Improving machine learning-based weather forecast 1 post-processing with clustering and transfer learning

Xiaomeng Huang¹, Yuwen Chen¹, Yi Li¹, Yue Chen¹, Chi Yan Tsui¹, Xing Huang¹, Mingqing Wang¹, and Jonathon S Wright¹

¹Tsinghua University

November 21, 2022

Abstract

Machine learning has been widely applied in numerical weather prediction, but the incorporation of new observational sites into models trained on stations with long historical records remains a challenge. Here we propose a post-processing framework consisting of three machine learning methods: station clustering with K-means, temperature prediction based on decision trees, and transfer learning for newly-built stations. We apply this framework to post-processing forecasts of surface air temperature at 301 weather stations in China. The results show significant reductions (as much as $39.4\%^{-2}20.0\%$) in the root-mean-square error of operational forecasts at lead times as long as 7 days. Moreover, the use of transfer learning to incorporate new stations improves forecasts at the new site by 36.4% after only one year of data collection. These results demonstrate the potential for clustering and transfer learning to boost existing applications of machine learning techniques in weather forecasting.

Improving machine learning-based weather forecast post-processing with clustering and transfer learning

Yuwen Chen¹, Xiaomeng Huang^{1,2,3}, Yi Li^{1,3}, Yue Chen¹, Chi Yan Tsui³, Xing Huang^{1,3}, Mingqing Wang^{1,3}, Jonathon S. Wright¹

5	¹ Ministry of Education Key Laboratory for Earth System Modeling, Department of Earth System
6	Science, Tsinghua University, Beijing 100084, China
7	² Laboratory for Regional Oceanography and Numerical Modeling, Qingdao National Laboratory for
8	Marine Science and Technology, Qingdao, 266237, China
9	³ National Supercomputing Center in Wuxi, Wuxi, 214011, China

¹⁰ Key Points:

3 4

11	•	A post-processing framework comprising clustering, decision tree, and transfer learn-
12		ing methods is employed to improve weather forecasts.
13	•	This framework reduces the root-mean-square error by 27.9% (0.81°C) compared
14		to operational ECMWF forecasts.
15	•	Transfer learning improves forecasts by 36.4% at new stations with only one year
16		of data available, reducing barriers to network expansion.

Corresponding author: Xiaomeng Huang, hxm@tsinghua.edu.cn

17 Abstract

Machine learning has been widely applied in numerical weather prediction, but the incorporation of new observational sites into models trained on stations with long historical records remains a challenge. Here we propose a post-processing framework consisting of three machine learning methods: station clustering with *K*-means, temperature

prediction based on decision trees, and transfer learning for newly-built stations. We ap-

²³ ply this framework to post-processing forecasts of surface air temperature at 301 weather

stations in China. The results show significant reductions (as much as $39.4\% \sim 20.0\%$)

²⁵ in the root-mean-square error of operational forecasts at lead times as long as 7 days.

Moreover, the use of transfer learning to incorporate new stations improves forecasts at $1 + 26 + 10^{-10}$

the new site by 36.4% after only one year of data collection. These results demonstrate

the potential for clustering and transfer learning to boost existing applications of ma-

²⁹ chine learning techniques in weather forecasting.

³⁰ Plain Language Summary

Statistical approaches have been used for decades to enhance and interpret numer-31 ical weather forecasts. Artificial intelligence models have greatly advanced this field but 32 the extension of these models to newly-built sites remains a challenge. To address this, 33 we design a framework that combines three machine learning methods: clustering to group 34 similar stations, decision trees to classify the forecasts, and transfer learning to adapt 35 the model to new stations. We apply this framework to real forecasts and evaluate it against 36 measurements from hundreds of weather stations in China. Station clustering and trans-37 fer learning both substantially improve predictions for recently-built sites, demonstrat-38 ing how these tools can supplement existing artificial intelligence techniques in weather 39 forecasting. 40

41 **1** Introduction

The skill of numerical weather prediction (NWP) has improved significantly in re-42 cent decades due to advances in numerical models, data assimilation, and observation 43 systems (Bauer et al., 2015). Nevertheless, the accuracy of NWP is still limited by im-44 perfect model physics, numerical schemes, and initial/boundary conditions (Bauer et al., 45 2015; Lynch, 2008). Following the pioneering work of Glahn and Lowry (1972), Model 46 Output Statistics (MOS) have been used operationally for over forty years. Raw model 47 forecasts are post-processed using statistical relationships between observations and NWP 48 results. However, the volume and variety of observational and model output data are in-49 creasingly overwhelming conventional implementations of these methods (e.g., Agapiou, 50 2017; Overpeck et al., 2011). 51

The emergence of machine learning (ML) techniques has provided new perspectives 52 in this field (e.g., Reichstein et al., 2019). The climate community has increasingly turned 53 to such techniques for applications such as improving subgrid-scale parameterizations 54 in numerical models (e.g., Gentine et al., 2018; Rasp et al., 2018; Schneider et al., 2017; 55 Jiang et al., 2018), improving forecasts at very short or very long lead times (e.g., Shi 56 et al., 2015; Ham et al., 2019; B. Pan et al., 2019), detecting extreme weather (Hwang 57 et al., 2019), and identifying complex teleconnection patterns (e.g., Runge et al., 2019; 58 Boers et al., 2019). ML techniques could also substantially improve the accuracy of NWP 59 results (McGovern et al., 2017; Rasp & Lerch, 2018; Scher, 2018). 60

The success of ML relies heavily on the quality and quantity of training data. Unfortunately, observations are usually sparse, especially for newly-built weather stations. Essential questions therefore arise regarding whether and by what means models trained on data-rich stations can be reliably extended to newly-built stations with limited data records.

