questions. If**
10 **the model is not reliable, I cannot be persuaded the further results which are based on the**
11 **modelling results is reliable. My suggestion is still to improve the performance of the**
12 **modelling or use more models to make the results robust.**
13 ((Reply)) As commented in the previous reply, the parameters of the VIC model were calibrated
14 individually for the seven historical gridded climate datasets (i.e., ANUSPLIN, Alberta Township,
15 PNWNAmet, CaPA, NARR, and two hybrid climate datasets) using an auto- calibration method
16 (dynamic dimensional search algorithm) under the same optimization constraints (e.g., 100
17 maximum iteration) for a fair comparison of the performance. For Christina, the performance of
18 the VIC model was not satisfactory for PNWNAmet and NARR during the calibration period and
19 for most cases for the validation period. For Firebag, the VIC model's performance was acceptable
20 except for NARR during both calibration and validation periods. The poor performance of NARR,
21 as commented in the previous reply, can be attributed to discontinued assimilation of observed
22 precipitation data over Canada since 2003, which is evident for the fact that the NSE values during
23 the validation period were much lower than those for the calibration (1985 to 1997) for all of the
24 hydrometric stations. In addition, such an underperformance at lower part of the Athabasca River
25 basin (i.e., lower Athabasca River basin) may be attributed to 1) relatively poor forcing datasets
26 within the drainage area of each hydrometric station, caused by the lack of observational stations
27 in the northern part of Alberta (refer to Figure 1) and 2) anthropogenic activities that were not
28 reflected in the VIC model simulations during the validation period when land cover changes and
29 water withdrawals mainly induced by Oil-Sand development have occurred. The table below
30 (Table 7 in the revised manuscript) shows the intercomparison of the performance of hydrologic
31 models applied for the Athabasca River basin in literature. Only one study (Shrestha et al., 2017)
32 has conducted hydrologic modeling at a sub-watershed level. The results of the current study were
33 from the VIC simulation forced by the hybrid climate dataset. The VIC model's performance was
34 better or comparable to the literature. In particular, this study significantly improved the
35 performance of streamflow simulation for the Firebag catchment from 0.28 to 0.56. Comparing to
36 the NSE values in Table 6, the NSE values of all cases for Firebag and Christina were better (or
37 comparable) than those of the literature. Both Table 6 and Table 7 also showed clearly that
38 regardless of climate forcing datasets and hydrologic models (i.e., VIC or SWAT), the hydrologic
39 performance measured by NSE for Christina and Firebag was not improved above 0.55 and 0.65,
40 respectively. Based on the intercomparison of performance between this study and the literature,

1

the quality of the hydrologic simulations in this study was improved (or comparable) considerably, in particular at the Firebag station, compared to the results of the literature. Furthermore, the reviewer needs to note that the main aim of this study is not to improve the performance hydrologic simulations, but to provide more reliable climate dataset (i.e., hybrid climate dataset) through REFRES suggested in this study. To validate the applicability of the hybrid climate dataset, a proxy validation approach was employed by comparing simulated streamflows derived from the generated hybrid climate data and other available climate datasets to recorded streamflows at the selected hydrometric stations in the Athabasca River basin.

Table 7. NSE values between the current study and literature for the Athabasca River basin. The NSE values were obtained for calibration and validation periods. For comparison of the current study to the literature, the NSE values of the current study were obtained from the VIC simulation for the hybrid climate dataset ($R_{ind}$).

| Stations | Current study/ VIC[1] | Literature/Hydrologic model | | | | |
|---|---|---|---|---|---|---|
| | | Shrestha et al. (2017b)/ SWAT[2] | Faramarzi et al. (2017)/ SWAT | Faramarzi et al. (2015)/ SWAT | Betrie et al. (2015)/ SWAT | Leong and Donner (2015) /IBIS-THMB[3] |
| Hinton | 0.80 | 0.87 | - | - | - | - |
| Pembina | 0.64 | 0.69 | - | - | - | - |
| Athabasca | 0.78 | 0.90 | - | - | | 0.50 |
| Fort McMurray | 0.77 | 0.89 | - | - | 0.41 | 0.35 |
| **Christina** | **0.52** | **0.49** | - | - | - | - |
| **Firebag** | **0.56** | **0.28** | - | - | - | - |
| Average for all stations | 0.58 | 0.57 | 0.21 | 0.11 | - | - |

[1] Variable Infiltration Capacity

[2] Soil and Water Assessment Tool

[3] Integrated BIosphere Simulator - Terrestrial Hydrology Model with Biogeochemistry

Betrie, G. D., Deng, B. and Wang, J.: Integrated modeling of the Athabasca River Basin using
    SWAT, Proceedings of Science and Technology Innovations, 27–38, 2015.
Faramarzi, M., Srinivasan, R., Iravani, M., Bladon, K. D., Abbaspour, K. C., Zehnder, A. J. B.
    and Goss, G. G.: Setting up a hydrological model of Alberta: Data discrimination

analyses prior to calibration, Environmental Modelling & Software, 74, 48–65, doi:10.1016/j.envsoft.2015.09.006, 2015.

Faramarzi, M., Abbaspour, K. C., Adamowicz, W. L. (Vic), Lu, W., Fennell, J., Zehnder, A. J. B. and Goss, G. G.: Uncertainty based assessment of dynamic freshwater scarcity in semi-arid watersheds of Alberta, Canada, Journal of Hydrology: Regional Studies, 9, 48–68, doi:10.1016/j.ejrh.2016.11.003, 2017.

Leong, D. N. S. and Donner, S. D.: Climate change impacts on streamflow availability for the Athabasca Oil Sands, Climatic Change, 133(4), 651–663, doi:10.1007/s10584-015-1479-y, 2015.

Shrestha, N. K., Du, X. and Wang, J.: Assessing climate change impacts on fresh water resources of the Athabasca River Basin, Canada, Science of the Total Environment, 601, 425–440, 2017.

The authors also addressed the performance of the VIC model used in this study as below:

*"Over the five hydrometric stations, most of the climate datasets performed well with the exception of NARR in the Pembina catchment. Most of the NSE values in calibration for Christina and Firebag were above 0.50, which was considered as a threshold of satisfactory performance in hydrologic models as suggested by Moriasi et al. (2007). However, model performance is not satisfactory for Christina and Firebag during the validation period. Such an underperformance at the lower reach of the Athabasca River basin may be attributed to 1) relatively poor forcing datasets within the drainage area of each hydrometric station, caused by the lack of observational stations in the northern part of Alberta (refer to Figure 1) and 2) anthropogenic activities that were not reflected in the VIC model simulations especially during the validation period when land cover changes and water withdrawals mainly induced by Oil-Sand development have occurred. Table 7 shows the NSE values of hydrologic models applied for the Athabasca River basin in literature. All of NSE values were obtained from the simulations for calibration and validation periods. The NSE values of the current study were obtained from the VIC simulation forced by Hybrid ($R_{ind}$) for comparison to the literature. It needs to note that the VIC model was calibrated for the entire ARB watershed to simulate historical flow over the ARB. The results of the VIC simulation for the entire Athabasca River basin were included in the discussion section. The VIC model's performance in this study was better or comparable to the literature for all stations in ARB. In particular, this study improved considerably the performance of streamflow simulation for the Firebag catchment. Comparing to the NSE values presented in Table 6, in addition, the NSE values of all cases for Firebag and Christina were better (or comparable) than those of the literature. Overall, the quality of hydrologic simulations in this study was improved (or comparable) considerably, compared to the results of the literature. Consequently, the VIC model performance is acceptable at all of hydrometric stations for the proxy validation."* (P22L3-L23)

1 **<<Reviewer 2>>**

2 **The authors had shown a tremendous effort in addressing the comments made from the**

3 **previous round of review. I am satisfied with the revised version. There are two minor**

4 **issues to be addressed before publication.**

5

6 **(1) It would be appreciated if the authors could add the results (Table in the response letter**

7 **and perhaps a hydrograph as a sub-figure in Figure 12) to support the argument that**

8 **hybrid dataset performed better a whole-basin scale in the revised manuscript [P23L24-**

9 **P24L1-2].**

10 ((Reply)) The authors added a table (Table 8 in the revised manuscript) and included the

11 following text in the manuscript to address the added-value of the hybrid climate dataset for the

12 whole ARB as below;

13 *"To further validate the utility of the hybrid climate dataset, the VIC model was calibrated for*

14 *the entire ARB to produce a long-term historical hydrologic simulation for the ARB. Table 8*

15 *presents the NSE values of hydrologic simulations forced by ANUSPLIN and Hybrid ($R_{ind}$) at the*

16 *hydrometric stations in the main stream of the ARB. The result shows that as the size of watershed*

17 *increases, the hybrid climate dataset starts performing better than ANUSPLIN used in Eum et al.,*

18 *(2017). In other words, the hybrid climate dataset improved the historical hydrologic simulation*

19 *for the ARB. This is mainly due to the fact that as the watershed area increases, the derived hybrid*

20 *climate dataset is no longer dominated by a single gridded climate dataset."* (P25L4 – L11)

21

22 **(2) A better explanation to the reasons of selecting the five climate datasets (especially NARR)**

23 **is needed to fully address the comment made by Fuad Yassin. First of all, Wong et al. (2017)**

24 **inter-compared multiple climate datasets at daily time step not monthly scale. Secondly, the**

25 **GPCC and CRU data mentioned by Fuad Yassin are actually referring to WFDEI [GPCC]**

26 **and WFDEI [CRU] as abbreviated in Wong et al. (2017). He tried to argue that WFDEI**

27 **[GPCC] and WFDEI [CRU] had already been shown to perform much better than NARR**

28 **across Canada. However, the authors still picked NARR as one of their candidates but not**

29 **WFDEI datasets. I suggest a better justification should be provided and included in the main**

30 **text (i.e. Section 2.2).**

31 ((Reply)) The WATCH Forcing Data methodology applied to the ERAInterim (WFDEI) dataset

32 provides reanalysis data from 1979 to 2016 globally at 0.5° ( ~ 50km), which are bias-corrected

33 by the Climatic Research Unit (CRU) and the Global Precipitation Climatology Centre (GPCC)

34 monthly precipitation data (Weedon et al., 2014). Another representative reanalysis data in North

35 America is the North American Regional Reanalysis (NARR) that provides a long-term set of

36 dynamically consistent 3-hourly climate data from 1979 to present at a regional scale (0.3°= ~

37 32km) for the North America domain. Wong et al. (2017) found that WFDEI performed better than

38 NARR over Canada. However, their study focused on only precipitation at the Canada-wide scale.

39 In addition, WFDEI is not an operational system but is updated when GPCC and CRU are available

40 for the bias-correction of monthly values. Furthermore, WFDEI provides rain and snow separately,

1   which requires another process to obtain daily total precipitation. On the contrary, the NARR data
2   provides total precipitation rate and is available from 1979 to present with ½ month delay as an
3   operational system. In other words, NARR is regularly updated every ½ month. Therefore, this
4   study selected NARR to provide a more recent climate dataset through the REFRES. The authors
5   added the justification of why NARR was included in this study to section 2.2.3.

