A robust error correction method for numerical weather

2 prediction wind speed based on Bayesian optimization,

3 Variational Mode Decomposition, Principal Component

4 Analysis, and Random Forest: VMD-PCA-RF (version 5 1.0.0)

6 Shaohui Zhou¹, Chloe Yuchao Gao^{2*}, Zexia Duan¹, Xingya Xi³, and Yubin Li¹

7 ¹Collaborative Innovation Centre on Forecast and Evaluation of Meteorological Disasters, Key

8 Laboratory for Aerosol-Cloud-Precipitation of China Meteorological Administration, School of

9 Atmospheric Physics, Nanjing University of Information Science and Technology, Nanjing, 210044,

10 China.

²Department of Atmospheric and Oceanic Sciences and Institute of Atmospheric Sciences, Fudan

12 University, Shanghai, 200438, China.

13 ³School of Atmospheric Sciences, Sun Yat-sen University, and Southern Marine Science and

14 Engineering Guangdong Laboratory (Zhuhai), Zhuhai, 519082, China

15 Correspondence to: gyc@fudan.edu.cn

16 Abstract. Accurate wind speed prediction is crucial for the safe and efficient utilization of wind 17 resources. However, current single-value deterministic numerical weather prediction methods 18 employed by wind farms do not adequately meet the actual needs of power grid dispatching. In this 19 study, we propose a new hybrid forecasting method for correcting 10-meter wind speed predictions 20 made by the Weather Research and Forecasting (WRF) model. Our approach incorporates Variational 21 Mode Decomposition (VMD), Principal Component Analysis (PCA), and five artificial intelligence 22 algorithms: Deep Belief Network (DBN), Multilayer Perceptron (MLP), Random Forest (RF), eXtreme 23 Gradient Boosting (XGBoost), light Gradient Boosting Machine (lightGBM), and the Bayesian 24 Optimization Algorithm (BOA). We first construct WRF-predicted wind speeds using the Global 25 Forecast System Global Prediction System (GFS) model output based on prediction results. We then 26 perform two sets of experiments with different input factors and apply BOA optimization to tune the 27 four artificial intelligence models, ultimately building the final models. Furthermore, we compare the 28 aforementioned five optimal artificial intelligence models suitable for five provinces in southern China 29 in the wintertime: VMD-PCA-RF in December 2021 and VMD-PCA-lightGBM in January 2022. We 30 find that the VMD-PCA-RF evaluation indices exhibit relative stability over nearly a year: correlation coefficient (R) is above 0.6, Forecasting Accuracy (FA) is above 85 %, mean absolute error (MAE) is 31

below 0.6 m s⁻¹, root mean square error (RMSE) is below 0.8 m s⁻¹, relative mean absolute error (rMAE)
is below 60 %, and relative root mean square error (rRMSE) is below 75 %. Thus, for its promising
performance and excellent year-round robustness, we recommend adopting the proposed
VMD-PCA-RF method for improved wind speed prediction in models.

36 1 Introduction

37 Sustainable energy plays a vital role in reducing carbon footprint and increasing system reliability 38 (Hanifi et al., 2020). As renewable energy sources have a negligible carbon footprint, they have 39 become the preferred choice for many industries in the power sector (Dhiman and Deb, 2020). Among 40 these sources, wind energy is a crucial low-carbon energy technology with the potential to become a 41 sustainable energy source (Tascikaraoglu and Uzunoglu, 2014). In 2022, the global wind power 42 capacity reached 906 GW, with a 9 % year-on-year increase due to a newly installed capacity of 77.6 43 GW. The global onshore wind market increased by 68.8 GW while facing a 5 % growth decline 44 compared to the previous year. Such change is attributed to a slowdown in China and the U.S., the 45 world's two largest wind markets that account for over two-thirds of the world's onshore wind farm 46 installations (Joyce and Feng, 2023). The instability and unpredictability of wind power generation can 47 lead to instability in the power system. In addition, the decline of the wind energy market also makes it 48 more challenging to improve the accuracy of wind speed forecasts. An accurate wind speed prediction 49 method is needed to reduce the instability risk of power system and the economic loss of wind power 50 enterprises (Huang et al., 2019). Therefore, accurate and stable wind speed prediction (WSP) is very 51 important for the safe and stable operation of the power grid system and for improving the utilization 52 rate of wind energy and economic development (Guo et al., 2021; Xiong et al., 2022; Tang et al., 53 2021).

54 Current WSP algorithms are primarily categorized into physical algorithms (Zhao et al., 2016), 55 statistical algorithms (Wang and Hu, 2015; Barthelmie et al., 1993), machine learning (ML) algorithms 56 (Huang et al., 2019; Salcedo-Sanz et al., 2011; Ma et al., 2020), and hybrid algorithms (Deng et al., 57 2020; Xu et al., 2021; Zhao et al., 2019; Xiong et al., 2022; Tang et al., 2021). Physical methods, such 58 as numerical weather prediction (NWP), are commonly used in wind speed forecasting. NWP, which 59 accounts for atmospheric processes and physical laws, solves discrete mass, momentum, and energy 60 conservation equations along with other fundamental physical principles, establishing itself as a widely 61 adopted and reliable physical method. Currently, the High-resolution Limited Area Model (HIRLAM) 62 (Služenikina and Männik, 2016), the European Center for Medium-Range Weather Forecast (ECMWF) 63 model, the fifth-generation mesoscale model (MM5) (Salcedo-Sanz et al., 2009), and the Weather 64 Research and Forecasting Model (WRF) (Skamarock et al., 2021) are extensively utilized for wind 65 speed prediction. However, NWP modeling faces challenges due to the selection of parameterization 66 schemes, such as model microphysics and systematic errors, which exhibit temporal and spatial 67 differences and uncertainties. These uncertainties hinder the accuracy of NWP models in wind speed 68 prediction, making it difficult to meet the rising demands of the grid system (Zhao et al., 2019; Xu et 69 al., 2021).

70 Studies have demonstrated that enhancing the accuracy of numerical weather prediction (NWP) 71 models and correcting prediction errors can effectively minimize the errors associated with wind speed 72 prediction. These research endeavors have typically sought to optimize the physical and dynamic 73 parameters of the NWP model (Cheng et al., 2013), refine the model structure (Jiménez and Dudhia, 74 2012), or improve the accuracy of model inputs through preprocessing and denoising techniques (Xu et 75 al., 2015). Additionally, improving initial field error through methods, such as target observation and 76 data assimilation (Williams et al., 2013), can also minimize wind speed errors predicted by NWP 77 models.

78 Physical methods are generally more appropriate for long-term wind speed prediction, such as 79 those 48-72 hours in advance, while their practical application in short-term forecasting is limited 80 (Zhao et al., 2019; Deng et al., 2020; James et al., 2018). In contrast, statistical methods utilize 81 historical data to establish a relationship between input and output variables and are therefore 82 well-suited for short-term wind speed prediction. They are usually time series models, such as 83 Autoregressive Moving Average (ARMA) (Erdem and Shi, 2011) and Autoregressive Integrated 84 Moving Average (ARIMA) (Wang and Hu, 2015). Whereas filtering models (Cassola and Burlando, 85 2012; Chen and Yu, 2014), machine learning models (Hu et al., 2013), and hybrid models (Huang et al., 86 2019) have been gradually developed to further improve wind speed prediction accuracy.

With purely statistical models becoming less suitable for wind speed predictions beyond 6 hours,
the use of a combination of physical and statistical methods has gained growing interest (Zjavka, 2015;

89 Xu et al., 2021). The error correction model improves the accuracy of the NWP model by training on 90 the relationship between the NWP predictor variables and the observed correlation variables (Sun et al., 91 2019). However, traditional error prediction models rely solely on historical wind speed sequences as 92 input factors (Deng et al., 2020; Guo et al., 2021) and do not incorporate the characteristic 93 meteorological factors forecasted by the WRF model. Studies have shown that considering all relevant 94 historical meteorological factors can lead to more accurate predictions compared to only taking into 95 account historical wind speed (Zhang et al., 2019c). Therefore, it is crucial to include meteorological 96 characteristic factors as input in the prediction model.

97 For an error prediction model, wind speed is the most important input factor. Traditionally, the 98 error prediction model uses historical wind speed data as input, without any feature selection. Feature 99 selection methods, such as filtering methods, are commonly used in time series analysis. Currently, 100 empirical mode decomposition (EMD) (Liu et al., 2018; Guo et al., 2012), ensemble empirical mode 101 decomposition (EEMD) (Wang et al., 2017), wavelet decomposition (WD) (Zhang et al., 2019b), 102 variational mode decomposition (VMD) (Hu et al., 2021; Zhang et al., 2019a), and other filtering 103 methods are used to select key features in the wind speed data. As mentioned above, studies have 104 shown that these feature selection methods can effectively extract the hidden features in the wind speed 105 series to improve wind speed prediction accuracy. However, despite the effectiveness of wind speed 106 filtering methods in wind speed prediction, only a few studies have applied these methods to the 107 correction of wind speed errors in NWP forecasting (Xu et al., 2021; Li et al., 2022).

