

1 **Journal**

2 Climatic Change

3

4 **Title**

5 Changes of diurnal temperature range over East Asia from 1901 to 2018 and its relationship  
6 with precipitation

7

8 **Supplementary Material**

9 1. Figure S1

10 2. CMA-LSAT v1.1 dataset

11 3. Classification method of urban and rural stations in East Asia (Figure S1)

12

13 **Authors**

14 Xiubao Sun<sup>1 2 3</sup>, Chunzai Wang<sup>1 2 3 \*</sup>, Guoyu Ren<sup>4 5</sup>

15

16 <sup>1</sup> *State Key Laboratory of Tropical Oceanography, South China Sea Institute of Oceanology,*  
17 *Chinese Academy of Sciences, Guangzhou, China 510301*

18 <sup>2</sup> *Southern Marine Science and Engineering Guangdong Laboratory (Guangzhou),*  
19 *Guangzhou, China 511458*

20 <sup>3</sup> *Innovation Academy of South China Sea Ecology and Environmental Engineering, Chinese*  
21 *Academy of Sciences, Guangzhou, China 510000*

22 <sup>4</sup> *Department of Atmospheric Science, School of Environmental Studies, China University of*  
23 *Geosciences, Wuhan, China 430074*

24 <sup>5</sup> *Laboratory for Climate Studies, National Climate Center, China Meteorological*  
25 *Administration, Beijing, China 100081*

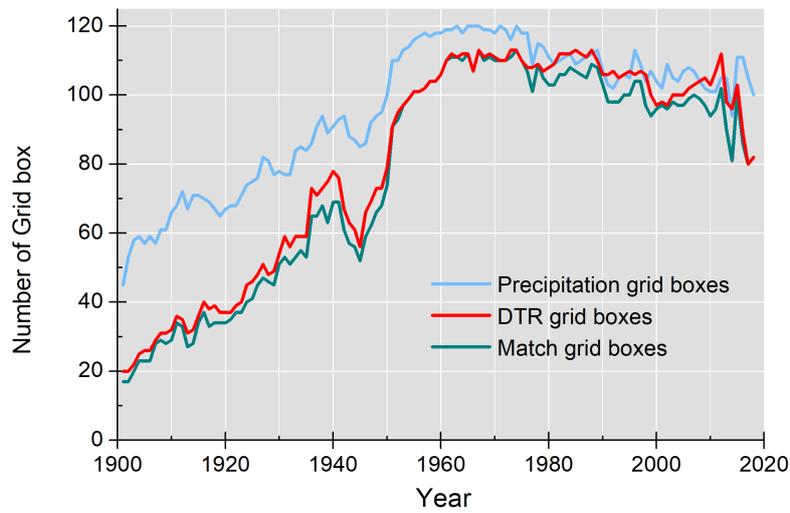
26

27 \*Corresponding author: Dr. Chunzai Wang (E-mail: cwang@scsio.ac.cn)

28

29 **1. Figure S1**

30



31

32 **Figure S1.** Long-term changes in the number of precipitation grid boxes (blue curve), DTR  
33 grid boxes (red curve), and match grid boxes (green curve).

34

35

36 **2. China Meteorological Administration - global land surface air temperature dataset**  
37 **(CMA-LSAT v1.1)**

38 The original data sources of CMA-LSAT data include 3 global datasets (Global Historical  
39 Climatology Network-V3 dataset (GHCN-V3), Climatic Research Unit Temperature 4.0  
40 dataset (CRUTEM4.0), and Berkeley Earth Surface Temperature dataset (BEST)); 2 regional  
41 datasets (European Climate Assessment & Dataset (ECA&D), and Historical Instrumental  
42 Climatological Surface time series of the greater Alpine region database (HISTALP)); and 8  
43 national datasets (China, USA, Russia, Canada, Australia, Korea, Japan, Vietnam) (Sun et al.  
44 2017; Xu et al. 2014). Compared with the other existing global land surface air temperature  
45 datasets, CMA-LSAT dataset shows similar ability in describing global temperature changes,  
46 despite there are some differences when describing regional temperature changes (Xu et al.  
47 2017). Compared with other global datasets, the obvious improvement of CMA-LSAT data  
48 spatial coverage is mainly in South America, Africa, East Asia especially in mainland China  
49 and its neighboring countries and regions (Sun et al. 2017).

