



**HAL**  
open science

# A century and a half precipitation oxygen isoscape for China generated using data fusion and bias correction

Jiacheng Chen, Jie Chen, Xunchang Zhang, Peiyi Peng, Camille Risi

## ► To cite this version:

Jiacheng Chen, Jie Chen, Xunchang Zhang, Peiyi Peng, Camille Risi. A century and a half precipitation oxygen isoscape for China generated using data fusion and bias correction. *Scientific Data*, 2023, 10 (1), pp.185. 10.1038/s41597-023-02095-1 . hal-04234610

**HAL Id: hal-04234610**

**<https://hal.science/hal-04234610v1>**

Submitted on 13 Oct 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License



OPEN

DATA DESCRIPTOR

# A century and a half precipitation oxygen isoscape for China generated using data fusion and bias correction

Jiacheng Chen<sup>1,2</sup>, Jie Chen<sup>1,2</sup>✉, Xunchang J. Zhang<sup>3</sup>, Peiyi Peng<sup>4</sup> & Camille Risi<sup>5</sup>

The precipitation oxygen isotopic composition is a useful environmental tracer for climatic and hydrological studies. However, accurate and high-resolution precipitation oxygen isoscapes are currently lacking in China. In this study, a precipitation oxygen isoscape in China for a period of 148 years is built by integrating observed and iGCMs-simulated isotope compositions using an optimal hybrid approach of three data fusion and two bias correction methods. The temporal and spatial resolutions of the isoscape are monthly and 50–60 km, respectively. Results show that the Convolutional Neural Networks (CNN) fusion method performs the best (correlation coefficient larger than 0.95 and root mean square error smaller than 1‰), and the other two data fusion methods perform slightly better than the bias correction methods. Thus, the isoscape is generated by using the CNN fusion method for the common 1969–2007 period and by using the bias correction methods for remaining years. The generated isoscape, which shows similar spatio-temporal distributions to observations, is reliable and useful for providing strong support for tracking atmospheric and hydrological processes.

## Background & Summary

Taking advantage of the fact that isotope composition varies sensitively with environmental conditions, environmental isotopes play an important role in the identification and characterization of the Earth's systems processes<sup>1</sup>. The study of the hydrologic cycle is one of the most important applications of stable isotopes. Firstly, the isotope composition of water provides an effective tracking method of water sources. In the process of moisture transport, the isotope composition changes with atmospheric processes, which can reflect moisture contribution<sup>2–4</sup>. In addition, for surface runoff, soil water and groundwater, the isotope composition can also reflect the water source, infiltration mechanism and evaporation consumption of each system<sup>5–8</sup>. Secondly, isotope composition can also reveal hydrological processes that cannot be achieved by other methods<sup>9</sup>. For example, evaporation processes can be better diagnosed by dual hydrogen-oxygen or triple oxygen isotope, which can be used to quantify the raindrop re-evaporation<sup>10,11</sup>. Thirdly, isotopes can be incorporated into surface hydrology models as diagnostic tools. The isotope composition of evapotranspiration, soil moisture, and runoff can be predicted by incorporating the isotope cycle, thus the distribution of isotopic variation in evapotranspiration and runoff can be better understood<sup>12</sup>. What's more, isotope composition can quantify evaporation rates, which is useful for understanding water balance and climate change from catchment to continental scales<sup>1</sup>. Precipitation isotopes can also be used to estimate the precipitation isotopic lapse rate by establishing relationships with climatic elements or elevation, so as to study paleoclimate and paleoelevation<sup>13,14</sup>.

The international observation of stable isotopes in precipitation began in the 1950s, and the Global Network of Isotopes in Precipitation (GNIP) established in 1961 provides first-hand data for the study of stable isotopes in precipitation. Since then, Austria (Austrian Network of Isotopes in Precipitation, ANIP)<sup>15</sup>, the United States

<sup>1</sup>State Key Laboratory of Water Resources & Hydropower Engineering Science, Wuhan University, Wuhan, 430072, China. <sup>2</sup>Hubei Key Laboratory of Water System Science for Sponge City Construction, Wuhan University, Wuhan, 430072, China. <sup>3</sup>USDA-ARS Oklahoma and Central Plains Agricultural Research Center, 7207W. Cheyenne St., El Reno, OK, 73036, USA. <sup>4</sup>Chongqing Southwest Research Institute for Water Transport Engineering, Chongqing Jiaotong University, Chongqing, 400016, China. <sup>5</sup>Laboratoire de Meteorologie Dynamique, IPSL, CNRS, Ecole Normale Supérieure, Sorbonne Université, PSL Research University, Paris, France. ✉e-mail: [jiechen@whu.edu.cn](mailto:jiechen@whu.edu.cn)

(United State Network of Isotopes in Precipitation, USNIP)<sup>16</sup>, Switzerland (Swiss National Network for the Observation of Isotopes in the Water Cycle, NISOT)<sup>17</sup>, Canada (Canadian Network of Isotopes in Precipitation, CNIP)<sup>18</sup> and other countries have also established their national networks, which provide strong data support for promoting and deepening the study of stable precipitation isotope.

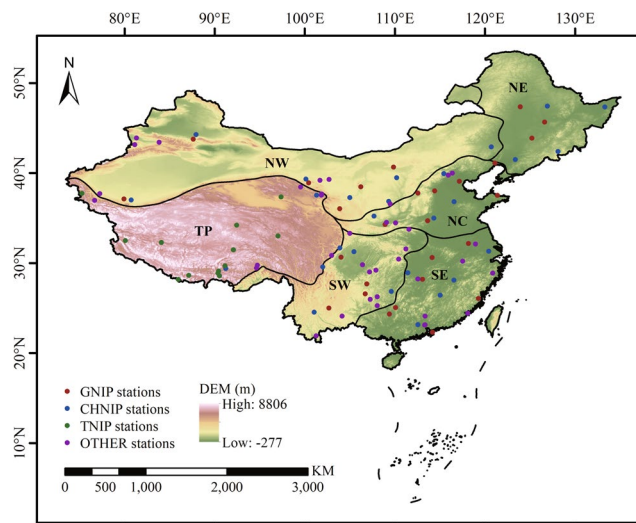
The establishment of the isotope observation network in China was relatively late. Before 1985, GNIP had only one station in Hong Kong of China, and it was not until 1985 that more stations were selected for inclusion in GNIP. Due to the scarcity of stations on the Tibetan Plateau, the Chinese Academy of Sciences (CAS) launched the Tibetan Plateau Network of Isotopes in Precipitation (TNIP) in 1991<sup>19</sup>. However, most Chinese stations in GNIP stopped monitoring in the early 2000s<sup>20</sup>, and until 2004 only one station remained. In order to continue the systematic study, the CAS established the Chinese Network of Isotopes in Precipitation (CHNIP) based on the Chinese Ecosystem Research Network (CERN) in 2004<sup>21</sup>. Due to the difficulty and high cost of measuring precipitation isotope ratios<sup>22</sup>, most of the observed data are short in length. The spatial distribution of observation stations is uneven, with few stations in inaccessible areas<sup>23</sup>.

The stable isotopes in precipitation can also be simulated by isotope-equipped general circulation models (iGCMs). In contrast to observations, iGCMs can provide time-continuous and space-regular isotope data<sup>24</sup>. Jousaume, *et al.*<sup>25</sup> incorporated the fractionation process of water stable isotope into GCM for the first time. They used the GCM of the Laboratoire de Météorologie Dynamique (LMD) to simulate the distribution of global water stable isotope, and the relationship between simulated precipitation oxygen isotope and meteorological elements was in good agreement with the measured results. Since then, an increasing number of GCMs have incorporated isotope cycles, for example, the ECHAM4 developed by the Max Planck Institute for Meteorology (MPI) in Germany<sup>26</sup>, the GISS E developed by the NASA Goddard Institute for Space Studies (GISS) in the United States<sup>27,28</sup>, the HadAM3 developed by the Hadley Centre for Climate Prediction and Research in the United Kingdom<sup>29</sup>, the LMDZ4 developed by the Laboratoire de Météorologie Dynamique in France<sup>30</sup>, and the MIROC32 developed by the Center for Climate System Research (CCSR) of the University of Tokyo in Japan<sup>31</sup>, etc.

On the basis of these, the comparison and evaluation of iGCMs in simulating isotopes have been conducted in many studies. Yoshimura, *et al.*<sup>32</sup> indicated that due to the limitation of spatial and temporal resolution, iGCMs are poor in simulating the short-term (days) variability of stable isotopes in precipitation, while they are good at the monthly or annual scale. Conroy, *et al.*<sup>33</sup> evaluated the spatio-temporal pattern of precipitation isotope variability in the tropical Pacific for iGCM simulations, and found that nudging models by reanalysis wind has a certain effect on precipitation isotope values, and the performance of models varies with regions. Zhang, *et al.*<sup>34</sup> selected four iGCMs to evaluate the average precipitation isotopic composition in East Asia. The results showed that the characteristics of measured values were well reproduced by iGCM simulation, but the simulated values were all lower in the inland areas at middle and high latitudes, and the amount effect in arid areas was incorrectly simulated. Wang, *et al.*<sup>23</sup> verified iGCM-simulated stable isotopes in precipitation in arid Central Asia. In general, the seasonality of stable isotopes in precipitation could be well simulated, but the values of oxygen isotopes were higher in summer and lower in winter, lower in the eastern section and higher in the western section. Che, *et al.*<sup>35</sup> concluded that nudged simulation by LMDZ has the best comprehensive performance by comparing the simulated values of different models with the measured values of GNIP in China. In terms of altitude effects, CAM and GISS E perform better, while in terms of continental effects, the free simulations by GISS E and LMDZ perform better.

To comprehensively consider the error characteristics and advantages of different sources of data to reduce uncertainty, data fusion is usually used. One of the common methods for data fusion is to use *in-situ* observations as baselines to correct estimates from other sources. Several data fusion methods such as cokriging<sup>36</sup>, probability matching<sup>37</sup>, statistical objective analysis<sup>38</sup>, Bayesian correction<sup>39</sup>, probability density function–optimal interpolation<sup>40</sup>, and variational<sup>41</sup> are usually used to fuse *in-situ* observation information. The key of these methods is to deal with the estimation errors directly based on weighted average, regression analysis, filtering analysis and other mathematical approaches. In contrast, neural network methods have stronger learning and generalization abilities, and have advantages in discovering complex relationships in data and processing large amounts of data<sup>42</sup>. So far, the neural network was mainly applied to precipitation data fusion in the field of hydrology but very little in isotopic hydrology. For example, Turlapaty, *et al.*<sup>43</sup> used Artificial Neural Network to fuse various satellite precipitation products, and found the fusion performance was statistically superior to each individual dataset for all seasons. Sun and Tang<sup>44</sup> combined information from satellite precipitation products and reanalysis data in Central Texas, U.S., by using an attention-based deep convolutional neural network (AU-Net), and found the Au-net models have achieved varying degrees of success under different climatic conditions. Wu, *et al.*<sup>45</sup> combined Convolutional Neural Network with Long Short-Term Memory Network to fuse the TRMM satellite data, thermal infrared images of Gridded satellite, rain gauge data and elevation data. The results showed that this method can improve the accuracy of original TRMM data in China, even for regions with different precipitation intensities or sparse gauges.

Overall, both observations and iGCM simulations have advantages and disadvantages. The effort of constructing a database by taking advantages and circumventing disadvantages of both becomes a challenge. With the motivations of resolving the lack and uneven distribution of observations, as well as the coarse and biased iGCM simulations, this study aims to take a hybrid approach that makes full use of observations to integrate the advantages of various iGCMs by using the optimal combination of data fusion and bias correction methods. In order to determine the best scheme to build the dataset, two bias correction methods (BCMs) and three neural network data fusion methods (DFMs) are first compared in terms of bias correcting and fusing iGCM simulations. The new isoscape in monthly temporal and approximately 0.5° spatial resolutions is produced by combining the optimal data fusion and bias correction methods for the 1870–2017 period. The spatial and temporal distribution characteristics of oxygen isotopes in precipitation are then analysed for China.



**Fig. 1** Map of the station locations and topography in the mainland of China. The dots indicate the distribution of isotope observation stations, with different colours representing different sources. The six sub-regions are plotted (NE – Northeast China, NC – North China, SE – Southeast China, SW – Southwest China, TP – Tibetan Plateau, NW – Northwest China).

## Methods

**Study area.** China is located in the east of Eurasia and on the west coast of the Pacific Ocean. The topography of China generally presents three steps descending to the east. The climate is complex and diverse in China. Heavily influenced by the continents and oceans, the monsoon climate is significant, especially for the east of China. The spatiotemporal variation of precipitation stable isotopes is very complex due to the significant changes in winter and summer circulation<sup>46</sup>. The mainland of China can be geographically classified into three sub-regions, the eastern monsoon region, the arid northwest region and the Qinghai-Tibet Plateau region, according to topography, climate, soil and vegetation. The eastern monsoon region is further divided into four sub-regions by taking terrain and climatic conditions into account, as well as ensuring the sufficient number of stations and data volumes in a sub-region. To sum up, the study area is divided into six sub-regions for our analysis: Northeast China (NE), North China (NC), Southeast China (SE), Southwest China (SW), Qinghai-Tibet Plateau (TP) and Northwest China (NW), as shown in Fig. 1.

**Datasets.** There are 107 oxygen isotope observation stations in the study area (Fig. 1), including 29 GNIP stations<sup>47</sup> (available at <https://nucleus.iaea.org/wiser>), 27 CHNIP stations<sup>46</sup>, 13 TNIP stations<sup>48,49</sup> and 38 stations from other sources (mainly from references). Monthly oxygen isotope composition of precipitation ( $\delta^{18}\text{O}_p$ ) is used for analysis. For a few references providing event isotope data, monthly precipitation weighted data are used. The time span of GNIP data mostly ranges from 1980 to 2000, and the length of the time period is basically 5–15 years. For CHNIP, most time periods are about 2–5 years ranging from 2005 to 2010. Most TNIP data are between 1995 and 2005, with varying lengths. More details about the observation stations can be found in Table S1.

