

1 **Prediction of hourly wind speed time series at unsampled**
2 **locations using machine learning**

3
4
5 Freddy Houndekindo^{1*}, Taha B.M.J. Ouarda¹

6
7 ¹Canada Research Chair in Statistical Hydro-Climatology, Institut national de la recherche
8 scientifique, Centre Eau Terre Environnement, INRS-ETE, 490 de la Couronne, Québec, QC,
9 G1K 9A9, Canada

10
11
12
13
14
15
16 *Corresponding author: Freddy Houndekindo
17 490, Couronne street, Québec, QC, G1K 9A9, Canada
18 Tel: +1 418-654-3842
19 E-mail: freddy.houndekindo@inrs.ca

22 Abbreviations

a.g.l	Above ground level
CDF	Cumulative distribution function
DEM	Digital elevation model
ECCC	Environment and Climate Change Canada
ERA5-WSQ	Wind speed quantiles extracted from the ERA5 dataset (m/s)
GWA	Global wind atlas
GWA-ERA5	Bias-corrected ERA5 using GWA (m/s)
IAV	Interannual variability
IDW	Inverse distance weighting
LGBM	Light gradient-boosting machine
LGBMQR	Lightgbm for quantile regression
LGBMSI	LGBM for spatial interpolation
LGBMSI-ERA5	LGBMSI using the ERA5 wind data as covariates
LGMBQR-ERA5	LGMBQR using ERA5-WSQ as covariates
MAE	Mean absolute error (m/s)
ME	Mean error (m/s)
MRMR	Minimum redundancy maximum relevancy algorithm
OP	Overlap percentage (%)
PC	Pearson correlation
PD	Probability distribution
QM	Quantile mapping
QM-ERA5	Quantile mapping bias correction of ERA5 wind data
QR	Quantile regression
R²	Coefficient of determination
RCov	Robust coefficient of variation
RFSI	Random forest for spatial interpolation
RMSE	Root-mean-squared error (m/s)
TS	time series
WDC	Wind Duration Curve method
WRA	Wind resource assessment
WS	Wind speed
WSD	Wind speed distribution
WSNEP	Wind speed non-exceedance probabilities
WSQ	Wind speed quantiles
WSTS	Wind speed time series

23

24

25 **Abstract**

26 Various models for wind speed mapping have been developed, with increasing attention on models
27 focusing on mapping wind speed distribution. This study extends these models to predict hourly
28 wind speed time series at unsampled locations. A model based on the quantile mapping (QM)
29 procedure was compared to a traditional and machine-learning model to interpolate wind speed
30 spatially. These proposed models were also used with inputs from the ERA5 reanalysis dataset,
31 enabling them to consider local variation in orography and large-scale wind fields. A widely used
32 procedure for mean bias correction of reanalysis based on the Global Wind Atlas (GWA) was
33 implemented and compared to the proposed models. It was found that the QM and machine learning
34 model, both using input from ERA5, significantly outperformed GWA bias correction in terms of
35 time series correlation and probability distribution. Despite being more computationally intensive
36 than GWA bias correction, both models are recommended due to their significantly (in a statistical
37 sense) superior performance.

38 **Keywords:** Bias-correction, ERA5, Light gradient-boosting machine, Quantile regression,
39 Reanalysis, Wind resource assessment

40

41 **1. Introduction**

42 The past decades have witnessed a significant uptake of wind energy in various parts of the world
43 [1]. This growth reflects a global shift toward more renewable energy sources, with wind power
44 playing a prominent role in energy supply [2]. The intermittent nature of wind speed still poses
45 some challenges to the development of the renewable energy source [3]. Due to the cubic
46 relationship between wind speed and power output, inaccuracies in estimating wind speed are

47 amplified when estimating the energy production, leading to suboptimal design of wind energy
48 infrastructure and jeopardizing the profitability and sustainability of the project [4].