In addition, traditional error correction methods generally adopt linear regression (Dong et al., 2013), multiple linear regression (Liu et al., 2016), machine learning (Salcedo-Sanz et al., 2011), and deep learning algorithms (Zhang et al., 2019c). However, the efficacy of machine learning and deep learning algorithms is highly dependent on the selection of model parameters (Guo et al., 2021; Xiong et al., 2022). The Bayesian optimization algorithm (Li and Shi, 2010; Guo et al., 2021) is considered a relatively advanced algorithm for optimizing model parameters and has been widely used in MATLAB and Python packages.

In this study, we investigate a multi-step wind speed forecasting model that combines NWP simulation and an error correction strategy. We present two sets of experiments divided into three steps: (1) we use the first group of experiments to extract hidden features from various meteorological

118 elements forecasted by NWP; The second group of experiments mainly focuses on the wind speed 119 forecast of NWP, and the VMD-PCA algorithm is used to extract the hidden features in the forecasted 120 wind speed; each set of experimental input factors is matched with the actual 10-meter wind speed data 121 of 410 stations in time and space; (2) we employ four advanced machine learning algorithms optimized 122 by the BOA algorithm, and DBN deep learning algorithm to train the two groups of experiments and perform 5-fold cross-validation; and (3) we analyze six distinct wind speed error indicators to compare 123 124 and identify the most suitable wind speed error correction schemes for the five southern provinces 125 (Yunnan, Guizhou, Guangxi, Guangdong, Hainan) in winter and throughout most of the year. The 126 remainder of this paper is organized into sections discussing the effects of the BOA-VMD-PCA 127 approach, the interpretability of RF feature importance, and the stability analysis of the proposed 128 models.

129 2 Data and methods

130 2.1 Data

131 The observed data comes from the China Meteorological Administration land data assimilation 132 system (CLDAS-V2.0) real-time product data set. According to the description of the documents on the 133 official website (https://data.cma.cn/data/cdcdetail/dataCode/NAFP CLDAS2.0 RT.html), the dataset 134 is constructed through the integration of multiple sources, including ground and satellite data, and is 135 refined using advanced techniques such as multi-grid variational assimilation, physical inversion, and 136 terrain correction. This dataset exhibits superior quality in comparison to other products, offering 137 higher spatial and temporal resolutions. The target observation data includes 2-m air temperature, 2-m 138 specific humidity, 10-meter wind speed, surface pressure, and precipitation. These data are processed 139 by the China Meteorological Public Service Center to equivalent latitude and longitude grid scale, 140 covering a geographical range of 15-32.97°N and 94-120.97°E. The spatial resolution of the grid is $0.03^{\circ} \times 0.03^{\circ}$ (3km by 3km) and the temporal resolution is 1 hour. China Meteorological Public 141 142 Service Center applied the nearest neighbor interpolation for precipitation and bilinear interpolation for 143 the other four meteorological elements with downscaling from 3km to 410 sites. We select the 144 10-meter wind speed data of 410 sites, as illustrated in Fig. 1.





147 Figure 1÷._WRF model simulation area elevation diagram. (d02 represents the nested area of the second
148 layer of the WRF model, and the black triangles represent the meteorological sites).

149

150 **2.2 Methods**

151 2.2.1 WRF simulation

152 The WRF 4.2 model (Skamarock et al., 2021), developed by the National Center for 153 Atmospheric Research (NCAR), represents a new generation of mesoscale numerical models with 154 numerous applications in research forecasting. When forecasting meteorological elements, the WRF 155 model normally uses the GFS data developed by the National Centers for Environmental Prediction 156 (NCEP). We use the WRF model in combination with daily GFS data resolution of $0.25^{\circ} \times 0.25^{\circ}$. The 157 GFS data used by us is released at 06:00 UTC with forecasting every 3 hours for a total duration of 90 158 <u>hours</u>Using the WRF model in combination with daily GFS data resolution of $0.25^{\circ} \times 0.25^{\circ}$, the GFS-159 data updates at 06:00 UTC and generates forecasting every 3 hours for a total duration of 90 hours. We

160 selected the 24-h forecasting data from the WRF-resulted file after a spin-up time of 18 hours. The 161 GFS data as the initial field and lateral boundary conditions for the WRF model. Surface static data, 162 such as terrain, soil data, and vegetation coverage, are derived from the Moderate Resolution Imaging 163 Spectroradiometer (MODIS) satellite with a resolution of 15 seconds (approximately 500 meters). 164 Incorporating a two-layer grid nesting configuration, the forecast area is illustrated in Fig. 1. The WRF 165 configuration process is detailed in Table 1. Given that the time scale of the meteorological station data 166 in the study area is 1 hour, the forecast data time interval of the WRF model is also set to 1 hour. As a 167 widely used numerical weather forecast model, the WRF model is suitable for weather studies from a 168 few meters to several thousand kilometers. Therefore, this paper uses the WRF model to predict 169 10-meter wind speed as the input factor for the error correction model.

- 170
- 171

Table 1+. WRF configuration scheme

Model (Version)	WRF (V4.2)			
Domains	D1	D2			
Horizontal grid points	600*500	967*535			
$\Delta x (km)$	9	3			
Vertical layers	58				
Longwave radiation	RRTMG (Iacono et al., 2008)				
Shortwave radiation	RRTMG (Iacono et al., 2008)				
Land surface	Noah LSM (Chen et al., 1997)				
Surface layer	MYJ (Janjić, 1994)				
Microphysics	Thompson (Thomp	oson et al., 2008)			
Boundary layer	MYJ (Janj	ić, 1994)			
Cumulus	Tiedtke (Tiedtke, 1989	; Zhang et al., 2011)			

172

173 2.2.2 Variational mode decomposition

As a new filtering method, VMD is robust in feature selection. The VMD algorithm decomposes a time series signal into several intrinsic mode functions (Isham et al., 2018). The sum of the modes equals the original signal, and the sum of the bandwidths is the smallest. The analysis signal is calculated using the Hilbert transform to estimate the modal bandwidth. The optimization model is described as

179
$$\left\{\min_{\{u_k\},\{\omega_k\}}\left\{\sum_{k=1}^{K} \left\|\partial_t\left[\left(\delta(t) + \frac{j}{\pi t}\right)u_k(t)\right]e^{-j\omega_k t}\right\|_2^2\right\} s.t. \quad \sum_{k=1}^{K} u_k = v \quad (1.1)$$

180 where *K* is the total number of modes, u_k is the decomposed *K*-th mode, w_k is the corresponding 181 center frequency, and *v* is the time-series signal, representing the wind speed sequence predicted by the 182 WRF model in this study.

183 The above-constrained problem can be transformed into an unconstrained problem using the 184 Lagrangian function:

185
$$L(\lbrace u_k \rbrace, \lbrace \omega_k \rbrace, \lambda) = \omega_k^{n+1} = \frac{\int_0^\infty \omega |\hat{u}_k(\omega)|^2 d\omega}{\int_0^\infty |\hat{u}_k(\omega)|^2 d\omega} \sum_{k=1}^K ||\partial_t \left[\left(\delta(t) + \frac{j}{\pi t} \right) u_k(t) \right] d\omega$$

186

187
$$\times]e^{-j\omega_{k}t} \parallel_{2}^{2} + \parallel v(t) - \sum_{k=1}^{K} u_{k}(t) \parallel_{2}^{2} + \left\langle \lambda(t), v(t) - \sum_{k=1}^{K} u_{k}(t) \right\rangle$$
(1.2)

188 where α is the penalty parameter and $\lambda(t)$ is the Lagrange multiplier.

189 Then we update u_k , w_k , and λ using the alternating direction method of the multiplier:

190
$$\hat{u}_{k}^{n+1}(\omega) = \frac{\hat{v}(\omega) - \sum_{i \neq k} \hat{u}_{i}(\omega) + \frac{\hat{\lambda}(\omega)}{2}}{1 + 2\alpha \left(\omega - \omega_{k}\right)^{2}}$$
(1.3)

191
$$\omega_k^{n+1} = \frac{\int_0^\infty \omega |\hat{u}_k(\omega)|^2 d\omega}{\int_0^\infty |\hat{u}_k(\omega)|^2 d\omega}$$
(1.4)

192
$$\hat{\lambda}^{n+1}(\omega) = \hat{\lambda}^{n}(\omega) + \tau \left[\hat{v}(\omega) - \sum_{k=1}^{K} \hat{u}_{k}^{n+1}(\omega) \right]$$
(1.5)

193 where τ is the update parameter.

194 When the accuracy (left side of the following expression) meets the following condition, u_k , w_k 195 and λ would stop updating:

196
$$\sum_{k=1}^{K} \frac{\left\| \hat{u}_{k}^{n+1} - \hat{u}_{k}^{n} \right\|_{2}^{2}}{\left\| \hat{u}_{k}^{n} \right\|_{2}^{2}} < \varepsilon$$
(1.6)

197 where ε is the tolerance of the convergence criterion.