Some physical-based ancillary data are introduced in the fusion methods, including elevation and meteorological data, to enrich the climate and terrain information in the process of data fusion. For regions with high spatial autocorrelation of isotope values, and regions with large variability in topography or climate, it is very necessary to introduce ancillary data<sup>50,51</sup>. Moreover, the neural network methods have the ability to integrate simulated data, observed data and ancillary data to extract enough spatial variability to generate more accurate oxygen isoscapes<sup>52</sup>. The Digital Elevation Model (DEM) with a spatial resolution of 90 m is derived from the USGS/NASA Shuttle Radar Topographic Mission<sup>53</sup> (SRTM, <https://srtm.csi.cgiar.org/>). The monthly temperature and precipitation gridded datasets on 0.5° spatial resolution in China developed by the National Meteorological Information Centre (<http://www.nmic.cn/>) are also used. The data are generated by the Thin Plate Spline interpolation method based on *in-situ* temperature and precipitation data in China. It should be noted the DEM ancillary data is not used for data fusion in TP, because the observation stations are distributed at lower altitudes, some isotope simulations of the grid points with higher altitudes are unreasonable in the fusion methods.

Seven  $\delta^{18}\text{O}_p$  spatio-temporal fields simulated by five iGCMs (CAM2, GISS E, HadAM3, LMDZ4 and MIROC32) are used, which are selected from the SWING2 archive (available at <https://data.giss.nasa.gov/swing2/>). Five of eight simulations are free-running, performed following the Atmospheric Model Intercomparison Project (AMIP) protocol, using prescribed sea surface temperatures (SST) and sea ice<sup>30,54</sup>. The remaining three (GISS E, IsoGSM2 and LMDZ4) are nudged to constrain large-scale atmospheric circulation, so that the dynamical fields in simulations are close to the observations<sup>54</sup>. In addition to SWING2 simulations, a zoomed simulation by LMDZ4, with the horizontal resolution of 50–60 km<sup>30,55</sup> and nudged to reanalyses, is used. The zoomed LMDZ4 simulations used the HPC resources of IDRIS under the allocation 0292 made by GENCI.



GCM	Simulation method	Horizontal resolution (longitude × latitude)	Time period	Key references
CAM2	Free-running	2.8° × 2.8°	1958–2003	Lee, <i>et al.</i> <sup>111</sup>
GISS E	Free-running and nudged by NCEP	2.5° × 2°	1969–2009	Schmidt, <i>et al.</i> <sup>27</sup>
HadAM3	Free-running	3.75° × 2.5°	1870–2001	Tindall, <i>et al.</i> <sup>29</sup>
IsoGSM2	nudged by NCEP	1° × 1°	1979–2017	Yoshimura, <i>et al.</i> <sup>54</sup>
LMDZ4	Free-running and nudged by ECMWF	3.75° × 2.5°	1979–2007	Risi, <i>et al.</i> <sup>30</sup>
	Zoomed (nudged by ECMWF)	50–60 km	1979–2017	Gao, <i>et al.</i> <sup>55</sup>
MIROC32	Free-running	2.8° × 2.8°	1979–2007	Kurita, <i>et al.</i> <sup>31</sup>

**Table 1.** Time periods and basic outputs information of selected iGCMs.

The isoGSM version 2 nudged to reanalysis with the horizontal resolution of 1° is also used<sup>54,56,57</sup>. Totally, nine simulations from six iGCMs are used, and detailed information about these iGCMs can be found in Table 1. In general, this study makes the maximum extent to use isotope observations and simulation in data fusion.

The time span of the isoscape built in this study covers the union set of all simulations ranging from 1870 to 2017. Since the temporal lengths of nine iGCM are not identical, the number of iGCM simulations used to build the isoscape varies. Specifically, for 1979–2001, a total of nine simulations from all six iGCMs in Table 1 are used; for 2002–2007, seven simulations from four iGCMs (GISS E, IsoGSM2, LMDZ4 and MIROC32) are used; and for 1969–1978, four simulations from three iGCMs (CAM2, GISS E and HadAM3) are used. For the remaining periods, there is only one simulation or two: for 1958–1968, CAM2 and HadAM3 are used; for 1870–1957, HadAM3 is used; and for 2008–2017, IsoGSM2 and zoomed LMDZ4 are used.

The stable isotope composition of precipitation is expressed in the relative permillage (‰) derived from the standard sample<sup>58</sup> as:

$$\delta = \left( \frac{R_{\text{sample}}}{R_{\text{V-SMOW}}} - 1 \right) \times 1000\text{‰} \quad (1)$$

where  $R$  is the ratio of heavier isotope to common isotope ( $^{18}\text{O}/^{16}\text{O}$ ), and the subscripts sample and V-SMOW represent standard sample and Vienna Standard Mean Ocean Water, respectively.

**Generation of isoscape.** Generally, the generation of isoscape can be divided into five steps.

- (1) Prior to generating the dataset, the inverse distance weighting (IDW) method is used to interpolate all iGCM simulations and ancillary data to observation stations.
- (2) Three neural network data fusion and two bias correction methods are trained using observations and iGCM simulations for all months within a season and all stations within a sub-region. In other words, observed and simulated monthly isotopes within a season and all stations within a region are pooled to train the data fusion and bias correction methods to ensure that the model is well trained with enough samples. For the fusion methods, ancillary data are also included in the training process, which is not necessary for bias correction methods.
- (3) The performance of each model is evaluated for the validation period by the cross-validation method to find the optimal data fusion and bias correction methods. Correlation coefficient (CC) and root mean square error (RMSE) are used as metrics to validate these methods for the common period of 1969–2007.
- (4) All iGCM simulations and ancillary data are interpolated to the LMDZ4 zoomed grid with a spatial resolution of approximately 50 km by the IDW method.
- (5) The optimal trained model and bias correction methods are applied to all grid points within a sub-region and all months within a season. Since the length of iGCM simulations is not identical, the optimal combination of data fusion and bias correction methods are used to generate the isoscape for a long period. In other words, for the common period of observations and iGCM simulations, the optimal data fusion method is used, while for the period with no observations, the bias correction methods are used.

**Neural network data fusion.** The neural network is a kind of mathematical model, which imitate the behaviour characteristics of human neural network and carry out distributed parallel computing<sup>59,60</sup>. Performing calculations and spreading information through large numbers of interconnected neurons, neural networks are often used to describe complex relationships between inputs and outputs, or to explore the internal structure and patterns of data<sup>61,62</sup>. In this study, Back Propagation Neural Network (BP), Long Short-Term Memory (LSTM) Neural Network and Convolutional Neural Network (CNN) are adopted for data fusion, considering BP's simplicity and practicality, LSTM's advantages in time series prediction and CNN's outstanding performance in various fields.

The structure and hyperparameters of these neural network methods are carefully considered and validated. Considering that previous studies<sup>63–66</sup> on hyperparameter sensitivity of neural networks have shown similar results, the hyperparameter selection scheme is determined based on these studies using a hierarchical stepwise search method to determine hyperparameter values. Specifically, the hyperparameters are divided into three parts, structural hyperparameters, sensitive algorithm hyperparameters and other algorithm hyperparameters, which are determined step by step. At each step, the performance of all hyperparameter combinations is tested using the grid search method. Referring to previous studies<sup>45,67–69</sup>, some conventional hyperparameter settings (such as filters are usually set to the power of 2) are considered. The details of hyperparameter selection

Model	Steps	Hyperparameters	Range tested	Selected
BP	Step 1. Structural hyperparameters	Hidden layers	2, 3, 4	3
		Learning rate	0.0001–0.005	0.005
	Step 2. Sensitive algorithm hyperparameters	Batch size	10–50	20
		Dense neurons	8, 16, 32, 64	16/32/64
Step 3. Other algorithm hyperparameters	Activation	ReLU, TanH	ReLU	
LSTM	Step 1. Structural hyperparameters	LSTM layers	2, 3, 4	3
		Learning rate	0.0001–0.005	0.001
	Step 2. Sensitive algorithm hyperparameters	Batch size	10–50	50
		LSTM neurons	8, 16, 32, 64	32
	Step 3. Other algorithm hyperparameters	Time steps	1–5	2
		Dropout rate	0–0.3	0.1
Activation		ReLU, TanH	TanH	
CNN	Step 1. Structural hyperparameters	Convolutional layers	1, 2, 3	2
		Dense layers	1, 2	1
	Step 2. Sensitive algorithm hyperparameters	Learning rate	0.0001–0.005	0.0005
		Batch size	10–50	50
		Kernel size	3, 4, 5	4
	Step 3. Other algorithm hyperparameters	Filters	8, 16, 32, 64	8/32
		Dense neurons	8, 16, 32, 64	16
		Dropout rate	0–0.3	0
		Activation	ReLU, TanH	ReLU

**Table 2.** Details of the hyperparameter selection in three neural network fusion methods. The complete hyperparameter setting of the neural network fusion methods can be seen in Table S2.

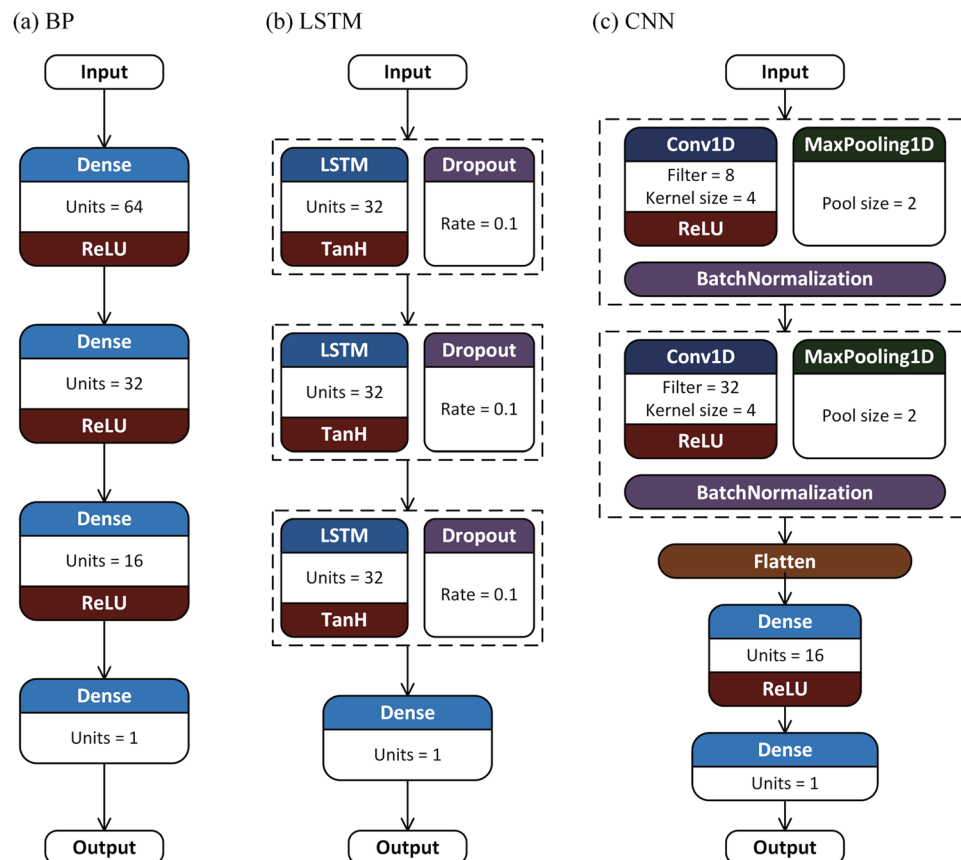
are shown in Table 2. Furthermore, when structural hyperparameter values produce similar performance, the simpler structure (i.e., the one with fewer hyperparameters) is chosen to avoid overfitting. When algorithm hyperparameter values produce similar performance, the one that is more efficient for computing is chosen. The structures of these three neural network DFMs are presented in Fig. 2.

BP, first proposed by Rumelhart, *et al.*<sup>70</sup>, is a multilayer feed-forward network trained by the error back-propagation algorithm. BP is one of the most widely used neural network models with high simulation accuracy for nonlinear functions. The main characteristic of BP is that the input signal is processed layer by layer from the input layer to the hidden layer and then to the output layer, and each neuron carries out the weighted sum of the input signal through the activation function. If the error between the actual output and the expected output is larger than the set value, the weight and bias of the network are continuously corrected by the back-propagation to minimize the loss function. The isotope simulations of iGCMs are selected as the input of BP and the observations as the expected output to calculate the loss function. The input layer is the corresponding input parameter, and the output layer is the fusion isotope value.

LSTM is very efficient for sequential data and is derived from Recurrent Neural Network (RNN) with memory function. RNN has a sequential feed-forward connection, so that the information of the past moment can affect the output of the present moment<sup>71</sup>. The traditional RNN has the problems of vanishing gradient and exploding gradient<sup>72</sup>. To solve these problems, the LSTM Neural Network was proposed<sup>73</sup>. A basic LSTM neuron usually consists of a memory cell and three gates (i.e. input gate, forget gate and output gate). Memory cells are used to store past information, realizing long-distance dependent learning of sequence features. The input gate determines which inputs are saved to the cell; the forget gate determines what information is retained from the previous moment; output gate determines what information needs to be output. In this study, a fully connected layer is added after three LSTM layers to generate fusion results. Dropout layers are applied to the three LSTM layers of the network to make the model more robust.

