



# Deep learning tool: Reconstruction of long missing climate data based on multilayer perceptron(MLP)

Zhang Yan<sup>2</sup>; Xu Tianxin<sup>1</sup>; Zhang Chenjia<sup>1\*</sup>; Ma Daokun<sup>1</sup>

<sup>1</sup>College of Information and Electrical Engineering, China Agricultural University Haidian District, Beijing 100089, China

<sup>2</sup>School of Yi Language and Culture, Xichang University, Xichang, Sichuan 615000, China

\* Correspondence to: Zhang Chenjia, Email: zcj136578646@outlook.com

1 **Abstract:** Long-term monitoring of climate data is significant for grasping the law and  
2 development trend of climate change and guaranteeing food security. However, some weather  
3 stations lack monitoring data for even decades. In this study, 62 years of historical monitoring  
4 data from 105 weather stations in Xinjiang were used for missing sequence prediction,  
5 validating proposed data reconstruction tool. First of all, study area was divided into three  
6 parts according to the climatic characteristics and geographical locations. A deep learning tool  
7 based on multilayer perceptron (MLP) was established to reconstruct meteorological data with  
8 three time scales ( Short term, cycle and long term ) and one spatio dimension as inputing,  
9 filling in long sequence blank data. By designing an end-to-end model to autonomously detect



10 the locations of missing data and make rolling predictions,we obtained complete  
11 meteorological monitoring data of Xinjiang from 1961 to 2022. Seven kinds of parameter  
12 reconstructed include maximum temperature (Max\_T), minimum temperature (Min\_T), mean  
13 temperature (Ave \_ T), average water vapor pressure (Ave \_ WVP), relative humidity (Ave \_  
14 RH), average wind speed (10 m Ave \_ WS), and sunshine duration (Sun\_H). The quality of  
15 reconstructed data was evaluated by calculating correlation coefficient with the monitored  
16 sequences of nearest station. Results show that,proposed model reached satisfied average  
17 correlation coefficient for Max\_T, Min\_T, Ave \_ T and Ave \_ WVP parameters are 0.969,  
18 0.961, 0.971 and 0.942 respectively. The average correlation coefficient of Sun\_H and Ave \_  
19 RH are 0.720 and 0.789. Although it is difficult to predict extreme values, it can still capture  
20 the period and trend; the reconstruction effect of 10 m Ave \_ WS is poor, with the average  
21 similarity of 0.488. Finally, we published the trained parameter files and prediction codes as a  
22 micro service on the Agricultural Smart Brain platform, which provides firstly a deep learning  
23 tool for rapid and reliable reconstruction of meteorological monitoring data.

24 **Keywords:**Deep learning; Meteorological monitoring; Data reconstruction; MLP

## 25 **1 Introduction**

26 Agriculture, as the most fundamental industry for human, is facing a serious threat of  
27 climate change.Meteorological disasters account for over 70% of the natural disasters  
28 globally,have caused serious economic losses (Qin et al. 2002). More severely, global climate  
29 is changing dramatically,due to amount of greenhouse gases human activities produced.  
30 Because of climate warming, the frequency and intensity of drought and flood are increasing,  
31 and the harm to agricultural production is increasingly intensified (IPCC. 2012). From 2010 to



32 2017, global average annual economic losses due to drought reached US \$23.125 billion, with  
33 annual grain production cuts ranging from millions of tons to more than 30 million tons (Buda  
34 et al. 2018). In past 40 years, flooding events caused over a trillion dollars damage as well  
35 (UNDRR. 2020). Monitoring and analysis of meteorological data are able to reduce food and  
36 economic losses (Ziolkowska & Zubillaga. 2018).

37 In order to effectively guide agricultural production, meteorological monitoring usually  
38 includes meteorological parameters such as temperature, humidity, water vapor pressure, wind  
39 speed and sunshine. What is more, conducting research on climate forecasting requires  
40 long-term, large-scale and comprehensive climate data (Bonnet et al. 2020). Governments and  
41 scientific communities have been committed to the construction of meteorological databases  
42 (Anderson et al. 2008), a large number of professional meteorological monitoring stations  
43 have been established around the world, including China. However, there is much missing  
44 historical data due to temporal differences of monitoring station establishment, sensor failure,  
45 and other reasons. Thus, it is crucial to reconstruct the complete meteorological monitoring  
46 data.

47 Usually, researchers use interpolation method combined with manual correction to  
48 reconstruct missing meteorological data. Which not only consumes a lot of manpower, but also,  
49 due to the spatial variability of geographical conditions, the data results reconstructed by the  
50 traditional method are too smooth and inaccurate (Yao et al. 2023). Machine learning is a  
51 better interpolation tool (Li et al. 2020), but which performed poorly when deal with  
52 long-sequence missing data scenarios. A simple and efficient method for data reconstruction,  
53 deep neural network, has a great potential in meteorological data reconstruction tasks. The



54 neurons in the hidden layers of the neural network can constantly update the weights under  
55 supervision of true value, learning high-dimensional association among different data, and  
56 more accurately complement missing data (Rajae et al. 2019). In the task of reconstructing  
57 monitoring data of turbomechanical particle flow. Deep learning method was more accurate  
58 compared with six commonly used interpolation methods (Ghasem&Nader. 2022). In fact, in  
59 the field of weather forecasting, some available deep learning models have been published.  
60 Training based on large amounts of data, FourCastNet2 can calculate the next 24 hours of  
61 climate for 100 sites in just 7s (Jaideep et al. 2022),orders of magnitude faster than the  
62 numerical weather prediction (NWP). The Pangu model proposed by Huawei team can  
63 accurately and quickly predict the global climate by learning global meteorological  
64 monitoring data of past four decades (Bi et al. 2023).

65 Selection and design of neural networks is a key step in reconstructing climate data.  
66 Ghose selected recursive neural network (RNN) for groundwater level prediction (Ghose et al.  
67 2018).Vu reconstructed 50 years groundwater level data in Normandy (France) based on the  
68 long and short-term memory (LSTM) (Vu et al. 2020). Differently, meteorological parameters  
69 are greatly spatially correlated, as a typical spatiotemporal sequence.Nature Subissue  
70 Geoscience published related research, which using image restoration technology combined  
71 with HadCRUT4 global historical temperature grid dataset, reconstructed complete global  
72 monthly grid temperature, and the reconstructed data sequence has extremely high correlation  
73 with the non-reconstructed data (Christopher et al. 2020). Continuity of time and spatial  
74 correlation must be considered simultaneously in the data reconstruction.Most of the frontier  
75 studies of spatio-temporal prediction are modeling based on graph neural network (GNN) and



76 Transformer (Pan&Li. 2021). But they have high computational complexity and memory  
77 overhead. Although MLP is a relatively simple deep learning model, the ability of  
78 spatiotemporal prediction is not inferior to complex models in recent studies.The MLPST  
79 model shows that, compared with RNN, GNN and Transformer, it can be very accurate even  
80 completely based on MLP (Zhang et al. 2023).Usually, different type of time series data have  
81 different characteristics, and screening some obvious characteristics can significantly improve  
82 the model performance (Tang et al. 2024).Therefore, in specific tasks, feature engineering and  
83 special model design need to be carried out to improve the prediction performance of the  
84 model.