Clustering techniques are widely used to extract information hidden in complex spatio-66 temporal data (Bador et al., 2015). Stations classified within the same cluster often share 67 similar meteorological features. This type of feature-based classification provides a nat-68 ural foundation for transfer learning, a technique by which knowledge gained in completing one task is repurposed for a different but related task (S. J. Pan & Yang, 2010). These 70 methods may permit models trained for data-rich stations to be rapidly fine-tuned for 71 application to data-poor stations. To take full advantage of these techniques, we pro-72 pose a new framework that combines three different ML methods: Clustering, Decision 73 trees, and Transfer learning, or CDT for short. We apply CDT to surface air temper-74 ature forecasts as an illustrative validation of this framework and its applicability. 75

76 **2** Data

NWP data are provided by The International Grand Global Ensemble (TIGGE) 77 project of the European Centre for Medium-Range Weather Forecasts (ECMWF) (e.g., 78 Bougeault et al., 2010; Swinbank et al., 2016). The numerical forecasts are initialized 79 twice per day at 00 and 12 UTC with lead times ranging from 6 to 168 hours at 6-hour 80 increments (for a total of 28 lead times). We use data for the period from 1 January 2013 81 to 31 December 2018. The sample size is therefore 4384 for each weather station and lead 82 time. Five variables are selected: temperature and dew point temperature at 2 m height, 83 surface pressure, and the zonal and meridional wind components at 10 m height. 84

Observations from weather stations in China are obtained from www.meteomanz.com 85 for the same period (1 January 2013 through 31 December 2018). As too few data are 86 available in Xizang and Qinghai, we omit these areas from the analysis. We select 301 87 weather stations with data covering at least half of the year 2018 (the testing period as 88 introduced below). Four variables (surface air temperature, surface pressure, surface air 89 relative humidity, and near-surface wind speed) are provided every 3 hours (00, 03, 06, 90 09, 12, 15, 18, and 21 UTC). Static information for each station is also used, including 91 latitude, longitude, and elevation. Missing values are filled via linear interpolation in the 92 time dimension. 93

The historical observations are processed to generate feature vectors with shapes 94 defined by $(n_{\text{samples}}, n_{\text{steps}}, n_{\text{features}})$, where n_{samples} is the number of records for a spec-95 ified station, n_{steps} is the number of time steps used for temporal pattern mining, and 96 n_{features} is equal to 4 (i.e., the number of measurements to match at each time step). For 97 example, the shape of the input vector for the Beijing station is (4384, 25, 4) when three 98 days of past observations are used. NWP data are interpolated to each station location 99 using an inversion-distance weighting iDW; (Myers, 1994) applied to forecast data from 100 the four nearest model grid cells. The observational and NWP data are combined for 101 input to the CDT framework. 102

103 3 Methods

The CDT framework consists of three individual ML modules: clustering, decisiontree, and transfer learning. The clustering module classifies the 301 stations into groups using the traditional *K*-means technique. Separate decision-tree-based post-processing modules are then developed for each cluster and each lead time. Each newly-built station is assigned to the best-fit existing cluster. The transfer learning module is then used to produce the final results.

110

3.1 Clustering Stations with K-means

The traditional K-means (Hastie et al., 2009) clustering technique is often used for climate data analysis (e.g., Bador et al., 2015; Bernard et al., 2013). Stations with similar features are categorized into K individual clusters by calculating the feature distance between them. The features used in this study are the annual averages and standard deviations of surface air temperature, surface air relative humidity, near-surface wind speed,
surface pressure, latitude, longitude, and elevation. Models are established and trained
for each cluster instead of for each station to reduce the computational cost and enlarge
the training sample for each model.

The clustering result is highly sensitive to the value of K. We use the Silhouette Coefficient (Rousseeuw, 1987) to identify the optimal value of K. This metric measures the consistency of samples within each cluster as the ratio between cluster tightness and cluster dissociation. A larger Silhouette Coefficient indicates an increase in the inter-cluster distance relative to the intra-cluster distance. The maximum coefficient thus marks the optimal clustering result according to this metric.

The average Silhouette Coefficient (ASC; Text S1 in the supporting information) 125 varies with the number of clusters K (Fig. 1a). We use the ASC to reduce the number 126 of candidate K values so that we do not need to train ML models for all possible val-127 ues of K. Although the ASC is useful for identifying potential optimal values of K, a 128 larger ASC does not necessarily translate to a better ML model result. We test clusters 129 based on K = 2, K = 4, and K = 8, which each produce climatologically coherent 130 station groups. The result for K = 2 divides stations into two main groups correspond-131 ing to northern and southern China (Fig. 1b), while that for K = 4 produces clusters 132 corresponding to the Northeast, North, and South regions along with some scattered sta-133 tions (Fig. 1c). The scattered stations in cluster 3 are grouped because they experience 134 much larger wind speeds than their geographic neighbors. The result for K = 8 fur-135 ther distinguishes some sub-regions with distinct climatological characteristics, such as 136 the northwestern region and Yunnan Province (Fig. 1d). 137

138

3.2 Temperature post-processing based on LightGBM

After clustering, we apply a decision-tree model (Quinlan, 1986) to characterize re-139 lationships between the NWP forecasts and observations, correct biases, and identify how 140 different features affect the prediction results. Decision trees are tree-like graph mod-141 els. Information is passed from the root (representing the raw data) and split into branches 142 at each level. The splitting rule is typically set by the variable that best discriminates 143 among the samples along each branch. Decision trees produce naturally explainable out-144 puts and can provide valuable insight into hidden relationships uncovered by the algo-145 rithm. This method has been successfully employed in a wide variety of weather appli-146 cations (McGovern et al., 2017). 147