CNN was first proposed by LeCun<sup>74</sup>, for the problem of handwritten digit recognition. CNN combines three advantages of local connectivity, weight sharing and pooling. On one hand, it reduces the number of weights, making the network easy to optimize. On the other hand, it reduces the complexity of the model and alleviates the overfitting problem. It is one of the most widely used neural networks with the best performance. The convolutional layer and pooling layer of the hidden layer are the core modules of CNN. The function of the convolutional layer is to extract features of the input data by convolutional kernels. The pooling layer performs feature selection and information filtering on the feature map output by the convolution layer. In this study, CNN is mainly composed of two convolutional layers and pooling layers. Two fully connected layers are added at the end to remove the spatial topology and output the results. The convolutional layer and the fully connected layer are connected by flattening the output of the convolutional layer through the flatten layer. Batch normalization layers are inserted into the model to improve the speed, performance and stability of the neural network.

**Bias correction methods.** BCs aim to correct the mean, variance and/or quantile of the climate model time series, so that the corrected model time series can better match those of the observations<sup>75</sup>. In this study, two typical methods (i.e. linear scaling (LS) and distribution translation (DT)) are used to correct the bias of iGCM at



**Fig. 2** Network structure of BP, LSTM and CNN fusion method. The complete hyperparameter setting of the neural network fusion methods can be seen in Table S2.

the monthly timescale. These two BCMs can be classified into mean-based scaling (i.e. LS) and distribution-based correction (i.e. DT) approaches<sup>76</sup>. The mean-based scaling uses a constant correction factor for the entire time series, while the distribution-based approach uses correction factors that vary with the quantiles of the distribution<sup>77</sup>.

The LS method is the simplest bias correction method. The differences between observations and raw iGCM simulations are applied to simulations to obtain the bias-corrected isotope time series for each season and sub-region. Specifically, for a particular sub-region, the differences (defined as correction factors) in mean values between observed and simulated isotopes are first calculated at the seasonal basis using Eq. (2). The calculated correction factors are then applied to simulated isotopes for the entire period using Eq. (3).

$$R_{LS,s,sr} = \overline{\delta O}_{obs,s,sr}^{ref} - \overline{\delta O}_{raw,s,sr}^{ref} \quad (2)$$

$$\delta O_{cor,s,sr} = \delta O_{raw,s,sr} + R_{LS,s,sr} \quad (3)$$

where  $R_{LS}$  is the correction factor;  $\overline{\delta O}$  is the mean value of isotope composition; the superscript *ref* represents the reference period; the subscripts *obs*, *raw* and *cor* represent observations, raw simulations and corrected simulations, respectively; and *s* and *sr* represent a specific season and a sub-region, respectively.

The implementation of the DT method is similar to the LS method. However, the differences (i.e. correction factors) between observed and simulated isotopes are calculated for each of 100 integral percentiles as shown in Eqs. (4–6), to represent the distribution for each season in each sub-region. The correction factors of grid points are obtained by interpolating or extrapolating the factors of observation stations using Eq. (5).

$$R_{DT,s,sr}^{ref} = \delta O_{obs,q,s,sr}^{ref} - \delta O_{raw,q,s,sr}^{ref} \quad (4)$$

$$R_{DT,s,sr}^{ref} \xrightarrow{\text{Interpolation/Extrapolation}} R_{DT,s,sr} \quad (5)$$

$$\delta O_{cor,s,sr} = \delta O_{raw,s,sr} + R_{DT,s,sr} \quad (6)$$

HadAM3 LS & DT ensemble average	CAM2 & HadAM3 LS & DT ensemble average	CAM2, GISS E, HadAM3 & ancillary data CNN fusion	CAM2, GISS E, HadAM3, IsoGSM2, LMDZ4, MIROC32 & ancillary data CNN fusion	GISS E, IsoGSM2, LMDZ4, MIROC32 & ancillary data CNN fusion	IsoGSM2 & LMDZ4 zoomed LS & DT ensemble average
1870-1957	1958-1968	1969-1978	1979-2001	2002-2007	2008-2017

**Fig. 3** The generation mode of dataset in each period.

where the subscript  $q$  is a percentile for a specific season in a sub-region. Other superscripts and subscripts are the same as Eqs. (2, 3).

**Model performance.** Correlation coefficient (CC) and root mean square error (RMSE) are used as metrics to quantify model performance.

$$CC = \frac{1}{N-1} \sum_{i=1}^N \left( \frac{O_i - \mu_O}{\sigma_O} \right) \left( \frac{S_i - \mu_S}{\sigma_S} \right) \quad (7)$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{i=1}^N (S_i - O_i)^2} \quad (8)$$

where  $S$  and  $O$  are the simulated and observed value, respectively;  $N$  is the number of samples; and  $\mu$  and  $\sigma$  are the mean and standard deviation, respectively.

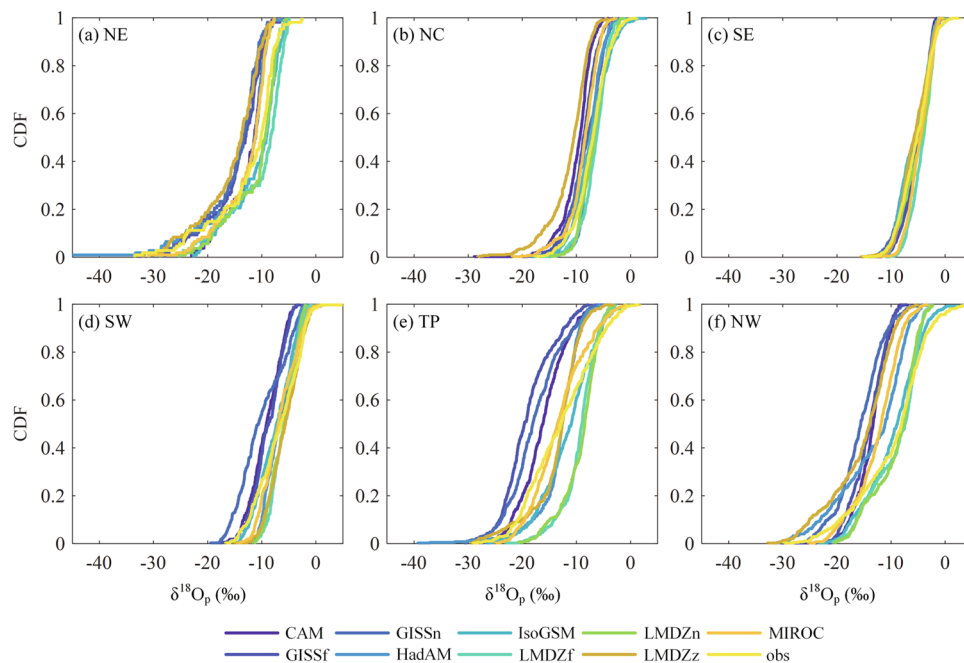
**Cross-validation experiments.** In order to make full use of the data and reduce the variation of model accuracy caused by the difference between the training set and the test set, K-fold cross-validation is adopted. In the K-fold ( $K=5$  in this study) experiment, the data set is randomly divided into  $K$  groups, and one of them is used as the test set each time, leaving  $K-1$  groups as the training set. To fully consider the variations of random division, K-fold cross-validation is repeated 100 times.

## Data Records

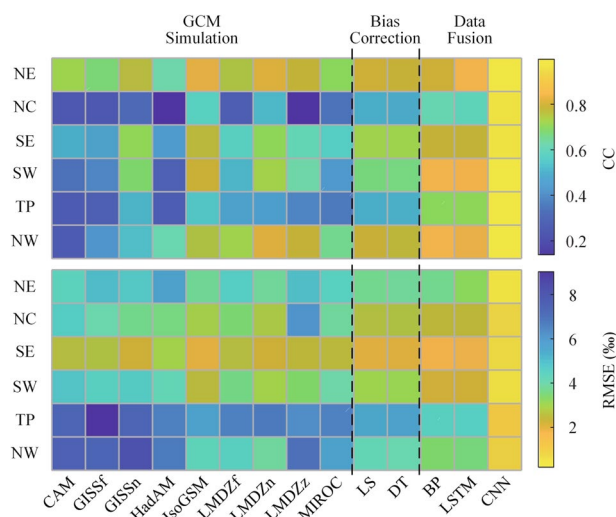
The dataset includes the stable oxygen isotope of precipitation for the mainland of China over the 1870–2017 period, at a spatial resolution of 50–60 km and a monthly temporal resolution. In order to make full use of observations to integrate the advantages of various iGCMs, the combination of data fusion and bias correction methods are used (as shown in Fig. 3). Specifically, (1) for the 1979–2001 period, nine simulations from six iGCMs (CAM2, GISS E, HadAM3, IsoGSM2, LMDZ4, and MIROC32) and ancillary data are fused with observations by using CNN fusion method; (2) for the 2002–2007 period, seven simulations from four iGCMs (GISS E, IsoGSM2, LMDZ4, and MIROC32) and ancillary data are fused by using CNN fusion method; (3) for the 1969–1978 period, four simulations from three iGCMs (CAM2, GISS E, and HadAM3) and ancillary data are fused by using CNN fusion method; (4) for the 1958–1968 and 2008–2017 periods, two iGCM simulations (CAM2 and HadAM3 for 1958–1968 and IsoGSM2 and LMDZ4 zoomed for 2008–2017) are corrected by using two BCs, and ensemble mean (mean of four simulations) is then calculated; (5) for the 1870–1957 period, one iGCM simulation (HadAM3) is corrected by using two BCs, and ensemble mean (mean of two simulations) is then calculated. The dataset<sup>78</sup> is freely available in Zenodo repository (<https://doi.org/10.5281/zenodo.7306199>) with the format of netCDF4.

## Technical Validation

**Evaluation of bias correction and data fusion methods.** Prior to applying BCs and DFMs to build the isoscape, the performance of iGCM simulations is evaluated by comparing gauged observations for the common period of 1969–2007. Figure 4 shows the cumulative distribution functions (CDFs) of  $\delta^{18}\text{O}_p$  for observations and iGCMs simulations in each sub-region. Generally, the CDFs of observed  $\delta^{18}\text{O}_p$  can be well represented by iGCM simulations for each sub-region, as the observed CDFs distribute in the centre of simulated ones. For specific regions, the envelope of CDFs is the narrowest for SE, indicating that iGCMs perform consistently better for this region. For NE, NC, SW and NW, the CDFs of IsoGSM2 and LMDZ4 (free and nudged) simulated  $\delta^{18}\text{O}_p$  are relatively close to the observations, while other iGCM simulations generally overestimate the  $\delta^{18}\text{O}_p$ . For TP, the differences between CDFs of observed and simulated  $\delta^{18}\text{O}_p$  are the largest, indicating the prominent variability of  $\delta^{18}\text{O}_p$  simulations. This is expected, as climate models generally perform worse for TP than other regions<sup>79,80</sup>. For NW, the variability of  $\delta^{18}\text{O}_p$  simulations is also large. This is because, on the one hand, the sparse coverage of stations coupled with complex topography over northwest China cannot well represent the full range of precipitation isotope conditions. This can lead to biases in the distribution of observations. On the other hand, the arid northwest region is one of the most sensitive regions to climate change due to its fragile ecosystem, which affects sub-cloud evaporation and local moisture re-cycling, leading to the large uncertainty in isotope simulation between different iGCMs<sup>23,81</sup>.



**Fig. 4** Cumulative distribution functions of  $\delta^{18}\text{O}_p$  for eight iGCM simulations in six sub-regions.



**Fig. 5** Average correlation coefficient (CC) and root mean square error (RMSE) metrics of raw, bias-corrected and fused  $\delta^{18}\text{O}_p$  in six sub-regions.

Root mean square error (RMSE) and correlation coefficient (CC) are also calculated to evaluate the accuracy of iGCM simulated  $\delta^{18}\text{O}_p$ . DFMs that introduce ancillary data generally perform better than DFMs that do not introduce that, with larger CC and smaller RMSE. Therefore, ancillary data are introduced into all DFMs in this study. Figure 5 presents the RMSE and CC for raw iGCM simulations, bias corrected and fused simulations for six sub-regions over the validation periods (1979–2001). Generally, DFMs perform better than BCMs, and both perform better than raw iGCM simulations. In addition, all the simulations are correlated with the observations with CC ranging between 0.12 and 0.99.

The CC and RMSE vary considerably for raw iGCM simulations, with CC ranging from 0.12 to 0.84 and RMSE ranging from 1.9‰ to 9.1‰. The simulations of IsoGSM2, nudged GISS E and nudged LMDZ4 have the strongest correlation with the observations, and their CC is basically above 0.5, ranging from 0.30–0.84. The error of IsoGSM2 and nudged LMDZ4 is the smallest, and their RMSE ranges between 1.9‰ and 6.8‰.

Generally, the performance of LS and DT is similar, even though the LS method performs slightly better than the DT method for some regions. The CCs range from 0.49 to 0.81 for the LS and DT corrected simulations, with an average increase from 0.53 to 0.66 relative to raw simulations; the RMSEs are between 1.9‰ and 5.7‰, with an average decrease of 23.7%.



For DFMs, BP and LSTM show similar performance, while CNN consistently performs the best. The CC of CNN-generated simulations is all greater than 0.97, increasing from 0.53 to 0.99 on average compared with the raw simulations. The RMSE of CNN-generated simulations is all smaller than 1.1%, showing an 84.3% reduction relative to the raw simulation on average.

Figure 5 shows that CC and RMSE of simulations in different sub-regions are quite different. For all raw, fused, and bias corrected simulations, CC is smaller for NC and TP while RMSE is larger for TP and NW than other sub-regions. For DFMs, especially CNN, the differences of CC and RMSE between different regions are smaller.