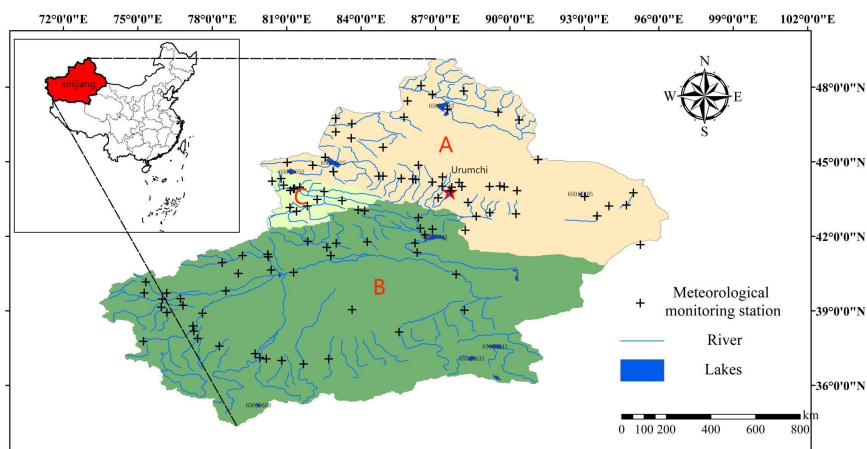
85 To meet the demand of meteorological data reconstruction in agricultural productions,  
86 we designed a neural network models as a reconstruction tool based on MLP. A total of 143  
87 missing data (43 weather stations) were reconstructed obtained from three divided study areas  
88 in Xinjiang.The parameters reconstructed include Max\_T, Min\_T, Ave \_ T, Ave \_ WVP, Ave \_  
89 RH, 10m Ave \_ WS, and Sun\_H. Inputs make up with short-term, cyclical, long-term trends  
90 and the same time data of weather stations with the highest sequence similarity, length of  
91 filled sequences ranged from one month to 38 years.The confidence of the results is measured  
92 by the correlation with the most adjacent station. Finally, datasets automatic construction  
93 module, automatic training module, missing positions automatic query module, and automatic  
94 rolling prediction module are integrated, realizing end-to-end data reconstruction and  
95 published as a micro service.

## 96 **2 Study area and data**

97 The study area is located in Xinjiang, northwest of China. Which is one of the most



198 important cotton production bases in China and most developed drought agricultural  
199 technology region(Liu. 2022). Located in the hinterland of Eurasia, due to complex terrain  
200 and frequent weather system activity, drought is the main climatic feature of this region (Mao  
201 et al. 2008). The Tianshan Mountains crosses the central region, divides Xinjiang into  
202 northern Xinjiang and southern Xinjiang. The water vapor could enter northern Xinjiang but  
203 hardly reach southern Xinjiang, so the drought degree difference of drought between the north  
204 and the south is obvious (Wang.2023). Yili River Valley located in west of the Tianshan  
205 Mountains in Xinjiang, surrounded by mountains on three sides, with abundant precipitation,  
206 forming special climate (Yan et al. 2017).



107  
108 **Figure 1. Study area division and location of weather stations**

109 105 weather stations in this study are distributed in three areas: northern Xinjiang (A),  
110 southern Xinjiang (B) and Yili (C), having 48, 44, 13 sites respectively, recorded  
111 meteorological data for nearly 62 years (Figure 1). Among those, 44 stations exist data  
112 missing in varying degrees, with missing parameter types and missing duration varied. Table  
113 1 listed codes and parameters of weather stations with missing data. Of these, 16 stations exist

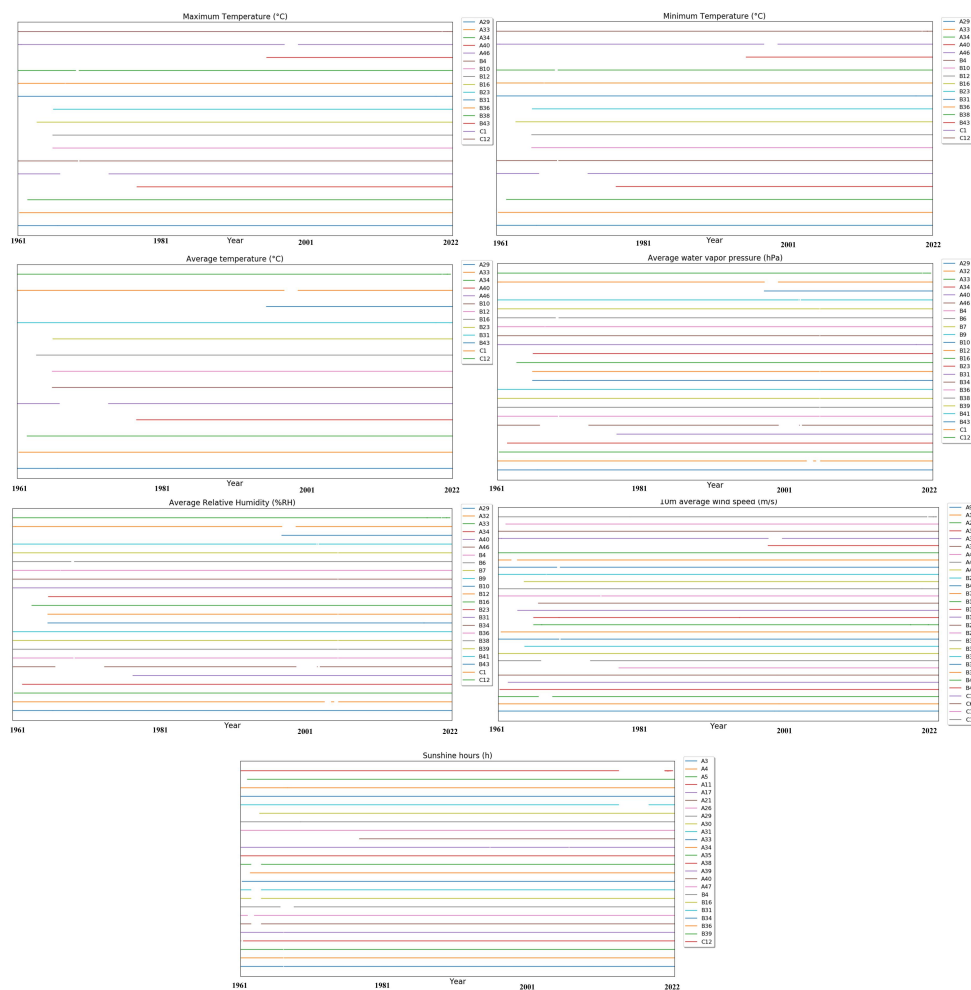


114 Max\_T and Min\_T data missing. The number of Ave\_T is 13, Ave\_WVP and Ave\_RH is 23,  
 115 10m Ave\_WS is 28, Sun\_H is 24. In totally, we need to reconstruct a total of 143  
 116 sequences,with time spans from 1961 to 2022. Figure 2 corresponding to Table 1, shows that  
 117 the specific missing period, these missing lengths are long, the missing location is different,  
 118 increased the difficulty of the reconstruction.