The VMD algorithm is implemented to decompose the wind speed signal predicted by the WRF model. When using multiple sub-signals instead of the original signal, more features of the wind speed can be obtained. Therefore, it is beneficial to improve the prediction accuracy when using the sub-signal as input to the error correction model (Xu et al., 2021; Li et al., 2022).

202 2.2.3 Principal Component Analysis

Subsequences obtained by VMD usually have several illusory components. Using PCA to extract the principal components of subsequences increases the number of features input to the model and reduces the dimension of the data decomposed by VMD. When principal components (PCs) are used as the input of the error prediction algorithm, the PCs fully reflect the characteristics of the subsequence and reduce the model complexity. The PCs y_k , k=1, 2, ..., K of the subsequence matrix U and the cumulative contribution rate η_n of the first *n* principal components are expressed as:

$$y_k = c'_k U \tag{1.7}$$

210
$$\eta_n = \frac{\sum_{k=1}^n \lambda_k}{\sum_{k=1}^K \lambda_k}$$
(1.8)

211 where c_k is the corresponding characteristic unit vector, with k=1, 2, ..., K; λ_k is the characteristic 212 root, with $\lambda_1 \ge \lambda_2 \ge ... \ge \lambda_K$.

213 2.2.4 Evaluation indicators

There are many commonly used predictive effect evaluation indicators. This article uses the following evaluation indicators: correlation coefficient (R), root mean square error (RMSE), mean absolute error (MAE), relative root mean square error (rRMSE), relative mean absolute error (rMAE), Forecasting Accuracy (FA). Six error indicators are used to evaluate the correction results of short-term wind speed forecasts of wind farms. The formula for calculating the error index is as follows:

219
$$R = \frac{\sum_{i}^{n} (y_{i} - \overline{y}) (\hat{y}_{i} - \overline{\hat{y}})}{\sqrt{\sum_{i=1}^{n} (y_{i} - \overline{y})^{2}} \sqrt{\sum_{i=1}^{n} (\hat{y}_{i} - \overline{\hat{y}})^{2}}}$$
(1.9)

220
$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2}$$
(1.10)

221
$$MAE = \frac{1}{n} \sum_{i=1}^{n} |\hat{y}_i - y_i|$$
(1.11)

222
$$rRMSE = \left[\sqrt{\frac{1}{n}\sum_{i=1}^{n} (\hat{y}_{i} - y_{i})^{2}} / \left(\frac{1}{n}\sum_{i=1}^{n} y_{i}\right)\right] \times 100\%$$
(1.12)

223
$$rMAE = \left(\frac{1}{n}\sum_{i=1}^{n} |\hat{y}_{i} - y_{i}| / \left(\frac{1}{n}\sum_{i=1}^{n} y_{i}\right)\right) \times 100\%$$
(1.13)

$$FA = N_r / N_f \tag{1.14}$$

Among them, *n* represents the number of samples, \hat{y}_i represents the *i*-th predicted value, y_i represents the *i*-th actual value; N_r represents the number of wind speed absolute errors not greater than 1 m s⁻¹, and N_f represents the number of research samples.

228

229 2.2.5 Proposed hybrid forecasting algorithms

230 This study used five machine learning algorithms to conduct ten experiments following two main 231 paths. The first path involves increasing the meteorological variables possibly related to wind speed in 232 the forecast field. The correlation between the WRF-predicted 10-m wind speed and the observed wind 233 speed is the highest. The purpose of the second experimental path is using the VMD-PCA algorithm to 234 dig out the hidden wind speed characteristics of the 10-meter forecast wind speed, reduce the input of 235 other meteorological factors such as WD_{10} and D_2 , and further prove that the VMD-PCA algorithm is 236 effective before correcting the WRF-predicted wind speed. The overarching goal is to achieve accurate 237 correction of the forecast field wind speed. The flowchart of the artificial intelligence models used to 238 correct the WRF predicted wind speed for the two main experimental paths is illustrated in Fig. 2 and 239 comprises the following three steps:

240 Step 1. Data fusion, cleaning, and standardization: As depicted in Fig. 2, this paper proposes two 241 distinct experimental paths, with the primary difference being the selection of input variables. In 242 Experiment 1, as shown in Fig. 2, 12 sets of data are selected from the WRF forecast field, including 243 altitude (HGT), 10-meter wind speed (WS₁₀), latitude (LAT), longitude (LON), surface pressure (PRS), 244 relative humidity (RH), 10-meter meridional wind (V10), 10-meter zonal wind (U10), 2-meter 245 temperature (T₂), 2-meter dew point temperature (D₂), 10-meter wind direction (WD₁₀), and hourly 246 precipitation (PRE). Experiment 2 derives 8 sets of data by reducing the selected WRF field forecast 247 data to include only altitude, 10-meter wind speed, latitude, longitude, surface pressure, relative 248 humidity, 2-meter temperature, and hourly precipitation. The focus is on unearthing hidden 249 characteristic information of forecast wind speed. In this experiment, the wind speed is decomposed 250 into 9 Intrinsic Mode Functions (IMFk, k=0, 1, 2, ..., 8) using VMD. Subsequently, a low-dimensional 251 wind speed vector is extracted from the 9 IMF components via PCA dimensionality reduction (pca0, 252 pcal, pcal, pcal, and all data are concatenated to construct the input factors for the model in Experiment 2. 253 The time points in the dataset where missing values are located are eliminated. Experiment 1 254 (Experiment 2) standardizes 12 sets of meteorological elements (8 sets of meteorological elements in 255 Fig. 3, 9 IMF components, and three PCA vectors in Fig. 4) and wind speed observation data, 256 respectively. Standardization addresses the issue of varying meteorological factor values during 257 training, which may result in different contributions. In this paper, the 24-hour forecast data correspond 258 to the observation data of the subsequent 24 hours. The dataset spans from 00:00 on December 1, 2021, 259 to 23:00 on February 28, 2022, totaling 2160 hours and encompassing 410 weather stations. 260 Consequently, the original dataset comprises 2160*410 samples, with each sample containing 12 261 meteorological features in Experiment 1 and 20 input features in Experiment 2. While similar past 262 studies for wind speed correction from NWP models usually use several years for training and at least 263 one year for testing whereas our periods are shorter, the size of our data set is sufficient. For example, 264 Sun et al. 2019 used a data set that contained 1827 days, from January 2012 to December 2016, using 265 143 grid points with a resolution of $0.5^{\circ}*0.5^{\circ}$ predicted by ECMWF, followed by 24 features for each 266 sample, with a training set size of 1827*143*24 for each prediction time. Meanwhile, the size of our 267 training set is about 2160*410*12. Therefore, even though it only took us a month to train, for this 268 project, we trained millions of data; Second, the training data we used here was obtained through daily 269 operational runs of numerical weather forecasting, so it would have taken several years to get an equal 270 amount of training data.

Step 2. BOA optimization of AI models and cross-validation: In this study, the dataset is partitioned into training, validation, and test sets in accordance with the time series. February 2022 serves as the training and validation sets, while December 2021 and January 2022 constitute the test set. The training and validation sets are divided based on five-fold cross-validation. Both experiments employ five machine learning algorithms (DBN, MLP, RF, XGBoost, and LightGBM) to construct distinct machine learning models. Concurrently, this paper utilizes the BOA algorithm to tune the parameters of all models, except for DBN, resulting in the optimal hyperparameters for each model.

Step 3. Model evaluation and error analysis: The trained machine learning models are applied to the test set to obtain the revised wind speed data, and ultimately, the accuracy of all models is assessed through the wind speed evaluation index. The ultimate goal here is to identify the best wind speed correction model suitable for the entire year. Accordingly, the generalization of all models is evaluated









Figure 3:-._Daily average hourly rainfall (a), surface pressure (b), 2-meter temperature (c), 2-meter relative
humidity (d), 10-meter wind speed (e), 2-meter dew point temperature (f), and 10-meter wind direction (g)
which are located at Guangdong Lechang Station from December 1, 2021, to February 28, 2022. (February
2022 represents the training and verification sets, and December 2021 to January 2022 represents the
testing set).



Figure 4:-._Three-dimensional view of 12 wind speed components after VMD and PCA processing of the
10-meter forecast wind speed at Lechang Station in Guangdong from December 1, 2021, to February 28,
2082.