Generally, all simulations perform worse in NC, TP and NW than in other sub-regions. The poor simulation performance in NC may be due to complex air mass movements<sup>82,83</sup>, which are difficult to be accurately simulated by iGCMs. The poor simulation performance in TP and NW may be due to the fact that GCMs cannot accurately describe the atmospheric physical process and simulate precipitation and other meteorological factors in these regions<sup>84–86</sup>.

The performance of BCMs and DFMs is also evaluated for three periods (1969–1978, 1979–2001 and 2002–2007) and six sub-regions on a seasonal basis. The seasonal average CC and RMSE for two BCMs and three neural network DFMs are presented in Figs. 6, 7. The seasons are divided into spring (SPR), summer (SUM), autumn (AUT) and winter (WIN), respectively defined as March–May, June–August, September–November, and December–February. Generally, all BCMs and DFMs perform very similarly for all three periods.

As for CC, the simulations in the northern region (NE, NC, and NW) show strong correlations with observations in spring, while those in the southern region (SE, SW) show strong correlations in summer and autumn. The correlation of CNN fusion simulations is significantly higher than that of the other methods, with CC being mostly above 0.95. The correlation of BP and LSTM fusion simulations is slightly higher than that of LS and DT corrected simulations, with CC being mostly between 0.3 and 0.7, varying with sub-regions and seasons. The BP and LSTM fusion methods perform slightly worse in NE, but better in other sub-regions, compared with BCMs.

As for RMSE, the errors of simulations in NE, TP and NW are relatively large, with an average error of about 5‰, while those in other sub-regions have an average error of about 3‰ or less. The northern region shows a small error of about 2‰ in summer, except for NW with a larger error of 3‰; while the southern region shows relatively small seasonal differences in error, mostly ranging between 1‰ and 3‰. On the whole, the DFMs perform better than the two BCMs. The errors of CNN simulations are the smallest in all regions and seasons, which are mostly smaller than 1‰. The errors of BP and LSTM simulations are slightly smaller than those of LS and DT simulations.

The plotted  $\pm$  one standard deviation shows the dispersion degree of CC and RMSE for the simulated results of all bias correction and data fusion methods over 100 trials. The standard deviations of CC and RMSE calculated by LS and DT corrected simulations are smaller, while those calculated by BP, LSTM and CNN fused simulations are relatively larger. Generally, the standard deviation of CC and RMSE calculated by the CNN fused simulations is the smallest among fusion methods. It can be considered that the correction methods show smaller uncertainties than the fusion methods in terms of CC and RMSE. This is as expected, since the simulations of DFMs show uncertainties brought by the neural network itself, in addition to the uncertainties brought by cross-validation. Furthermore, CNN fusion methods show smaller uncertainties than the other two fusion methods.

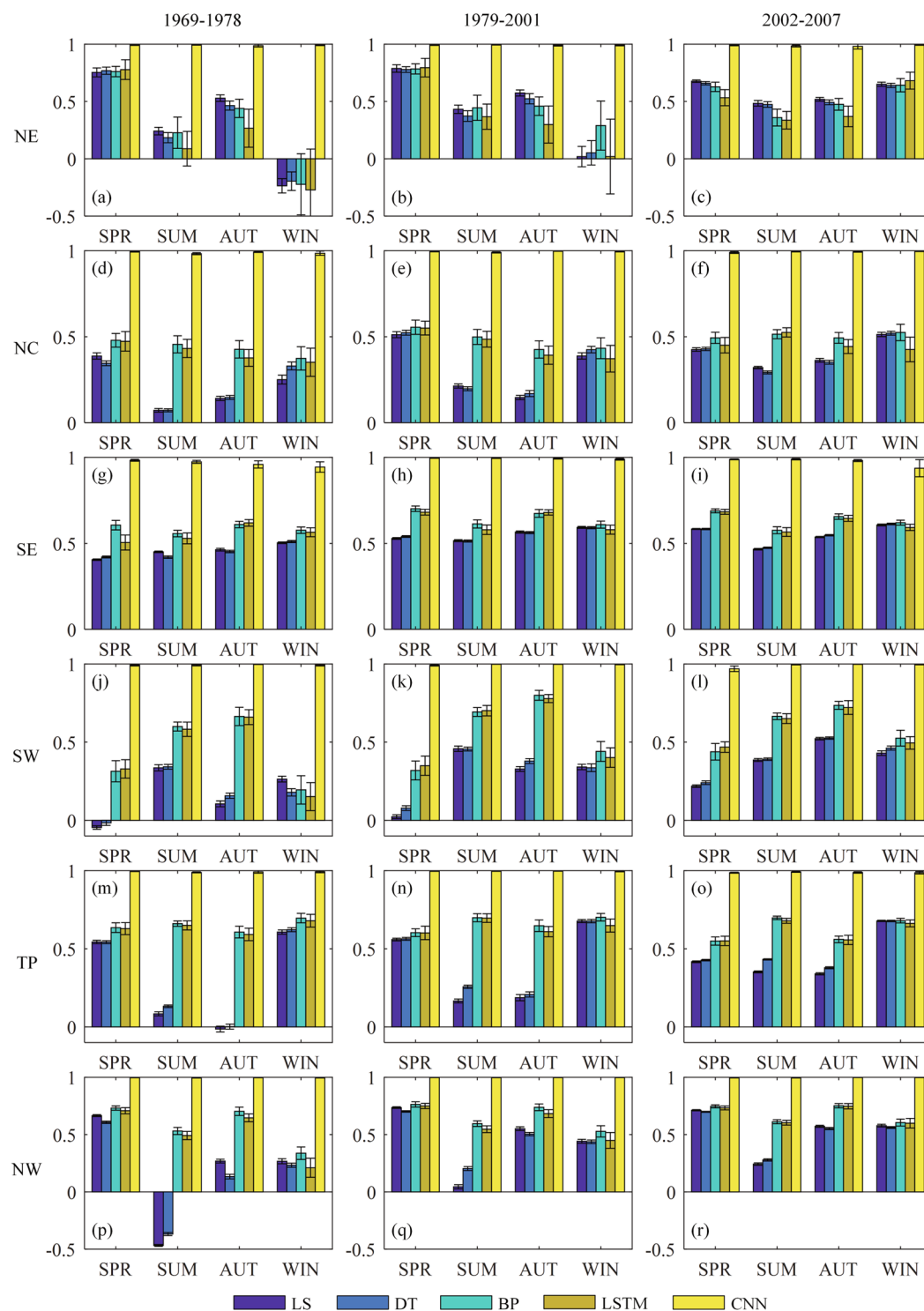
The above results show that the CNN fusion method consistently performs better than the other methods. To further confirm this conclusion, scatter plots of fused and corrected against observed  $\delta^{18}\text{O}_p$  are presented in Fig. 8 for the 1979–2001 period. The overall mean CC and RMSE corresponding to the figure are shown in Table 3. It can be seen that the CNN fusion method shows a stronger correlation with the observations than the other fusion methods and BCMs. The CNN fused  $\delta^{18}\text{O}_p$  consistently shows the largest CC and smallest RMSE, showing a strong positive linear correlation with the observations with CC being almost all larger than 0.99.

Based on the above performances of the correction and fusion methods, the best combination to build the dataset can be determined. Since the CNN fusion method consistently performs better than other methods, the CNN was used for the common period of all climate simulation and observations, while the bias correction methods were used for the periods with only one or two climate simulations, and with no observations. The specific generation mode of the dataset for each period can be found in the Data Records section.

To evaluate the dataset for all stations, CC and RMSE are calculated for  $\delta^{18}\text{O}_p$  series between observations and raw iGCM simulations, and between observations and the fused isoscape for all stations over the common period (Tables S3–4). Also, the distributions of these CCs and RMSEs are shown in the form of histogram (Fig. S1). The results show that the built isoscape performs excellent for the vast majority of stations, with much larger CCs and smaller RMSEs than iGCM simulations. Specifically, the CCs between the isoscape simulations and observations are larger than 0.8 for 77.6% of stations and larger than 0.9 for 57.0% of stations. The RMSEs between the isoscape simulations and observations are smaller than 3‰ for 80.4% of stations and less than 2‰ for 56.1% of stations.

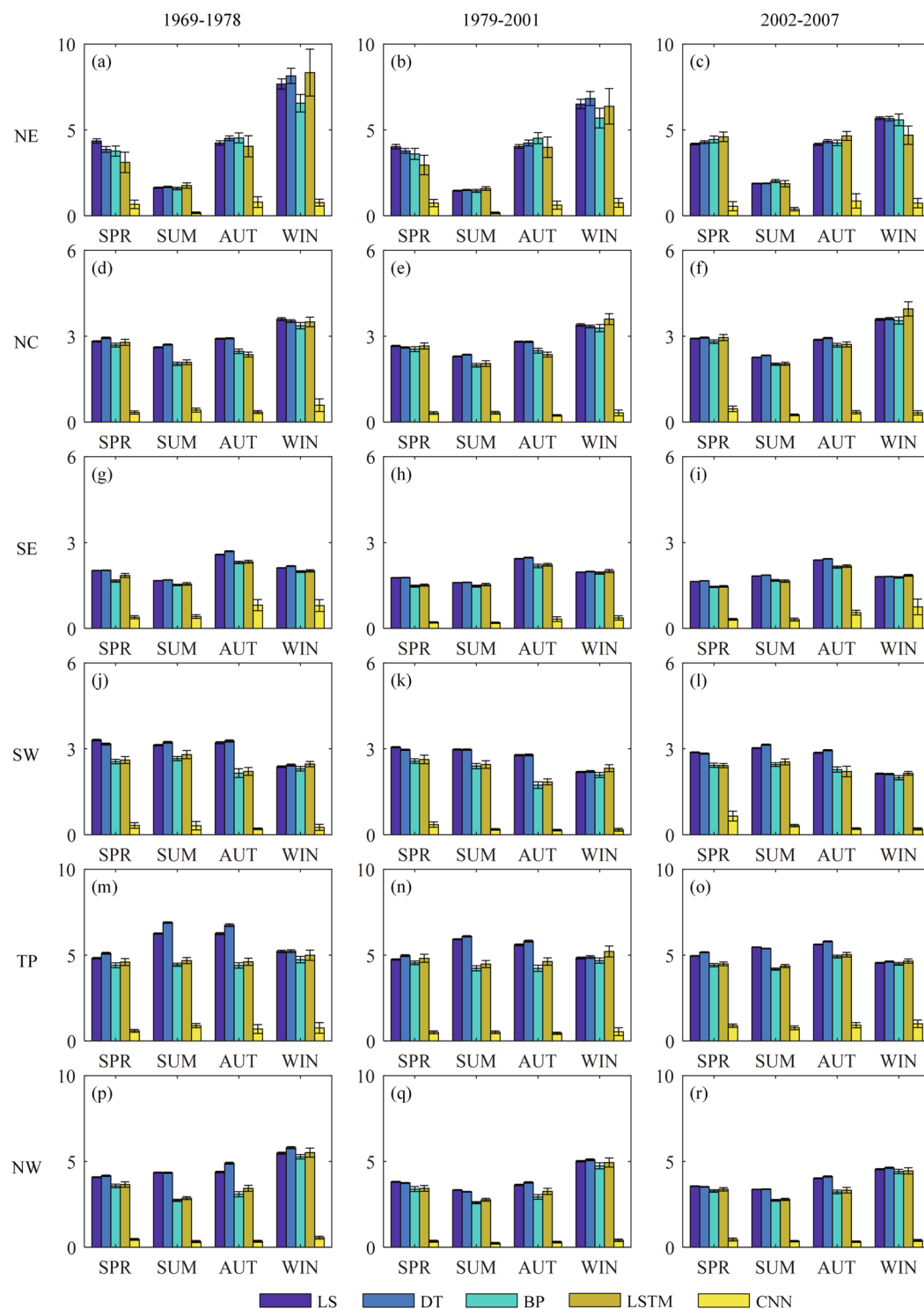
To further demonstrate the dataset quality, twenty stations with appropriate lengths of observation are selected with two stations randomly selected from each sub-region. The time series of  $\delta^{18}\text{O}_p$  are plotted for observations, iGCM simulations, and the generated isoscape (Fig. 9). Figure 9 shows that the variations of fused  $\delta^{18}\text{O}_p$  time series show consistent patterns with observations, and it also performs much better than raw iGCM simulations. In particular, for the period before 2007, the CNN model integrates the advantages of various simulations and captures most features of the observed data. These results generally prove the reliability of fused isoscape.

Fig. S2 shows the spatial distribution of monthly mean observed, fused, and best-performing iGCM raw  $\delta^{18}\text{O}_p$  for their common period (i.e. 1979–2007). The spatial pattern presented by the isoscape shows the best consistency with observations. The strength of the CNN model has been demonstrated, which can make good



**Fig. 6** Seasonal average results of correlation coefficient (CC) metrics for BCMs and DFMs in six sub-regions. The whiskers denote  $\pm$  one standard deviation.