119 **Table 1. Meteorological parameter types of missing data and weather station number of subtasks**

Meteorological parameters	Weather station number of subtasks
Max_T(°C)	A29,A33,A34,A40,A46,B4,B10,B12,B16,B23,B31,B36,B38,B43,C1,C12
Min_T(°C)	A29,A33,A34,A40,A46,B4,B10,B12,B16,B23,B31,B36,B38,B43,C1,C12
Ave_T(°C)	A29,A33,A34,A40,A46,B10,B12,B16,B23,B31,B43,C1,C12
Ave_WVP(hPa)	A29,A32,A33,A34,A40,A46,B4,B6,B7,B9,B10,B12,B16, B23,B31,B34,B36,B38,B39,B41,B43,C1,C12
Ave_RH(%RH)	A29,A32,A33,A34,A40,A46,B4,B6,B7,B9,B10,B12,B16, B23,B31,B34,B36,B38,B39,B41,B43,C1,C12
10mAve_WS(m/s)	A9,A16,A29,A33,A34,A35,A40,A46,A48,B2,B4,B7,B10,B12,B16, B23,B26,B31,B34,B36,B38,B39,B40,B43,C1,C6,C11,C12
Sun_H (h)	A3,A4,A5,A11,A17,A21,A26,A29,A30,A31,A33,A34,A35,A38, A39,A40,A47,B4,B16,B31,B34,B36,B39,C12

120



121

122 **Figure 2. Measurement time-window at 105 weather stations over 62 years from 1961 to 2022**

123 **3 Methodology and model design**

124 **3.1 MLP**

125 MLP is the most classic deep neural network, widely used to solve the classification and  
126 regression problems of nonlinearity. Whose structure (Figure 3 lower-right) includes input  
127 layer, hidden layer and output layer, each layer contains several neurons, and neurons in the  
128 upper and lower layers are connected to each other for information exchange (Benedict. 1988).





129 When training it, the weight parameters of the neurons are constantly updated until a good fit  
130 is achieved. Forward propagation and backpropagation are required to complete each time the  
131 weights are updated (Rumelhart et al. 1986). Forward propagation takes the outputs of the  
132 previous layer as the inputs of the next layer, calculates the outputs of the next layer according  
133 to the weight. Consider the layer1 and layer2 as examples, outputs of the layer2 is:

$$134 \quad a_i^{layer2} = \sigma(b_i + \sum w_j a_j^{layer1})$$

135 Where  $\sigma$  is the activation function, which is the key for MLP to achieve nonlinear fitting.  
136 The most commonly used activation function, ReLU, is selected in this study (Glorot et al.  
137 2011):

$$138 \quad f(x) = \max(0, x)$$

139 Prediction errors is measured by cost function LOSS. The back propagation process is  
140 based on the chain conduction law, calculating gradients each layer parameters in the  
141 network to represent the influence of the parameters on the prediction errors, and updating  
142 the weight through multiply by learning rate  $\alpha$ , until the loss value no longer drops. Which  
143 can be considered that the MLP model fitting has reached the optimal solution. The initial  
144 learning rate selected for this study was 0.001. The backpropagation process is:

$$145 \quad w_{jnew} = w_j - \alpha \bullet \frac{\partial LOSS(y, \hat{y})}{\partial w_j}$$

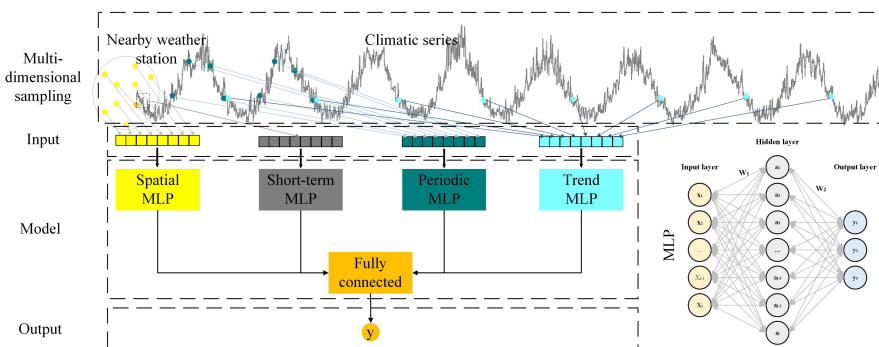
### 146 **3.2 Spatiotemporal MLP**

147 Past studies have proved that, climate shows a short-term dependence, and which is  
148 cyclical and shows a trend in the long term (Liu et al. 2020), and closely associated with the  
149 adjacent site data. According to these experiences, we designed four modules based on the  
150 MLP (Figure 3): Spatial MLP, Short-term MLP, Periodic MLP, Trend MLP. Time series, with



151 different time scales resampled, were entered separately Short-term MLP, Periodic MLP,  
152 Trend MLP models, extracting short-term trends, cyclical and long-term trend characteristics  
153 of historical data respectively. Monitoring values from nearby stations were fed into the  
154 Spatial MLP module to obtain spatial associations between them. Results of the four modules  
155 are combined as inputs of predictive header, two fully connected layers. Which enable  
156 spatio-temporal association in the sequence is captured.

157 Dataset size and input sequence length are negatively correlated, need to be balanced  
158 when designing the inputs. In our model, all of the inputs length were set to 8, ensuring input  
159 format is unified. Inputs of short-term MLP are values last 8 days. Inputs of Periodic MLP  
160 and Trend MLP are values resampled according to 90 days intervals and 365 days respectively.  
161 Inputs of Spatial MLP were the monitoring values of eight stations with highest similarity to  
162 the target sequence within the study region. Using pyramid structure, the number of neurons  
163 in each layer is half the number of neurons in the previous layer, and which is usually able to  
164 extract features at different scales more effectively (Yang et al. 2020).



166 **Figure 3. MLP and Model Framework**

167 All meteorological data were standardized according to the following formula:



168 
$$x_i^{normal} = \frac{x_i - x_{min}}{x_{max} - x_{min}}$$

169 Where  $x_i^{normal}$  is normalized value,  $x_i$  is actual value,  $x_{max}$  and  $x_{min}$  are maximum and  
170 minimum value of sequence, respectively. This normalization method standardizes the values  
171 to between 0-1, be able to eliminate the effect of dimension and negative values for model  
172 fitting.