Gradient Boosting Decision Tree (GBDT; e.g., Chen & Guestrin, 2016) is a pop-148 ular decision tree approach that involves an ensemble of sequentially-trained decision trees 149 and gains knowledge by fitting negative gradients. In this work we use LightGBM (Ke 150 et al., 2017), a highly efficient and scalable GBDT algorithm, to explore the relationships 151 between NWP forecasts and observations in each cluster. LightGBM has been applied 152 to sorting, classification, and regression tasks in a number of big-data studies (e.g., Cao 153 & Gui, 2019; Ju et al., 2019). Adopting a leaf-wise growth strategy with depth limita-154 tion and gradient-based one-side sampling, LightGBM seldom overfits on small train-155 ing datasets (Ke et al., 2017). More details on the LightGBM model and its implemen-156 tation in this study are provided in Text S2 and Fig. S2 of the supporting information. 157

158

3.3 Transfer Learning for Newly-built Stations

In practice, ML models may malfunction due to data deficiencies or over-fitting.
 Transfer learning helps to reduce the likelihood of these types of failures by transferring knowledge from a previously trained model. The transferred model is then fine-tuned using newly-added data. This approach has been widely applied, including for the pre-



Figure 1. The effect of the number of clusters (K) on the clustering results. (a) The average Silhouette Coefficient (ASC, Text S1 in SI) as a function of K. Local maxima occur at K = 2, K = 4, and K = 8. (b) The spatial distribution of clusters for K = 2. (c) Same as (b) but for K = 4. (d) Same as (b) but for K = 8.

diction of wind speed (e.g., Hu et al., 2016; Qureshi & Khan, 2019). The LightGBM model 163 for each cluster is taken as a pre-trained model, transferred and further trained on ob-164 servations from newly-built stations identified as belonging to that cluster. The cluster 165 to which each new station belongs is determined by static geolocation information along 166 with the estimated annual means and standard deviations of key meteorological features 167 (surface air temperature, pressure, wind speed, and relative humidity). The latter are 168 IDW-interpolated from gridded NWP forecasts to accommodate the limited observational 169 records at these stations. The refined LightGBM model is then applied to surface air tem-170 perature forecasts at the newly-built station. 171

172 **4 Results**

Data spanning the six-year period from 2013 to 2018 are divided into three parts. 173 Data from 2013 to 2017 are used for training (80% of the data) and validation (the re-174 maining 20%). All data for 2018 are used for testing. We construct a separate model to 175 post-process ECMWF forecasts at each lead time (28 in all; Sect. 2) in each cluster. The 176 benefits are most significant at short lead times, with error reductions as large as 39.4%177 $(1.02^{\circ}C)$ for 1-day forecasts (6~24 h lead times; Table 1). Improvements decrease steadily 178 to 20.0% (0.68°C) for 7-day forecasts (144~168 h lead times). The average RMSE across 179 all lead times is reduced by 0.81°C, corresponding to a 27.9% increase in accuracy. Clus-180

tering improves the effectiveness of the decision tree algorithm, with the greatest error reduction achieved when stations are grouped into four clusters. Compared to models without clustering (i.e., a single model trained on all stations), the RMSE is reduced by 0.54% when two clusters are used (K = 2), 0.62% when K = 4, and 0.41% when K =8. Since the K = 4 result produces the smallest RMSE, we adopt this model for all subsequent experiments. In addition to improving the overall forecast quality, clustering reduces the RMSE at 296 out of 301 individual stations (98.3%) when K = 4 (Fig. 2a).

Table 1 and Fig. 2 also show results for three alternative ML algorithms that are 188 also widely used in meteorological applications (e.g., Gensler et al., 2017; Akram & El, 189 2016; Qing & Niu, 2018; Cao & Gui, 2019): linear regression (LR), artificial neural net-190 work (ANN), and long short-term memory (LSTM) with a fully-connected network (FCN). 191 LR, ANN, and LSTM-FCN are used as control models to predict temperature using iden-192 tical inputs. Detailed descriptions of the ANN and LSTM-FCN models are given in Text S3 193 and Figs. S3–S4 in the supporting information. The overall RMSE is reduced by 0.49°C 194 (16.8%) under LR, 0.71°C (24.7%) under ANN, and 0.71°C (24.7%) under LSTM-FCN 195 in the K = 4 scenario, including RMSE reductions at 211 stations under LR (Fig. 2b), 196 270 stations under ANN (Fig. 2c), and 272 stations under LSTM-FCN (Fig. 2d). Light-197 GBM outperforms all three models, providing a further reduction of the RMSE for sur-198 face air temperature forecasts of 14.2% relative to LR, 3.8% relative to ANN, and 2.6%199 relative to LSTM-FCN, indicating that LightGBM is more effective for this application. 200 LightGBM also takes less time for training ($\sim 10 \text{ minutes}$) than ANN ($\sim 20 \text{ minutes}$) or 201 LSTM-FCN (~ 40 minutes). 202

Table 1. RMSE of surface air temperature based on five different models for seven different lead times (Unit: °C). See text for details and definitions.

FCMWF	LightCBM	IB	ANN	I STM FCN
LOWIWI	LightGDM	LIU	AININ	LSIM-PON
2.59	1.57	1.94	1.63	1.60
2.72	1.83	2.22	1.91	1.89
2.83	2.00	2.37	2.10	2.09
2.93	2.15	2.48	2.25	2.23
3.05	2.30	2.60	2.40	2.39
3.21	2.49	2.76	2.61	2.61
3.41	2.73	2.95	2.85	2.95
	ECMWF 2.59 2.72 2.83 2.93 3.05 3.21 3.41	ECMWFLightGBM2.59 1.57 2.72 1.83 2.83 2.00 2.93 2.15 3.05 2.30 3.21 2.49 3.41 2.73	ECMWFLightGBMLR2.59 1.57 1.942.72 1.83 2.222.83 2.00 2.372.93 2.15 2.483.05 2.30 2.603.21 2.49 2.763.41 2.73 2.95	ECMWFLightGBMLRANN2.59 1.57 1.941.632.72 1.83 2.221.912.83 2.00 2.372.102.93 2.15 2.482.253.05 2.30 2.602.403.21 2.49 2.762.613.41 2.73 2.952.85