301 3 Results

302 **3.1 Experiment 1 evaluation**

303 In Experiment 1, the BOA optimization algorithm was applied to five AI models to correct the 304 10-meter wind speed forecasted by WRF. There were 12 meteorological element features to establish 305 five different AI models (see Table 2-S1 for the hyper-parameters of the five AI models). The training, 306 validation, and testing results for 10-meter wind speed are shown in Figs. S1-5 in the supplementary 307 material. It is clear from Table 3-2 that all models, except the DBN model, can fit the training set data 308 well. The DBN model exhibits the weakest performance on both the training and validation sets. 309 Alternatively, the LightGBM and XGBoost models demonstrate superior prediction performance on 310 the training set compared to the validation set. The scatters of the training sets of these two models 311 accumulate on the 1:1 diagonal, indicating slight overfitting. As shown in Figs. S1-5d, f, considering 312 different evaluation indices, the revision effects of the five models in two months demonstrate that 313 RMSE in January 2022 is generally lower than in December 2021; FA in January 2022 is generally 314 higher than in December 2021; R in January 2022 is generally lower than in December 2021. Overall, 315 the prediction performance of the five models in January 2022 surpassed that in December 2021. 316 Furthermore, the LightGBM and RF models exhibited the best performance among the five models in 317 the two-month test sets, while the DBN model had the least effective correction effect.

As illustrated in Fig. 5a, b, WS_{10} showed the strongest positive correlation with WS_{obs} , with the highest R of 0.51, which was consistent with the highest variable importance value of 31 % (23 %) in experiment 1 (experiment 2). In addition to WS_{10} , experiment 1 (experiment 2) also had another three dominant variables namely, LAT, HGT, and LON, with importance values of 16 % (14 %), 15 % (15 %), and 15 % (13 %), respectively. Meanwhile, in experiment 2, IMF0 and pca0 generated by the VMD-PCA algorithm have a good importance value of 9 % and 4 %, and the R values of them with WS_{obs} are as high as 0.47 and 0.45.

Concerning the importance of RF characteristics (Fig. 5a, c), it is indisputable that the 10 m wind speed predicted by WRF plays a dominant role in correcting the actual wind speed. The ones following are latitude, longitude, and topographic height, which represent spatial geographic information, and the 328 actual wind speed is closely related to geographic information. Subsequently, relative humidity is of 329 lesser importance. The distribution of the humidity field typically correlates with the movement of the 330 atmosphere, which is also closely related to wind speed. Certain meteorological elements, such as

- 331 rainfall, 2 m dew-point temperature, and 2 m temperature, contribute less importance.
- 332

333

Table 2. The best hyper-parameters of the models-

Model	parameters
VMD-PCA-lightGBM	<pre>'max_depth': 28, 'min_child_samples': 30, 'n_estimators': 436,</pre>
	<u>'num_leaves': 287</u>
VMD-PCA-XGBoost	<pre>'gamma': 1, 'max_depth': 19, 'min_child_weight': 1, 'n_estimators':</pre>
	408
VMD-PCA-RF	<pre>'max_depth': 31, 'max_features': 14, 'min_samples_leaf': 28,</pre>
	<pre>'min_samples_split': 3, 'n_estimators': 371</pre>
VMD-PCA-DBN	<pre>'input_length': 20, 'output_length': 1, 'loss_function':-</pre>
	'MSE', 'optimizer' : 'Adam', 'hidden_units' : [400,-
	200], 'batch_size' :20000, 'epoch_pretrain' : 100, 'epoch_finetune' :-
	200
VMD-PCA-MLP	<pre></pre>
lightGBM	<pre>'max_depth' : 21, 'min_child_samples' : 19, 'n_estimators' : 312,</pre>
	<pre>'num_leaves': 297</pre>
XGBoost	<pre>'gamma': 0, 'max_depth': 21, 'min_child_weight': 9, 'n_estimators':-</pre>
	299
RF	<pre>'max_depth': 40, 'max_features': 12, 'min_samples_leaf': 23,</pre>
	<pre>'min_samples_split': 2, 'n_estimators': 440</pre>
DBN	<pre>'input_length': 12, 'output_length': 1, 'loss_function':</pre>
	'MSE', 'optimizer' : 'Adam', 'hidden_units' : [400, 200], 'batch_size' :
	20000, 'epoch_pretrain': 100, 'epoch_finetune': 200
MLP	<pre> 'batch_size': 10232, 'hidden_layer_sizes': 494, 'max_iter': 311 </pre>

334 **Table 3. Table of evaluation indices of wind speed error trained and verified by 10 models in February 2022**

	training set			validation set		
Model –	R	RMSE(m s ⁻¹)	FA	R	RMSE(ms ⁻¹)	FA
VMD-PCA-lightGBM	0.96	0.33	0.99	0.88	0.53	0.9 4
VMD-PCA-XGBoost	0.96	0.31	1.00	0.87	0.54	0.94
VMD-PCA-RF	0.89	0.52	0.94	0.86	0.57	0.93
VMD-PCA-DBN	0.74	0.75	0.87	0.74	0.75	0.87
VMD-PCA-MLP	0.84	0.60	0.91	0.81	0.66	0.90
lightGBM	0.93	0.41	0.98	0.88	0.54	0.94
XGBoost	0.96	0.31	0.99	0.87	0.56	0.93
RF	0.89	0.52	0.94	0.86	0.57	0.93
DBN	0.76	0.73	0.88	0.76	0.73	0.88
MLP	0.85	0.59	0.92	0.83	0.62	0.91

 335
 Table 2. Table of evaluation indices of wind speed error trained and verified by 10 models in February 2022

M- 1-1	training set			validation set		
<u>Model</u> –	<u>R</u>	<u>RMSE(m s⁻¹)</u>	<u>FA</u>	<u>R</u>	<u>RMSE(m s⁻¹)</u>	<u>FA</u>
VMD-PCA-lightGBM	<u>0.96</u>	<u>0.33</u>	<u>0.99</u>	<u>0.88</u>	<u>0.53</u>	<u>0.94</u>
VMD-PCA-XGBoost	<u>0.96</u>	<u>0.31</u>	<u>1.00</u>	<u>0.87</u>	<u>0.54</u>	<u>0.94</u>
VMD-PCA-RF	<u>0.89</u>	<u>0.52</u>	<u>0.94</u>	<u>0.86</u>	<u>0.57</u>	<u>0.93</u>
VMD-PCA-DBN	<u>0.74</u>	<u>0.75</u>	<u>0.87</u>	<u>0.74</u>	<u>0.75</u>	<u>0.87</u>
VMD-PCA-MLP	<u>0.84</u>	<u>0.60</u>	<u>0.91</u>	<u>0.81</u>	<u>0.66</u>	<u>0.90</u>
<u>lightGBM</u>	<u>0.93</u>	<u>0.41</u>	<u>0.98</u>	<u>0.88</u>	<u>0.54</u>	<u>0.94</u>
<u>XGBoost</u>	<u>0.96</u>	<u>0.31</u>	<u>0.99</u>	<u>0.87</u>	<u>0.56</u>	<u>0.93</u>
<u>RF</u>	<u>0.89</u>	<u>0.52</u>	<u>0.94</u>	<u>0.86</u>	<u>0.57</u>	<u>0.93</u>
DBN	<u>0.76</u>	<u>0.73</u>	<u>0.88</u>	<u>0.76</u>	<u>0.73</u>	<u>0.88</u>
MLP	0.85	0.59	0.92	0.83	0.62	0.91



Figure 5: Schematic diagram of correlation coefficients (represented the WS₁₀ and input variables) and
feature importance (calculated by the scikit-learn python package) for two sets of experiments. (a) and (c)
represent experiment 1, and (b) and (d) represent experiment 2.

337

342 **3.2 Experiment 2 evaluation**

Experiment 2 builds upon Experiment 1, concentrating on the predicted 10-meter wind speed by the WRF model. We use the VMD algorithm to decompose the predicted wind speed into 9 components and use the PCA algorithm to extract the main 3 principal components. In the RF feature importance analysis (Fig. 5b, d), it is evident that the VMD algorithm can decompose IMF0 and IMF1, with contributions surpassing those of 2-meter temperature and precipitation, respectively. The importance of the pca0 component, after PCA principal component extraction, reaches up to 8%. What 349 is particularly interesting is that in the correlation analysis, the correlation values between the IMF0 350 and pca0 components and the actual wind speed are 0.50 and 0.51, which are second only to the 351 forecasted wind speed.

352 From the indices (RMSE, FA, R) of the training set and validation set shown in Table 32, in comparison to the above five artificial intelligence methods, the training results of VMD-PCA-DBN 353 354 are relatively inferior. VMD-PCA-lightGBM and VMD-PCA-XGBoost models still train the processed 355 data effectively. According to the scatter density figure (Fig. 6a, Fig. 7a), the scatters are relatively 356 concentrated on the 1:1 line. From the indicators (RMSE, FA, R) of the testing set shown in Figs. 357 S6-8d, f and Figs. 6-7d, f, the test results of the five models in Experiment 2 in December 2021 and 358 January 2022 show that the error indices of RMSE and FA of each model exhibit a minimal difference 359 in two months. Nonetheless, disregarding the R results, the performance of the five models in 360 December 2021 is inferior to that in January 2022. The diurnal variation scatter plot of two months is 361 tested. As is shown in Figs. S6-8d, f and Figs. 6-7d, f, the red scatters represent the nighttime wind 362 speed, which is more concentrated on the 1:1 line. In contrast, the blue scatters represent the afternoon 363 wind speed, which is slightly away from the 1:1 line. This suggests that the correction effect of the five 364 models (VMD-PCA-lightGBM, VMD-PCA-XGBoost, VMD-PCA-RF, VMD-PCA-DBN, and 365 VMD-PCA-MLP) exhibits a noticeable diurnal variation.