173 When predicting, we restore the results, and output the dimensional results:

174 
$$y_i^{pred} = \hat{y}_i \cdot (x_{max} - x_{min}) + x_{min}$$

175 Where  $y_i^{pred}$  is predicted value,  $\hat{y}_i$  is output value of neural network.

### 176 3.3 Assessment methods

177 Meteorological similarity is measured by Euclidean distance of the sequence  
178 commonly, but there is a big difference between the different parameters. In order to  
179 standardize this index to 0-1, we define similarity based on Euclidean distance of two  
180 sequences:

181 
$$SM_{mn} = \frac{1}{e^{\sum_{i=1}^n |y_i - y'_i| / 100 * n}}$$

182  $SM_{mn}$  represents the similarity of m sequence and n sequence,  $y_i, y'_i$  are the value of two  
183 sequences at the same time, respectively. n is the number of non-missing value. SM is closer  
184 to 1, the more similar, the closer to 0, the lower similarity. SM was used to select the inputs  
185 for the Spatial MLP module.

186 We used two common indicators, mean squared error (MSE) and mean absolute error  
187 (MAE), to evaluate the quality of the prediction:

188 
$$MAE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)$$



189 
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

190 where  $y_i$  is the real measure of GWL;  $\hat{y}_i$  is the estimated value of GWL; and  $\bar{y}_i$  is

191 the mean of  $y_i$ . MAE and MSE were used to select the best number of hidden layers.

192 When evaluating the credibility of the reconstructed data, we used the correlation

193 coefficient as the evaluation index:

194 
$$corr_{m-n} = \frac{\sum_{i=1}^n (h_i^m - \bar{h}^m)(h_i^n - \bar{h}^n)}{\sqrt{\sum_{i=1}^n (h_i^m - \bar{h}^m)^2} \sqrt{\sum_{i=1}^n (h_i^n - \bar{h}^n)^2}}$$

195  $h_i^m$  and  $h_i^n$  are value of m sequence and n sequence, respectively,  $\bar{h}^m$  and  $\bar{h}^n$  is the average

196 value of m sequence and n sequence, respectively. Correlation coefficient is closer to 1, the

197 reconstructed data is more credible, and correlation coefficient closer to 0, the more unreliable

198 it is.

## 199 **4 Reconstruction of missing climate data**

### 200 **4.1 Sub-task division**

201 The reconstruction task was divided into 21 scenarios, 143 sub-tasks depending on the

202 region and the parameters. Climate region A consists of 53 sub-tasks, among these, Max\_T,

203 Min\_T and Ave\_T take up 5 sub-tasks respectively; Ave\_WVP and Ave\_RH take up 6

204 sub-tasks respectively; Ave\_WS and Sun\_H take up 9 and 17 sub-tasks respectively. Climate

205 region B has 75 sub-tasks, Max\_T, Min\_T, Ave\_T, Ave\_WVP, Ave\_RH, Ave\_WS, Sun\_H take

206 up 9, 9, 6, 15, 15, 15, 6 sub-tasks respectively. The numbers in climate region C are 2, 2, 2, 2,

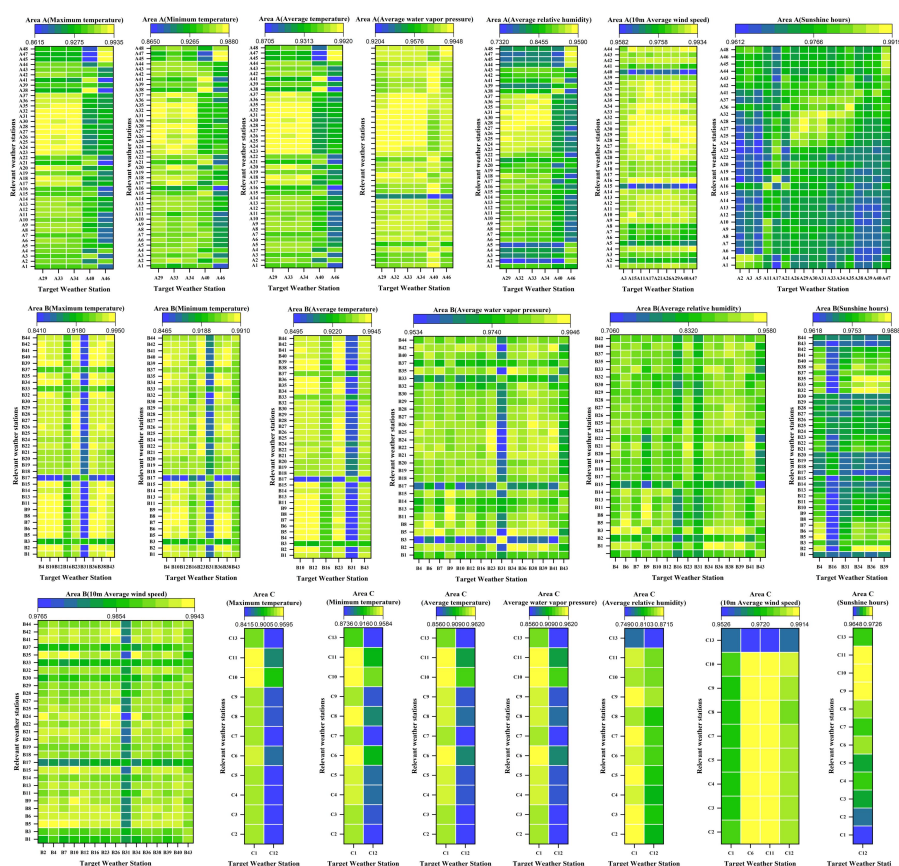
207 2, 4, 1 respectively, total sub-tasks numbers were 15.

208 The SM of each sequence was calculated and derived as SM table. We show the SM of

209 the target station and stations in the same region in Figure 4 due to the large amount of table



210 data. The more pronounced the yellow color, the higher the similarity, and the more  
211 pronounced blue the lower similarity, medium similarity shows green. Under all of the  
212 reconstruction scenarios, overall, the SM of these sequences ranged between 70.6% and  
213 99.48%.Based on SM, 8 stations with the highest similarity to each target task, 143 groups in  
214 total. Following the spatiotemporal sampling method shown in Figure 3, build model inputs.