50 The CMA-LSAT dataset has been processed by quality-control and homogenization. The  
51 method of quality-control was based on that used for GHCN dataset (Menne et al. 2010;  
52 Lawrimore et al. 2011), which mainly includes climatological outlier test, spatial  
53 inconsistency test, and internal inconsistency test (Sun et al. 2017; Xu et al. 2018). Data from  
54 non-homogenized stations in the CMA-LSAT data were also detected for temporal  
55 inhomogeneities by applying the RHtest-V3 system developed by Environment Canada  
56 (Wang et al. 2007; Wang and Fang 2011). The RHtest-V3 model consists of two parts: the  
57 penalized maximal t-test (PMT), and the penalized maximal F test (PMF, Wang and Fang  
58 2011). Firstly, the station with breakpoint is detected as the target station, and then the  
59 reference stations of the breakpoint station are established according to the method invented  
60 by Peterson et al (1998). By comparing the target series with reference time series of  
61 neighboring stations, the background climate signal could be removed from the target station  
62 series. Finally, the PMF algorithm in RHtest-V3 system is used to correct the breakpoints in  
63 the temperature series of the target station.

64 The dataset contains 11,036 stations, including maximum, minimum and mean temperature  
65 data. There are 2,051 stations in the East Asian region ( $80^{\circ}$ – $150^{\circ}$ E;  $0^{\circ}$ – $60^{\circ}$ N).

### 66 **3. Classification method of urban and rural stations**

67 The station classification model used in this paper is a machine learning method developed  
68 by Zhang (2020) and Zhang et al (2020). The United States National Climate Data Center  
69 (NCDC) has built the United States climate reference network (USCRN) since 1997, and it  
70 can be considered that the observation environment around the USCRN will not be affected  
71 by urbanization in the next 50 years (Diamond et al. 2013). The USCRN network is  
72 recognized to be able to accurately characterize regional climate change in the United States  
73 (Diamond et al. 2013; Ren and Chu 2019). Therefore, we use USCRN data to train machine  
74 learning algorithm model. Our goal is to select a global rural station network similar to the  
75 USCRN from all global stations. Based on this global rural station network, we can also select  
76 the rural station network in East Asia.

77 The “isolated forest” algorithm in machine learning is used to select rural stations, which is  
78 often used to identify outliers or abnormal data from all data (Liu et al. 2012). The first step of

79 the calculation method is to use USCRN fit model, and then to excavate the land use data  
80 around global stations through the model, and finally determine the dense parts (rural stations)  
81 and outliers (urban stations) in global stations network (Zhang 2020).

82 The isolated forest algorithm does not use any method based on distance or density  
83 measurement. Instead, it divides the data randomly by recursion, and stops after calculating  
84 that all the data are isolated. In this random segmentation strategy, outliers often have short  
85 cutting paths (isolated by fewer cuts), while the points with higher density are cut multiple  
86 times isolated. It can be divided into three stages:

87 (1) The first step is to construct the classification model of urban and rural stations. The  
88 data excavated by the station classification model is the percentage of land use around the  
89 station in 2018 (data available from [www.esa-cci.org](http://www.esa-cci.org)). The land use data around the station  
90 can reflect the intensity of the station affected by urbanization, and the classification of station  
91 types based on the land use data has been applied in some studies (Hu et al. 2009). The reason  
92 for using only the land use data around the stations in 2018 is to ensure that the selected rural  
93 stations can represent the stations that have not been affected by urbanization so far. If the  
94 impact of station migration is ignored, it can be considered that the selected rural stations are  
95 similar to the USCRN. The disadvantage of the classification method is that the selected rural  
96 stations may include stations that have moved from urban areas to rural areas. However, the  
97 data used in our study have been processed by homogenization process, which can be  
98 considered that the impact of urbanization in relocated stations is removed through the  
99 homogenization correction process.

100 Determining the station buffer range is another key point in the model construction. Taking  
101 China as an example, the observation specifications of China Meteorological Administration  
102 (2003) have strict restrictions on the building height and surrounding observation  
103 environment within 2 km of the national basic station. However, these observational  
104 specifications only consider the direct impact of the station surrounding environment on the  
105 detection, and do not fully consider the impact of urbanization effect. A large number of  
106 studies have found that the influence area of urban heat island effect is not only limited to the  
107 local stations, but also can reach the outskirts of city, and it increases with the size of city  
108 (Zhou and Shu 1994; Souch and Grimmond 2006). Therefore, it is difficult to determine the

109 urbanization impact range for a large number of stations, and there has never been a  
110 consensus on the impact range of urbanization. The station buffer range we used in this study  
111 is 12 km around the station in the current definition method, mainly referring to the station  
112 buffer range used in the previous study of Tysa et al. (2019).