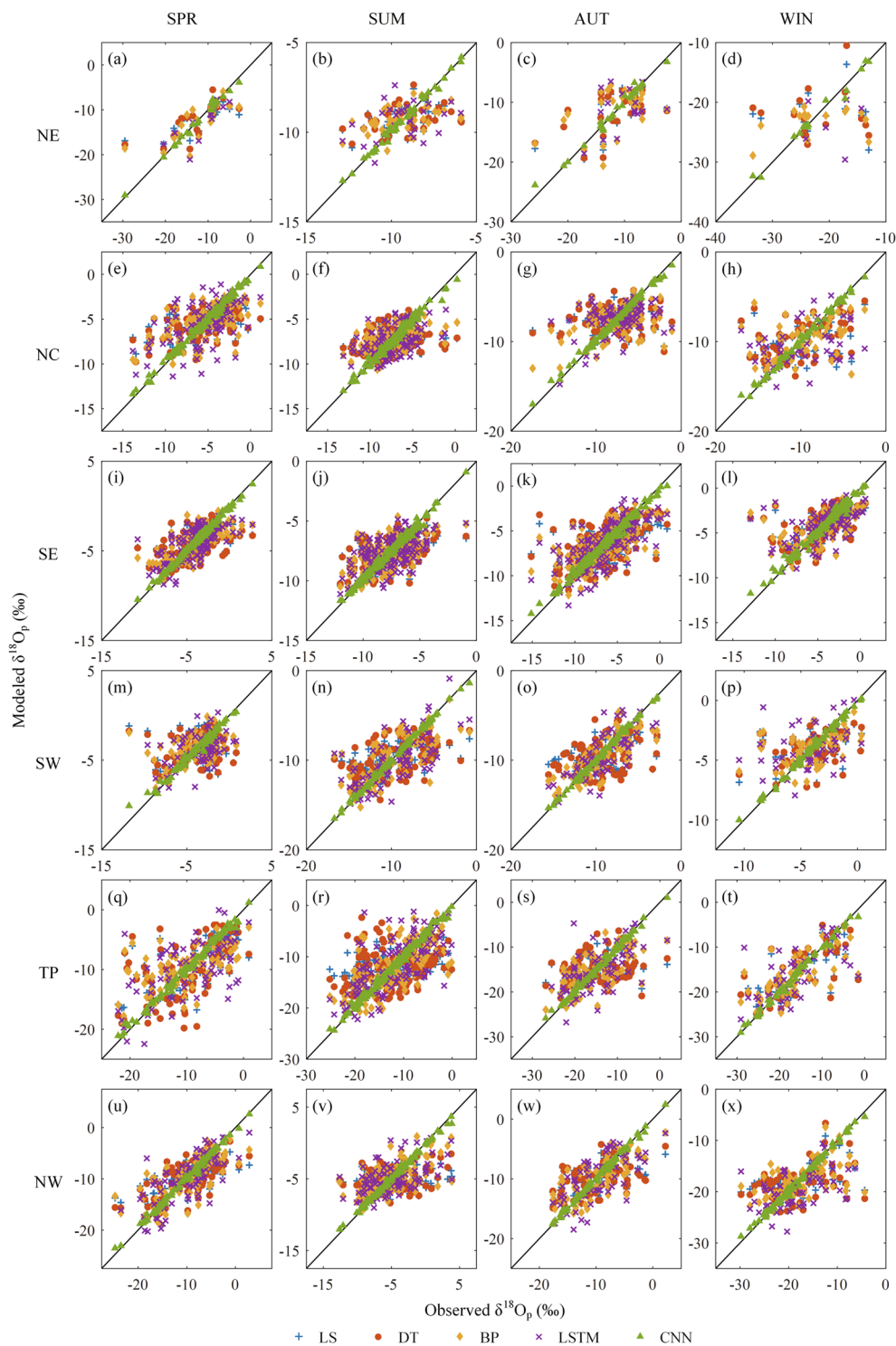
use of the advantages of each simulation to accurately capture the characteristics of observations. For example, nudged LMDZ4 model shows a strong ability to reproduce the spatial distribution of  $\delta^{18}\text{O}_p$  for the eastern region in summer and autumn, but a slightly poor performance in the Qinghai-Tibet Plateau. The built isoscape combines nudged LMDZ4, GISS E, IsoGSM2 and other simulations, which show reasonable performance for the Qinghai-Tibet Plateau, and well reproduces the spatial distribution of  $\delta^{18}\text{O}_p$  for the mainland of China. In addition, the spatial and seasonal variability presented in the built isoscape is consistent with the monthly isoscape (C-Isoscape) created by Wang, *et al.*<sup>87</sup> based on regionalized fuzzy clustering, with only minor differences that



**Fig. 7** Seasonal average results of root mean square error (RMSE) metrics (%) for BCMs and DFMs in six sub-regions. The whiskers denote  $\pm$  one standard deviation.

may be due to different data or methods. Compared with C-Isoscape, the  $\delta^{18}\text{O}_p$  of our isoscape is generally lower in the Qinghai-Tibet Plateau, and the seasonal variation of  $\delta^{18}\text{O}_p$  in northwest China is smaller.

**Spatial variability of precipitation oxygen isotope.** To further evaluate the isoscape fused by the CNN method, the observations and the isoscape simulations are compared at the seasonal scale for analysing the spatial variability. Figure 10 presents the spatial variability of mean  $\delta^{18}\text{O}_p$  for observed and fused data for the 1969–2007 period. Generally, the CNN fused  $\delta^{18}\text{O}_p$  shows similar spatial distribution to observations for all four seasons.



**Fig. 8** Scatter plots of seasonal  $\delta^{18}\text{O}_p$  from bias-corrected and fused output against observations in six sub-regions.

In NE,  $\delta^{18}\text{O}_p$  decreases with increasing latitude for all seasons, and its spatial variation is basically parallel to the latitude, which reflects the latitude effect<sup>88</sup>. Indeed, most of the water vapour in the atmosphere is formed at low latitudes, and Rayleigh distillation continuously depletes the residual water vapour as air masses move toward higher latitudes, thus depleting the  $\delta^{18}\text{O}_p$  of the residual water vapour and thus of the rain forming in clouds.

In SE and SW,  $\delta^{18}\text{O}_p$  decreases from the southeast coast to inland. This phenomenon is consistent with the continental effect. As water vapour transfers from the ocean to the interior of the continent, precipitation is formed along the way. The separation process of heavy isotopes takes place preferentially than that of light isotopes, which leads to the gradual dilution of heavy isotopes in the cloud, and thus makes the proportion of heavy isotopes in the subsequent precipitation lower.

Sub-region	Season	LS		DT		BP		LSTM		CNN	
		CC	RMSE	CC	RMSE	CC	RMSE	CC	RMSE	CC	RMSE
NE	SPR	0.791	4.058	0.774	3.769	0.784	3.642	0.799	3.030	0.996	0.628
	SUM	0.433	1.457	0.379	1.497	0.448	1.445	0.361	1.507	0.995	0.156
	AUT	0.578	4.030	0.530	4.203	0.468	4.513	0.331	4.059	0.993	0.623
	WIN	0.024	6.508	0.067	6.834	0.277	5.696	0.079	6.289	0.992	0.770
NC	SPR	0.510	2.665	0.521	2.612	0.556	2.550	0.556	2.654	0.995	0.315
	SUM	0.213	2.297	0.199	2.348	0.504	1.979	0.485	1.991	0.991	0.329
	AUT	0.151	2.806	0.167	2.800	0.422	2.509	0.396	2.382	0.997	0.237
	WIN	0.386	3.388	0.425	3.337	0.434	3.302	0.376	3.571	0.997	0.310
SE	SPR	0.529	1.771	0.540	1.777	0.701	1.474	0.683	1.519	0.995	0.208
	SUM	0.520	1.596	0.516	1.607	0.616	1.472	0.576	1.547	0.995	0.197
	AUT	0.568	2.432	0.567	2.468	0.674	2.182	0.680	2.225	0.994	0.355
	WIN	0.597	1.958	0.593	1.980	0.607	1.938	0.579	1.999	0.990	0.372
SW	SPR	0.021	3.051	0.078	2.953	0.325	2.568	0.346	2.560	0.993	0.364
	SUM	0.462	2.970	0.456	2.967	0.696	2.387	0.704	2.466	0.999	0.186
	AUT	0.330	2.771	0.380	2.791	0.801	1.728	0.774	1.778	0.998	0.176
	WIN	0.343	2.189	0.339	2.206	0.442	2.075	0.413	2.319	0.997	0.199
TP	SPR	0.560	4.742	0.563	4.977	0.619	4.476	0.643	4.495	0.997	0.467
	SUM	0.169	5.926	0.260	6.090	0.709	4.171	0.703	4.483	0.996	0.514
	AUT	0.189	5.585	0.211	5.791	0.714	3.884	0.667	4.216	0.998	0.427
	WIN	0.676	4.836	0.676	4.887	0.697	4.705	0.645	5.307	0.996	0.639
NW	SPR	0.738	3.799	0.702	3.746	0.763	3.389	0.749	3.426	0.998	0.402
	SUM	0.044	3.341	0.205	3.231	0.596	2.594	0.545	2.748	0.998	0.282
	AUT	0.553	3.628	0.506	3.766	0.734	2.957	0.684	3.250	0.997	0.323
	WIN	0.441	5.015	0.437	5.090	0.533	4.730	0.449	5.030	0.997	0.448

**Table 3.** Correlation coefficient (CC) and root mean square error (RMSE, %) metrics corresponding to the Fig. 8.

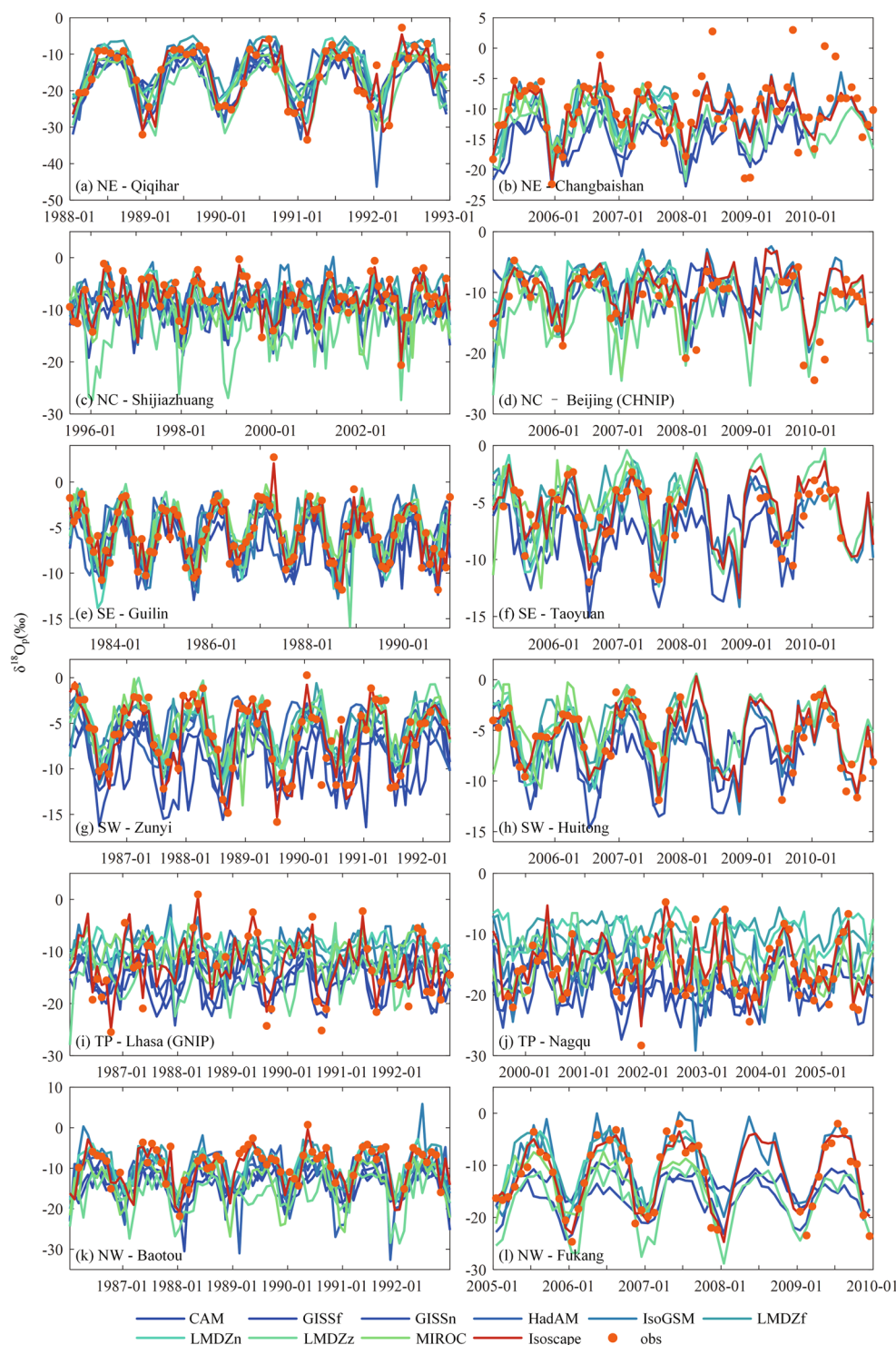
The  $\delta^{18}\text{O}_p$  in TP is low in general except for the southeast corner, which is mainly due to the special effect of large landforms<sup>89</sup>. The low  $\delta^{18}\text{O}_p$  in TP is mainly due to its high altitude, with an average altitude being above 4,000 m. The moisture in the air mass is gradually removed during the orographic uplift, with heavy isotopes preferred to be removed during the condensation process, which leads to the dilution of heavy isotopes in water vapour<sup>88</sup>. Higher  $\delta^{18}\text{O}_p$  in the southeast corner of TP indicates closer vapour sources such as the Bay of Bengal and the Arabian Sea in the Indian Ocean. Due to the terrain barrier of the Himalayas, most of the water vapour can only pass through its southeast corner, along the valley of the major rivers (Nujiang River, Jinsha River, etc.) into the plateau, or through the Yarlung Zangbo River valley into the plateau<sup>90</sup>.