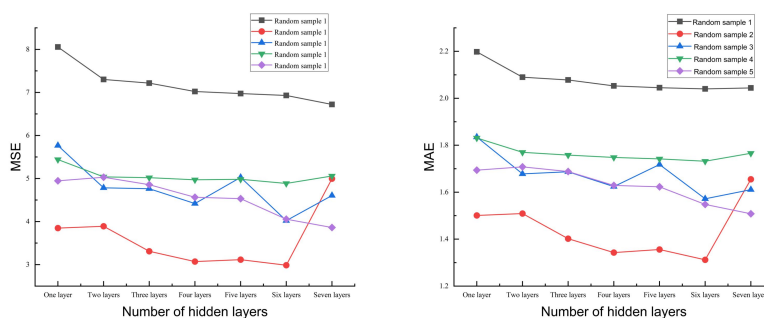
215  
216 **Figure 4. The similarity between the target weather station and the related weather station (under**  
217 **the sequence reconstruction scenario of different regions and different parameters)**

### 218 4.2 Determine the number of MLP hidden layers

219 The number of layers of the MLP hidden layer is one of the most important parameters  
220 of model. Generally, increasing hidden layers can enhance the fitting ability of model, but



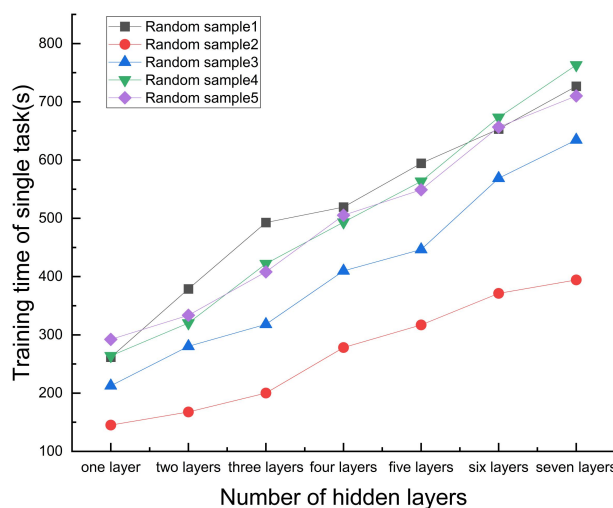
221 enhancement is limited, while largely increasing the number of model parameters, and  
222 causing slower run speed of model. To determine the number of hidden layers, we randomly  
223 picked five datasets, testing MSE, MAE and training time of the model prediction. Numbers  
224 of hidden layers is seted from 1 to 7(training 1000 epochs). Figure 5 displays, when setting  
225 1-4 hidden layers, MAE and MSE of the model did have a significant downward trend as the  
226 hidden layer increasing. But when number of hidden layers greater than 4, MAE and MSE  
227 showed little improvement, even rose on some datasets. This may be related to the appearance  
228 of gradient explosion when model is too deep.  
229



230

231 **Figure 5. MAE and MSE trend with the number of hidden layers increases**

232 Figure 6 shows, time consumption to complete training increased significantly with the  
233 increase of the number of hidden layers, on the five randomly selected datasets, which  
234 displays almost linear. Considering the results of the above trials, the number of hidden layers  
235 chosed for MLP is 4, and the longest time consumption for a single task is 519.35s.



236

237

**Figure 6. Time-consuming trend of training model with the number of hidden layers increases**

238

### **4.3 Train model and reconstruct missing data**

239

240

241

242

243

244

245

246

247

248

249

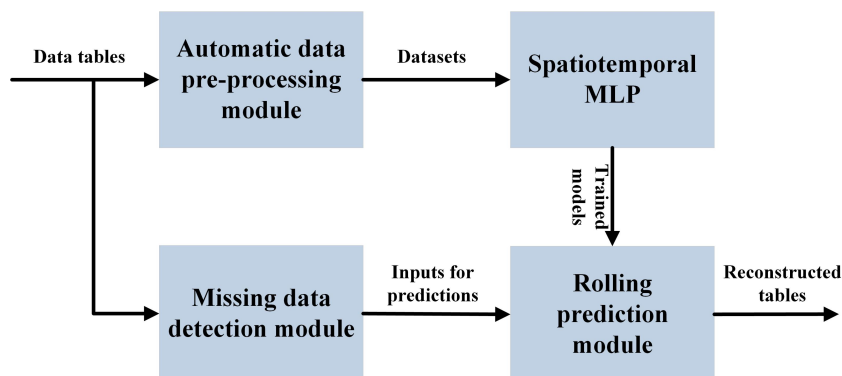
250

Figure 7 shows our prediction process. Using pycharm2022.1 for our programming, we integrated multiple modules to implement end-to-end programs, with pre-processing data automatically, training model automatically, detecting missing data automatically, and reconstructing data automatically. The situation of missing data is complex, and the task of manually constructing datasets is large. Automatic data pre-processing model could generate datasets by reading data sheets according to task list, and complete normalization. Data were disrupted the order before entering. Automatic training model could complete multiple tasks and save as different parameter files. When predicting, missing data detection model could detect location of missing data. Later, according to detecting results, rolling prediction model automatically forward or backward predicting.

Tensorflow (Martín et al. 2016) was selected to be development framework in this study. Using Adaptive Moment Estimation (Adam) optimizer to improve learning efficiency



251 (Kingma&Ba. 2014), it can adjust automatically the learning rate according to historical  
252 gradient information. At the beginning of training, the larger learning rate helps the model to  
253 converge quickly. While later, learning rate adjusts smaller to improve accuracy of  
254 model. Meanwhile, Adam normalized the weight parameters, which also alleviates overfitting.  
255 MSE was choosed to be loss function. As a skill of training, Dropout layer can effectively  
256 prevent model overfitting (Nitish et al. 2014). In this study, the super-parameter of Dropout  
257 layers was set to 0.5.