113 Since the station buffer area is determined as 12km, the land use data of 1-12 km in the  
114 station buffer zone can be defined as 12 dimensions in the model. Among these 12 dimensions,  
115 a cutting point (or a hyperplane) can be generated randomly in a random dimension, and the  
116 data space can be cut into two subspaces by cutting points. The data larger than the cutting  
117 point and the data smaller than the cutting point are divided into two parts:

$$118 \quad \min(x_{i,j} = q, x_{ij} \in X') < p < \max(x_{i,j} = q, x_{ij} \in X')$$

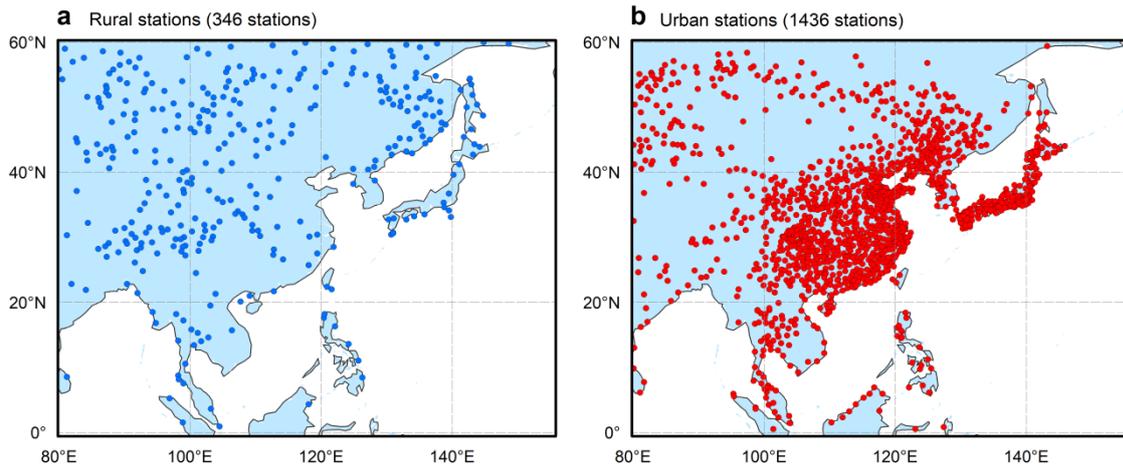
119 where  $p$  represents a stochastic cutting point,  $q$  represents 12 dimensions,  $X'$  represents all  
120 data in the classification model,  $x_{i,j}$  represents a station data in the model. All the data points  
121 can be divided into two categories: one is the dense part points which can be isolated after  
122 many times of cutting (rural stations), and the other is the outliers which are isolated after few  
123 cuts (urban stations).

124 USCRN stations are used as training data to fit the model. Then, the algorithm iterates until  
125 each station on the global land belongs to an independent space. In other words, according to  
126 the calculation results of the model, the global land stations can be divided into dense part  
127 (rural network similar to the USCRN) and outliers (urban stations).

128 (2) Verify the reliability of the model and training data, and determine the pollution  
129 parameters in the model. The pollution parameters in the model represent the degree of station  
130 affected by urbanization, which is a key parameter in the model. 70% of the station data are  
131 randomly selected from USCRN station network to train the model, and the remaining 30%  
132 data are used as the verification dataset to verify the accuracy of the model results. The  
133 pollution parameters were set to 0-0.5 (interval 0.05). The results show that the abnormal ratio  
134 of the validation dataset and the training dataset is very close. This also shows that USCRN  
135 data training model is reliable.

136 (3) All the USCRN data are used to fit the model. Then the 1-12 km land use data of global  
137 stations are put into the model for calculation. Finally, the global stations are divided into two

138 categories, dense stations (rural stations) and outlier stations (urban stations). When the  
139 pollution parameter is set at 0.3, the selected rural stations are similar to those of USCRN  
140 stations. With the increase of pollution parameter, the number of rural stations gradually  
141 decreases. Based on the above method, East Asian stations can be divided into 1436 urban  
142 stations and 346 rural stations (Figure S2a-b). The urbanization land in the buffer range of  
143 1-12 km around the East Asian rural stations is less than 5%.