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    Hyung-Il Eum[1], ·and Anil Gupta[1,2]

5

6

7    **For submission to Hydrology and Earth System Sciences (HESS)**

8

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 content, 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

## 1. Introduction

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

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

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

3

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

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

4

1 region (Faramarzi et al., 2015; Shrestha et al., 2017a). In addtion, no systematic framework has been

2 developed yet that could be employed by climatic and hydrologic studies.

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

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.

21

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 $2^{nd}$ 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).

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.

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

## 2.2 Historical gridded climate datasets

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

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

### 2.2.1 Station-based datasets

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

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.

**2.2.2 Multiple source-based dataset**

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

**2.2.3 Reanalysis-based dataset**

Reanalysis products are another common type of gridded dataset used in climate and hydrologic studies. The WATCH Forcing Data methodology applied to the ERAInterim (WFDEI) dataset provides reanalysis data from 1979 to 2016 globally at 0.5° ( ~ 50km), which are bias-corrected by the Climatic Research Unit (CRU) and the Global Precipitation Climatology Centre (GPCC) monthly precipitation data (Weedon et al.,

2014). Another representative reanalysis data in the North America is the North American Regional Reanalysis (NARR) that has been developed to create a long-term set of dynamically consistent 3-hourly climate data from 1979 to 2003 at a regional scale (0.3°= ~ 32km) for the North America domain (Mesinger et al., 2006). By utilizing advanced land-surface modeling and data assimilation through the Eta Data Assimilation System (EDAS), NARR improved the National Centers for Environmental Prediction/National Center for Atmospheric Research (NCEP/NCAR) global reanalysis data. NARR cycled every 3 hours to produce a climate dataset from 1979 to present. Choi et al. (2009) tested the applicability of NARR for hydrologic modeling in Manitoba for a region with a poor monitoring network density. However, the NARR dataset after 2004 is not consistent with that of prior years (i.e., 1979 to 2003) because assimilation of precipitation observations was discontinued in 2003 (Eum et al., 2014a). Wong et al. (2017) found that WFDEI performed better than NARR over Canada. However, their study focused on only precipitation at the Canada-wide scale. In addition, WFDEI is not an operational system but is updated when GPCC and CRU are available for the bias-correction of monthly values. Furthermore, WFDEI provides rain and snow separately, which requires another process to obtain total precipitation. On the contrary, the NARR data provides total precipitation rate and is available from 1979 to the current with ½ month delay as an operating system. In other words, NARR is operationally updated every ½ month. Therefore, this study selected NARR to provide more recent climate dataset through the REFRES. Using the 3-hr NARR climate data, daily precipitation and minimum and maximum temperature were calculated by adjusting the time zone to MT from the original NARR dataset (UTC zone).

**3. Methodology**

**3.1 Reference Reliability Evaluation System (REFRES)**

This study suggests a **REF**ference **R**eliability **E**valuation **S**ystem (REFRES) that consists of three main modules (refer to Figure 2): 1) a performance measure module (PMM) to evaluate various performance measures for each climate dataset, 2) a ranking module (RM) to identify the most reliable

1    climate data for a target grid cell using a multi-criteria decision-making technique based on the performance

2    measures provided by PMM, and 3) a data generation module (DGM) to produce a hybrid climate dataset

3    by selecting the most reliable climate dataset based on the ranking provided by the RM (ranking model).

4    These three modules are seamlessly integrated and exchange the required data and information to generate

5    a hybrid climate dataset. The next section provides further details on each module.

6    Figure 2. Structure of REFRES comprised of three modules; 1) Performance Measure Module (PMM), 2)

7                 Ranking Module (RM), and 3) Data Generation Module (DGM)

8

9    **3.1.1 Performance Measure Module (PMM)**

10      AHCCD is a point (station) dataset while the other climate datasets used in this study (refer to Table

11    1) are regularly spaced gridded datasets with varying time period, spatial resolution, and coverage (i.e.,

12    domain). Therefore, the inverse distance squared weighting method was applied to obtain the values at the

13    AHCCD stations from all the gridded climate datasets. Then, performance measures were calculated by

14    comparing the interpolated values with the data collected at AHCCD stations. The choice of the

15    performance measures is vital in REFRES, as the ranking of climate datasets entirely depends on included

16    performance measures. In this study, performance measures were selected based on three criteria: 1)

17    distribution, 2) sequencing, and 3) spatial pattern. Distribution-related performance is assessed by the

18    Kolmogorov-Smirnov $D$ statistic ($D_{KS}$) and standard deviation ratio ($\sigma_{\text{ratio}}$). Sequence-related performance

19    is assessed by the percentage of bias ($P_{\text{bias}}$), root mean square error (RMSE), and temporal correlation

20    coefficient (TCC). Spatial pattern-related performance is evaluated by the pattern correlation coefficient

21    (PCC) as shown in Eq. (1) to Eq. (5). The equations of TCC and PCC are identical but TCC is calculated

22    with the daily time series of climate variables and PCC is obtained by the mean annual precipitation and

23    temperature of the AHCCD stations over a target domain. Therefore, PCC varies with the user specified

24    target domain.

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

$$\sigma_{\text{ratio}} = \{(\sigma_G/\sigma_O) - 1\} \tag{2}$$

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

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

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

where $\sigma_G$ and $\sigma_O$ are the standard deviation of gridded and observed climate datasets, $G_i$ and $O_i$ represent gridded and observed climate datasets at $i$th time step, respectively; $F$ is the empirical distribution function of a climate dataset; $\sigma$ is standard deviation; $\bar{G}$ and $\bar{O}$ represent the mean of gridded and observed climate datasets, respectively and $N$ is a total number of data points. These six performance measures were calculated for all the selected climate datasets and variables at each AHCCD station. Figure 2 (blue box in PMM) shows an example of 6 PMs calculated for the precipitation variable using the ANUSPLIN gridded data. Thus, 15 tables (5 climate datasets × 3 variables) were generated by PMM and transferred to the RM.

### 3.1.2 Ranking Module (RM)

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

As multiple performance measures are employed in this study, there are situations when a climate dataset may perform well for some measures but not for others. Therefore, a multi-criteria decision-making

1   (MCDM) technique is required to systematically rank all of the climate datasets while considering multiple

2   performance measures. This study applied a multi-criteria decision-making technique called the Technique

3   for Order of Preference by Similarity to Ideal Solution (TOPSIS, Hwang and Yoon 1981) to systematically

4   determine the order of preference for all climate datasets at each target grid cell. TOPSIS calculates the

5   geometric distance between alternatives and an ideal solution defined by the best performance on each

6   criterion from the alternatives, and then determines the best and worst alternatives based on the distance.

7   TOPSIS has been successfully applied to watershed management for multi-criteria problems (Jun et al.,

8   2013; Lee et al., 2013). TOPSIS starts with the averaged performance measures, $(x_{ij})_{m \times n}$ for the $i^{th}$ alternative

9   (climate dataset in this study) and $j^{th}$ criterion (i.e., a performance measure). A weighted normalized decision

10  matrix, $(t_{ij})_{m \times n}$ is given by

11
$$\left(t_{ij}\right)_{m \times n} = \left(w_j n_{ij}\right)_{m \times n} \quad i = 1,2,\cdots,m; \quad j = 1,2,\cdots,n \tag{6}$$

12
$$n_{ij} = \frac{x_{ij}}{\Sigma_{i=1}^m x_{ij}^2} \tag{7}$$

13  where, $m$ and $n$ are the total number of alternatives and criteria, respectively, $n_{ij}$ is normalized matrix by Eq.