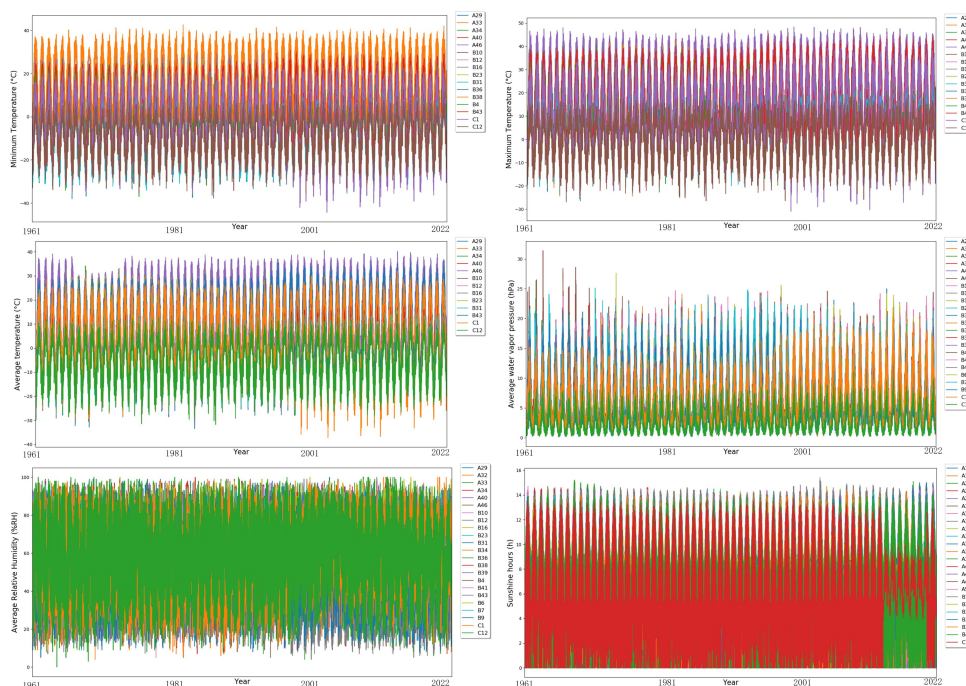


258

259

Figure 7. Structure of ensemble forecasting





260

261

**Figure 8. Results of reconstructed sequence(except Ave\_WS)**

262

Figure 8 shows reconstruction results of Max\_T, Min\_T, Ave\_T, Ave\_WVP, Ave\_RH and

263

Sun\_H. Except Sun\_H, from the figure, reconstructed sequences of other five parameters are

264

indistinguishable from the real sequences. Sun\_H usually represents time length, that the solar

265

radiation above certain intensity. Which influenced by all kinds of meteorological factors,

266

especially the change of the clouds. Our prediction values almost no 0 while measured exists

267

some 0 values, can not mining to the occurrence of 0 values. But we can clearly see that, even

268

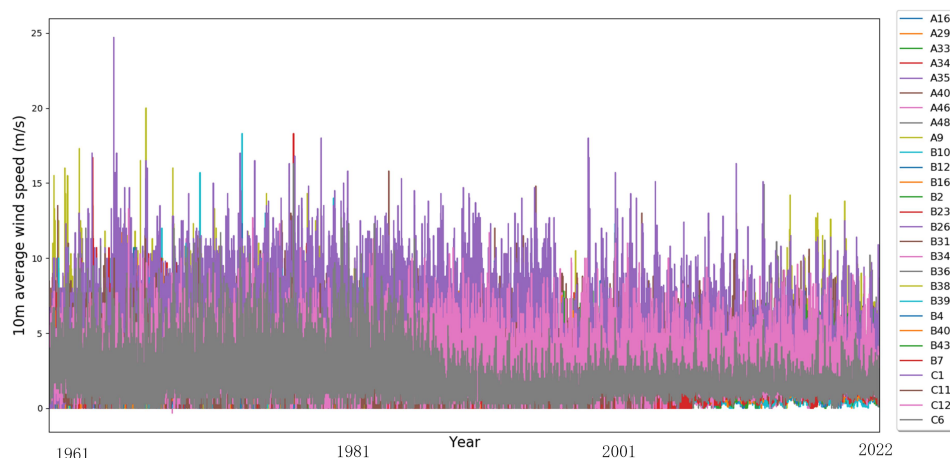
without filtering, proposed model could dig out the cycle and trend laws. These reconstructed

269

Sun\_H data still very useful. For Ave\_WS (Figure 9), model is difficult to predict the

270

suddenly extreme wind speed. Which also shows shortage when predicting extreme values.



271

272

Figure 9. Result of reconstructed sequence(Ave\_WS)

273

#### 4.4 Evaluate quality of reconstructed data

274

Compared with the visual evaluation, the evaluation index can give more information

275

about the reconstruction results. Assessing the quality of reconstructed data is very difficult

276

due to the difficulty in tracing the real data of the past. By calculating correlation coefficient

277

of sequences between reconstructed data and the nearest weather station data, credibility of

278

reconstructed meteorological data was scientifically evaluated. A higher correlation

279

coefficient indicates a higher confidence.

280

Table 2. Correlation of reconstructed sequences and nearest neighbor sequences

Weather station	Meteorological parameter Types						
	Max_T (°C)	Min_T (°C)	Ave_T (°C)	Ave_WVP (hPa)	Ave_RH (%RH)	10m Ave_WS (m/s)	Sun_H (h)
A11-A10							0.838
A16-A10						0.521	
A17-A23							0.878
A21-A43							0.755
A26-A27							0.858
A29-A27	0.998	0.992	0.998	0.984	0.948	0.568	0.943
A3-A2							0.818
A30-A37							0.850
A31-A27							0.903
A32-A37				0.967	0.926		
A33-A41	0.937	0.952	0.955	0.961	0.450	0.314	0.774
A34-A37	0.987	0.982	0.989	0.959	0.927	0.371	0.879
A35-A36						0.681	0.944
A38-A37							0.780



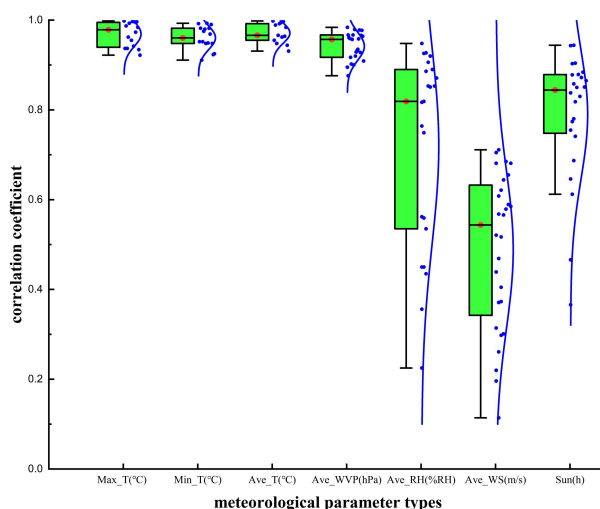
A39-A27							0.830
A4-A1							0.741
A40-A1	0.937	0.911	0.948	0.917	0.225	0.220	0.466
A46-A36	0.962	0.952	0.966	0.902	0.562	0.196	
A47-A45							0.872
A48-A43						0.261	
A5-A2							0.646
A9-A8						0.517	
B10-B11	0.993	0.975	0.992	0.957	0.764	0.566	
B12-B11	0.994	0.980	0.994	0.967	0.817	0.579	
B16-B17	0.955	0.948	0.962	0.900	0.559	0.373	0.612
B2-B35						0.705	
B23-B22	0.997	0.993	0.997	0.979	0.906	0.608	
B26-B25						0.711	
B31-B33	0.942	0.949	0.964	0.895	0.450	0.589	0.687
B34-B32				0.920	0.819	0.621	0.884
B36-B32	0.996	0.982		0.929	0.886	0.644	0.904
B38-B37	0.973	0.969		0.932	0.535	0.439	
B39-B32				0.935	0.853	0.685	0.851
B4-B5	0.996	0.990		0.958	0.851	0.655	0.865
B40-B21						0.469	
B41-B22				0.978	0.890		
B43-B14	0.984	0.948	0.984	0.928	0.749	0.405	
B6-B5				0.966	0.853		
B7-B5				0.977	0.920	0.681	
B9-B15				0.964	0.871		
C1-C10	0.934	0.923	0.944	0.876	0.435	0.301	
C11-C8						0.298	
C12-C3	0.922	0.925	0.931	0.909	0.356	0.114	0.366
C6-C8						0.585	

281 Correlation coefficients of all reconstructed sub-tasks are shown in Table 2. From the  
 282 reconstruction effect of temperature, we very approach the results of Christopher (0.9941)  
 283 (Christopher et al.2020), even exceed theirs in 4 tasks of temperature reconstruction (45 in  
 284 total). More importantly, the data we reconstructed are of daily scale, smaller than their time  
 285 granularity (monthly). Our work demonstrates that MLP with special spatiotemporal design  
 286 can better reconstruct climate data.