144  
145 **Figure S2.** Spatial locations of urban (a) and rural (b) stations. The classification method of  
146 urban and rural stations is the isolated forest algorithm in machine learning.

147

## 148 **References**

- 149 China Meteorological Administration (2003) Criterion of surface meteorological observation.  
150 China Meteorological Press, Beijing. (in Chinese)
- 151 Diamond HJ, Karl TR, Palecki MA, et al (2013) U.S. Climate Reference Network after one  
152 decade of operations: status and assessment. *Bulletin of the American Meteorological*  
153 *Society* 94(4): 485–498
- 154 Hu Y, Jia G, Guo H, (2009) Linking primary production, climate and land use along an  
155 urban–wildland transect: a satellite view. *Environmental Research Letters* 4(4):044009
- 156 Makowski K, Jaeger EB, Chiacchio M, et al (2009) On the relationship between diurnal  
157 temperature range and surface solar radiation in Europe. *Journal of Geophysical*  
158 *Research* 114: D00D07. doi:10.1029/2008JD011104
- 159 Lawrimore JH, Menne MJ, Gleason BE, et al (2011) An overview of the global historical

160 climatology network monthly mean temperature data set, version 3. *Journal of*  
161 *Geophysical Research: Atmospheres* 116:5454-5466

162 Liu FT, Ting KM, Zhou ZH (2012) Isolation-Based Anomaly Detection. *ACM Transactions*  
163 *on Knowledge Discovery from Data*, 2012, 6(1): 1–39.

164 Menne MJ, Williams CN, Vose RS, (2010) The U.S. Historical Climatology Network monthly  
165 temperature data, version 2. *Bulletin of the American Meteorological Society*  
166 90:993-1007

167 Peterson TC, Easterling DR, Karl TR, et al (1998) Homogeneity adjustments of in situ  
168 atmospheric climate data: a review. *International Journal of Climatology*  
169 18:1493-1517

170 Ren GY, Chu ZY (2019) A brief introduction to the U.S. Climate Reference Network.  
171 *Advances in Meteorological Science & Technology* 9(04): 56–61 (in Chinese)

172 Souch C, Grimmond S (2006) Applied climatology: urban climate. *Progress in Physical*  
173 *Geography* 30(2): 270-279

174 Sun XB, Ren GY, Xu WH, et al (2017) Global land-surface air temperature change based on  
175 the new CMA GLSAT data set. *Science Bulletin* 62(4):236-238

176 Sun XB, Ren GY, You QL, et al (2019) Global diurnal temperature range (DTR) changes  
177 since 1901. *Climate Dynamics* 52(5):3343-3356

178 Tysa SK, Ren GY, Qin Y, et al (2019) Urbanization effect in regional temperature series  
179 based on a remote sensing classification scheme of stations. *Journal of Geophysical*  
180 *Research: Atmospheres* 124(20): 10646–10661

181 Wang XL, Feng Y (2011) RHtestsV3 user manual. Climate Research Division Atmospheric  
182 Science and Technology Directorate Science and Technology Branch, Environment  
183 Canada Toronto, Ontario, Canada

184 Wang XL, Wen QH, Wu Y, (2007) Penalized maximal t test for detecting undocumented mean  
185 change in climate data series. *Journal of Applied Meteorology & Climatology* 46:916-931

186 Xu WH, Li QX, Jones PD, et al (2017) A new integrated and homogenized global monthly  
187 land surface air temperature dataset for the period since 1900. *Climate Dynamics*  
188 (15):1-24. doi: 10.1007/s00382-017-3755-1

189 Xu WH, Li QX, Yang S, et al (2014) Overview of global monthly surface temperature data in

- 190 the past century and preliminary integration. *Advances in Climate Change Research*  
191 5:111-117
- 192 Zhang PF (2020) Observed trend changes in extreme temperature over the global land,  
193 1951-2018, China University of Geosciences (Wuhan). (in Chinese)
- 194 Zhang PF, Ren GY, Qin Y, et al. (2021) Urbanization effects on estimates of global trends in  
195 mean and extreme air temperature. *Journal of Climate*, 2021, 34(5):1923–1945.
- 196 Zhou SZ, Shu J (1994) *Urban Climatology*. China Meteorological Press, Beijing. (in Chinese)