14  (7), and $w_j$ represents weighting on the $j^{th}$ criterion. Under the assumption that all performance measures

15  are important, this study used an equal weighting. Then, Euclidean distances ($d_{ib}$ and $d_{iw}$) of climate datasets

16  from the best ($A_b$) and worst ($A_w$) conditions were calculated respectively by Eq. (8) to Eq. (11)

17
$$A_w = \left\{ \left\langle \max(t_{ij}|i = 1,2,\cdots,m) | j \in J_- \right\rangle, \left\langle \min(t_{ij}|i = 1,2,\cdots,m) | j \in J_+ \right\rangle \right\} \equiv \left\{ t_{wj}|j = 1,2,\cdots,n \right\} \tag{8}$$

18
$$A_b = \left\{ \left\langle \min(t_{ij}|i = 1,2,\cdots,m) | j \in J_- \right\rangle, \left\langle \max(t_{ij}|i = 1,2,\cdots,m) | j \in J_+ \right\rangle \right\} \equiv \left\{ t_{bj}|j = 1,2,\cdots,n \right\} \tag{9}$$

19
$$d_{iw} = \sqrt{\Sigma_{j=1}^n (t_{ij} - t_{wj})^2} \quad i = 1,2,\cdots,m \tag{10}$$

20
$$d_{ib} = \sqrt{\Sigma_{j=1}^n (t_{ij} - t_{bj})^2} \quad i = 1,2,\cdots,m \tag{11}$$

1    Where, $t_{bj}$ and $t_{wj}$ are the best and worst decision matrices determined by Eq. (8) and (9), respectively, and

2    $J_+$ and $J_-$ represent criteria that have a positive and a negative impact on performance. For example, TCC

3    and PCC are in $J_+$ while $D_{KS}$, $\sigma_{\mathrm{ratio}}$, $P_{\mathrm{bias}}$, and RMSE are in $J_-$. Using the Euclidean distances, the order of

4    preference for all climate datasets was determined by the similarity ($S_{iw}$) to the worst condition in Eq. (15).

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

6    $s_{iw} = 1$ when the alternative is equal to the best condition ($A_b$) and $s_{iw} = 0$ if the alternative is equal to the

7    worst condition ($A_w$). In other words, a higher $s_{iw}$ represents higher preference among alternatives. As we

8    evaluate the performance measures (criteria) for individual climate variables, TOPSIS can be applied to

9    decide the preference of climate datasets considering the performance measures for either individual or

10    multiple variables. In this study, TOPSIS provides two types of ranking information by using performance

11    measures from i) individual climate variable and ii) all climate variables. That is, one is the ranking for

12    precipitation and temperature separately ($R_{ind}$) and the other is the ranking for multiple variables ($R_{mul}$). For

13    example, in this study, $R_{ind}$ was determined by a 5×6 decision matrix (5 climate datasets and 6 performance

14    measures) for precipitation and temperature individually, while $R_{mul}$ was determined by a 4×18 decision

15    matrix (4 climate datasets excluding CaPA that provides only precipitation by 18 performance measures

16    from three variables). To alleviate the erroneous output that minimum temperature is higher than maximum

17    temperature on a certain day when producing the hybrid climate dataset by the ranking of temperature

18    values individually, the performance measures of both minimum and maximum temperature are employed

19    together to rank the climate datasets for temperature.

20

### 3.1.3 Data Generation Module (DGM)

22      DGM extracts the most reliable climate data for a user-specified target region based on the ranking

23    information obtained from the RM. The tool is flexible enough to provide output in various common

24    formats, i.e., NetCDF, ASCII (text) or in the specific format of a numerical model. As all of the historical

1  gridded climate datasets have been tested and employed in numerous climatic and hydrologic studies, an

2  assumption was made in generating the hybrid climate dataset that all of the climate datasets are equally

3  qualified for inclusion but the final selection can be determined by the proven superiority evaluated through

4  the performance measures. Under this assumption, the available datasets can be combined systematically

5  based on the rank (performance) of each dataset at target grid cells. As each climate dataset has different

6  data periods shown in Table 1, the first ranked dataset cannot fully cover a whole target period to be

7  extracted from a set of climate data candidates. DGM provides a systematic procedure to identify the most

8  reliable dataset for a target region and extracts the data from the inventory of climate datasets considering

9  the ranking and availability of each dataset for a desired period. For instance, if CaPA and ANUSPLIN

10  ranked first and second for precipitation and the desired period is 1950 to 2016, DGM starts searching for

11  the availability of precipitation in 1950. As CaPA is only available between 2002 to 2016, DGM reorders

12  the rank to select ANUSPLIN as the best climate dataset available in 1950. In this way, a hybrid dataset

13  over the period 1950 to 2016 is generated by extracting from ANUSPLIN from 1950 to 2001 and CaPA

14  from 2002 to 2016 in this particular case. Once the best climate datasets are extracted over all the target

15  grid cells (study domain), the hybrid climate dataset is produced in a user-defined format. This study

16  generated the hybrid climate datasets in the form of the VIC forcing input format to be directly employed

17  into the hydrologic model.

18

19  **3.2 Proxy validation**

20  Although the AHCCD dataset has been adjusted to provide better estimates of actual precipitation and

21  temperature, it contains statistical artifacts that include inevitable errors from sequential data processes that

22  can be propagated in the derived hybrid climate dataset. Given that the AHCCD stations, the reference

23  dataset for the performance measures, are not regularly distributed and have especially poor density in the

24  northern parts of the study area (refer to Figure 1), it is questionable if the hybrid climate dataset can

25  represent a historical climate better than the individual gridded climate dataset. Utilizing a proxy validation

approach (Klyszejko, 2007), this study applied streamflow records to validate the utility of the derived hybrid climate dataset over other existing climate datasets in hydrologic simulations. In this study, the proxy validation was conducted using an existing hydrologic model (Eum et al., 2017), Variable Infiltration Capacity (VIC, Liang et al., 1994), for the Athabasca River basin (ARB). The VIC model was further refined at $1/32°$ (2~3 km) for a finer spatial resolution and to better simulate the complex river network in the Lower Athabasca River basin. Five of the catchment areas listed in Table 2 were selected for the proxy validation based on three criteria: i) hydrometric record length, ii) location defined by upper, middle and lower reaches (Northern River Basin Study, 2002), and iii) the number of gridded climate datasets used to generate a hybrid climate dataset for the catchment area of the selected hydrometric station. 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. Hinton is located near the headwaters of ARB, which are characterized by mountainous topography and snow- and glacier-ice melt dominated hydrologic regimes. Pembina is one of the major rivers in the middle reach. The other three stations (Christina, Clearwater above Christina and Firebag) are located in the lower reach, which is a water-limited (dry) region due to a higher amount of evapotranspiration (Eum et al., 2014b). The sub-basins of Hinton, Firebag, and Clearwater include a partial area outside of the Township data domain, thus inducing a higher or lower number of climate datasets in the derived hybrid dataset.

A total of seven climate datasets (five individual and two hybrid climate datasets from the $R_{ind}$ and $R_{mul}$) are available to calibrate the VIC hydrologic model parameter set related to soil properties and routing. 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. More details on calibration can be found in Eum et al. (2017). 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

16

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 × 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.

**4. Results**

**4.1 Precipitation performance measures in Alberta**

1       Although the performance measures were calculated for 190 AHCCD stations in western Canada, the

2       target area of this study is in Alberta, where only 45 stations are located. Therefore, the results for the 45

3       AHCCD stations are given in this study. Table 3 shows spatially-averaged performance measures for

4       precipitation. The Township data outperformed other climate datasets for all performance measures except

5       $P_{bias}$. ANUSPLIN is the second best climate dataset for Alberta. All climate datasets underestimate the

6       standard deviation of observed daily precipitation (i.e., negative $\sigma_{ratio}$), especially PNWNAmet and CaPA

7       which underestimated by 34 % and 39 %, respectively. Interestingly, two station-based gridded climate

8       datasets, ANUSPLIN and Township, show negative $P_{bias}$ while PNWNAmet, CaPA, and NARR datasets

9       have positive $P_{bias}$. This indicates that ANUSPLIN and Township may underestimate extreme precipitation,

10     as they employed the raw station data instead of the adjusted precipitation data which is higher than the raw

11     station data by 5%-20%. In contrast, other climate datasets (especially multiple sources and reanalysis data)

12     overestimate extreme precipitation. These results are consistent with findings in Eum et al. (2014a) that

13     CaPA and NARR overestimate extreme precipitation events by overly reflecting the orographic effects on

14     precipitation in western Alberta.

15      Figure 4 shows the temporal correlation coefficient (TCC) data averaged over the AHCCD stations in

16     Alberta to investigate the similarity between historical precipitation datasets employed in this study. As

17     expected, station-based climate datasets (i.e., ANUSPLIN, PNWNAmet, and Township) showed better

18     TCCs than CaPA and NARR. The TCC between ANUSPLIN and Township was the highest among climate

19     datasets except for the observations (i.e., OBS), even though they incorporated different interpolation

20     techniques. PNWNAmet showed the highest TCC with ANUSPLIN because they both are based on thin

21     plate spline interpolation. TCCs between CaPA and other climate datasets are similar, as CaPA is produced

22     from multiple sources such as GEM's outputs and weather radar networks of Canada and US. NARR, the

23     reanalysis-based climate dataset, showed higher TCC with CaPA than with other datasets, as it is assimilated

24     with multiple sources of observations.

1    Maps of each performance measure are shown in Figure 5. It is evident from the spatial variability that

2    the ANUSPLIN and Township datasets outperformed the other datasets in $D_{KS}$ throughout Alberta. In the

3    mountainous region of southwest Alberta, most of the climate datasets performed poorly in $P_{bias}$, $\sigma_{ratio}$,

4    RMSE, and PCC, resulting mainly from the sparse observation network and inconsistent observations near

5    the Canada-US border. PNWNAmet highly overestimates the mean annual precipitation in the mountainous

6    area (e.g., 300 mm/year higher than that observed at station ID 3050519), which may considerably affect

7    simulated streamflows originating in mountainous headwaters and further downstream.

8

9    **4.2 Air temperature performance measures in Alberta**

10    The performance measures for air temperature averaged over 37 AHCCD stations in Alberta are

11    presented in Table 4. As CaPA provides only precipitation, it was excluded in the assessment for temperature.