Based on these findings, we conclude that LightGBM in combination with four clus-203 ters presents a substantial improvement over both the original operational forecasts and 204 other ML-learning post-processing products. We therefore apply transfer learning to fine-205 tune the LightGBM model for extension to data-poor stations. To replicate the oper-206 ational scenario, we randomly select 20% of the stations to serve as synthetic newly-built 207 stations, using the remaining 80% stations to produce pre-trained models for each of the 208 four clusters. We then fine-tune the pre-trained models using data covering between zero 209 and 24 months at 2-month increments. The use of zero months of data corresponds to 210 applying the pre-trained model directly without fine-tuning. We then evaluate the cor-211 rected forecasts for the 'new' stations using testing data from the year 2018. To validate 212 the transfer learning results, we select seven lead times ranging from 24 h to 168 h at 24-213 h increments. The pre-trained models outperform the original NWP by $0.56^{\circ}C$ (16.8%) 214 even without fine-tuning (Fig. 3). The RMSE reduction continues to improve as the data 215 span used for fine-tuning is extended, reaching 36.4% (1.23°C) when 12 months of data 216 are used. Further improvements are negligible, indicating that the fine-tuning benefits 217 plateau once the annual cycle is fully represented. 218



Figure 2. Model assessment for test data. (a) Spatial distribution of relative error reduction by the LightGBM model with four clusters. Blue colors indicate improvement; red colors indicate deterioration. (b) Same as (a) but for LR. (c) Same as (a) but for ANN. (d) Same as (a) but for LSTM-FCN.

LightGBM, as a GBDT variant, is a 'grey box' AI algorithm. Information gain, split 219 times, and coverage rate can be calculated for each feature and used to explain the re-220 sults (Gilpin et al., 2019). For example, the raw (NWP) surface air temperature fore-221 cast contributes the most information for most lead times and cluster members when K =222 4 (Fig. 4). Temperature observations are the second most influential feature, but make 223 only marginal contributions in most cases. For clusters where the RMSE of the opera-224 tional ECMWF forecasts is already relatively small, such as cluster 2, the NWP fore-225 casts account for a larger proportion of the overall influence. Conversely, observed tem-226 peratures play a larger role for clusters with larger RMSEs in the operational forecasts, 227 such as cluster 4. The importance of the operational forecasts also increases as lead time 228 increases, with concomitant reductions in the importance of the direct observations. 229



Figure 3. Results of transfer learning for the 60 sites randomly selected to serve as synthetic newly-built stations. The time span of training data used to fine-tune the model ranges from zero to 24 months, where zero months means the pre-trained model is used directly without fine-tuning. (a) RMSE values at seven different lead times using pre-trained models based on four clusters. (b) RMSE of the ECMWF forecasts and LightGBM post-processed results at seven different lead times. The LightGBM results reflect average RMSEs for training data time spans ranging from zero to 24 months.

230 5 Conclusion

ML algorithms show great potential for post-processing numerical weather fore-231 casts, but their application is often restricted by the amount of available observations. 232 In this paper we propose the CDT framework, based on clustering, decision tree, and trans-233 fer learning, and assess its performance in post-processing ECMWF forecasts of surface 234 air temperature at lead times ranging from 6 to 168 h for 301 weather stations in China. 235 The stations are first divided into two, four, and eight clusters, as these classifications 236 produce climatologically and geographically meaningful station groupings. The CDT frame-237 work reduces the average RMSE of temperature forecasts at the 301 stations by up to 238 $0.81^{\circ}C$ (27.9%). These benefits are seen for all clustering scenarios and at all lead times, 239 but the greatest improvements are for the 4-cluster scenario at 6-24 h lead times. Trans-240 fer learning aids the extension of models trained on data-rich stations to data-sparse sta-241 tions within the same cluster. The RMSE at new stations is reduced by 16.8% (0.56° C) 242 relative to the raw ECMWF forecasts even without fine-tuning, rising to 36.4% (1.23°C) 243 once one year of observations is available for fine-tuning the algorithm. These improve-244 ments illustrate the great potential of the CDT framework for operational model post-245 processing, since newly-built sites typically suffer from short data records that restrict 246 the application of AI techniques. 247

An attractive feature of decision tree-based models is that the results can be explained in terms of the contributions from each input feature. Here the main contribution is from the raw ECMWF forecast, especially at longer lead times. However, the station temperature observations are most important contributor for short lead times at stations in cluster 4, where the operational forecasts are less accurate than in other clusters. Overall, the CDT framework can help to correct prediction biases between NWP and observations, especially for newly-built stations or sites with sparse data records.

255 Acknowledgments

This work is based on data provided by the TIGGE project and Meteomanz.com. TIGGE

²⁵⁷ (The Interactive Grand Global Ensemble) is an initiative of the World Weather Research