287 The consistency of evaluation indicators in different tasks is also one of the goals we  
 288 pursue. Figure 10 shows the distribution of its correlation coefficient index. Which can be  
 289 easily see, Max\_T, Min\_T, Ave\_T, Ave\_WVP shows excellent performance, with average  
 290 correlation coefficient is over 0.9 and distribution is very concentrated; Ave\_RH and Sun\_H  
 291 performance unstable, although average correlation coefficient is over 0.7 but dispersed;



292 Ave\_WS shows poor results, average correlation coefficient is around 0.5 and distributing is  
293 dispersed.



294

295 **Figure 10. The evaluation indicators distribution of different type parameters reconstruction**

#### 296 **4.5 Release model**

297 In order to provide convenient services for people in the agriculture field, the model  
298 accomplished in this study will published on the Agricultural Smart Brain platform as a  
299 tool. Which developed by Beijing Lianchuang Siyuan Measurement and control Technology  
300 Co., LTD, providing scientific research data, computing power and publishing AI  
301 micro-services for scientific researchers. Users can obtain it by purchasing access authority  
302 to the platform. The link of our microservices as follow: [http://192.168.50.201:15000/  
303 app/services/visionarytech/test-1/alg-26bc86f42a137f8f](http://192.168.50.201:15000/app/services/visionarytech/test-1/alg-26bc86f42a137f8f)

#### 304 **5 Conclusion**

305 In order to reconstruct the long-term missing meteorological monitoring data in  
306 agricultural field, proposed an end-to-end rapid reconstruction method based on MLP.



307 Spatio-temporal datasets were built according to the similarity indicator SM,  
308 standardized data to ensure good performance of the model. Backbone was designed four  
309 MLP modules with 4-hidden layers to jointly learn short-term trends, periodicity, long-term  
310 trends, and spatial associations. Predictive head consisted with two fully connected layers.  
311 After that, the automatic preprocessing, automatic detection of missing data location,  
312 automatic model training and rolling prediction modules are coded and integrated to realize  
313 end-to-end long sequence reconstruction. Our model is able to complete a single  
314 reconstruction task within 10 minute.

315 Daily meteorological monitoring data of 44 meteorological stations (143 tasks) in  
316 Xinjiang from 1961 to 2022 were reconstructed using our method. The evaluation indexes  
317 show that, average correlation coefficient of Max\_T, Min\_T, Ave\_T and Ave\_WVP are  
318 0.969,0.961,0.971, and 0.942 respectively, showing high consistency and high credibility;  
319 average correlation coefficient of Ave\_RH and Sun\_H are 0.720 and 0.789 respectively,  
320 showing low consistency and general credibility; average correlation coefficient of Ave\_WS  
321 is 0.488, showing low consistency and low credibility. The results demonstrate that MLP was  
322 useful and reliable in the task of rapidly reconstructing missing meteorological data, which  
323 will provide an important solution to solve the problem of missing data in agrometeorological  
324 field.

325 Finally, we released our model on Agricultural Smart Brain platform, provided users a  
326 tool of data reconstruction, in the form of micro service.

**Conflict of interest:** None.



**Code/data availability:** Both the code and data are freely available by contacting the corresponding authors.

**Author contributions:**

ZhangYan: methodology;essay writing;proofreading of dissertations

XuTianxin: data processing;method validation;essay writing;proofreading of dissertations

Zhangchenjia: methodology;method validation;coding;proofreading of dissertations

MaDaokun: financial support;project management

327 **Thanks**

328 Thanks to financial support of Silk Road Economic Belt Innovation-Driven  
329 Development Pilot Zone, WuChangShi National Independent Innovation Demonstration Zone  
330 project(2022LQ04001) .

331 Thanks to Beijing Lianchuang Siyuan Measurement and Control Technology Co., Ltd.  
332 and Beijing Zhiyu Chuangyi Co., Ltd. for their help in the release of micro service.

333 **Reference**

334 Anderson, S.P., Bales, R.C., Duffy, C.J., 2008. Critical Zone Observatories: building a  
335 network to advance interdisciplinary study of Earth surface processes. Mineral. Mag. 72  
336 (1), 7 - 10.

337 Benedict A K. 1988. Learning in the multilayer perceptron. Journal of Physics A:  
338 Mathematical and General.21(11).

339 Bonnet, R., Bo´ e, J., Habets, F., 2020. Influence of multidecadal variability on high and low  
340 flows: the case of the Seine basin. Hydrol. Earth Syst. Sci. 24, 1611 - 1631.

341 Bi Kaifeng, Xie Lingxi, Zhang Hengheng, Chen Xin, Gu Xiaotao, Tian Qi. Accurate



342 medium-range global weather forecasting with 3D neural networks. *Nature* 619,  
343 533–538 (2023). <https://doi.org/10.1038/s41586-023-06185-3>

344 Buda Su, Jinlong Huang, T. Fischer, Yanjun Wang, Z. Kundzewicz, J. Zhai, Hemin  
345 Sun, Anqian Wang, X. Zeng, Guojie Wang, H. Tao, M. Gemmer, Xiucang Li, T. Jiang.  
346 2018. Drought losses in China might double between the 1.5 ° C and 2.0 °  
347 C warming. *Proc Natl Acad Sci USA*, 115 (42) : 10600-10605

348 Christopher Kadow, David Matthew Hall & Uwe Ulbrich. 2020. Artificial intelligence  
349 reconstructs missing climate information. *Nat. Geosci.* 13, 408 - 413 (2020).  
350 <https://doi.org/10.1038/s41561-020-0582-5>

351 Ghose, D., Das, U., Roy, P., 2018. Modeling response of run off and evapotranspiration for  
352 predicting water table depth in arid region using dynamic recurrent neural network.  
353 *Groundwater Sustainable Dev.* 6, 263 - 269.