The  $\delta^{18}\text{O}_p$  in NW is lower than that in the southern region, but higher than that in NE and TP. This is because NW is far away from the ocean and has a dry climate, so the amount of heavy isotope in water vapour from the ocean is limited. However, a large part of water vapour to generate precipitation in NW comes from terrestrial evaporation<sup>91</sup>. The  $\delta^{18}\text{O}_p$  in surface water in the arid area is high, resulting in high  $\delta^{18}\text{O}_p$  in evaporation water vapour and heavy isotope enrichment in precipitation. Another process is the re-evaporation of raindrop in arid climate, enriching heavy isotopes in precipitation water<sup>92</sup>. Under the joint control of both, the  $\delta^{18}\text{O}_p$  in this sub-region varies greatly. In the southern part of the Taklimakan Desert, Xinjiang,  $\delta^{18}\text{O}_p$  is obviously higher. This is because, the Taklimakan Desert is located in the heart of Eurasia, and it is surrounded by high mountains and has extremely low rainfall<sup>93</sup>. In the southern part of the desert, there is more precipitation in the Kashi-Hotan line, and the water vapour mainly comes from the evaporation of local lakes and rivers<sup>94</sup>, so the ratio of isotopes in precipitation is high.

The seasonality of  $\delta^{18}\text{O}_p$  varies in sub-regions and is influenced by various factors. For NE and NC,  $\delta^{18}\text{O}_p$  is lower in winter than in other seasons. NE and NC are influenced by westerly wind and polar continental air mass constantly, with no convergence or strong convection with isotope-depleted air mass. Compared with Pacific air mass, westerly wind and polar air mass are drier and have higher  $\delta^{18}\text{O}_p$ <sup>95</sup>. The seasonal distribution pattern of  $\delta^{18}\text{O}_p$  in NE and NC is consistent with the temperature effect<sup>96</sup>. Although the amount effect is not significant at the annual scale in these regions, it cannot be ignored in the wet season<sup>97</sup>. In particular, the maximum value of  $\delta^{18}\text{O}_p$  is observed in spring for NC, when a large part of the precipitation vapour comes from local re-evaporation. The temperature effect is also reflected in NW. Due to the long-term influence of continental air mass, the temperature difference between winter and summer is large, and the  $\delta^{18}\text{O}_p$  changes synchronously with temperature in these two seasons.

For SE and SW,  $\delta^{18}\text{O}_p$  is lower in summer than in other seasons. The climate features of SE and SW are related to the deep convection driven by the East Asian monsoon<sup>98</sup>, which brings water vapour from the Pacific Ocean to eastern China and dominates the sub-regions in summer. Air masses from the Pacific Ocean are more isotopically depleted than those from SE and SW<sup>83,89</sup>, so convergence with the Pacific depleted air masses will dilute the isotopic content of precipitation in SE and SW. Therefore, although the temperature in summer is generally

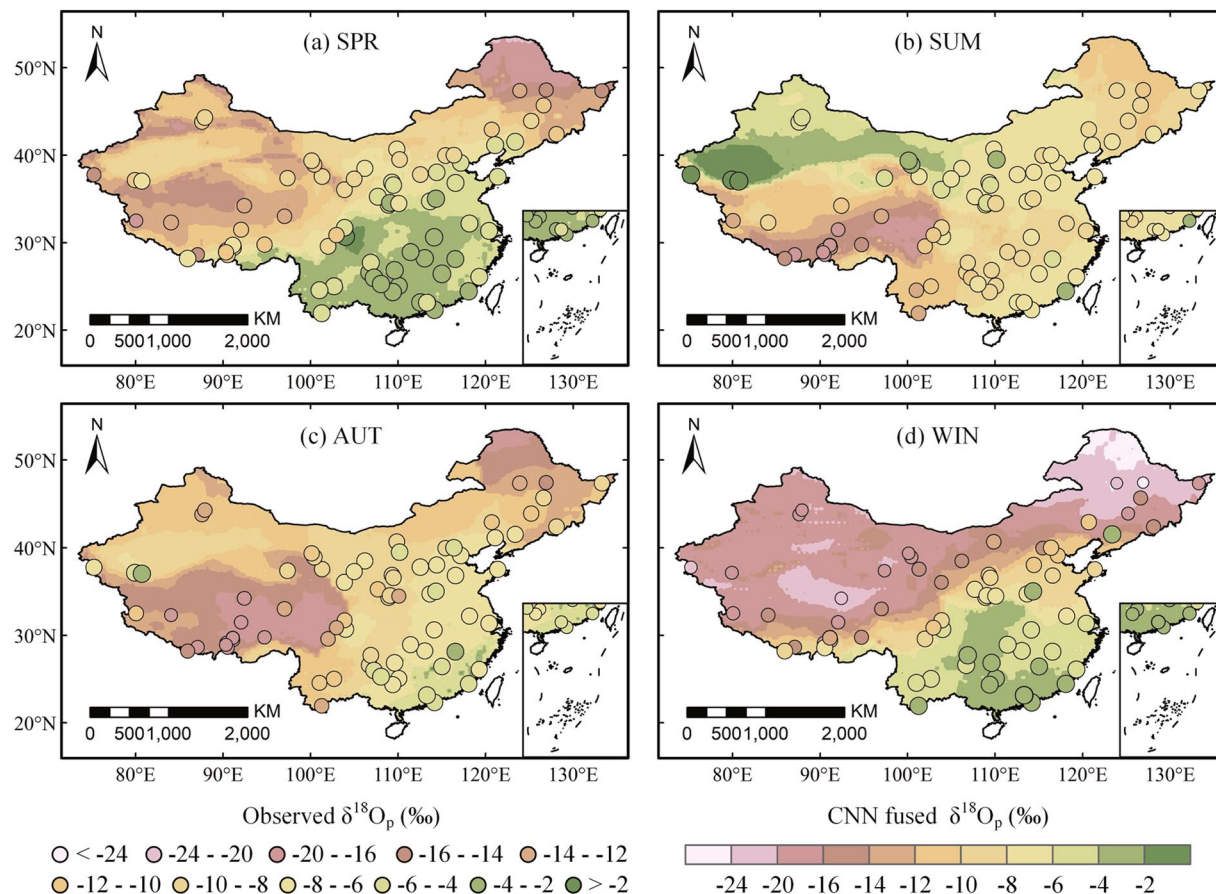




**Fig. 9** Time-series comparisons of  $\delta^{18}\text{O}_p$  among the built isoscape, iGCM simulations, and *in-situ* observations at selected stations in each sub-region.

higher, depletion of  $\delta^{18}\text{O}_p$  is usually larger during the monsoon season than winter season. The effect of surface temperature on isotopic fractionation during precipitation is masked by the effect of precipitation amount<sup>90</sup>. The temporal distribution pattern of  $\delta^{18}\text{O}_p$  in SE and SW is influenced by heavy monsoon precipitation and follows the amount effect<sup>88,96</sup>.

In TP,  $\delta^{18}\text{O}_p$  is positively correlated with temperature in the non-monsoon region (northern part of the plateau), with high  $\delta^{18}\text{O}_p$  in summer and low in winter, reflecting the temperature effect. For the monsoon region (southern part of the plateau),  $\delta^{18}\text{O}_p$  is high in winter and spring and low in summer and autumn, which is



**Fig. 10** Seasonal averaged observations and CNN fused simulations of  $\delta^{18}\text{O}_p$  in the mainland of China.

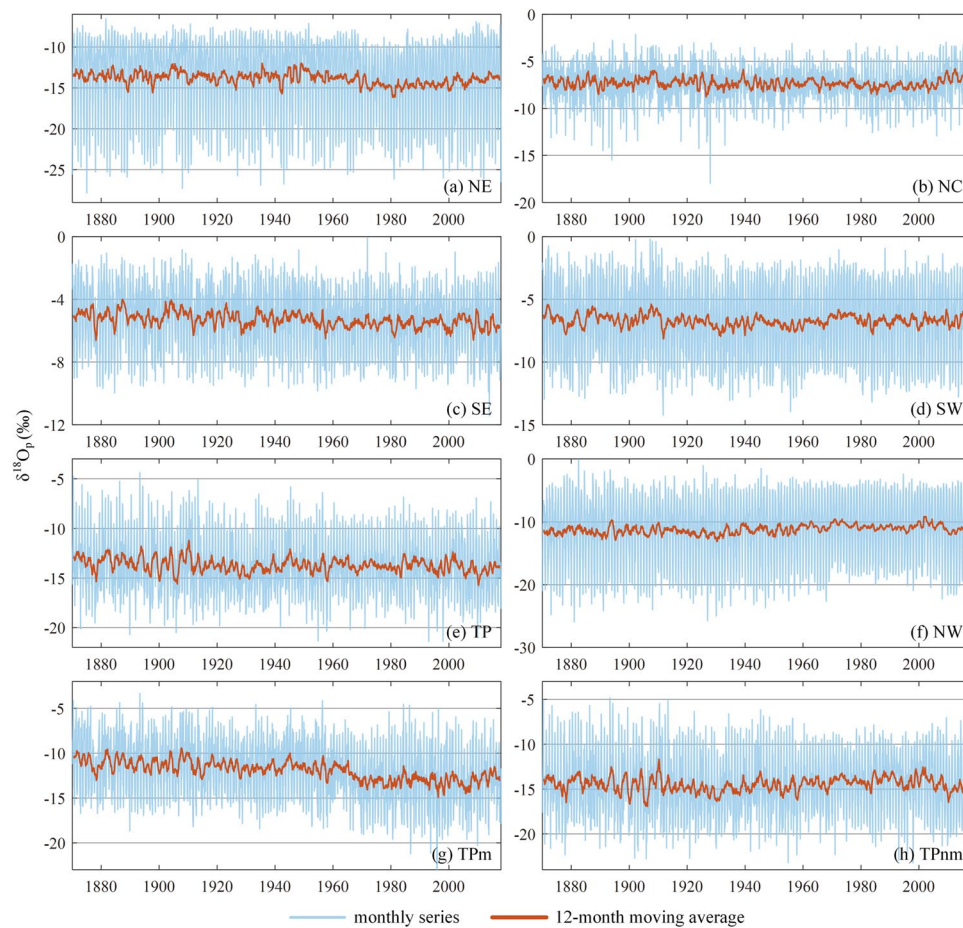
obviously influenced by marine air mass and shows obvious amount effect. These results are similar to previous studies<sup>48,99</sup>.

**Temporal variability of precipitation oxygen isotope.** As mentioned earlier, the isoscape is generated by a combination of bias correction and data fusion methods for the 1870–2017 period. The CNN fusion is used for 1969–2007, and BCMs (mean of LS and DT) are used for the rest of the period. Figure 11 shows the monthly time series of generated  $\delta^{18}\text{O}_p$  and their 12-month moving averages for eight sub-regions over the 1870–2017 period. TP is divided into monsoon and non-monsoon regions, according to the research of Yu, *et al.*<sup>99</sup>. In our study, the region with significant correlations between  $\delta^{18}\text{O}_p$  and temperature is the non-monsoon region, while the rest is the monsoon region.

The Mann-Kendall tests show that the  $\delta^{18}\text{O}_p$  significantly increased in NE and NC for the recent 40 years at the  $P = 0.01$  level. These two regions are consistent with the temperature effect, and it can be inferred that the temperature of these sub-regions has a rising trend during this period. This has been proved in many studies. For example, the studies of Ren, *et al.*<sup>100</sup> and Ding, *et al.*<sup>101</sup> have shown that the temperature in China has had a rising trend in recent years, especially in northeast, north and northwest China. A slight upward trend is also observed in NW from the 1930s to the 1970s (significant at the 0.1 level). The temperature effect is more significant in inland areas at middle and high latitudes. In winter, NW is mainly controlled by the westerlies, and the amount effect can be ignored<sup>102,103</sup>. Therefore, the temperature effect in NW is more significant in winter. From 2001 to 2012, the  $\delta^{18}\text{O}_p$  values in NW showed a decreasing trend, mainly in winter. Some studies have shown that the temperature in NW during this period is consistent with the global land warming hiatus phenomenon, and even shows obvious cooling, especially in winter<sup>104,105</sup>.

The  $\delta^{18}\text{O}_p$  in SE presents a gentle decline trend for the past 80 years (significant at the 0.05 level), indicating that precipitation has an upward trend in this period, since the  $\delta^{18}\text{O}_p$  conforms to the rule of amount effect in this region. While the  $\delta^{18}\text{O}_p$  in SW shows no significant trend in recent years, which indicates no significant trend in precipitation. These trends are consistent with existing researches that precipitation increased in the east coast and northwest of China, decreased in the north and northeast, and showed no significant changes in the southwest<sup>106,107</sup>.

There is no significant trend for  $\delta^{18}\text{O}_p$  in TP. However,  $\delta^{18}\text{O}_p$  shows a significant decreasing trend for the monsoon region of TP from the 1950s to the 2000s (significant at the 0.01 level), while it shows a significant increasing trend for the non-monsoon region (significant at the 0.05 level). This is because the non-monsoon region shows the temperature effect, while the monsoon region shows the amount effect, which is consistent



**Fig. 11** Monthly time series of the generated  $\delta^{18}\text{O}_p$  (‰) and their 12-month moving average in eight sub-regions from 1870 to 2017. (g) shows the monsoon region of TP. (h) shows the non-monsoon region of TP.

with the increasing trend of temperature and precipitation in the Qinghai-Tibet region in recent years<sup>48,99</sup>. Overall, the changing trend of temperature and precipitation derived from the isotope effect analysis is consistent with that analysed directly using temperature and precipitation data in the mainland of China. This further proved the reasonable performance of the isoscape built in this study.