						(a)	Cluste	r 1					
	150 - 168 h	0. 72%	0.36%	4.87%	0. 73%	0.85%	0. 22%	89.15%	1.30%	1. 40%	0. 29%	0.11%	
me	126 - 144 h	0.69%	0.35%	4.86%	0.64%	0.70%	0. 20%	89.58%	1.21%	1.39%	0.29%	0.10%	- 0. 7
	102 <mark>-</mark> 120 h	0. 66%	0.34%	4.98%	0. 57%	0.60%	0. 18%	89.82%	1.10%	1. 36%	0. 28%	0. 09%	- 0. 6
Ч	78 - 96 h	0. 64%	0. 33%	5.26%	0. 53%	0.52%	0.17%	89.94%	0.95%	1. 31%	0. 27%	0. 09%	- 0. 4
Lea	54 - 72 h	0. 60%	0. 30%	5.76%	0. 48%	0.44%	0.17%	89.97%	0.71%	1. 24%	0.26%	0. 08%	- 0. 3
	30 - 48 h	0. 54%	0. 27%	6.74%	0. 45%	0.38%	0.16%	89.71%	0.30%	1. 12%	0.25%	0. 08%	-0.1
	6 - 24 h	0. 44%	0. 20%	8.32%	0. 42%	0.30%	0.16%	88.81%	0.15%	0. 90%	0.23%	0. 07%	
Lat_Lon ELE OBS_Temp OBS_RH OBS_Press OBS_WS ECD_Temp EC_Dew EC_Press EC_WS_V EC_WS_U Feature Name													
						(b)	Cluste	r 2					_
	150-168 h	0.68%	0. 38%	5.71%	0. 50%	0.39%	0.16%	90. 53%	0.17%	1. 25%	0.11%	0. 12%	- 0. 8
	126-144 h	0.68%	0. 38%	5.65%	0. 45%	0.34%	0. 14%	90. 70%	0.16%	1.26%	0.11%	0. 12%	
ime	102 - 120 h	0.67%	0. 38%	5.67%	0. 41%	0.29%	0. 13%	90. 80%	0.15%	1.26%	0.11%	0. 12%	- 0. 6
T De	78 - 96 h	0.66%	0. 38%	5.79%	0.39%	0.27%	0. 12%	90. 78%	0.14%	1. 25%	0.11%	0. 12%	- 0, 4
Le	54 - 72 h	0.64%	0.37%	6.04%	0.36%	0.24%	0.11%	90.68%	0.12%	1.21%	0.11%	0.12%	
	30 - 48 h	0.60%	0.35%	6.58%	0. 33%	0.22%	0. 11%	90. 33%	0.10%	1.15%	0.10%	0. 12%	- 0. 2
	6 - 24 h	0.52%	0. 30%	8.06%	0. 30%	0.21%	0.11%	89.24%	0.07%	1.00%	0.10%	0.11%	
		Lat_Lon	ELE	UBS_Temp	082 <u>k</u> H	UBS_Press Fea	OBS_WS ature Na Cluste	r 3	EG_Dew	EU_Press	EC_WS_V	EC_WS_U	
	150 - 168 h	0. 55%	2.26%	11.95%	0.44%	0.34%	0.19%	80.95%	2.54%	0.18%	0.46%	0.13%	- 0. 7
	126 - 144 h	0. 52%	2.17%	12.90%	0. 43%	0.33%	0.19%	80. 28%	2.41%	0. 18%	0.46%	0.13%	-0.6
me	102 - 120 h	0. 48%	2.05%	14. 21%	0. 41%	0.31%	0. 18%	79.35%	2.25%	0.17%	0.46%	0.13%	0.0
Ξ	78 - 96 h	0. 43%	1.88%	16.11%	0. 42%	0.30%	0.19%	77.91%	2.02%	0.16%	0.45%	0.13%	- 0. 4
Lea	54 - 72 h	0. 37%	1.65%	19.16%	0. 40%	0. 28%	0. 18%	75.56%	1.68%	0.14%	0.44%	0.13%	- 0. 3
	30 - 48 h	0. 29%	1.31%	24. 27%	0.37%	0.25%	0.17%	71.59%	1.09%	0.12%	0.42%	0.12%	- 0. 1
	6 - 24 h	0.19%	0.83%	34. 50%	0. 31%	0. 20%	0.14%	63.02%	0.24%	0. 08%	0.40%	0. 09%	
		Lat_Lon	ELE	OBS_Temp	OBS_RH	OBS_Press Fea	OBS_WS ature Na	ECD_Temp ame	EC_Dew	EC_Press	EC_WS_V	EC_WS_U	
						(d)	Cluste	r 4					0 -
	150 - 168 h	2. 64%	0. 12%	20. 10%	0. 53%	0.46%	0.14%	75.30%	0.30%	0.15%	0.16%	0. 09%	- 0. /
	126 - 144 h	2. 57%	0. 12%	21. 75%	0. 48%	0.38%	0. 13%	73.88%	0.29%	0.15%	0.17%	0. 09%	- 0. 6
ime	102 - 120 h	2. 47%	0. 12%	24. 12%	0. 42%	0.30%	0. 11%	71.79%	0.27%	0.15%	0.17%	0. 08%	- 0 4
Ч	78 - 96 h	2. 30%	0.11%	27. 92%	0. 38%	0.26%	0. 10%	68.28%	0.26%	0.14%	0.17%	0. 08%	5.
Lea	54 - 72 h	2. 04%	0.11%	34. 08%	0.36%	0. 22%	0. 09%	62.49%	0.23%	0.14%	0.16%	0. 08%	- 0. 3
	30 - 48 h	1. 58%	0.10%	46. 30%	0.35%	0. 20%	0. 09%	50.81%	0.20%	0.13%	0.16%	0. 08%	- 0. 1
	6 - 24 h	0. 82%	0. 08%	69. 05%	0. 32%	0.17%	0. 09%	29.00%	0.14%	0.11%	0.15%	0. 07%	
		Lat_Lon	ELE	OBS_Temp	OBS_RH	OBS_Press	OBS_WS	ECD_Temp	EC_Dew	EC_Press	EC_WS_V	EC_WS_U	

Feature Name

Figure 4. The relative importance of features at different lead times and for different clusters. The "EC" prefix indicates variables from the original ECMWF forecasts, while the "OBS" prefix indicates direct observations. Temp stands for temperature; RH for relative humidity; Press for surface pressure; WS for wind speed; dew for dew point temperature; WS_U and WS_V for the zonal and meridional components of wind speed, respectively; Lat_Lon for the latitude and longitude of the station; and ELE for the elevation of the station. The cluster numbers correspond to the K = 4 clustering result (Fig. 1c).