12    All of the performance measures for temperature are better than those for precipitation except $P_{bias}$. NARR

13    is highly biased as it underestimates minimum and maximum temperatures, which might be an attribute of

14    discontinuation of observation assimilation since 2003 (Eum et al., 2014a). ANUSPLIN and Township

15    showed an almost perfect linear relationship (TCC) with the observations (i.e., > 0.97 for all of the climate

16    datasets). The performance measures for maximum temperature are better than those for minimum

17    temperature as maximum temperature is dominated by mainly large-scale heat waves while minimum

18    temperature is affected by local physical processes, e.g., topography and surface conditions (Eum et al.,

19    2012). NARR showed less skill in capturing these local effects due to the coarse spatial resolution (~32km)

20    compared to other station-based climate datasets. As with precipitation, the maps of performance measures

21    for minimum and maximum temperature presented in Figure 6 and Figure 7 showed that data from the

22    mountainous areas performed poorly in most of the performance measures. NARR showed positive and

23    negative $P_{bias}$ for minimum and maximum temperature, respectively, in the mountainous region, indicating

24    that NARR has a warm bias in extreme cold temperatures and a cold bias in extreme warm temperatures.

25

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

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

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

For temperature, ANUSPLIN was ranked first (in 62% grid cells) for the whole ARB, followed by Township (31%) and PNWNAmet (7%). In the upper and middle reaches, i.e., Hinton and Pembina, PNWNAmet and Township were mostly ranked first, respectively, while ANUPLIN outperformed other climate datasets for the subbasins in the lower reach. When considering the performance measures for multiple variables simultaneously, the Township dataset was ranked first, followed by ANUSPLIN for 64% and 36% of the grid cells for the whole ARB. Figure 9 shows maps of the first-ranked climate datasets for each case in Table 5, i.e., individual variable (Case A and B) and multi-variables (Case C). Due to the limited spatial coverage of the Township dataset, other climate datasets were ranked first in the headwaters of the ARB and the area of the river basin in Saskatchewan. For instance, ANUSPLIN and PNWNAmet were ranked first in the headwaters, while no specific climate dataset dominated in Saskatchewan for

1　precipitation (refer to Figure 9A). For temperature, ANUSPLIN outperformed in the northern part (middle

2　and lower reaches of the ARB) due to outstanding performance of the $P_{bias}$ performance measure for

3　minimum temperature as shown in Table 4 and Figure 6(b). For multi-variables, Township was mostly

4　ranked first within its domain and ANUSPLIN was ranked first outside the Township dataset domain and

5　also for a small part of lower reach area in the ARB.

6　　　Figure 10 shows the percentage of each climate dataset at each rank for the three cases (e.g. A, B, and

7　C in Table 5). For precipitation (Case A), Township overwhelmed other climate datasets. The second

8　alternative was ANUSPLIN in the majority of grid cells in the ARB. PNWNAmet, NARR and CaPA were

9　mostly ranked 3rd, 4th and 5th, respectively. For temperature (Case B), ANUSPLIN was ranked mostly first

10　and Township was a distinct second choice in the majority of grid cells, followed by PNWNAmet and

11　NARR. For multi-variables (Case C), Township and ANUSPLIN were the first and second choices in the

12　majority of grid cells in the ARB, respectively.

13　　　As two different hybrid climate datasets were generated using the ranking information from single-

14　and multi-variable approaches, i.e., Hybrid ($R_{ind}$) and Hybrid ($R_{mul}$), further investigation is required to

15　identify which hybrid climate dataset may provide better performance and consequently will be

16　recommended for future climate-related studies. A proxy validation approach was applied using both

17　generated hybrid climate datasets to validate the utility of one dataset over the other.

18

19　**4.4 Proxy validation of generated hybrid climate datasets**

20　　　In addition to the five gridded climate datasets, the two hybrid climate datasets were implemented for

21　proxy validation using the VIC model. In contrast to the station-based climate datasets, both CaPA and

22　NARR were produced from climate models and multiple sources of observations, consequently showing a

23　higher correlation with each other as shown in Figure 4. Since CaPA also provides only precipitation, this

24　study combined precipitation of CaPA with the NARR temperature to prepare the CaPA climate forcing

25　dataset for the proxy validation. Table 6 presents the Nash-Sutcliffe Efficiency (NSE) for the calibration

21

and validation periods at the selected hydrometric stations (Hinton, Pembina, Christina, Clearwater, and Firebag) in the ARB to assess the suitability of each climate dataset as a climate forcing input data for hydrologic simulations. Over the five hydrometric stations, most of the climate datasets performed well with the exception of NARR in the Pembina catchment. Most of the NSE values in calibration for Christina and Firebag were above 0.50, which was considered as a threshold of satisfactory performance in hydrologic models as suggested by Moriasi et al. (2007). However, model performance is not satisfactory for Christina and Firebag during the validation period. Such an underperformance at the lower reach of the Athabasca River basin may be attributed to 1) relatively poor forcing datasets within the drainage area of each hydrometric station, caused by the lack of observational stations in the northern part of Alberta (refer to Figure 1) and 2) anthropogenic activities that were not reflected in the VIC model simulations especially during the validation period when land cover changes and water withdrawals mainly induced by Oil-Sand development have occurred. Table 7 shows the NSE values of hydrologic models applied for the Athabasca River basin in literature. All of NSE values were obtained from the simulations for calibration and validation periods. The NSE values of the current study were obtained from the VIC simulation forced by Hybrid ($R_{ind}$) for comparison to the literature. It needs to note that the VIC model was calibrated for the entire ARB watershed to simulate historical flow over the ARB. The results of the VIC simulation for the entire Athabasca River basin were included in the discussion section. The VIC model's performance in this study was better or comparable to the literature for all stations in ARB. In particular, this study improved considerably the performance of streamflow simulation for the Firebag catchment. Comparing to the NSE values presented in Table 6, in addition, the NSE values of all cases for Firebag and Christina were better (or comparable) than those of the literature. Overall, the quality of hydrologic simulations in this study was improved (or comparable) considerably, compared to the results of the literature. Consequently, the VIC model performance is acceptable at all of hydrometric stations for the proxy validation. The two hybrid climate datasets performed well, with comparably good and better NSE values than other climate datasets, especially at Pembina, Clearwater, and Firebag, located in the middle and lower reaches.

1    Figure 11 presents the boxplots of NSEs obtained through the multiset-parameter VIC simulations.

2    The NSE ranges were obtained from multiple VIC simulations, with each climate dataset used as climate

3    forcing for all the plausible model parameter sets, which were calibrated with seven climate datasets,

4    individually. The values above each boxplot represent the averaged value of the NSEs over the multiset-

5    parameter hydrologic simulations. A narrower range of NSE values represents a higher precision for a

6    climate dataset and a higher averaged NSE value means higher accuracy. Therefore, a climate dataset

7    showing both a higher averaged NSE and a narrow range of NSEs indicates that it is a relatively more

8    appropriate and reliable climate forcing dataset for hydrologic simulations.

9    At Hinton, all of the climate datasets showed satisfactory NSE values for accuracy, while ANUSPLIN,

10   Hybrid($R_{ind}$), and Hybrid($R_{mul}$) showed better precision. The validation period of CaPA is only six years

11   from 2010 to 2016, as CaPA data are only available between 2002 to 2016. This might be a reason why

12   CaPA produced the highest NSE (accuracy) among the climate datasets used in this study. Therefore, the

13   results of CaPA need to be considered carefully otherwise they might be misleading. In this context, the

14   CaPA dataset was excluded from further assessment of the precision and accuracy even though all of the

15   results of CaPA were included in Figure 11 for reference only. Hybrid($R_{mul}$) and ANUSPLIN showed the

16   highest accuracy as forcing data, followed by Hybrid($R_{ind}$), PNWNAmet, and NARR. In the Pembina and

17   Christina catchments, the Hybrid($R_{ind}$), Hybrid($R_{mul}$), and Township datasets had the highest precision and

18   accuracy. NARR produced negative NSEs at Pembina, indicating it is not reliable or suitable as a forcing

19   dataset. For Clearwater, Hybrid($R_{ind}$) is the top performer, followed by Hybrid($R_{mul}$), ANUSPLIN,

20   PNWNAmet, and NARR. Clearwater had the highest number of climate datasets combined in the hybrid

21   climate dataset within the basin for precipitation as shown in Figure 9. Interestingly, the precision of NARR

22   is similar to that of CaPA because they shared the temperature data from NARR. For Firebag, Hybrid($R_{ind}$)

23   also showed top performance in both precision and accuracy, followed by Hybrid($R_{mul}$), ANUSPLIN,

24   PNWNAmet, and NARR. Overall, Hybrid($R_{ind}$) showed the best accuracy and precision at all hydrometric

23

1    stations, indicating that it has the potential not only to improve historical hydrologic simulations but also

2    to be used as reference data for statistical downscaling of climate change projections in the province.

3

4    **5. Discussion**

5          Among the station-based gridded climate datasets, the Township dataset outperformed other station-

6    based gridded climate datasets. As PNWNAmet set a common period from 1945 to 2012 for all stations

7    included in the interpolation, many stations might be left out in the data generation processes. While

8    ANUSPLIN used the Canada-wide archive (raw) station data collected by only ECCC, the Alberta

9    Township data has been produced on the basis of the archive (raw) station data collected by ECCC, AEP,

10   and AF over Alberta. Therefore, one of the possible reason for outperformance of Township dataset might

11   be the difference in the numbers of stations (i.e. station density) employed to produce the gridded climate

12   datasets. In addition, PNWNAmet showed a positive $P_{bias}$ for precipitation, especially in the mountainous

13   areas, while ANUSPLIN, which employs similar thin plate spline interpolation, generated negative $P_{bias}$.