49 Prospective studies to evaluate the wind resource across a large region at a high spatial and
50 temporal resolution provide valuable sources of information for the expansion of wind energy [5,
51 6]. In-situ wind speed (WS) data are generally accepted as the most reliable data source for wind
52 resource assessment (WRA). However, measuring stations are often sparsely available in a given
53 region and have limited record length for WRA. Several publicly available datasets exist that give
54 access to wind data at the global scale with high temporal resolution and extensive record length.
55 The European Centre for Medium-Range Weather Forecasts (ECMWF) reanalysis v5 [ERA5: 7],
56 and the NASA's Modern Era Retrospective Analysis for Research and Applications-2 [MERRA-2:
57 8] have been used extensively to conduct WRA across large regions [9, 10]. Samal [11] evaluated
58 the adequacy of MERRA-2 for WRA in India. The author compared the wind data from the
59 reanalysis dataset with observed data collected at meteorological stations. The study found that the
60 reanalysis dataset was more suitable for long-term than short-term planning. In another study,
61 MERRA-2 was used to perform a preliminary evaluation of the wind resource in South Sudan [12].
62 The authors identified areas in the region with high wind potential. Five global reanalysis datasets
63 including ERA5 and MERRA-2 were evaluated for WRA by comparing them with measured WS
64 data from meteorological stations distributed worldwide [13]. The comparative study was based on
65 estimated mean WS, variability, and trends. From the study results, the ERA5 dataset was
66 recommended for wind energy applications.

67 Direct application of reanalysis datasets for WRA still has some drawbacks. Notably, the coarse
68 spatial resolution of reanalysis datasets renders them unable to resolve local variations in orography
69 and surface roughness influencing near-surface WS [14]. A review of the uncertainties associated

70 with the application of reanalysis data for WRA was presented by Gualtieri [9]. Several studies
71 endeavoured to increase the spatial resolution and bias-correct reanalysis datasets using ground
72 measurements and other datasets with higher spatial resolution. The Global Wind Atlas (GWA) is
73 a popular dataset used to correct the bias in reanalysis WS data [15]. In this procedure, the mean
74 WS from the reanalysis dataset is corrected to match the GWA mean WS by applying a correction
75 factor estimated during the overlapping period of both datasets.

76 Alternatively, to reanalysis datasets, spatial interpolation and machine learning models have been
77 used to map wind data at a high spatial resolution using in-situ observations. The main advantage
78 of this approach over the use of reanalysis data is its ability to account for the rapid change in the
79 topography and surface roughness by using covariates extracted from DEM and land use maps. A
80 comparative analysis of several spatial interpolation methods for hourly WS mapping was
81 performed by Collados-Lara, et al. [16]. The authors found that the regression kriging model produced
82 the best results and was selected to generate hourly wind speed time series (WSTS) between 1996
83 and 2016 in The Granada province, Spain. In another study, Cellura, et al. [17] developed a machine-
84 learning model to interpolate mean WS in Sicily, Italy. The author recommended the approach for
85 its ease of application and transferability to other regions. A similar study was conducted in
86 Venezuela [18] to create a regional mean WS map. It should be noted that wind speed distribution
87 (WSD) is often skewed, and the mean is not a good representative of the most typical value of the
88 distribution.

89 In recent studies, authors have been interested in mapping the entire WSD, allowing a better
90 evaluation of the wind resource variability at unsampled locations of interest. For example,
91 Veronesi, et al. [19] mapped the parameters of the Weibull distribution fitted to WS data across the
92 United Kingdom (UK). Jung [20] mapped the parameter of the Wakeby distribution fitted to WS

93 data to estimate the annual wind energy yield with a high spatial resolution in Germany. In another
94 study, Jung and Schindler [21] developed a global model that estimates the parameters of the Kappa
95 and Wakeby distribution for WS variability assessment using estimated L-moments. Houndekindo
96 and Ouarda [22] recently proposed a nonparametric approach for WSD mapping. The approach does
97 not restrict the region to a single WSD distribution family. The availability of methods to map the
98 entire WSD is a crucial step forward compared to past studies where only aggregated values of WS
99 were estimated. However, For the evaluation of WS variability at different temporal resolutions
100 (ex., daily, seasonal, annual), WSTS with a high temporal resolution (e.g., ten min. or one hour)
101 are still required.

102 This study proposes expanding upon previously developed techniques for mapping WSD to predict
103 hourly WSTS at unsampled locations. The proposed method named the Wind Duration Curve
104 (WDC) is inspired by an approach commonly used for environmental variables (see, for instance,
105 Castellarin, et al. [23] and Requena, et al. [24] for application to streamflow data and Ouarda, et al. [25]
106 for application to daily river temperature) and can be seen as an adaptation of the quantile mapping
107 (QM) technique often used to downscale global circulation models and regional climate model
108 outputs [26, 27]. A comprehensive evaluation of the WDC method is performed and the approach
109 is compared to other methods for WSTS estimation at unsampled locations.

110 The paper is structured as follows: Section 2 describes the study area and the datasets. The
111 methodology employed is presented in section 3. The results of the comprehensive evaluation of
112 the different approaches are presented in section 4. The discussion follows in section 5, and section
113 6 gives the conclusions of the study.