**Summary.** Long-time sequences of  $\delta^{18}\text{O}_p$  are of great significance for hydrological and meteorological studies. In view of the lack of long and reliable  $\delta^{18}\text{O}_p$  datasets in China, this study generates a new dataset by integrating multi-iGCM data to overcome the limitations of short duration and uneven distribution of observed data. This dataset contains monthly  $\delta^{18}\text{O}_p$  over the mainland of China for the 1870–2017 period with a spatial resolution of 50–60 km. The dataset from 1969 to 2007 is generated by using the CNN fusion method, when the observed time series and multiple iGCM simulations are available. For other periods, it is generated by bias correcting iGCMs simulations. Two BCMs (i.e. LS and DT) with similar performances are used to produce the ensemble mean. Prior to building the isoscape, the performance of two BCMs (LS and DT) and three DFMs (BP, LSTM and CNN) is evaluated using RMSE and CC as criteria. The results show that the CNN fusion method consistently performs the best for all sub-regions in China, and BP and LSTM fusion methods perform slightly better than LS and DT (BCMs). The performance of the LS and DT methods is similar. In terms of spatial distribution and temporal variability of  $\delta^{18}\text{O}_p$ , the generated data show very similar spatial distributions to observations, and the temporal trend of  $\delta^{18}\text{O}_p$  is consistent with the observed changes in precipitation and temperature for different regions in China. All these show that the built isoscape is reliable and useful to extend the time and space of observations in China.

## Usage Notes

**Advantages and limitations.** The generated isoscape dataset has high spatio-temporal resolution and a long series covering 1870–2017. Compared with the existing iGCMs, the isoscape has high quality and stability for a large region in China at the monthly scale. Benefiting from the characteristics of optimal neural network and bias correction methods, the isoscape makes full use of observations to integrate the advantages of various iGCMs. In other words, by using the combination of data fusion and bias correction methods, all observations and iGCM simulations are used to the utmost extent to ensure the highest accuracy throughout the entire time period. Studies have shown that the CNN model has strong abilities for generalization and information synthesis<sup>51,108</sup>, while the bias correction methods have commonly used in climate change studies. Moreover, the hybrid



generation method of the isoscape has the characteristics of high accuracy and simplicity, which can be easily extended to the generation of isoscape datasets in other regions. Even though the methods used in this study have been widely used, this is the first time to generate a high-quality isoscape with a long time period in China. The generated isoscape would be very used for hydro-meteorological studies. However, it should be noted that the isoscape may be more reliable for the common periods of most iGCMs (1969–2007), but mediocre for other periods. What's more, affected by the data quality and representativeness of observation stations, the accuracy of the isoscape in some regions still needs to be improved. It is believed that this problem will be solved as observed data become more abundant.

**Data applications.** Based on this built isoscape, the physical mechanisms driving the spatio-temporal variation of  $\delta^{18}\text{O}_p$  can be deeply explored. This dataset is useful for tracing atmospheric and hydrological processes. It can be used to study the effect of meteorological variables and air mass trajectory on stable isotope distribution, and quantify the source and fate of moisture<sup>55,109,110</sup>. For example, over East Asia, where the length of observed isotope data is short, or over the Tibetan Plateau, where data are unevenly distributed, the influence of climate change on moisture source and contribution can be studied based on this long series precipitation isoscape. The isoscape can also be used with regional climate models through data assimilation. For example, the precipitation isoscape can be combined with the physical constraints of regional climate models to reconstruct hydrological and climatic elements such as water vapour and precipitation. It can be a useful attempt to advance the study of climatic and hydrological data.

### Code availability

The codes for two bias correction methods (LS and DT) and three neural network data fusion methods (BP, LSTM and CNN) are available at <https://doi.org/10.5281/zenodo.7306199>. The codes were programmed using MATLAB version 2022a and Python 3.8.

Received: 10 November 2022; Accepted: 22 March 2023;

Published online: 06 April 2023

### References

1. Bowen, G. J. Isoscapes: Spatial pattern in isotopic biogeochemistry. *Annu. Rev. Earth Planet. Sci.* **38**, 161–187 (2010).
2. Gibson, J. J. *et al.* Progress in isotope tracer hydrology in Canada. *Hydrol. Process.* **19**, 303–327 (2005).
3. Galewsky, J. *et al.* Stable isotopes in atmospheric water vapor and applications to the hydrologic cycle. *Rev. Geophys.* **54**, 809–865 (2016).
4. Ansari, M. A., Noble, J., Deodhar, A. & Saravana Kumar, U. Atmospheric factors controlling the stable isotopes ( $\delta^{18}\text{O}$  and  $\delta^2\text{H}$ ) of the Indian summer monsoon precipitation in a drying region of Eastern India. *J. Hydrol.* **584**, 124636 (2020).
5. Zhang, Y., Jones, M., Zhang, J., McGowan, S. & Metcalfe, S. Can  $\delta^{18}\text{O}$  help indicate the causes of recent lake area expansion on the western Tibetan Plateau? A case study from Aweng Co. *J. Paleolimnol.* **65**, 169–180 (2020).
6. McGuire, K., DeWalle, D. & Gburek, W. Evaluation of mean residence time in subsurface waters using oxygen-18 fluctuations during drought conditions in the mid-Appalachians. *J. Hydrol.* **261**, 132–149 (2002).
7. Gazis, C. & Feng, X. A stable isotope study of soil water: evidence for mixing and preferential flow paths. *Geoderma* **119**, 97–111 (2004).
8. Chen, J. S. *et al.* Groundwater maintains dune landscape. *Nature* **432**, 459–460 (2004).
9. Bowen, G. J., Cai, Z., Fiorella, R. P. & Putman, A. L. Isotopes in the water cycle: regional- to global-scale patterns and applications. *Annu. Rev. Earth Planet. Sci.* **47**, 453–479 (2019).
10. Worden, J., Noone, D. & Bowman, K. Tropospheric Emission Spectrometer Science, T. & Data, c. Importance of rain evaporation and continental convection in the tropical water cycle. *Nature* **445**, 528–532 (2007).
11. Froehlich, K. *et al.* Deuterium excess in precipitation of Alpine regions - moisture recycling. *Isotopes Environ. Health Stud.* **44**, 61–70 (2008).
12. Fekete, B. M., Gibson, J. J., Aggarwal, P. & Vörösmarty, C. J. Application of isotope tracers in continental scale hydrological modeling. *J. Hydrol.* **330**, 444–456 (2006).
13. Rowley, D. B. & Garzione, C. N. Stable isotope-based paleoaltimetry. *Annu. Rev. Earth Planet. Sci.* **35**, 463–508 (2007).
14. Johnson, K. R. & Ingram, B. L. Spatial and temporal variability in the stable isotope systematics of modern precipitation in China: implications for paleoclimate reconstructions. *Earth Planet. Sci. Lett.* **220**, 365–377 (2004).
15. Kralik, M., Papesch, W. & Stichler, W. Austrian Network of Isotopes in Precipitation (ANIP): Quality assurance and climatological phenomenon in one of the oldest and densest networks in the world. *Isotope hydrology and integrated water resources management*, 146–149 (2003).
16. Lynch, J., Grimm, J. & Bowersox, V. Trends in precipitation chemistry in the United States: A national perspective, 1980–1992. *Atmospheric Environ.* **29**, 1231–1246 (1995).
17. Schürch, M., Kozel, R., Schotterer, U. & Tripet, J.-P. Observation of isotopes in the water cycle? The Swiss National Network (NISOT). *Environ. Geol.* **45**, 1–11 (2003).
18. Fritz, P., Drimmie, R., Frapet, S. & O'Shea, K. The isotopic composition of precipitation and groundwater in Canada. In *Isotope techniques in water resources development. Proc. IAEA symposium, Vienna, 1987*. 539–550 (1987).
19. Yu, W. *et al.* Stable isotope variations in precipitation over Deqin on the southeastern margin of the Tibetan Plateau during different seasons related to various meteorological factors and moisture sources. *Atmos. Res.* **170**, 123–130 (2016).
20. Zhang, M. & Wang, S. A review of precipitation isotope studies in China: Basic pattern and hydrological process. *J. Geogr. Sci.* **26**, 921–938 (2016).
21. Song, X. *et al.* Establishment of Chinese Network of Isotopes in Precipitation (CHNIP) based on CERN. *Advances in Earth Science* **22**, 738–747 (2007).
22. Allen, S. T., Kirchner, J. W. & Goldsmith, G. R. Predicting spatial patterns in precipitation isotope ( $\delta^2\text{H}$  and  $\delta^{18}\text{O}$ ) seasonality using sinusoidal isoscapes. *Geophys. Res. Lett.* **45**, 4859–4868 (2018).
23. Wang, S. *et al.* Comparison of GCM-simulated isotopic compositions of precipitation in arid central Asia. *J. Geogr. Sci.* **25**, 771–783 (2015).
24. Hoffmann, G., Jouzel, J. & Masson, V. Stable water isotopes in atmospheric general circulation models. *Hydrol. Process.* **14**, 1385–1406 (2000).
25. Joussaume, S., Sadourny, R. & Jouzel, J. A general-circulation model of water isotope cycles in the atmosphere. *Nature* **311**, 24–29 (1984).

26. Hoffmann, G., Werner, M. & Heimann, M. Water isotope module of the ECHAM atmospheric general circulation model: A study on timescales from days to several years. *J. Geophys. Res. Atmos.* **103**, 16871–16896 (1998).
27. Schmidt, G. A., LeGrande, A. N. & Hoffmann, G. Water isotope expressions of intrinsic and forced variability in a coupled ocean-atmosphere model. *J. Geophys. Res. Atmos.* **112**, D10103 (2007).
28. Schmidt, G. A., Hoffmann, G., Shindell, D. T. & Hu, Y. Modeling atmospheric stable water isotopes and the potential for constraining cloud processes and stratosphere-troposphere water exchange. *J. Geophys. Res. Atmos.* **110**, D21314 (2005).
29. Tindall, J. C., Valdes, P. J. & Sime, L. C. Stable water isotopes in HadCM3: Isotopic signature of El Niño–Southern Oscillation and the tropical amount effect. *J. Geophys. Res. Atmos.* **114**, D04111 (2009).
30. Risi, C., Bony, S., Vimeux, F. & Jouzel, J. Water-stable isotopes in the LMDZ4 general circulation model: Model evaluation for present-day and past climates and applications to climatic interpretations of tropical isotopic records. *J. Geophys. Res. Atmos.* **115**, D12118 (2010).
31. Kurita, N. *et al.* Intraseasonal isotopic variation associated with the Madden-Julian Oscillation. *J. Geophys. Res. Atmos.* **116**, D24101 (2011).
32. Yoshimura, K., Oki, T., Ohte, N. & Kanae, S. A quantitative analysis of short-term  $\delta^{18}\text{O}$  variability with a Rayleigh-type isotope circulation model. *J. Geophys. Res. Atmos.* **108** (2003).
33. Conroy, J. L., Cobb, K. M. & Noone, D. Comparison of precipitation isotope variability across the tropical Pacific in observations and SWING2 model simulations. *J. Geophys. Res. Atmos.* **118**, 5867–5892 (2013).
34. Zhang, X. *et al.* GCM simulations of stable isotopes in the water cycle in comparison with GNIP observations over East Asia. *Acta Meteorol. Sin.* **26**, 420–437 (2012).
35. Che, Y. *et al.* Stable water isotopes of precipitation in China simulated by SWING2 models. *Arab. J. Geosci.* **9**, 732 (2016).
36. Krajewski, W. F. Cokriging radar-rainfall and rain gage data. *J. Geophys. Res. Atmos.* **92**, 9571–9580 (1987).
37. Rosenfeld, D., Wolff, D. B. & Amitai, E. The window probability matching method for rainfall measurements with radar. *J. Appl. Meteorol. Climatol.* **33**, 682–693 (1994).
38. Pereira Fo, A. J., Crawford, K. C. & Hartzell, C. L. Improving WSR-88D hourly rainfall estimates. *Weather Forecast.* **13**, 1016–1028 (1998).
39. Todini, E. A Bayesian technique for conditioning radar precipitation estimates to rain-gauge measurements. *Hydrol. Earth Syst. Sci.* **5**, 187–199 (2001).
40. Shen, Y., Zhao, P., Pan, Y. & Yu, J. A high spatiotemporal gauge-satellite merged precipitation analysis over China. *J. Geophys. Res. Atmos.* **119**, 3063–3075 (2014).
41. Bianchi, B., van Leeuwen, P. J., Hogan, R. J. & Berne, A. A variational approach to retrieve rain rate by combining information from rain gauges, radars, and microwave links. *J. Hydrometeorol.* **14**, 1897–1909 (2013).
42. LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).
43. Turlapaty, A. C., Anantharaj, V. G., Younan, N. H. & Joseph Turk, F. Precipitation data fusion using vector space transformation and artificial neural networks. *Pattern Recognit. Lett.* **31**, 1184–1200 (2010).
44. Sun, A. Y. & Tang, G. Downscaling satellite and reanalysis precipitation products using attention-based deep convolutional neural nets. *Frontiers in Water* **2** (2020).
45. Wu, H., Yang, Q., Liu, J. & Wang, G. A spatiotemporal deep fusion model for merging satellite and gauge precipitation in China. *J. Hydrol.* **584** (2020).
46. Liu, J., Song, X., Yuan, G., Sun, X. & Yang, L. Stable isotopic compositions of precipitation in China. *Tellus B Chem. Phys. Meteorol.* **66**, 22567 (2014).
47. IAEA/WMO. Global network of isotopes in precipitation. *The GNIP Database* <https://nucleus.iaea.org/wiser> (2022).
48. Yao, T. *et al.* A review of climatic controls on  $\delta^{18}\text{O}$  in precipitation over the Tibetan Plateau: Observations and simulations. *Rev. Geophys.* **51**, 525–548 (2013).
49. Gao, J. Data set of  $\delta^{18}\text{O}$  stable isotopes in precipitation from Tibetan Network for Isotopes(1991–2008). *National Tibetan Plateau Data Center* <https://doi.org/10.11888/Geogra.tpcd.270940> (2020).
50. Shang, K. *et al.* DNN-MET: A deep neural networks method to integrate satellite-derived evapotranspiration products, eddy covariance observations and ancillary information. *Agric. For. Meteorol.* **308–309**, 108582 (2021).
51. Hengl, T., Nussbaum, M., Wright, M. N., Heuvelink, G. B. M. & Graler, B. Random forest as a generic framework for predictive modeling of spatial and spatio-temporal variables. *PeerJ* **6**, e5518 (2018).
52. Jing, Y., Lin, L., Li, X., Li, T. & Shen, H. Cascaded downscaling–calibration networks for satellite precipitation estimation. *IEEE Geosci. Remote Sens. Lett.* **19**, 1–5 (2022).
53. Jarvis, A., Reuter, H. I., Nelson, A. & Guevara, E. Hole-filled seamless SRTM data V4. *International Centre for Tropical Agriculture (CIAT)* <https://srtm.csi.cgiar.org> (2008).
54. Yoshimura, K., Kanamitsu, M., Noone, D. & Oki, T. Historical isotope simulation using reanalysis atmospheric data. *J. Geophys. Res. Atmos.* **113**, D19108 (2008).
55. Gao, J. *et al.* Precipitation water stable isotopes in the south Tibetan Plateau: observations and modeling. *J. Clim.* **24**, 3161–3178 (2011).
56. Chiang, J. C. H., Herman, M. J., Yoshimura, K. & Fung, I. Y. Enriched East Asian oxygen isotope of precipitation indicates reduced summer seasonality in regional climate and westerlies. *Proc. Natl. Acad. Sci. USA* **117**, 14745–14750 (2020).
57. Chiang, J., Herman, M., Yoshimura, K. & Fung, I. Data from: Enriched East Asian oxygen isotope of precipitation indicates reduced summer seasonality in regional climate and westerlies. *Dryad* <https://doi.org/10.6078/D1MM6B> (2020).
58. Clark, I. D. & Fritz, P. *Environmental Isotopes in Hydrogeology*. 6–7 (CRC Press, 1997).
59. Rumelhart, D. E., Widrow, B. & Lehr, M. A. The basic ideas in neural networks. *Commun. ACM* **37**, 87–92 (1994).
60. Krenker, A., Bešter, J. & Kos, A. Introduction to the artificial neural networks. *Artificial Neural Networks: Methodological Advances and Biomedical Applications. InTech*, 1–18 (2011).
61. Hsu, K.-I., Gupta, H. V. & Sorooshian, S. Artificial neural network modeling of the rainfall-runoff process. *Water Resour. Res.* **31**, 2517–2530 (1995).
62. French, M. N., Krajewski, W. F. & Cuykendall, R. R. Rainfall forecasting in space and time using a neural network. *J. Hydrol.* **137**, 1–31 (1992).
63. Zhang, Y. & Wallace, B. A sensitivity analysis of (and practitioners' guide to) convolutional neural networks for sentence classification. In *Proceedings of the Eighth International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. 253–263 (Asian Federation of Natural Language Processing, 2017).
64. Taylor, R., Ojha, V., Martino, I. & Nicosia, G. Sensitivity analysis for deep learning: ranking hyper-parameter influence. In *2021 IEEE 33rd International Conference on Tools with Artificial Intelligence (ICTAI)*. 512–516 (IEEE, 2021).
65. Mboga, N., Persello, C., Bergado, J. R. & Stein, A. Detection of informal settlements from VHR images using convolutional neural networks. *Remote Sens.* **9**, 1106 (2017).
66. Bengio, Y. Practical recommendations for gradient-based training of deep architectures. In *Neural networks: Tricks of the trade* 437–478 (Springer, 2012).
67. Xue, M., Hang, R., Liu, Q., Yuan, X.-T. & Lu, X. CNN-based near-real-time precipitation estimation from Fengyun-2 satellite over Xinjiang, China. *Atmos. Res.* **250** (2021).