14   PNWNAmet overestimated precipitation over the mountainous area, which considerably affects simulated

15   low flows at Hinton in the ARB. Figure 12 shows the observed and simulated hydrographs from gridded

16   climate datasets at (a) Hinton and (b) Pembina. It clearly shows that PNWNAmet highly overestimated the

17   low and high, which is caused by overestimated precipitation in the drainage area of the sub-basins. As with

18   PNWNAmet, NARR also overestimated the low and high flows, which is induced by the combined effects

19   of overestimating precipitation and warm biases in cold temperature. The temperature bias of NARR is thus

20   further confirmed and is consistent with the earlier finding of Eum et al., (2014) and Islam and Dery (2016).

21        In Figure 12, the hybrid climate datasets underestimated the peak flows (in 2009, 2010, 2014, and

22   2015) at Hinton, and hydrograph is similar to the hydrograph produced by ANUSPLIN data set that

23   dominantly ranked first in this watershed. On the contrary, the hydrograph of the hybrid climate datasets at

24   Pembina is similar to that of Township that is dominantly ranked first in Pembina (refer to Table 5). These

25   results indicate that the hybrid climate dataset has the intrinsic limitation that the performance of the hybrid

1    dataset for a basin may closely resemble that of the climate dataset that is dominantly ranked first for the

2    basin. However, the utility of the hybrid climate dataset can be clearly found at a whole-basin scale for a

3    large watershed, as the added values of the hybrid climate dataset in sub-basins can be cumulated to the

4    main stem at the downstream in the watershed. To further validate the utility of the hybrid climate dataset

5    the VIC model was calibrated for the entire ARB to produce a long-term historical hydrologic simulation

6    for the ARB. Table 8 presents the NSE values of hydrologic simulations forced by ANUSPLIN and Hybrid

7    ($R_{ind}$) at the hydrometric stations in the main stream of the ARB. The result shows that as the size of

8    watershed increases, the hybrid climate dataset starts performing better than ANUSPLIN used in Eum et

9    al., (2017). In other words, the hybrid climate dataset improved the historical hydrologic simulation for the

10   ARB. This is mainly due to the fact that as the watershed area increases, the derived hybrid climate dataset

11   is no longer dominated by a single gridded climate dataset.

12       Among the station-based gridded climate datasets, ANUSPLIN and Township employed a different

13   number of stations depending on their periods of record. Therefore, there is an inconsistency in these climate

14   datasets over time. For example, the Township dataset employed only 300~400 stations in the 1960s, but

15   has increased to 400~500 since 1970. A change-point analysis of these datasets may provide some useful

16   information to end-users with respect to when and where changes occurred, which will help in establishing

17   spatial and temporal accuracies of these datasets (Eum et al. 2014a). Further, PNWNAmet employed the

18   same number of stations over time to avoid the above mentioned inconsistency, but this study found that it

19   induced overestimation of precipitation in data-poor regions such as mountainous regions in Alberta. As

20   the hybrid climate datasets are generated from the multiple historical gridded datasets, they may also have

21   the same inconsistencies identified in other datasets. The proxy validation, however, demonstrated that the

22   generated hybrid climate datasets can improve the performance of hydrologic simulations.

23       This study identified the preference order of all gridded climate datasets based on the performance

24   measures evaluated at the AHCCD stations, therefore the ranking somewhat relies on the spatial distribution

25   of the AHCCD stations. As shown in Figure 1, the density of AHCCD stations varies across western Canada,

1    and is low in the cold climates of mountainous and northern areas. Therefore, the ranking could further be

2    improved with a more uniform density of AHCCD stations over western Canada.

3         Literature has demonstrated that NARR, a reanalysis-based climate dataset, can be an alternative as a

4    climate forcing dataset for hydrologic simulations in data sparse regions (Choi et al., 2009; Praskievicz and

5    Bartlein, 2014; Islam and Dery, 2016). In this study, the NARR dataset performed quite well in high-

6    elevation regions (Hinton in this study) while it did not perform so well in the middle and lower reaches,

7    i.e., lower-elevation watersheds. NARR performed especially poorly in the Pembina sub-basin, a region

8    where hydrologic simulations are highly sensitive to model parameters (Eum et al., 2014b). In Figure 11

9    (b), however, the NARR parameter set produced fair NSE values in hydrologic simulations forced by the

10   other climate datasets except for CaPA and PNWNAmet. Such result indicates that 1) all of parameter sets

11   used in this study were calibrated reasonably and 2) climate forcing input data plays a more crucial role in

12   hydrologic simulations as any parameter sets did not produce a fair NSE value from NARR in Pembina.

13   CaPA was more suitable than NARR for the selected sub-basins in this study, which indicates that CaPA

14   might be a better alternative in low station-density regions such as the ARB. However, since the validation

15   period in this study is only 7 years from 2010 to 2016, a longer data period is necessary to validate the

16   suitability of CaPA as indicated in Eum et al. (2014a) and Wong et al. (2017).

17        In the proxy validation, Hybrid($R_{ind}$) performed well in the Clearwater sub-basin where the highest

18   number of climate datasets were combined in the generated hybrid climate datasets. The Township dataset,

19   which mostly ranked first within its spatial domain, partially covers the drainage area of Clearwater, so that

20   the generated hybrid climate dataset, Hybrid($R_{ind}$), is composed of many climate datasets in this sub-basin.

21   In a traditional approach to hydrological modelling for Clearwater, either the Township dataset might be

22   completely excluded (as it does not cover the entire Clearwater watershed), or potentially combined with

23   other gridded climate datasets to cover the entire watershed. However, combining different climate datasets

24   to construct the climate forcing for a larger region requires an evaluation of the datasets to identify the order

25   of preference for such aggregation when multiple choices are available. Therefore, this study suggested the

REFRES methodology to systematically compare all-available climate datasets for a region to produce a hybrid climate dataset that covers a desired period of record and spatial domain by considering the order of preference for combining various climate datasets at each grid cell. The proxy validation approach also confirmed the utility of a generated hybrid climate dataset over other data sets, especially in hydrologic simulations.

**6. Summary and concluding remarks**

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

- Among the five climate datasets, the station-based climate datasets performed better than multi-source- and reanalysis-based datasets. The Township dataset, in particular, outperformed other climate datasets in the selected performance measures over northern Alberta.

- Most of the climate datasets performed poorly in the mountainous areas of southwest Alberta, due to a sparse observation network, orographic effects, topographic complexity, and inconsistencies in observation between Canada and the US.

- As a result of REFRES' application for the ARB, the Township and ANUSPLIN datasets are mostly ranked the highest among the five climate datasets for precipitation and temperature, respectively.

1     -     In the proxy validation, two hybrid climate datasets, Hybrid($R_{ind}$) and Hybrid($R_{mul}$), performed

2         better in terms of precision and accuracy as forcing data for hydrologic simulations.

3     -     Hybrid($R_{ind}$) especially outperformed other climate datasets in the Clearwater sub-basin where the

4         highest number of climate datasets were combined in generating Hybrid($R_{ind}$) for precipitation. This

5         indicates that the hybrid climate dataset generated by REFRES may lead to more reliable

6         hydrologic simulations, resulting in improved hydrologic predictions.

7     This study provided the preference order of climate datasets available in Alberta, which may be useful

8 for modelers and decision-makers as to which climate dataset is the most suitable for their studies and

9 projects. Furthermore, this study demonstrated that the hybrid climate dataset produced by REFRES is more

10 representative of historical climatic conditions. Therefore, the hybrid climate dataset is recommended to be

11 used as a reference dataset for statistical downscaling and hydrologic model forcing, resulting in more

12 reliable high-resolution climatic and hydrologic projections.

13

14 *Code availability.* A package of REFRES is available by contacting at hyung.eum@gov.ab.ca when

15 requested. Variable Infiltration Capacity (VIC) is also freely downloaded at https://github.com/UW-

16 Hydro/VIC.

17

18 *Data availability.* ANUSPLIN can be access via ftp://ftp.nrcan.gc.ca/pub/outgoing/canada_daily_grids and

19 PNWNAmet is downloaed at https://data.pacificclimate.org/portal/gridded_observations/map/. The Alberta

20 Township data can be downloaded at http://agriculture.alberta.ca/acis/township-data-viewer.jsp. The

21 archives of CaPA can be access via http://collaboration.cmc.ec.gc.ca/science/outgoing/capa.grib/ and

22 http://collaboration.cmc. ec.gc.ca/science/outgoing/capa.grib/hindcast/ and the last 30 days of CaPA data is

23 available at http://dd.weather.gc.ca/analysis/precip/rdpa/grib2/polar_stereographic. The NARR dataset is

24 available at https://www.esrl.noaa.gov/psd/data/gridded/data.narr.monolevel.html. The hybrid climate

25 dataset for Alberta is also available by contacting at hyung.eum@gov.ab.ca when requested.

14

15 **References**

16 Asong, Z. E., Khaliq, M. N. and Wheater, H. S.: Regionalization of precipitation characteristics in the

17    Canadian Prairie Provinces using large-scale atmospheric covariates and geophysical attributes,

18    Stochastic Environmental Research and Risk Assessment, 29(3), 875–892, 2015.

19 Betrie, G. D., Deng, B. and Wang, J.: Integrated modeling of the Athabasca River Basin using SWAT,

20    Proceedings of Science and Technology Innovations, 27–38, 2015.

21 Choi, W., Kim, S. J., Rasmussen, P. F. and Moore, A. R.: Use of the North American Regional Reanalysis

22    for hydrological modeling in Manitoba, Can. Water Resour. J., 34, 13–36, 2009.