367Figure 62: The scatter density map compared with the actual 10-meter wind speed: (a) 10-fold368cross-validation training set of VMD-PCA-RF model in February 2022, (b) 10-fold cross-validation369validation set of VMD-PCA-RF model in February 2022. The 24-hour scatter map compared with the370actual 10-meter wind speed: (c) WRF forecasts in December 2021, (d) VMD-PCA-RF model forecasts in371December 2021, (e) WRF forecasts in January 2022, and (f) VMD-PCA-RF model forecasts in January 2022.372(The time is UTC + 08:00.)





Figure 7:-.__The scatter density map compared with the actual 10-meter wind speed: (a) 10-fold cross-validation training set of VMD-PCA-lightGBM model in February 2022, (b) 10-fold cross-validation validation set of VMD-PCA-lightGBM model in February 2022. The 24-hour scatter map compared with the actual 10-meter wind speed: (c) WRF forecasts in December 2021, (d) VMD-PCA-lightGBM model forecasts in December 2021, (e) WRF forecasts in January 2022, and (f) VMD-PCA-lightGBM model forecasts in January 2022. (The time is UTC + 08:00.)

382 **3.3 Comparison of the two experiments**

383 Firstly, all 10 models effectively corrected the 10-meter wind speed forecasted by WRF. Table 3-384 S2 and Table 4-S3 represent the evaluation indices of wind speed errors predicted by 10 models in 385 December 2021 and January 2022. From the two tables, it is evident that the VMD-PCA-RF and 386 VMD-PCA-lightGBM models have the best performance in December 2021 and January 2022, 387 respectively, with the most comprehensive performance of the forecast indicators. The MAE, RMSE, 388 rMAE, rMAE, and FA for the two models VMD-PCA-RF (VMD-PCA-lightGBM) were 0.46 m s⁻¹ (0.45 m s⁻¹), 0.62 m s⁻¹ (0.63 m s⁻¹), 37.36 % (34.75 %), 50.39 % (48.65 %), and 91.79 % (91.49 %) in 389 390 December 2021 (January 2022). Additionally, based on the analysis of the Taylor chart (Fig. 8e, f) of 391 10 models in Fig. 8, it can also be seen that the scatter distance of VMD-PCA-RF and 392 VMD-PCA-lightGBM models is closest to the observed black dotted line and the black triangle 393 position. The two models show that the standard deviation is close to the observed wind speed, with the 394 lowest RMSE and the highest R. Secondly, in the comparison of cumulative probability distributions, 395 all models passed Kolmogorov's 5 % confidence interval test when the interval of wind speed is 0.5 m s⁻¹ (Fig. 8a, d). However, when the interval of wind speed is 0.2 m s^{-1} (Fig. 8b, e), the 396 397 VMD-PCA-lightGBM model deviated from Kolmogorov's 5 % confidence interval detection in 398 December 2021. This indicates that the VMD-PCA-RF model has a better predictive effect than the 399 VMD-PCA-lightGBM model in December 2021 when the actual wind speed is within the range of 0.4 400 m s⁻¹-0.8 m s⁻¹.

401

Table 3. Table of evaluation indices of wind speed error predicted by 10 models in December 2021

Model	MAE (m s ⁻¹)	RMSE(m s ⁻¹)	rMAE-(%)	rRMSE (%)	FA (%)	R
VMD-PCA-lightGBM	0.47	0.63	37.67	51.25	91.13	0.81
VMD-PCA-XGBoost	0.49	0.68	39.8 4	54.82	89.22	0.78
VMD-PCA-RF	0.46	0.62	37.36	50.39	91.79	0.82
VMD-PCA-DBN	0.53	0.75	4 3.32	61.13	87.93	0.71
VMD-PCA-MLP	0.53	0.72	4 3.0 4	58.47	87.2	0.75
lightGBM	0.49	0.67	39.59	54.16	89.68	0.79
XGBoost	0.51	0.70	41.51	56.64	87.9	0.77
RF	0.48	0.65	38.80	52.32	90.64	0.81

DBN	0.56	0.77	4 5.25	62.46	86.74	0.71
MLP	0.55	0.74	44 .65	60.1	86.08	0.75
402						
403 Table 4. Table	e of evaluation ind	ices of wind speed	error predicted	y 10 models in Ja	nuary 2022	
Model	MAE (m s ⁻¹)	RMSE(m s ⁻¹)	rMAE (%)	rRMSE (%)	FA (%)	R
VMD-PCA-lightGBM	0.45	0.63	34.75	4 8.65	91.49	0.78
VMD-PCA-XGBoost	0.47	0.66	36.31	51.01	90.23	0.76
VMD-PCA-RF	0.46	0.64	35.06	49.00	91.57	0.78
VMD-PCA-DBN	0.53	0.75	4 0.96	57.49	87.61	0.67
VMD-PCA-MLP	0.50	0.69	38.46	53.16	88.94	0.73
lightGBM	0.46	0.64	35.24	4 9.3 4	91.11	0.77
XGBoost	0.48	0.67	36.68	51.38	89.88	0.75
RF	0.46	0.64	35.18	4 9.13	91.36	0.78
DBN	0.53	0.74	4 0.97	56.86	87.71	0.68
MLP	0.49	0.68	37.83	52.26	89.57	0.74





406 | Figure 8:-._The cumulative distribution probability scatter plots of the actual wind speed and the predicted
407 wind speed of 10 models in wind speed intervals of 0.5 m s⁻¹ ((a) represents December 2021, (d) represents
408 January 2022) and 0.2 m s⁻¹ ((b) represents December 2021, (e) represents January 2022) respectively;
409 Taylor distribution map ((c) represents December 2021, (f) represents January 2022).

412 **3.4 Spatial-temporal variations in the best models**

413 Based on our comparative analysis results, we conclude that the best performing combination 414 models in December 2021 and January 2022 are VMD-PCA-RF and VMD-PCA-lightGBM 415 respectively. Fig. 9 and Fig. S9 shows the diurnal variation corrections of the two best models for a 416 given month, as well as the diurnal variation of wind speed in the original WRF forecast. The wind 417 speed of the original WRF numerical weather forecast shows a noticeable overestimation, which is 418 confirmed in Fig. 7c and 7e. The scatters of WRF forecast predominantly deviate towards the upper left 419 corner, with relatively low correlation coefficients, 0.56 and 0.23, respectively. Furthermore, the wind 420 speed forecast by WRF displays obvious diurnal variation traits, characterized by large errors between 421 afternoon and evening, specifically between 11:00 and 20:00 (Fig. 9a, Fig. S9ab). Moreover, the actual 422 average wind speed in January 2022 deviates from the range of one standard deviation of the WRF 423 forecast wind speed at 17:00 and 18:00 (Fig. S9a). This demonstrates that the wind speed forecast by 424 WRF is inaccurate and exhibits substantial diurnal variation errors.

After the best model was corrected, the error of diurnal variation was significantly reduced (Fig. 9e9b, Fig. S9bd). First, the average wind speed corrected by the best model is essentially consistent with the actual average wind speed curve, with minimal error and no diurnal variation. Second, the one standard deviation range of the corrected and actual wind speeds is also well-matched, indicating that the corrected and actual wind speed distributions are consistent. The correction effect at 16:00 and 17:00 on January 2022 is suboptimal, which may be due to the insufficient generalization of the training model and the excessive fluctuation of the actual wind speed at these two-time points.

432 The FA (Fig. 10a, Fig. S10ab) and RMSE (Fig. 10e10b, Fig. S10bf) distribution of WRF forecast 433 10-meter wind speed at 410 stations in the five southern provinces shows that the 10-meter wind speed 434 prediction effect of the WRF model in Yunnan is superior to that in the other four provinces. In the 435 regions of Hainan, Guangxi, and Guangdong, the number of sites with a RMSE for 10-meter wind 436 speed forecast in December 2021 ranging from 5.6 to 6.0 m s⁻¹ was significantly higher than in January 437 2022, especially in coastal areas (Fig. 10b, Fig. S10b). In the Yunnan area, the FA of most WRF 438 forecast station 10-meter wind speeds exceeds 40 %, and the RMSE value is mostly below 2.4 m s⁻¹. 439 Conversely, in other regions, such as Guangxi, Guangdong, and Hainan, the terrain is relatively flat. 440 The FA of the 10-meter wind speed forecast by WRF is as low as 30 % at some stations, and the

RMSE reaches up to 5.4 m s⁻¹. However, after the VMD-PCA-RF and VMD-PCA-lightGBM models
are corrected, the FA of most stations in the five southern provinces is as high as 90 %, and the RMSE
is as low as 0.6 m s⁻¹. Moreover, in Guangxi, Guangdong, and Hainan, where the WRF forecast effect
is subpar, the accuracy of the corrected 10-meter wind speed by VMD-PCA-RF
(VMD-PCA-lightGBM) is significantly improved.