Programme (WWRP). Meteomanz.com collects observations released by official weather

- stations. The authors are sponsored by grants from the State's Key Project of Research
- and Development Plan (2016YFB0201100, 2017YFC1502200, 2018YFB0505000, 2018YFB1502800),
- the National Natural Science Foundation of China (41776010), and the Pilot National
- Laboratory for Marine Science and Technology (Qingdao)(QNLM2016ORP0108).

263 **References**

- Agapiou, A. (2017). Remote sensing heritage in a petabyte-scale: satellite data and 264 heritage Earth Engine© applications. International Journal of Digital Earth, 265 10(1), 85-102. doi: 10.1080/17538947.2016.1250829266 Akram, M., & El, C. (2016). Sequence to Sequence Weather Forecasting with Long 267 Short-Term Memory Recurrent Neural Networks. International Journal of 268 Computer Applications, 143(11), 7–11. doi: 10.5120/ijca2016910497 269 Bador, M., Naveau, P., Gilleland, E., Castellà, M., & Arivelo, T. (2015).Spatial 270 clustering of summer temperature maxima from the CNRM-CM5 climate 271 model ensembles & E-OBS over Europe. Weather and Climate Extremes, 9, 272 17-24. doi: 10.1016/j.wace.2015.05.003 273 Bauer, P., Thorpe, A., & Brunet, G. (2015).The quiet revolution of numerical 274 weather prediction. Nature, 525(7567), 47–55. doi: 10.1038/nature14956 275 Bernard, E., Naveau, P., Vrac, M., & Mestre, O. (2013). Clustering of maxima: Spa-276 tial dependencies among heavy rainfall in france. Journal of Climate, 26(20), 277 7929–7937. doi: 10.1175/JCLI-D-12-00836.1 278 Boers, N., Goswami, B., Rheinwalt, A., Bookhagen, B., Hoskins, B., & Kurths, J. 279 (2019). Complex networks reveal global pattern of extreme-rainfall teleconnec-280 tions. Nature, 566(7744), 373–377. doi: 10.1038/s41586-018-0872-x 281 Bougeault, P., Toth, Z., Bishop, C., Brown, B., Burridge, D., De Chen, H., ... 282 Worley, S. (2010).The thorpex interactive grand global ensemble. Bul-283 letin of the American Meteorological Society, 91(8), 1059–1072. doi: 284 10.1175/2010BAMS2853.1 285 Cao, Y., & Gui, L. (2019). Multi-Step wind power forecasting model Using LSTM 286 networks, Similar Time Series and LightGBM. 2018 5th International Confer-287 ence on Systems and Informatics, ICSAI 2018(Icsai), 192–197. doi: 10.1109/ 288 ICSAI.2018.8599498 289 Chen, T., & Guestrin, C. (2016, 3). XGBoost: A scalable tree boosting system. Pro-290 ceedings of the ACM SIGKDD International Conference on Knowledge Discov-291 ery and Data Mining, 13-17-Augu, 785-794. doi: 10.1145/2939672.2939785 292 Gensler, A., Henze, J., Sick, B., & Raabe, N. (2017).Deep Learning for so-293 lar power forecasting - An approach using AutoEncoder and LSTM Neu-294 ral Networks. 2016 IEEE International Conference on Systems, Man, and 295 Cybernetics, SMC 2016 - Conference Proceedings (April), 2858–2865. doi: 296 10.1109/SMC.2016.7844673 297 Gentine, P., Pritchard, M., Rasp, S., Reinaudi, G., & Yacalis, G. (2018). Could Ma-298 chine Learning Break the Convection Parameterization Deadlock? Geophysical 299 Research Letters, 45(11), 5742-5751. doi: 10.1029/2018GL078202 300 Gilpin, L. H., Bau, D., Yuan, B. Z., Bajwa, A., Specter, M., & Kagal, L. (2019).301 Explaining explanations: An overview of interpretability of machine learning. 302 Proceedings - 2018 IEEE 5th International Conference on Data Science and 303 Advanced Analytics, DSAA 2018, 80–89. doi: 10.1109/DSAA.2018.00018 304 Glahn, H. R., & Lowry, D. A. (1972). The Use of Model Output Statistics (MOS) in 305 Objective Weather Forecasting (Vol. 11) (No. 8). doi: 10.1175/1520-0450(1972) 306 011(1203:tuomos)2.0.co;2307 Ham, Y.-g., Kim, J.-h., & Luo, J.-j. (2019, 9). Deep learning for multi-year ENSO 308 forecasts. Nature. doi: 10.1038/s41586-019-1559-7 309
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). The Elements of Statistical Learn-