23 Christensen, N. S. and Lettenmaier, D. P.: A multimodel ensemble approach to assessment of climate

24    change impacts on the hydrology and water resources of the Colorado River Basin, Hydrology and

25    Earth System Sciences, 11(4), 1417–1434, 2007.

1    Côté, J., Desmarais, J.-G., Gravel, S., Méthot, A., Patoine, A., Roch, M. and Staniforth, A.: The

2        operational CMC–MRB global environmental multiscale (GEM) model. Part II: Results, Monthly

3        Weather Review, 126(6), 1397–1418, 1998a.

4    Côté, J., Gravel, S., Méthot, A., Patoine, A., Roch, M. and Staniforth, A.: The operational CMC–MRB

5        global environmental multiscale (GEM) model. Part I: Design considerations and formulation,

6        Monthly Weather Review, 126(6), 1373–1395, 1998b.

7    Deacu, D., Fortin, V., Klyszejko, E., Spence, C. and Blanken, P. D.: Predicting the Net Basin Supply to

8        the Great Lakes with a Hydrometeorological Model, Journal of Hydrometeorology, 13(6), 1739–

9        1759, doi:10.1175/JHM-D-11-0151.1, 2012.

10   Demaria, E.M, Nijssen, B., Wagener, T.: Monte Carlo sensitivity analysis of land surface parameters

11       using the variable infiltration capacity model, Journal of Geophysical Research, 112, D11113, 2007

12   Dibike, Y., Eum, H.-I. and Prowse, T.: Modelling the Athabasca watershed snow response to a changing

13       climate, Journal of Hydrology: Regional Studies, 15, 134–148, doi:10.1016/j.ejrh.2018.01.003,

14       2018.

15   Essou, G. R. C., Sabarly, F., Lucas-Picher, P., Brissette, F. and Poulin, A.: Can Precipitation and

16       Temperature from Meteorological Reanalyses Be Used for Hydrological Modeling?, Journal of

17       Hydrometeorology, 17(7), 1929–1950, doi:10.1175/JHM-D-15-0138.1, 2016.

18   Eum, H.-I. and Cannon, A. J.: Intercomparison of projected changes in climate extremes for South Korea:

19       application of trend preserving statistical downscaling methods to the CMIP5 ensemble,

20       International Journal of Climatology, 37(8), 3381–3397, doi:10.1002/joc.4924, 2017.

21   Eum, H.-I., Dibike, Y. and Prowse, T.: Climate-induced alteration of hydrologic indicators in the

22       Athabasca River Basin, Alberta, Canada, Journal of Hydrology, 544, 327–342,

23       doi:10.1016/j.jhydrol.2016.11.034, 2017.

Eum, H.-I., Dibike, Y. and Prowse, T.: Comparative evaluation of the effects of climate and land-cover changes on hydrologic responses of the Muskeg River, Alberta, Canada, Journal of Hydrology: Regional Studies, 8, 198–221, doi:10.1016/j.ejrh.2016.10.003, 2016.

Eum, H.-I., Dibike, Y., Prowse, T. and Bonsal, B.: Inter-comparison of high-resolution gridded climate data sets and their implication on hydrological model simulation over the Athabasca Watershed, Canada, Hydrological Processes, 28(14), 4250–4271, doi:10.1002/hyp.10236, 2014a.

Eum, H.-I., Dibike, Y. and Prowse, T.: Uncertainty in modelling the hydrologic responses of a large watershed: a case study of the Athabasca River basin, Canada, Hydrological Processes, 28(14), 4272–4293, doi:10.1002/hyp.10230, 2014b.

Eum, H.-I., Gachon, P., Laprise, R. and Ouarda, T.: Evaluation of regional climate model simulations versus gridded observed and regional reanalysis products using a combined weighting scheme, Climate Dynamics, 38(7-8), 1433–1457, doi:10.1007/s00382-011-1149-3, 2012.

Faramarzi, M., Abbaspour, K. C., Adamowicz, W. L. (Vic), Lu, W., Fennell, J., Zehnder, A. J. B. and Goss, G. G.: Uncertainty based assessment of dynamic freshwater scarcity in semi-arid watersheds of Alberta, Canada, Journal of Hydrology: Regional Studies, 9, 48–68, doi:10.1016/j.ejrh.2016.11.003, 2017.

Faramarzi, M., Srinivasan, R., Iravani, M., Bladon, K. D., Abbaspour, K. C., Zehnder, A. J. B. and Goss, G. G.: Setting up a hydrological model of Alberta: Data discrimination analyses prior to calibration, Environmental Modelling & Software, 74, 48–65, doi:10.1016/j.envsoft.2015.09.006, 2015.

Garand, L. and Grassotti, C.: Toward an objective analysis of rainfall rate combining observations and short-term forecast model estimates, Journal of Applied Meteorology, 34, 1962–1977, 1995.

Gutmann, E., Pruitt, T., Clark, M. P., Brekke, L., Arnold, J. R., Raff, D. A. and Rasmussen, R. M.: An intercomparison of statistical downscaling methods used for water resource assessments in the United States, Water Resources Research, 50(9), 7167–7186, doi:10.1002/2014WR015559, 2014.

1    Harpold, A. A., Kaplan, M. L., Klos, P. Z., Link, T., McNamara, J. P., Rajagopal, S., Schumer, R. and

2        Steele, C. M.: Rain or snow: hydrologic processes, observations, prediction, and research needs,

3        Hydrology and Earth System Sciences, 21(1), 1–22, doi:10.5194/hess-21-1-2017, 2017.

4    Hopkinson, R. F., McKenney, D. W., Milewska, E. J., Hutchinson, M. F., Papadopol, P. and Vincent, L.

5        A.: Impact of Aligning Climatological Day on Gridding Daily Maximum–Minimum Temperature

6        and Precipitation over Canada, Journal of Applied Meteorology and Climatology, 50(8), 1654–

7        1665, doi:10.1175/2011JAMC2684.1, 2011.

8    Hutchinson, M. F., McKenney, D. W., Lawrence, K., Pedlar, J. H., Hopkinson, R. F., Milewska, E. and

9        Papadopol, P.: Development and Testing of Canada-Wide Interpolated Spatial Models of Daily

10       Minimum–Maximum Temperature and Precipitation for 1961–2003, Journal of Applied

11       Meteorology and Climatology, 48(4), 725–741, doi:10.1175/2008JAMC1979.1, 2009.

12   Hwang, C.L. and Yoon, K.: Multiple attribute decision making: methods and applicasions. Springer, New

13       York, 1981.

14   Islam, S. U. and Déry, S. J.: Evaluating uncertainties in modelling the snow hydrology of the Fraser River

15       Basin, British Columbia, Canada, Hydrology and Earth System Sciences, 21(3), 1827–1847,

16       doi:10.5194/hess-21-1827-2017, 2017.

17   Jun, K. S., Chung, E.-S., Kim, Y.-G. and Kim, Y.: A fuzzy multi-criteria approach to flood risk

18       vulnerability in South Korea by considering climate change impacts, Expert Systems with

19       Applications, 40(4), 1003–1013, 2013.

20   Kay, A. L., Davies, H. N., Bell, V. A. and Jones, R. G.: Comparison of uncertainty sources for climate

21       change impacts: flood frequency in England, Climatic Change, 92(1-2), 41–63,

22       doi:10.1007/s10584-008-9471-4, 2009.

23   Kirchner, J. W.: Getting the right answers for the right reasons: Linking measurements, analyses, and

24       models to advance the science of hydrology: GETTING THE RIGHT ANSWERS FOR THE

25       RIGHT REASONS, Water Resources Research, 42(3), doi:10.1029/2005WR004362, 2006.

Klyszejko, E. S.: Hydrologic Validation of Real-Time Weather Radar VPR Correction Methods, University of Waterloo., 2007.

Lee, G., Jun, K.-S. and Chung, E.-S.: Integrated multi-criteria flood vulnerability approach using fuzzy TOPSIS and Delphi technique, Natural Hazards and Earth System Science, 13(5), 1293–1312, doi:10.5194/nhess-13-1293-2013, 2013.

Leong, D. N. S. and Donner, S. D.: Climate change impacts on streamflow availability for the Athabasca Oil Sands, Climatic Change, 133(4), 651–663, doi:10.1007/s10584-015-1479-y, 2015.

Lespinas, F., Fortin, V., Roy, G., Rasmussen, P. and Stadnyk, T.: Performance Evaluation of the Canadian Precipitation Analysis (CaPA), Journal of Hydrometeorology, 16(5), 2045–2064, doi:10.1175/JHM-D-14-0191.1, 2015.

Liang, X., Lettenmaier, D. P., Wood, E. F. and Burges, S. J.: A simple hydrologically based model of land surface water and energy fluxes for general circulation model, Joural of Geophysical Reserch, 99(D7), 14,415–14,428, doi:10.1029/94JD00483, 1994.

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

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

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

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

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

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

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

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

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

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

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

Shrestha, N. K., Du, X. and Wang, J.: Assessing climate change impacts on fresh water resources of the Athabasca River Basin, Canada, Science of the Total Environment, 601, 425–440, 2017b.

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

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

1  Vincent, L. A., Wang, X. L., Milewska, E. J., Wan, H., Yang, F. and Swail, V.: A second generation of

2      homogenized Canadian monthly surface air temperature for climate trend analysis:

3      HOMOGENIZED CANADIAN TEMPERATURE, Journal of Geophysical Research:

4      Atmospheres, 117(D18), n/a–n/a, doi:10.1029/2012JD017859, 2012.

5  Wang, X. L. and Lin, A.: An algorithm for integrating satellite precipitation estimates with in situ

6      precipitation data on a pentad time scale: BLENDED PENTAD PRECIPITATION DATA, Journal

7      of Geophysical Research: Atmospheres, 120(9), 3728–3744, doi:10.1002/2014JD022788, 2015.