Figure 9:-._VMD-PCA-lightGBM, VMD-PCA-RF, and WRF daily variation of predicted and actual wind
speeds in December 2021-and January 2022. (The shading areas represent an interval of 1 standard
deviation, which is a 68% confidence interval. <u>The time is UTC + 08:00</u>.)





Figure 10: FA ((a), (b), (c), and (d)) and RMSE ((eb), (fd), (g), and (h)) distribution maps of VMD-PCA-RF,
VMD-PCA-lightGBM and WRF models on 410 sites in five southern provinces in -((a), (c), (c), and (g)
represent December 2021; (b), (d), (f), and (h) represent January 2022).

457 4. Discussion

453

458 4.1 The effects of BOA-VMD-PCA

459 It is shown in Table 2-S1 that the hyper-parameters of the 10 models in the two experiments are 460 different. Since the DBN model is not added to the scikit-learn Python learning package, it is 461 challenging to call the BOA algorithm for tuning parameters. Apart from the DBN model, all the other 462 models are optimized using the BOA algorithm. From the various evaluation indicators in Table 3-S2 463 and Table 4<u>S3</u>, the DBN model, which does not use the BOA algorithm to adjust the model parameters 464 to obtain an optimal parameter configuration, yields the worst prediction results in December 2021 and 465 January 2022. Moreover, studies (Xiong et al., 2022) also have shown that BOA can further improve 466 the model's prediction accuracy by configuring optimal hyper-parameters. The hyper-parameters such 467 as the number of neurons and learning rate in the hidden layer, significantly impact the model's 468 performance. When the same model is applied to different data sets of two experiments, the BOA 469 adaptively obtains the optimal combination of hyper-parameters, overcoming the limitations of manual 470 parameter adjustment (Guo et al., 2021). This suggests that the selection of model hyper-parameters 471 introduces considerable uncertainty in our prediction results. Therefore, the choice of optimization 472 model parameters represents one source of uncertainty in the correction results, which entails the
473 complexity of parameter selection. However, a more advanced parameter tuning method, such as the
474 BOA tuning algorithm, is essential.

The VMD is used to obtain unknown but meaningful features hidden in the 10-meter wind speed sequences predicted using WRF models (Li et al., 2022). In addition, the PCA can extract important components of anemometer subsequences. When the stationary subsequence serves as an input to the error correction model, it contains more valuable information than the previous non-stationary wind speed sequences (Xu et al., 2021).

480 The complexity of the input factors in this study is one of the sources of uncertainty in the process 481 of correcting WRF prediction results. The input factors of the two experiments are not identical. In the 482 second set of experiments, the input of meteorological factors is reduced based on the first set of 483 experiments, while component information of the 10-meter wind speed predicted by WRF is increased. 484 Multiple wind speed components processed by VMD-PCA and noise reduction are introduced. Among 485 them, the importance of pca0 and IMF0 introduced is approximately 5 %. In the 4013-month test sets, 486 the correction accuracy of experiment 2 is no less than the results of experiment 1 (Figs. \$9\$11, 1012), 487 indicating that the 10-meter wind speed components introduced by the VMD-PCA contribute 488 positively to the correction results.

489

490 **4.2 RF feature importance**

491 To further understand the feature importance ranking of the RF models, we divided the model 492 prediction results and actual wind speeds of the 410 stations into 20 equal parts according to terrain 493 height above sea level (Fig. 11). First of all, the actual wind speed in December 2021 and January 2022 494 varies with the height of the station, showing that the lower the height of the station, the more 495 significant the change of wind speed. This relationship is associated with the wind speed profile of the 496 atmosphere, where wind speed increases as height decreases. Secondly, the wind speed during the day 497 is generally greater than the wind speed at night, which is related to the turbulent motion of the 498 atmosphere during the day. Solar radiation causes the atmosphere to mix, resulting in convective 499 movement. The 10-meter wind speed at night is affected by the cooling radiation of the surface, and the 500 atmosphere is relatively stable.

501 The 10-meter wind speed predicted by WRF has the highest feature importance in the correction 502 process of the RF models. Input factors with distinct geographic information, such as latitude, 503 longitude, and height, rank highly in feature importance. Similarly, when Sun et al. 2019 used machine 504 learning to correct the 10-meter wind speed predicted by the numerical weather prediction model 505 ECMWF, the characteristic weight of the 10-meter wind speed predicted by the model was the highest, 506 followed by the sea-land factor. Also, as the 10-meter wind speed forecast by WRF increases, the 507 instability of the 10-meter wind speed corrected by the 10 machine learning models gradually increased, 508 and the correction accuracy gradually decreased (Fig. 12). This partly explains the higher importance 509 of the 10-meter wind speed forecast by WRF.

510 With 1 km as the center, the measured 10-meter wind speed is more variable in areas where the 511 station terrain height increases or decreases. However, the pink box of the 10-meter wind speed 512 predicted by WRF becomes wider as the station terrain height decreases (Fig. 11). The distance 513 between the gray box and the pink box is greater as the station terrain height decreases. It shows that 514 the 10-meter wind speed predicted by WRF has less accuracy with the station terrain height decreases. 515 The VMD-PCA-RF and VMD-PCA-lightGBM models significantly reduce the variability of the 516 10-meter wind speed predicted by WRF. When the height of the station increases or decreases at 1 km, 517 the correction intensity tends to increase gradually. This further explains the higher importance of the 518 height factor in the RF model training.





520

521 Figure 11:-_The boxplots of the predicted wind speeds of the VMD-PCA-RF (yellow), VMD-PCA-lightGBM 522 (blue), and WRF (pink) models at 20 stations at different height intervals, and the boxplots of the actual 523 wind speeds (gray).



526 Figure 12:-_The prediction error boxplots of 10 models in different WRF prediction intervals.

525

528 4.3 Stability analysis of the proposed models

In order to identify the best model of the five southern provinces and assess the model's stability, we evaluated all 10 models over <u>10-13</u> different months. Fig. 13 shows the evaluation histogram of the 10-meter wind speed predicted by the 10 models in Experiment 1 and Experiment 2, as well as the actual wind speed in various months. Meanwhile, Fig. <u>S9-S11</u> and Fig. <u>S10-S12</u> can more effectively illustrate the daily changes of the revised results of 10 models in <u>10-13</u> different months. As shown in Fig. 13, the evaluation indices of the model trained in Experiment 2, after VMD-PCA processing, outperform those of the model trained in Experiment 1. The RF model demonstrates exceptional

536	robustness, while the MLP model exhibits the poorest performance. VMD-PCA-RF evaluation indices
537	are relatively stable across the $\frac{10-13}{10}$ months, with a correlation coefficient R above 0.6, FA above
538	85 %, MAE below 0.6 m s ⁻¹ , RMSE below 0.8 m s ⁻¹ , rMAE below 60 %, and rRMSE below 75 %.
539	However, the robustness of the VMD-PCA-lightGBM and VMD-PCA-XGBoost models is inferior to
540	that of the VMD-PCA-RF, with all six evaluation indices performing worse than the VMD-PCA-RF as
541	the seasons and months change. In general, VMD-PCA-lightGBM is the superior wind speed
542	correction model for the winter, and VMD-PCA-RF performs the best throughout the entire year in the
543	five southern provinces. In cases where ample machine CPU and other hardware resources, as well as
544	training time, are available, we recommend using VMD-PCA-lightGBM for modeling each season.
545	However, when dealing with limited resources such as a laptop and constrained training time, we
546	recommend using VMD-PCA-RF to train data for a single month, as this yields the most robust
547	correction results.







550 551

Figure 13: __Evaluation histograms of 10-meter wind speed predicted by 10 models in different months in Experiment 1 and Experiment 2 ((a), (b), (c), (d), (e), and (f) represent R, FA (%), MAE (m s⁻¹), RMSE (m s⁻¹), rMAE (%), and rRMSE (%) respectively).

552 553

554 **5.** Conclusions

In an effort to enhance the wind speed prediction performance for wind farms, this study developed a WRF-based multi-step wind speed prediction model. A hybrid error correction strategy combining BOA, VMD, PCA, and RF (LightGBM) is proposed to increase the accuracy of WRF simulations. The first group of experiments used various meteorological elements as input factors in a control experiment. In the second group of experiments, the wind speed sequence predicted by the WRF model was decomposed into multiple IMFs using the VMD algorithm for feature extraction. A principal component analysis method is used to extract meaningful principal components from these subsequence IMFs to improve computational efficiency. In the error correction model, RF (lightGBM) and other algorithms are used to train the relationship between different input factors and the actual wind speed error, respectively.