114 **2. Study area and dataset**

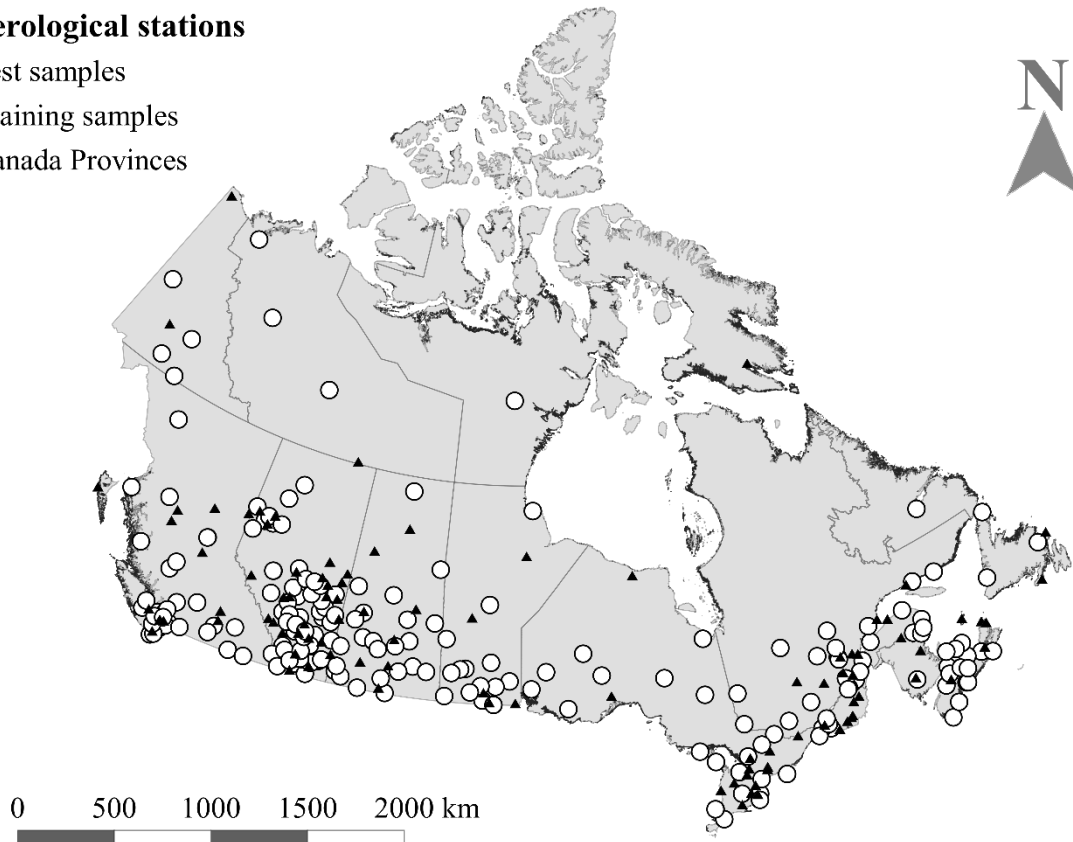
115 Experimental data for the study were obtained from Environment and Climate Change Canada
116 (ECCC) historical climate database (<https://climate.weather.gc.ca/>). Stations with less than 10%
117 missing values between 2011 and 2021 (11 years of mean hourly WS) were selected from the
118 database, resulting in 303 meteorological stations available for the study. WS data at the
119 meteorological stations were typically collected at 10 m above ground level according to ECCC.
120 The measured WS data was considered the most representative of the actual WS condition. Figure
121 1 illustrates the study area and the location of the 303 meteorological stations. In the figure, stations
122 represented with circles were used during the training of the models and those represented with
123 triangles were solely used as test samples.

124 Reanalysis WS data were obtained from ERA5 dataset. Wind speed data from ERA5 are provided
125 in a grid format with a temporal resolution of 1 h available between 1980 and the present. The
126 eastward and northward WS components at 10 m were obtained from the dataset
127 (<https://doi.org/10.24381/cds.adbb2d47>), and the 10 m horizontal WS was calculated and
128 interpolated at the 303 meteorological stations using nearest neighbor interpolation.

129 The WS covariates used in the study are presented in detail in Table S1 of the supporting material.
130 Topographical covariates were calculated from the Advanced Land Observing Satellite (ALOS)
131 Digital Elevation model (DEM) of 30m resolution [ALOS DEM: 28] obtained freely from the Japan
132 Aerospace Exploration Agency. The surface roughness length was estimated from a 2015 land use
133 map of Canada [29] obtained from Natural Resource Canada.