8  Weedon, G. P., Balsamo, G., Bellouin, N., Gomes, S., Best, M. J. and Viterbo, P.: The WFDEI

9      meteorological forcing data set: WATCH Forcing Data methodology applied to ERA-Interim

10     reanalysis data, Water Resources Research, 50(9), 7505–7514, 2014.

11  Werner, A. T. and Cannon, A. J.: Hydrologic extremes – an intercomparison of multiple gridded

12     statistical downscaling methods, Hydrology and Earth System Sciences, 20(4), 1483–1508,

13     doi:10.5194/hess-20-1483-2016, 2016.

14  Werner, A., Schnorbus, M., Shrestha, R., Cannon, A., Zwiers, F., Dayon, G., and Anslow, F.: A long-

15     term, temporally consistent, gridded daily meteorological dataset for northwestern North America,

16     Scientific Data, 6, 180299, 2019.

17  Wong, J. S., Razavi, S., Bonsal, B. R., Wheater, H. S. and Asong, Z. E.: Inter-comparison of daily

18     precipitation products for large-scale hydro-climatic applications over Canada, Hydrology and

19     Earth System Sciences, 21(4), 2163–2185, doi:10.5194/hess-21-2163-2017, 2017.

20  Zhao, K.: Validation of the Canadian Precipitation Analysis (CaPA) for hydrological modelling in the

21     Canadian Prairies, University of Manitoba (Canada)., 2013.

22

23

24

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) |

PRCP: precipitation, TMX: maximum temperature, TMN: minimum temperature, Tave: average

temperature, Tair: air temperature, WS: wind speed, RH: relative humidity, SR: solar radiation, GH:

Geopotential Height

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

7

8

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 |
| $\sigma_{ratio}$ | -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 |
| $\sigma_{ratio}$ | -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

Table 5. First ranked number of grid cells in the five sub-basins and the whole Athabasca River Basin

(ARB) and their percentages 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%) | - |
| | Firebag | 924 (97%) | 28 (3%) | 0 (0%) | 0 (0%) | - |
| (C) Multi-variables | 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%) | - |

5

6

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 SPLIN | 0.88 | 0.83 | 0.61 | 0.64 | 0.52 | 0.46 | 0.76 | 0.54 | 0.61 | 0.49 |
| Town ship | - | - | 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

Table 7. NSE values between the current study and literature for the Athabasca River basin. The NSE values were obtained for calibration and validation periods. For comparison of the current study to the literature, the NSE values of the current study were obtained from the VIC simulation for the hybrid climate dataset ($R_{ind}$).

| Stations | Current study/ VIC[1] | Literature/Hydrologic model | | | | |
|---|---|---|---|---|---|---|
| | | Shrestha et al. (2017b)/ SWAT[2] | Faramarzi et al. (2017)/ SWAT | Faramarzi et al. (2015)/ SWAT | Betrie et al. (2015)/ SWAT | Leong and Donner (2015) /IBIS-THMB[3] |
| Hinton | 0.80 | 0.87 | - | - | - | - |
| Pembina | 0.64 | 0.69 | - | - | - | - |
| Athabasca | 0.78 | 0.90 | - | - | | 0.50 |
| Fort McMurray | 0.77 | 0.89 | - | - | 0.41 | 0.35 |
| Christina | 0.52 | 0.49 | - | - | - | - |
| Firebag | 0.56 | 0.28 | - | - | - | - |
| Average for all stations | 0.58 | 0.57 | 0.21 | 0.11 | - | - |

[1] Variable Infiltration Capacity
[2] Soil and Water Assessment Tool
[3] Integrated BIosphere Simulator - Terrestrial Hydrology Model with Biogeochemistry

1

2    Table 8. Comparison of NSE values for hydrologic simulations forced by ANUSPLIN and the

3    hybrid climate datasets at the main stream of the ARB.

| 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 |

4
5
6

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

1

2    Figure 2. Structure of REFRES comprised of three modules; 1) Performance Measure Module (PMM), 2)

3                    Ranking Module (RM), and 3) Data Generation Module (DGM)

4

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

|          | OBS* | ANUSPLIN | PNWNAmet | 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        Figure 4. Temporal Correlation Coefficient (TCC) between historical precipitation data.

3         *: AHCCD data

4

(a) $D_{KS}$

(b) $P_{bias}$

(c) $\sigma_{ratio}$

(d) RMSE

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

1

2                              (e) TCC                           (f) Mean annual precipitation

3                                      Figure 5. Continued

4

(a) $D_{KS}$

(b) $P_{bias}$

(c) $\sigma_{ratio}$

(d) RMSE

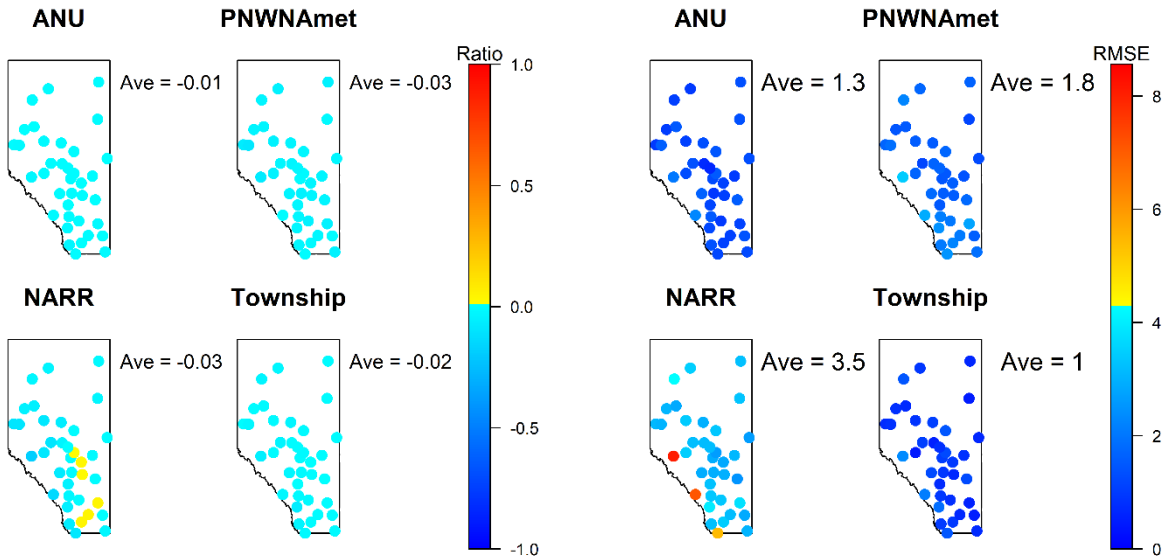Figure 6. Maps of performance measures for minimum temperature over the AHCCD stations in Alberta
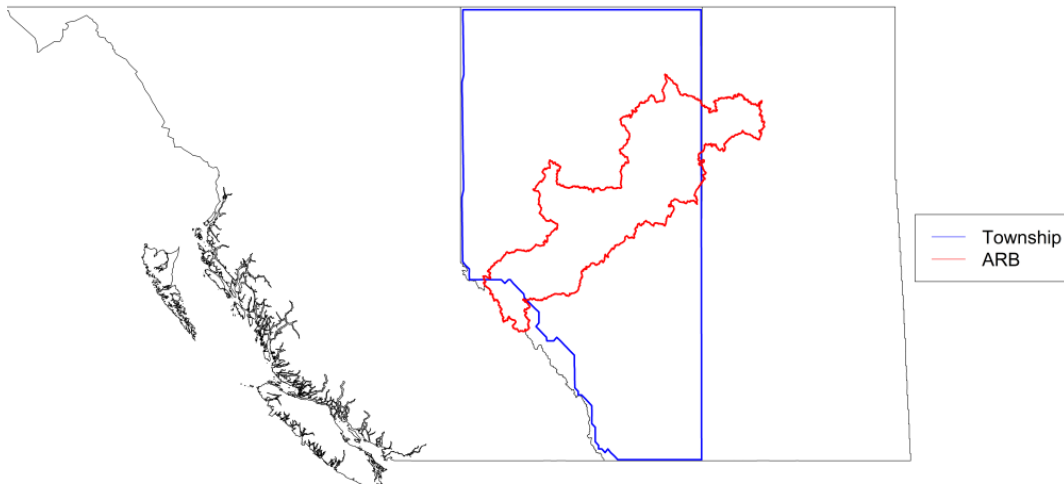
1

2                    (e) TCC                                (f) Mean annual minimum temperature

3                              Figure 6. Continued

4

1

2                          (a) D<sub>KS</sub>                                      (b) P<sub>bias</sub>

3

4                          (c) σ<sub>ratio</sub>                                   (d) RMSE

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

6

51

(e) TCC

(f) Mean annual maximum temperature

Figure 7. Continued

52

1

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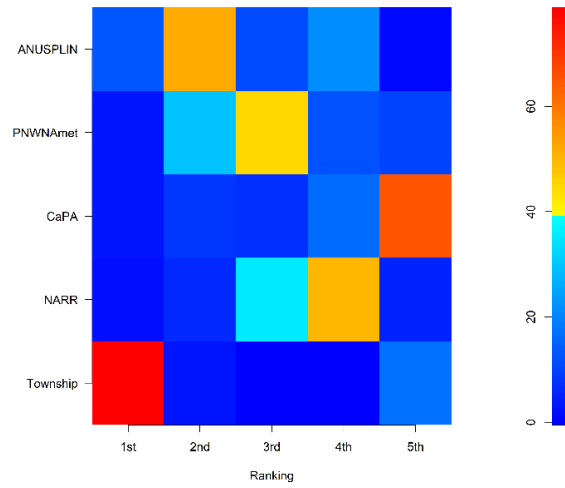
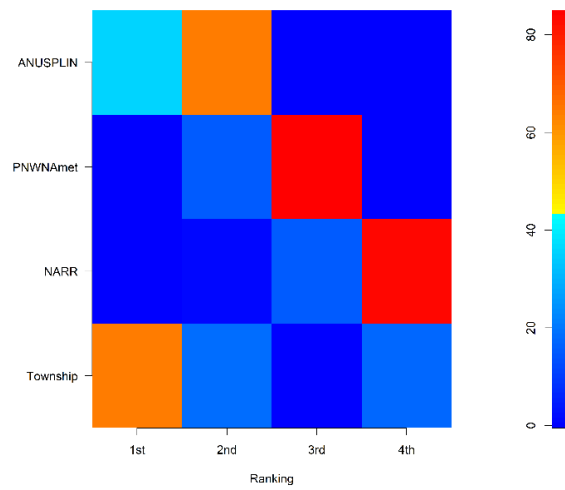3                                                              line)