565 Through a case analysis of 410 stations in five southern provinces in China, the following 566 conclusions can be drawn: (1) The machine learning models tuned by the BOA-VMD-PCA algorithm 567 exhibit a positive impact on wind speed error correction; (2) Feature importance analysis revealed that 568 the top eight contributing factors for correcting WRF forecasted wind speed include WRF forecast 569 10-meter wind speed (WS₁₀), latitude, longitude, altitude, pca0 (pca0 physically represents the lowest 570 frequency wind speed series after PCA treatment of all IMFk (k=0, 1, 2, ..., 8) sub-series with reduced 571 dimension), humidity, pressure, IMF0 (IMF0 physically represents the wind speed stationary series 572 with a specific lowest center frequency after the original wind speed series has been processed by 573 VMD); (3) VMD-PCA-RF and VMD-PCA-lightGBM are the most suitable wind speed correction 574 algorithms for December 2021 and January 2022, respectively. The MAE, RMSE, FA, rMAE, rRMSE, 575 and R of the corrected wind speed and the actual wind speed are 0.46 (0.45), 0.62 m s⁻¹ (0.63 m s⁻¹), 37.36 % (34.75 %), 50.39 % (48.65 %), 91.79 % (91.49 %), and 0.82 (0.78); and (4) The proposed 576 577 wind speed correction model (VMD-PCA-RF) demonstrates the highest prediction accuracy and 578 stability in the five southern provinces in nearly a year and at different heights. VMD-PCA-RF 579 evaluation indices for 10-13 months remain relatively stable: R is above 0.6, FA is above 85 %, MAE 580 is below 0.6 m s⁻¹, RMSE is below 0.8 m s⁻¹, rMAE is below 60 %, and rRMSE is below 75 %. In 581 future research, the proposed VMD-PCA-RF algorithm can be extrapolated to the 3 km grid points of 582 the five southern provinces to generate a 3km grid-corrected wind speed product.

585 Code availability

586 The code and model available are free-access repository on Zenodo at as а 587 https://doi.org/10.5281/zenodo.8108889 (Zhou, 2023).

588 Data Availability

589 The data is available as a free-access repository on Zenodo at https://doi.org/10.5281/zenodo.8108889
590 (Zhou, 2023).

591 Author contributions

592 SZ developed the software, visualized the data, and prepared the original draft. SZ and YG developed

593 the methodology and carried out the formal analysis. XX and SZ validated data. SZ, YG, XX, ZD, and

594 YL reviewed and edited the text. All authors have read and agreed to the published version of the 595 paper.

596 Competing interests

597 The authors declare that they have no conflict of interest.

598 Financial support

- 599 This research has been supported by the second batch of service public bidding projects for EHV
- 600 transmission companies in 2022 (2022-FW-2-ZB) (grant no. CG0100022001526556).

603 **References**

- Barthelmie, R. J., Palutikof, J. P., and Davies, T. D.: Estimation of sector roughness lengths and the
 effect on prediction of the vertical wind speed profile, Boundary-Layer Meteorol, 66, 19–47,
 https://doi.org/10.1007/BF00705458, 1993.
- Cassola, F. and Burlando, M.: Wind speed and wind energy forecast through Kalman filtering of
 Numerical Weather Prediction model output, Applied Energy, 99, 154–166,
 https://doi.org/10.1016/j.apenergy.2012.03.054, 2012.
- 610 Chen, F., Janjić, Z., and Mitchell, K.: Impact of Atmospheric Surface-layer Parameterizations in the
- 611 new Land-surface Scheme of the NCEP Mesoscale Eta Model, Boundary-Layer Meteorology, 85,
- 612 391–421, https://doi.org/10.1023/A:1000531001463, 1997.
- 613 Chen, K. and Yu, J.: Short-term wind speed prediction using an unscented Kalman filter based
 614 state-space support vector regression approach, Applied Energy, 113, 690–705,
 615 https://doi.org/10.1016/j.apenergy.2013.08.025, 2014.
- Cheng, W. Y. Y., Liu, Y., Liu, Y., Zhang, Y., Mahoney, W. P., and Warner, T. T.: The impact of
 model physics on numerical wind forecasts, Renewable Energy, 55, 347–356,
 https://doi.org/10.1016/j.renene.2012.12.041, 2013.
- 619 Deng, Y., Wang, B., and Lu, Z.: A hybrid model based on data preprocessing strategy and error
- 620 correction system for wind speed forecasting, Energy Conversion and Management, 212, 112779,
- 621 https://doi.org/10.1016/j.enconman.2020.112779, 2020.
- 622 Dhiman, H. S. and Deb, D.: A Review of Wind Speed and Wind Power Forecasting Techniques,
- 623 arXiv:2009.02279 [cs, eess], 2020.
- 624 Dong, L., Ren, L., Gao, S., Gao, Y., and Liao, X.: Studies on wind farms ultra-short term NWP wind
- 625 speed correction methods, in: 2013 25th Chinese Control and Decision Conference (CCDC), 2013 25th
- 626 Chinese Control and Decision Conference (CCDC), Guiyang, China, 1576–1579,
- 627 https://doi.org/10.1109/CCDC.2013.6561180, 2013.
- 628 Erdem, E. and Shi, J.: ARMA based approaches for forecasting the tuple of wind speed and direction,
- 629 Applied Energy, 88, 1405–1414, https://doi.org/10.1016/j.apenergy.2010.10.031, 2011.

- 630 Guo, X., Zhu, C., Hao, J., Zhang, S., and Zhu, L.: A hybrid method for short-term wind speed
- 631 forecasting based on Bayesian optimization and error correction, Journal of Renewable and Sustainable
- 632 Energy, 13, 036101, https://doi.org/10.1063/5.0048686, 2021.
- 633 Guo, Z., Zhao, W., Lu, H., and Wang, J.: Multi-step forecasting for wind speed using a modified
- 634 EMD-based artificial neural network model, Renewable Energy, 37, 241-249,
- 635 https://doi.org/10.1016/j.renene.2011.06.023, 2012.
- 636 Hanifi, S., Liu, X., Lin, Z., and Lotfian, S.: A Critical Review of Wind Power Forecasting
- 637 Methods—Past, Present and Future, Energies, 13, 3764, https://doi.org/10.3390/en13153764, 2020.
- 638 Hu, H., Wang, L., and Tao, R.: Wind speed forecasting based on variational mode decomposition and
- 639 improved echo state network, Renewable Energy, 164, 729-751,
- 640 https://doi.org/10.1016/j.renene.2020.09.109, 2021.
- 641 Hu, J., Wang, J., and Zeng, G.: A hybrid forecasting approach applied to wind speed time series,
- 642 Renewable Energy, 60, 185–194, https://doi.org/10.1016/j.renene.2013.05.012, 2013.
- 643 Huang, Y., Yang, L., Liu, S., and Wang, G.: Multi-Step Wind Speed Forecasting Based On Ensemble
- 644 Empirical Mode Decomposition, Long Short Term Memory Network and Error Correction Strategy,
- 645 Energies, 12, 1822, https://doi.org/10.3390/en12101822, 2019.
- 646 Iacono, M. J., Delamere, J. S., Mlawer, E. J., Shephard, M. W., Clough, S. A., and Collins, W. D.:
- 647 Radiative forcing by long-lived greenhouse gases: Calculations with the AER radiative transfer models,
- 648 J. Geophys. Res., 113, D13103, https://doi.org/10.1029/2008JD009944, 2008.
- 649 Isham, M. F., Leong, M. S., Lim, M. H., and Ahmad, Z. A.: Variational mode decomposition: mode
- 650 determination method for rotating machinery diagnosis, J VIBROENG, 20, 2604-2621,
- 651 https://doi.org/10.21595/jve.2018.19479, 2018.
- James, E. P., Benjamin, S. G., and Marquis, M.: Offshore wind speed estimates from a high-resolution
- 653 rapidly updating numerical weather prediction model forecast dataset, Wind Energy, 21, 264–284,
- 654 https://doi.org/10.1002/we.2161, 2018.
- 655 Janjić, Z. I.: The Step-Mountain Eta Coordinate Model: Further Developments of the Convection,
- 656 Viscous Sublayer, and Turbulence Closure Schemes, Monthly Weather Review, 122, 927–945,
- 657 https://doi.org/10.1175/1520-0493(1994)122<0927:TSMECM>2.0.CO;2, 1994.

- 558 Jiménez, P. A. and Dudhia, J.: Improving the Representation of Resolved and Unresolved Topographic
- Effects on Surface Wind in the WRF Model, Journal of Applied Meteorology and Climatology, 51,

660 300–316, https://doi.org/10.1175/JAMC-D-11-084.1, 2012.