311 312	<i>ing</i> (Vol. 27) (No. 2). New York, NY: Springer New York. doi: 10.1007/978-0 -387-84858-7
313	Hu, Q., Zhang, R., & Zhou, Y. (2016). Transfer learning for short-term wind speed
314	prediction with deep neural networks. <i>Renewable Energy</i> , 85, 83–95. doi: 10
315	.1016/j.renene.2015.06.034
316	Hwang, J., Orenstein, P., Cohen, J., Pfeiffer, K., & Mackey, L. (2019). Improving
317	subseasonal forecasting in the western U.S. With machine learning. <i>Proceedings</i>
318	of the ACM SIGKDD International Conference on Knowledge Discovery and
319	Data Mining, 2325–2335. doi: 10.1145/3292500.3330674
320	Jiang G Q Xu J & Wei J (2018) A Deep Learning Algorithm of Neural
321	Network for the Parameterization of Typhoon-Ocean Feedback in Typhoon
322	Forecast Models. Geophysical Research Letters, 45(8), 3706–3716. doi:
323	10.1002/2018GL077004
324	Ju. Y., Sun, G., Chen, Q., Zhang, M., Zhu, H., & Rehman, M. U. (2019). A model
325	combining convolutional neural network and lightsbm algorithm for ultra-
326	short-term wind power forecasting. <i>IEEE Access</i> , $7(c)$, 28309–28318. doi:
327	10.1109/ACCESS.2019.2901920
328	Ke. G., Meng. Q., Finley, T., Wang, T., Chen, W., Ma, W.,, Liu, T. Y. (2017).
329	LightGBM: A highly efficient gradient boosting decision tree. Advances in
330	Neural Information Processing Systems, 2017-Decem(Nips), 3147–3155.
331	Lynch, P. (2008). The origins of computer weather prediction and climate modeling.
332	Journal of Computational Physics, 227(7), 3431–3444, doi: 10.1016/j.jcp.2007
333	.02.034
334	McGovern, A., Elmore, K. L., Gagne, D. J., Haupt, S. E., Karstens, C. D.,
335	Lagerquist, R., Williams, J. K. (2017). Using artificial intelligence
336	to improve real-time decision-making for high-impact weather. Bul-
337	letin of the American Meteorological Society, 98(10), 2073–2090. doi:
338	10.1175/BAMS-D-16-0123.1
339	Myers, D. E. (1994). Spatial interpolation: an overview. <i>Geoderma</i> , 62(1-3), 17–28.
340	doi: 10.1016/0016-7061(94)90025-6
341	Overpeck, J. T., Meehl, G. A., Bony, S., & Easterling, D. R. (2011, 2). Climate Data
342	Challenges in the 21st Century. Science, 331(6018), 700–702. doi: 10.1126/
343	science.1197869
344	Pan, B., Hsu, K., AghaKouchak, A., & Sorooshian, S. (2019). Improving Precipi-
345	tation Estimation Using Convolutional Neural Network. Water Resources Re-
346	search, 55(3), 2301–2321. doi: 10.1029/2018WR024090
347	Pan, S. J., & Yang, Q. (2010, 10). A Survey on Transfer Learning. IEEE Trans-
348	actions on Knowledge and Data Engineering, 22(10), 1345-1359. doi: 10.1109/
349	TKDE.2009.191
350	Qing, X., & Niu, Y. (2018). Hourly day-ahead solar irradiance prediction using
351	weather forecasts by LSTM. Energy, 148, 461–468. doi: 10.1016/j.energy.2018
352	.01.177
353	Quinlan, J. R. (1986, 3). Induction of decision trees. Machine Learning, 1(1), 81-
354	106. doi: 10.1007/BF00116251
355	Qureshi, A. S., & Khan, A. (2019). Adaptive transfer learning in deep neural net-
356	works: Wind power prediction using knowledge transfer from region to region
357	and between different task domains. Computational Intelligence, 35(4), 1089-
358	1113. doi: 10.1111/coin.12236
359	Rasp, S., & Lerch, S. (2018, 11). Neural Networks for Postprocessing Ensemble
360	Weather Forecasts. Monthly Weather Review, $146(11)$, $3885-3900$. doi: 10
361	.1175/MWR-D-18-0187.1
362	Rasp, S., Pritchard, M. S., & Gentine, P. (2018, 9). Deep learning to represent sub-
363	grid processes in climate models. Proceedings of the National Academy of Sci-
364	ences, $115(39)$, 9684–9689. doi: 10.1073/pnas.1810286115
365	Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais,

366	N., & Prabhat. (2019, 2). Deep learning and process understanding
367	for data-driven Earth system science. <i>Nature</i> , 566(7743), 195–204. doi:
368	10.1038/s41586-019-0912-1
369	Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and vali-
370	dation of cluster analysis. Journal of Computational and Applied Mathematics,
371	20(C), 53-65. doi: $10.1016/0377-0427(87)90125-7$
372	Runge, J., Bathiany, S., Camps-Valls, G., Coumou, D., Deyle, E., Kretschmer,
373	M., Zscheischler, J. (2019). Inferring causation from time series
374	in Earth system sciences. Nature Communications, $10(1)$, 2553. doi:
375	10.1038/s41467-019-10105-3
376	Scher, S. (2018). Toward Data-Driven Weather and Climate Forecasting: Approxi-
377	mating a Simple General Circulation Model With Deep Learning. Geophysical
378	Research Letters, $45(22)$, 616–12. doi: 10.1029/2018GL080704
379	Schneider, T., Lan, S., Stuart, A., & Teixeira, J. (2017). Earth System Modeling 2.0:
380	A Blueprint for Models That Learn From Observations and Targeted High-
381	Resolution Simulations. Geophysical Research Letters, $44(24)$, 396–12. doi:
382	10.1002/2017 GL076101
383	Shi, X., Chen, Z., & Wang, H. (2015). Convolutional LSTM Network: A Machine
384	Learning Approach for Precipitation Nowcasting. Nips, 2–3. doi: 10.1007/978
385	$-3-319-21233-3_6$
386	Swinbank, R., Kyouda, M., Buchanan, P., Froude, L., Hamill, T. M., Hewson,
387	T. D., Yamaguchi, M. (2016). The TIGGE project and its achieve-
388	ments. Bulletin of the American Meteorological Society, $97(1)$, $49-67$. doi:

ments. Bulletin of the American Meteorological Society, 97(1), 49–67. 388 10.1175/BAMS-D-13-00191.1 389

Supporting Information for "Improving machine learning-based weather forecast post-processing with clustering and transfer learning"

Yuwen Chen¹, Xiaomeng Huang^{1,2,3}, Yi Li^{1,3}, Yue Chen¹, Chi Yan Tsui³,

Xing Huang^{1,3}, Mingqing Wang^{1,3}, Jonathon Wright¹

 1 Ministry of Education Key Laboratory for Earth System Modeling, Department of Earth System Science, Tsinghua University,

Beijing 100084, China

 $^{2} {\rm Laboratory \ for \ Regional \ Oceanography \ and \ Numerical \ Modeling, \ Qingdao \ National \ Laboratory \ for \ Marine \ Science \ and$

Technology, Qingdao, 266237, China

³National Supercomputing Center in Wuxi, Wuxi, 214011, China

Contents of this file

- 1. Text S1 to S3 $\,$
- 2. Figures S1 to S4

Introduction

Supporting information for the manuscript "Improving machine learning based weather forecast post-processing with clustering and transfer learning" is provided here. Text S1 introduces the calculation of Average Silhouette Coefficient (ASC), Text S2 introduces

details of the LightGBM as applied in this work, and Text S3 introduces the alternative machine learning frameworks used for comparison with the LightGBM results. Figure S1 shows a flow chart of our CDT framework, Figure S2 shows an example tree structure based on LightGBM, and Figure S3 and Figure S4 illustrate the structures of the LSTM-FCN and ANN models, respectively.