Meteorological stations

- ▲ Test samples
- Training samples
- Canada Provinces



135

136 Figure 1: Study area and location of the 303 meteorological stations used in the study.

137

138 3. Methods

139 3.1. Wind speed distribution mapping

140 In recent studies, different methodologies to map WSD were introduced. Most of these approaches
 141 relied on mapping the parameters of a distribution function fitted to WS data. More recently, a
 142 nonparametric method was developed by Houndekindo and Ouarda [22] to map hourly WSD. The
 143 approach starts by mapping hourly wind speed quantiles (WSQ) using a machine learning model
 144 and WS covariates. Then, the estimated WSQ are used as input of an asymmetric kernel function

145 to estimate the WS cumulative distribution function (CDF) at unsampled locations. The approach
146 is flexible and does not restrict the region to a unique WSD family. In their study, Houndekindo and
147 Ouarda [22] extracted 13 quantiles from observed WSTS and then built a regression model between
148 the covariates and each WSQ. The present study proposes a quantile regression (QR) model to
149 directly estimate 13 conditional WSQ. Although QR models have been used in previous studies for
150 WS forecasting [30] and for the estimation of other hydro-climatic variables at unsampled locations
151 [31], to the author's knowledge, it is the first time they are applied to estimate conditional WSQ at
152 unsampled locations. As done by Houndekindo and Ouarda [22], WSQ at the following 13 percentile
153 points were considered: 5.0% (P1), 12.5% (P2), 20.0% (P3), 27.5% (P4), 35.0% (P5), 42.5% (P6),
154 50.0% (P7), 57.5.0% (P8), 65.0% (P9), 72.5% (P10), 80.0% (P11), 87.5% (P12), and 95.0% (P13).

155 The Light Gradient-Boosting Machine [LGBM: 32] with the pinball loss function (Eq 1) was used
156 as the QR model (herein referred to as LGBMQR). The LGBM was adopted based on its efficiency,
157 scalability for large datasets, and proven high prediction accuracy [33-35]. The LGBM is a
158 histogram-based gradient-boosting model that sequentially builds additive decision trees to
159 minimize a loss function. By discretizing the continuous values of the covariates into a fixed
160 number of bins, the LGBM can significantly reduce the training time and memory usage for large
161 datasets (ex., $N > 10,000$) while maintaining good prediction accuracy. In addition, the LGMB
162 adopts a leaf-wise tree expansion with a fixed maximum depth, improving the model's training
163 performance. Table 1 shows the different model parameters that were tuned. Random search with
164 1000 iterations was used to select the best parameters for the QR model. Random search is not an
165 optimal algorithm for parameter tuning but can still find suitable parameters when allocated a
166 sufficient number of iterations [36]. LGBMQR is a single-output QR model. Thus, it needs to be
167 trained separately for each conditional WSQ of interest. Also, parameter searches can be performed

168 independently for each considered quantile. To reduce the computation burden associated with
 169 performing parameter tuning independently for every quantile of interest, the best parameters
 170 selected when training the model to predict the median (P7) were used for all quantiles.

$$171 \quad \rho_{\tau}(w - w_{\tau}) = \begin{cases} (\tau - 1)|w - w_{\tau}| & (w - w_{\tau}) < 0 \\ \tau|w - w_{\tau}| & (w - w_{\tau}) \geq 0 \end{cases} \quad (1)$$

172 where w_{τ} is the τ -quantile defined as follows:

$$173 \quad w_{\tau} = \inf\{w : F(w|X = x) \geq \tau\} \quad (2)$$

174 with $F(w|X = x)$ the conditional cumulative distribution function of the random variable w .

175 In addition to the covariates presented in Table S1 of the supporting material, hourly WSQ
 176 extracted from the ERA5 dataset (ERA5-WSQ) were assessed as covariates in the current study.
 177 As stated by Jung and Schindler [21], covariates from the ERA5 reanalysis dataset can represent the
 178 large-scale wind field unaffected by local surface properties. The LGBMQR that uses the ERA5-
 179 WSQ will be referred to as LGMBQR-ERA5, and the benefit of using the ERA5-WSQ as
 180 covariates will be evaluated and discussed in the following sections of the paper.

181 Furthermore, to select the optimal number of covariates to include in the model, the available
 182 covariates were ranked according to their relevance and redundancy using the minimum
 183 redundancy maximum relevancy algorithm [MRMR: 37]. Then, the number of covariates to use with
 184 LGBMQR and LGMBQR-ERA5 was treated as an additional hyperparameter during the
 185 implementation of the random search algorithm. The MRMR algorithm has already demonstrated
 186 good performance for WSQ mapping in a comparative study of covariate selection techniques [38].