- Joyce, L. and Feng Z.: Global Wind Report 2023, Global Wind Energy Council,
 https://gwec.net/globalwindreport2023 (last access: 9 May 2023), 2023.
- 663 Li, G. and Shi, J.: Application of Bayesian model averaging in modeling long-term wind speed
- distributions, Renewable Energy, 35, 1192–1202, https://doi.org/10.1016/j.renene.2009.09.003, 2010.
- Li, Y., Tang, F., Gao, X., Zhang, T., Qi, J., Xie, J., Li, X., and Guo, Y.: Numerical Weather Prediction
- 666 Correction Strategy for Short-Term Wind Power Forecasting Based on Bidirectional Gated Recurrent
- 667 Unit and XGBoost, Front. Energy Res., 9, 836144, https://doi.org/10.3389/fenrg.2021.836144, 2022.
- 668 Liu, H., Mi, X., and Li, Y.: An experimental investigation of three new hybrid wind speed forecasting
- models using multi-decomposing strategy and ELM algorithm, Renewable Energy, 123, 694-705,
- 670 https://doi.org/10.1016/j.renene.2018.02.092, 2018.
- Liu, Y., Wang, Y., Li, L., Han, S., and Infield, D.: Numerical weather prediction wind correction
 methods and its impact on computational fluid dynamics based wind power forecasting, Journal of
 Renewable and Sustainable Energy, 8, 033302, https://doi.org/10.1063/1.4950972, 2016.
- 674 Ma, Z., Chen, H., Wang, J., Yang, X., Yan, R., Jia, J., and Xu, W.: Application of hybrid model based
- on double decomposition, error correction and deep learning in short-term wind speed prediction,
- Energy Conversion and Management, 205, 112345, https://doi.org/10.1016/j.enconman.2019.112345,
- 677 2020.
- 678 Salcedo-Sanz, S., Ángel M. Pérez-Bellido, Ortiz-García, E. G., Portilla-Figueras, A., Prieto, L., and
- 679 Paredes, D.: Hybridizing the fifth generation mesoscale model with artificial neural networks for
- short-term wind speed prediction, Renewable Energy, 34, 1451–1457,
 https://doi.org/10.1016/j.renene.2008.10.017, 2009.
- 682 Salcedo-Sanz, S., Ortiz-García, E., Pérez-Bellido, Á., Portilla-Figueras, A., and Prieto, L.: Short term
- 683 wind speed prediction based on evolutionary support vector regression algorithms, Expert Syst. Appl.,
- 684 38, 4052–4057, https://doi.org/10.1016/j.eswa.2010.09.067, 2011.

- 685 Skamarock, W. C., Klemp, J. B., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., Wang, W., Powers, J. G.,
- Duda, M. G., Barker, D. M., and Huang, X.-Y.: A Description of the Advanced Research WRF Model
 Version 4, 2021.
- 688 Služenikina, J. and Männik, A.: Impact of the ASCAT scatterometer winds on the quality of HIRLAM 689 Proc. Estonian analysis in case of severe storms, Acad. Sci., 65, 177, 690 https://doi.org/10.3176/proc.2016.3.03, 2016.
- Sun, Q., Jiao, R., Xia, J., Yan, Z., Li, H., Sun, J., Wang, L., and Liang, Z.: Adjusting Wind Speed
 Prediction of Numerical Weather Forecast Model Based on Machine Learning Methods.
 Meteorological Monthly, 45(3): 426-436. https://doi.org/10.7519/j.issn.1000-0526.2019.03.012, 2019.
- 694 Tang, R., Ning, Y., Li, C., Feng, W., Chen, Y., and Xie, X.: Numerical Forecast Correction of
- 695 Temperature and Wind Using a Single-Station Single-Time Spatial LightGBM Method, Sensors, 22,
- 696 193, https://doi.org/10.3390/s22010193, 2021.
- Tascikaraoglu, A. and Uzunoglu, M.: A review of combined approaches for prediction of short-term
 wind speed and power, Renewable and Sustainable Energy Reviews, 34, 243–254,
 https://doi.org/10.1016/j.rser.2014.03.033, 2014.
- 700 Thompson, G., Field, P. R., Rasmussen, R. M., and Hall, W. D.: Explicit Forecasts of Winter
- 701 Precipitation Using an Improved Bulk Microphysics Scheme. Part II: Implementation of a New Snow
- 702 Parameterization, Monthly Weather Review, 136, 5095–5115,
- 703 https://doi.org/10.1175/2008MWR2387.1, 2008.
- Tiedtke, M.: A Comprehensive Mass Flux Scheme for Cumulus Parameterization in Large-Scale
 Models, Monthly Weather Review, 117, 1779–1800,
 https://doi.org/10.1175/1520-0493(1989)117<1779:ACMFSF>2.0.CO;2, 1989.
- $100 \qquad \text{hups.//doi.org/10.1175/1520-0495(1969)11/<1779. ACMIPSP 2.0.00, 2, 1989.}$
- Wang, C., Zhang, H., Fan, W., and Ma, P.: A new chaotic time series hybrid prediction method of wind
 power based on EEMD-SE and full-parameters continued fraction, Energy, 138, 977–990,
 https://doi.org/10.1016/j.energy.2017.07.112, 2017.
- 710 Wang, J. and Hu, J.: A robust combination approach for short-term wind speed forecasting and analysis
- 711 Combination of the ARIMA (Autoregressive Integrated Moving Average), ELM (Extreme Learning
- 712 Machine), SVM (Support Vector Machine) and LSSVM (Least Square SVM) forecasts using a GPR

713 (Gaussian Process Regression) model, Energy, 93, 41–56, https://doi.org/10.1016/j.energy.2015.08.045,

714 2015.

- Williams, J. L., Maxwell, R. M., and Monache, L. D.: Development and verification of a new wind
 speed forecasting system using an ensemble Kalman filter data assimilation technique in a fully
 coupled hydrologic and atmospheric model: Data Assimilation in a Coupled Forecasting System, J.
 Adv. Model. Earth Syst., 5, 785–800, https://doi.org/10.1002/jame.20051, 2013.
- 719 Xiong, X., Guo, X., Zeng, P., Zou, R., and Wang, X.: A Short-Term Wind Power Forecast Method via
- XGBoost Hyper-Parameters Optimization, Front. Energy Res., 10, 905155,
 https://doi.org/10.3389/fenrg.2022.905155, 2022.
- 722 Xu, Q., He, D., Zhang, N., Kang, C., Xia, Q., Bai, J., and Huang, J.: A Short-Term Wind Power
- 723 Forecasting Approach With Adjustment of Numerical Weather Prediction Input by Data Mining, IEEE
- 724 Trans. Sustain. Energy, 6, 1283–1291, https://doi.org/10.1109/TSTE.2015.2429586, 2015.
- Xu, W., Liu, P., Cheng, L., Zhou, Y., Xia, Q., Gong, Y., and Liu, Y.: Multi-step wind speed prediction
- by combining a WRF simulation and an error correction strategy, Renewable Energy, 163, 772–782,
- 727 https://doi.org/10.1016/j.renene.2020.09.032, 2021.
- 728 Zhang, C., Wang, Y., and Hamilton, K.: Improved Representation of Boundary Layer Clouds over the
- 729 Southeast Pacific in ARW-WRF Using a Modified Tiedtke Cumulus Parameterization Scheme*,
- 730 Monthly Weather Review, 139, 3489–3513, https://doi.org/10.1175/MWR-D-10-05091.1, 2011.
- 731 Zhang, D., Peng, X., Pan, K., and Liu, Y.: A novel wind speed forecasting based on hybrid
- 732 decomposition and online sequential outlier robust extreme learning machine, Energy Conversion and
- 733 Management, 180, 338–357, https://doi.org/10.1016/j.enconman.2018.10.089, 2019a.
- 734 Zhang, Y., Chen, B., Pan, G., and Zhao, Y.: A novel hybrid model based on VMD-WT and
- 735 PCA-BP-RBF neural network for short-term wind speed forecasting, Energy Conversion and
- 736 Management, 195, 180–197, https://doi.org/10.1016/j.enconman.2019.05.005, 2019b.
- 737 Zhang, Z., Ye, L., Qin, H., Liu, Y., Wang, C., Yu, X., Yin, X., and Li, J.: Wind speed prediction
- 738 method using Shared Weight Long Short-Term Memory Network and Gaussian Process Regression,
- 739 Applied Energy, 247, 270–284, https://doi.org/10.1016/j.apenergy.2019.04.047, 2019c.

- 740 Zhao, J., Guo, Z.-H., Su, Z.-Y., Zhao, Z.-Y., Xiao, X., and Liu, F.: An improved multi-step forecasting
- 741 model based on WRF ensembles and creative fuzzy systems for wind speed, Applied Energy, 162,
- 742 808–826, https://doi.org/10.1016/j.apenergy.2015.10.145, 2016.
- 743 Zhao, J., Wang, J., Guo, Z., Guo, Y., Lin, W., and Lin, Y.: Multi-step wind speed forecasting based on
- numerical simulations and an optimized stochastic ensemble method, Applied Energy, 255, 113833,
- 745 https://doi.org/10.1016/j.apenergy.2019.113833, 2019.
- 746 Zjavka, L.: Wind speed forecast correction models using polynomial neural networks, Renewable
- 747 Energy, 83, 998–1006, https://doi.org/10.1016/j.renene.2015.04.054, 2015.
- 748