354 Ghasem A ,Nader M .2022. Reconstruction of particle image velocimetry data using  
355 flow-based features and validation index: a machine learning approach. *Measurement*  
356 *Science and Technology*, 2022, 33(1)

357 Glorot, Xavier, Antoine Bordes, Yoshua Bengio.2011. “Deep Sparse Rectifier Neural  
358 Networks.” *International Conference on Artificial Intelligence and Statistics* (2011).

359 IPCC. 2012. Summary for policymakers // *Managing the Risks of Extreme*  
360 *Events and Disasters to Advance Climate Change Adaptation. A Special*  
361 *Report of Working Groups I and II of the Intergovernmental Panel on Climate Change.*  
362 *Cambridge: Cambridge University Press*, 1-19

363 Jaideep Pathak, Shashank Subramanian, Peter Harrington, Sanjeev Raja, Ashesh



- 364 Chattopadhyay, Morteza Mardani, Thorsten Kurth, David Hall, Zongyi Li, Kamyar  
365 Azizzadenesheli, Pedram Hassanzadeh, Karthik Kashinath, Animashree Anandkumar.  
366 FourCastNet: a global data-driven high-resolution weather model using adaptive Fourier  
367 neural operators. Preprint at <https://arxiv.org/abs/2202.11214> (2022).
- 368 Kingma P D ,Ba J . 2014. Adam: A Method for Stochastic Optimization. CoRR 2014,  
369 abs/1412.6980
- 370 Li C ,Ren X ,Zhao G .2023. Machine-Learning-Based Imputation Method for Filling Missing  
371 Values in Ground Meteorological Observation Data. Algorithms, 2023, 16 (9):
- 372 Liu Kai, Nie Gege, Zhang Sen. 2020. Study on the Spatiotemporal Evolution of Temperature  
373 and Precipitation in China from 1951 to 2018. Advances in Earth Science, 2020, 35(11):  
374 1113-1126 DOI:10.11867/j.issn.1001-8166.2020.102
- 375 Liu Yi. 2022. Build a national high-quality cotton production base. Xinjiang Daily,  
376 2022-09-05 (001). DOI:10.28887/n.cnki.nxjrb.2022.003680.(in chinese)
- 377 Mao Weiyi,Nan Qinghong,Shi Hongzheng. 2008. Research on climate change characteristics  
378 and climate zoning methods in Xinjiang. Meteorological Calendar, 2008, (10): 67-73.(in  
379 chinese)
- 380 Martín Abadi, Paul Barham, Jianmin Chen, Zhifeng Chen, Andy Davis, Jeffrey Dean,  
381 Matthieu Devin, Sanjay Ghemawat, Geoffrey Irving, Michael Isard, Manjunath Kudlur,  
382 Josh Levenberg, Rajat Monga, Sherry Moore, Derek Gordon Murray, Benoit Steiner,  
383 Paul A. Tucker, Vijay Vasudevan, Pete Warden, Martin Wicke, Yuan Yu, Xiaoqiang  
384 Zhang. 2016. TensorFlow: A system for large-scale machine learning. CoRR, 2016,  
385 abs/1605.08695





- 386 Nitish Srivastava, Geoffrey E. Hinton, Alex Krizhevsky, Ilya Sutskever, Ruslan  
387 Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting.  
388 Journal of Machine Learning Research, 2014, 15 (1): 1929-1958.
- 389 Pan Zhisong, Li Wei.2021. Survey of Spatio-temporal Series Prediction Methods Based on  
390 Deep Learning. Data Acquisition and Processing, 2021, 36 (03): 436-448.  
391 DOI:10.16337/j.1004-9037.2021.03.003. (in Chinese)
- 392 Qin D H, Ding Y H, Wang S W, Wang S M, Dong G R, Lin E D, Liu C Q, She Z X,  
393 Sun H N , Wang S R , Wu G H. 2002.  
394 Ecological and environmental change in West China and its response strategy.  
395 Adv Earth Sci, 17 (3) : 314-319 (in Chinese)
- 396 Rajae, T., Ebrahimi, H., Nourani, V.2019. A review of the artificial intelligence methods in  
397 groundwater level modelling. J. Hydrol. 572, 336–351.
- 398 Rumelhart, D., Hinton, G., Williams, R. 1986. Learning representations by back-propagating  
399 errors. Nature 323, 533 – 536 (1986). <https://doi.org/10.1038/323533a0>
- 400 Tang D ,Zhan Y ,Yang F. 2024. A review of machine learning for modeling air quality:  
401 Overlooked but important issues. Atmospheric Research, 2024, 300, 107261-.
- 402 UNDRR. Human Cost of Disasters: An Overview of the Last 20 Years 2000–2019; United  
403 Nations for Disaster Risk Reduction (UNISDR): Geneva, Switzerland, 2020.
- 404 Vu M.T., Jardani A., Massei N., Fournier M. 2020. Reconstruction of missing groundwater  
405 level data by using Long Short-Term Memory (LSTM) deep neural network. Journal of  
406 Hydrology, 2020, (prepublish): 125776-.
- 407 Wang Jiaoyan.2023. Distribution and evolution characteristics of drought under the



- 408 background of climate warming and humidification in Xinjiang. Arid Environment  
409 Monitoring, 2023, 37 (01): 15-21.(in chinese)
- 410 Yan Junjie, Yan Min, Cui Dong, Liu Haijun, Chen Chen, Xia Qianqian. 2017. Analysis of  
411 temperature and precipitation trends in the Ili River Valley of Xinjiang in the past 55  
412 years. Hydropower Energy Science, 2017, 35 (10): 13-16+12.
- 413 Yang C., Xu Y., Shi J., Dai B., Zhou B. 2020. Temporal pyramid network for action  
414 recognition. Proceedings of the IEEE Computer Society Conference on Computer Vision  
415 and Pattern Recognition, 2020, 588-597.
- 416 Yao Ziqiang, Zhang Tao, Wu Li, Wang Xiaoying, Huang Jianqiang. Physics-Informed Deep  
417 Learning for Reconstruction of Spatial Missing Climate Information in the  
418 Antarctic. Atmosphere, 2023, 14 (4)
- 419 Zhang Zijian, Huang Ze, Hu Zhiwei, Zhao Xiangyu, Wang Wanyu, Liu Zitao, Zhang Junbo,  
420 Qin S. Joe, Zhao Hongwei. 2023. MLPST: MLP is All You Need for Spatio-Temporal  
421 Prediction. In Proceedings of the 32nd ACM International Conference on Information  
422 and Knowledge Management (CIKM '23). Association for Computing Machinery, New  
423 York, NY, USA, 3381 - 3390. <https://doi.org/10.1145/3583780.3614969>
- 424 Ziolkowska Jadwiga R, Zubillaga Jesus. 2018. Importance of weather monitoring for  
425 agricultural decision-making - an exploratory behavioral study for Oklahoma Mesonet..  
426 Journal of the science of food and agriculture (13), 4945-4954.