4

**(A)Precipitation**  **(B)Temperature**  **(C)Multi-variables**

Legend: ● 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)

1

(a) Precipitation                    (b) Temperature

3



4

5                    (c) Multi-variables

6                    Figure 10. Percentage of climate datasets on each rank for $R_{ind}$ and $R_{mul}$
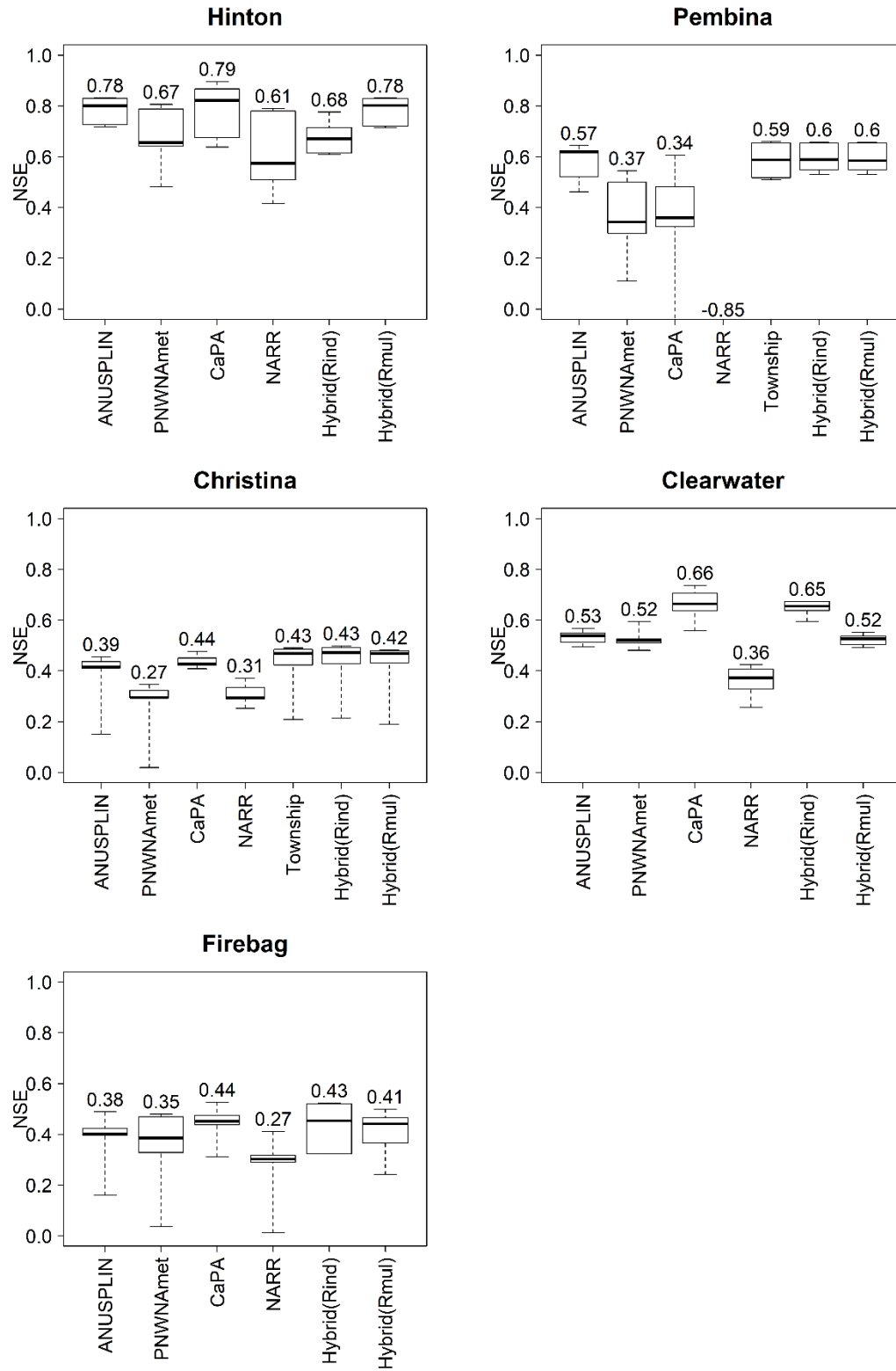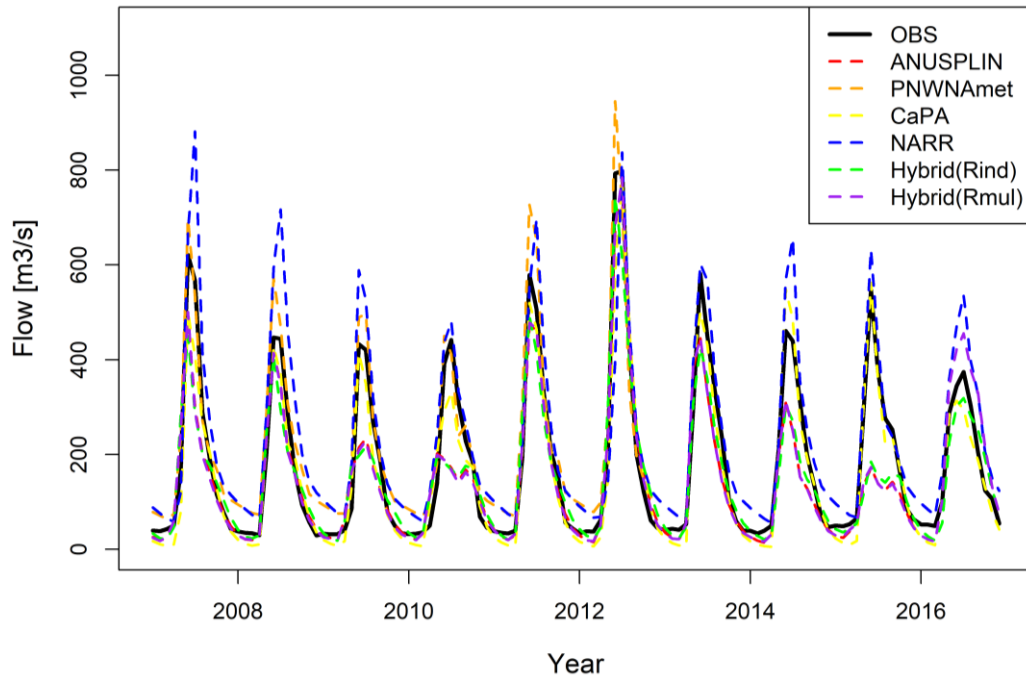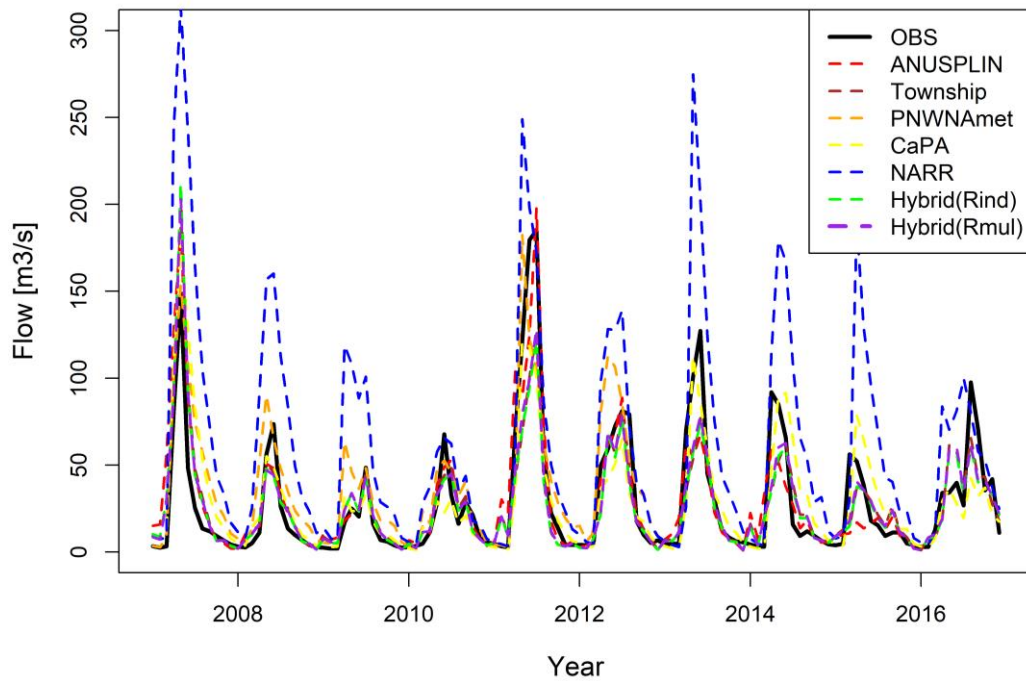
7

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

1

2                                    (a) Hinton



3

4                                    (a) Pembina

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