68. Langford, Z. L., Kumar, J. & Hoffman, F. M. Convolutional neural network approach for mapping arctic vegetation using multi-sensor remote sensing fusion. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*. 322–331 (IEEE, 2017).
69. Chen, H., Sun, L., Cifelli, R. & Xie, P. Deep learning for bias correction of satellite retrievals of orographic precipitation. *IEEE Trans. Geosci. Remote Sens.* **60**, 1–11 (2022).
70. Rumelhart, D. E., Hinton, G. E. & Williams, R. J. Learning representations by back-propagating errors. *Nature* **323**, 533–536 (1986).
71. Zhang, J., Zhu, Y., Zhang, X., Ye, M. & Yang, J. Developing a Long Short-Term Memory (LSTM) based model for predicting water table depth in agricultural areas. *J. Hydrol.* **561**, 918–929 (2018).
72. Hochreiter, S. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *Int. J. Uncertain. Fuzz.* **6**, 107–116 (1998).
73. Hochreiter, S. & Schmidhuber, J. Long short-term memory. *Neural Comput.* **9**, 1735–1780 (1997).
74. LeCun, Y. Generalization and network design strategies. *Connectionism in perspective* **19**, 143–155 (1989).
75. Maraun, D. Bias correction, quantile mapping, and downscaling: Revisiting the inflation issue. *J. Clim.* **26**, 2137–2143 (2013).
76. Chen, J., Brissette, F. P., Chaumont, D. & Braun, M. Finding appropriate bias correction methods in downscaling precipitation for hydrologic impact studies over North America. *Water Resour. Res.* **49**, 4187–4205 (2013).
77. Chen, J., St-Denis, B. G., Brissette, F. P. & Lucas-Picher, P. Using natural variability as a baseline to evaluate the performance of bias correction methods in hydrological climate change impact studies. *J. Hydrometeorol.* **17**, 2155–2174 (2016).
78. Chen, J., Chen, J., Zhang, X. J., Peng, P. & Risi, C. Precipitation oxygen isotope for mainland China from 1870 to 2017 generated based on data fusion and bias correction of iGCMs simulations (Version 2). *Zenodo* <https://doi.org/10.5281/zenodo.7306199> (2022).
79. Zhu, Y.-Y. & Yang, S. Evaluation of CMIP6 for historical temperature and precipitation over the Tibetan Plateau and its comparison with CMIP5. *Adv. Clim. Chang. Res.* **11**, 239–251 (2020).
80. Su, F., Duan, X., Chen, D., Hao, Z. & Cuo, L. Evaluation of the global climate models in the CMIP5 over the Tibetan Plateau. *J. Clim.* **26**, 3187–3208 (2013).
81. Pang, Z. *et al.* Processes affecting isotopes in precipitation of an arid region. *Tellus B Chem. Phys. Meteorol.* **63**, 352–359 (2011).
82. Yang, Q., Ma, Z. & Xu, B. Modulation of monthly precipitation patterns over East China by the Pacific Decadal Oscillation. *Clim. Change* **144**, 405–417 (2016).
83. Peng, P., John Zhang, X. & Chen, J. Bias correcting isotope-equipped GCMs outputs to build precipitation oxygen isotope for eastern China. *J. Hydrol.* **589**, 125153 (2020).
84. Miao, C., Duan, Q., Yang, L. & Borthwick, A. G. On the applicability of temperature and precipitation data from CMIP3 for China. *PLoS One* **7**, e44659 (2012).
85. Jiang, D., Tian, Z. & Lang, X. Reliability of climate models for China through the IPCC Third to Fifth Assessment Reports. *Int. J. Climatol.* **36**, 1114–1133 (2016).
86. Chen, L. & Frauenfeld, O. W. A comprehensive evaluation of precipitation simulations over China based on CMIP5 multimodel ensemble projections. *J. Geophys. Res. Atmos.* **119**, 5767–5786 (2014).
87. Wang, S. *et al.* Spatial and seasonal isotope variability in precipitation across China: Monthly isoscapes based on regionalized fuzzy clustering. *Journal of Climate* **35**, 3411–3425 (2022).
88. Rozanski, K., Araguás-Araguás, L. & Gonfiantini, R. Isotopic patterns in modern global precipitation. *AGU Geophys. Monogr.* **78**, 1–36 (1993).
89. Zhao, L. *et al.* Factors controlling spatial and seasonal distributions of precipitation  $\delta^{18}\text{O}$  in China. *Hydrol. Process.* **26**, 143–152 (2012).
90. Araguás-Araguás, L., Froehlich, K. & Rozanski, K. Stable isotope composition of precipitation over southeast Asia. *J. Geophys. Res. Atmos.* **103**, 28721–28742 (1998).
91. Li, Z. *et al.* Contributions of local terrestrial evaporation and transpiration to precipitation using  $\delta^{18}\text{O}$  and D-excess as a proxy in Shiyang inland river basin in China. *Global Planet. Change* **146**, 140–151 (2016).
92. Liu, J. *et al.* Characteristics of  $\delta^{18}\text{O}$  in precipitation over Eastern Monsoon China and the water vapor sources. *Chin. Sci. Bull.* **55**, 200–211 (2009).
93. Sun, C. *et al.* Spatial and temporal characteristics of stable isotopes in the Tarim River Basin. *Isotopes Environ. Health Stud.* **52**, 281–297 (2016).
94. Yao, J. *et al.* Climatic and associated atmospheric water cycle changes over the Xinjiang, China. *J. Hydrol.* **585**, 124823 (2020).
95. Tang, Y. *et al.* Using stable isotopes to understand seasonal and interannual dynamics in moisture sources and atmospheric circulation in precipitation. *Hydrol. Process.* **31**, 4682–4692 (2017).
96. Dansgaard, W. Stable isotopes in precipitation. *Tellus* **16**, 436–468 (1964).
97. Yamanaka, T., Tsujimura, M., Oyunbaatar, D. & Davaa, G. Isotopic variation of precipitation over eastern Mongolia and its implication for the atmospheric water cycle. *J. Hydrol.* **333**, 21–34 (2007).
98. Vuille, M., Werner, M., Bradley, R. S. & Keimig, F. Stable isotopes in precipitation in the Asian monsoon region. *J. Geophys. Res. Atmos.* **110**, D23108 (2005).
99. Yu, W. *et al.* Temperature signals of ice core and speleothem isotopic records from Asian monsoon region as indicated by precipitation  $\delta^{18}\text{O}$ . *Earth Planet. Sci. Lett.* **554** (2021).
100. Ren, G. *et al.* Recent progress in studies of climate change in China. *Adv. Atmos. Sci.* **29**, 958–977 (2012).
101. Ding, Y. *et al.* Detection, causes and projection of climate change over China: An overview of recent progress. *Adv. Atmos. Sci.* **24**, 954–971 (2007).
102. Yang, X., Yao, T., Yang, W., Yu, W. & Qu, D. Co-existence of temperature and amount effects on precipitation  $\delta^{18}\text{O}$  in the Asian monsoon region. *Geophys. Res. Lett.* **38**, L21809 (2011).
103. Yang, X., Davis, M. E., Acharya, S. & Yao, T. Asian monsoon variations revealed from stable isotopes in precipitation. *Clim. Dyn.* **51**, 2267–2283 (2017).
104. Wen, X., Wu, X. & Gao, M. Spatiotemporal variability of temperature and precipitation in Gansu Province (Northwest China) during 1951–2015. *Atmos. Res.* **197**, 132–149 (2017).
105. Ma, L., Li, H., Liu, T. & Liang, L. Abrupt temperature change and a warming hiatus from 1951 to 2014 in Inner Mongolia, China. *J. Arid Land* **11**, 192–207 (2019).
106. Qin, N., Chen, X., Fu, G., Zhai, J. & Xue, X. Precipitation and temperature trends for the Southwest China: 1960–2007. *Hydrol. Process.* **24**, 3733–3744 (2010).
107. Liu, B., Xu, M., Henderson, M. & Qi, Y. Observed trends of precipitation amount, frequency, and intensity in China, 1960–2000. *J. Geophys. Res. Atmos.* **110** (2005).
108. Oyebo, O. & Stretch, D. Neural network modeling of hydrological systems: A review of implementation techniques. *Nat. Resour. Model.* **32**, e12189 (2019).
109. Tian, L. *et al.* Stable isotopic variations in west China: A consideration of moisture sources. *J. Geophys. Res. Atmos.* **112**, D10112 (2007).
110. Peng, P., Zhang, X. J. & Chen, J. Modeling the contributions of oceanic moisture to summer precipitation in eastern China using  $\delta^{18}\text{O}$ . *J. Hydrol.* **581**, 124304 (2020).
111. Lee, J.-E., Fung, I., DePaolo, D. J. & Henning, C. C. Analysis of the global distribution of water isotopes using the NCAR atmospheric general circulation model. *J. Geophys. Res. Atmos.* **112**, D16306 (2007).

## Acknowledgements

This work was partially supported by the National Natural Science Foundation of China (grant nos. U2240201, 52109007 and 52079093), the Wuhan knowledge innovation project (grant no. 2022020801010106), and the Overseas Expertise Introduction Project for Discipline Innovation (111 Project) funded by the Ministry of Education and State Administration of Foreign Experts Affairs, P. R. China (grant no. B18037). The zoomed LMDZ simulation was performed using HPC resources from GENCI-IDRIS (Project 0292).

## Author contributions

Jiacheng Chen collected data, performed the analyses, and wrote the manuscript. Jie Chen designed the framework of the research, supervised the progress of research, revised the manuscript and provided resources and financial support. Xunchang J. Zhang designed the framework of the research and revised the manuscript. Peiyi Peng and Camille Risi provided data and methodological support and revised the manuscript.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41597-023-02095-1>.

**Correspondence** and requests for materials should be addressed to J.C.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2023