Text S1: Average Silhouette Coefficient

In this paper, we use the average Silhouette coefficient (ASC) as a guide to find viable candidates for the clustering number K. The ASC is calculated via the following steps: (1) For data point i in cluster C_i , define

$$a(i) = \frac{1}{|C_i| - 1} \sum_{j \in C_i, i \neq j} d(i, j)$$
(1)

as the average distance from point i to the other points in cluster C_i , where d(i, j) is the distance between point i and point j;

(2) Define

$$b(i) = \min_{k \neq i} \frac{1}{|C_k|} \sum_{j \in C_k} d(i, j)$$
(2)

as the smallest average distance of point i to all points in any other cluster;

(3) Define the Silhouette of point i as

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}};$$
(3)

(4) Define the ASC as the mean Silhouette of all points.

Text S2: LightGBM

We use the LightGBM model to post-process ECMWF temperature forecasts. Figure S2 provides an illustrative example of the LightGBM decision tree structure in the CDT July 8, 2020, 10:07am

framework. Split features include ECMWF-predicted temperature (EC) and observations of temperature (T), pressure(p), wind speed (WS) and relative humidity (RH) measured at previous time steps. The path selects different branches in sequence depending on the split conditions, with the leaf value (ovals) returned as the final result. For the illustrated decision tree (Fig. S2), the value of leaf 10 (-0.176) is calculated using a series of criteria from the leftmost node $EC \leq 6.445$ to the rightmost node $T(t-1) \leq -0.450$). LightGBM provides the added advantages of rapid training speed, low memory usage, high accuracy, and parallel learning capacity.

LightGBM only supports two-dimensional structured datasets. Therefore, the observations are converted from the three-dimensional shape $(n_{\text{samples}}, n_{\text{steps}}, n_{\text{features}})$ to the two-dimensional shape $(n_{\text{samples}}, n_{\text{steps}} \times n_{\text{features}})$. When combined with the ECWMF data, latitude, longitude and elevation of stations, the final training data is organized in the shape $(n_{\text{samples}}, n_{\text{steps}} \times n_{\text{features}} + 8)$. LightGBM predictions represent the results of boosting all trees. The RMSE between the predicted result and observations collected at the valid forecast time is used to evaluate the prediction. To control overfitting, we tune the maximum tree depth to eight in this paper. This hyper-parameter provides the maximum depth that each tree is allowed to have. A smaller value indicates a weaker predictor.

Text S3: LSTM-FCN and ANN

Long-Short-Term-Memory (LSTM) models are widely used in time series prediction. In this paper we exploit temporal auto-correlation in observational time series in constructing a two-layer LSTM. The input is formatted as a three-dimensional array of the shape

 $(n_{\text{samples}}, n_{\text{steps}}, n_{\text{features}})$, and the LSTM output is merged with ECMWF forecasts using a four-layer, fully connected network (FCN). This deep neural network is then used as a control model. Fig. S3 shows the structure of our LSTM-FCN model.

:

As the most basic neural network frameworks, Artificial Neural Networks (ANNs) are also widely used for time series prediction. In this paper, we use a four-layer ANN as the control model. The input data to the ANN are the same as those provided to the LightGBM. Fig. S4 shows the structure of our ANN model.



:

Figure S1. A flow chart of our CDT framework. Starting from the top, the NWP forecasts and observational data are downloaded and pre-processed. The training branch builds a 'cluster estimator' and groups the existing stations into clusters. Separate LightGBM models are then trained for each cluster and each lead time. Post-processed forecasts for existing stations are generated using the pre-trained model for the corresponding cluster. New stations are grouped into the best-fit existing cluster, after which the corresponding LightGBM model is fine-tuned to produce the final forecast.



Figure S2. An example LightGBM decision tree. The features comprise the ECMWF result and observations of temperature (T), relative humidity (RH), pressure (p), and wind speed (WS) from a specified number of preceding time steps. Split features marked EC refer to the ECMWF prediction; T(t-1) refers to the observed temperature one time step prior; p(t-24) refers to the observed pressure 24 time steps prior; WS(t-0) refers to the observed wind speed at the current time; and RH(t-21) refers to the observed relative humidity 21 time steps prior. A LightGBM model consists of multiple such decision trees, and a LightGBM prediction is the result of boosting the returned leaf values from all trees.





Figure S3. The structure of the LSTM-FCN model. The input consists of two parts, the time series of the observed data (left), and the inverse distance weighted (IDW) NWP result (right). The left input is in the shape of $(n_{\text{samples}}, n_{\text{steps}}, n_{\text{features}})$, where n_{steps} equals to 25, n_{features} equals to 4. The right input is in the shape of $(n_{\text{samples}}, 8)$, where the number 8 means 5 variables from ECMWF forecasts, and 3 location information of the stations (latitude, longitude and elevation).



Figure S4. The structure of the ANN model. The input shape is $(n_{\text{samples}}, n_{\text{steps}} \times n_{\text{features}} + 8)$. Four FCN(Dense) layer with batch normalization function and Relu activation function are used to build this ANN model.

input:

output:

Dense

(None, 32)

(None, 1)