187 The estimated conditional WSQs were used as input for the Birnbaum-Saunders asymmetric kernel
 188 estimator of CDF [39] to estimate the WS CDF at unsampled locations. For more details on fitting

189 the Birnbaum-Saunders kernel using the WSQ as input, the readers are referred to Houndekindo and
 190 Ouarda [22].

191 Table 1: Parameters of LGBMQR and LGBMQR-ERA5. The same set of randomly selected
 192 parameters was tested for LGBMQR and LGBMQR-ERA5 to implement the random search.

Model parameter	Description	Range
learning_rate	Learning rate	0.02-0.1
max_depth	Maximum depth of the regression trees	3-8
feature_fraction	Fraction of covariate to use to build each tree	0.1-0.9
bagging_fraction	Fraction of data to sample to build each tree	0.1-0.9
extra_trees	Use of extremely randomized trees [40]	True, False
lambda_l2	L2 regularization	0-1000
lambda_l1	L1 regularization	0-1000
num_leaves	maximum number of leaves per regression tree	2-50
max_bin	max number of bins for the discretization of the covariates	50-400
min_data_in_leaf	minimal amount of data in one leaf	100-20000
num_boost_round	Number of trees to build (boosting iteration)	90-400
n_features	Number of features to include in the model	5-30

193

194

195 3.2. Prediction of wind speed time series at unsampled locations

196 It is proposed to adapt the QM [41] procedure to predict WSTS at unsampled locations using the
 197 following general formula [26]:

$$198 \hat{w}_t(s_0) = \hat{F}_{S_0}^{-1}[\hat{F}(w_t)] \quad (3)$$

199 where: $\hat{w}_t(s_0)$ is the estimated WS at time t and unsampled location s_0 . \hat{F}_{s_0} is the estimated WS
 200 CDF at the unsampled location s_0 , and $\hat{F}_{s_0}^{-1}$ is its inverse. $\hat{F}(w_t)$ is the estimated wind speed non-
 201 exceedance probabilities (WSNEP) at time t . The methodology to estimate \hat{F}_{s_0} at any unsampled
 202 location in the region was described in section 3.1. For the estimation of $\hat{F}(w_t)$ two approaches
 203 have been put forward in previous studies:

- 204 1. Some authors [25, 42] proposed using information from nearby locations to estimate $\hat{F}(w_t)$
 205 at any unsampled location. This technique assumes that observed non-exceedance
 206 probabilities (or exceedance probabilities) between nearby locations are correlated. Thus, a
 207 spatial interpolation method could be applied to estimate the WSNEP at unsampled
 208 locations. The Inverse Distance Weighting (IDW) was used to interpolate the WSNEP. The
 209 method was named the flow duration curve and the temperature duration curve for
 210 streamflow and temperature modelling. Following this nomenclature, the technique was
 211 referred to as the Wind Duration Curve (WDC) in the context of WS modelling.
- 212 2. Jung and Schindler [43] derived the non-exceedance probabilities $\hat{F}(w_t)$ directly from a
 213 reanalysis dataset, thereby performing bias correction. This approach will be applied with
 214 ERA5 and named quantile mapping bias correction of ERA5 (QM-ERA5) in the following
 215 sections.

216

217 The Weibull plotting position was used to estimate the WSNEP from the WSTS as follows:

$$218 F_n(w_t) = i_t / n + 1 \quad (4)$$

219 where: $i_t = 1, 2, 3, \dots, n$ is the rank of the WS value observed at time t (w_t) after sorting the time
 220 series in ascending order.

221 **3.3. Spatial interpolation methods**

222 Two spatial interpolation methods were selected and evaluated to interpolate the WSTS directly.
223 The IDW technique was selected for its ease of application and set as the baseline method in the
224 study. The general formula of the IDW methods is:

225
$$\hat{w}_t(s_0) = \sum_{i=1}^k \lambda_i w_t(s_i) \tag{5}$$

226 where:

227
$$\lambda_i = \frac{d_i^{-p}}{\sum_{j=1}^k d_j^{-p}} \tag{6}$$

228 where: $w_t(s_{i=1:k})$ is the observed WS value at time t and the nearest location s_i , located at a
229 distance d_i from the target location s_0 . The parameters p and k are the exponents and the number
230 of nearest neighbours to consider. It should be noted that the IDW was used in the study to
231 interpolate observed WSTS and WSNEP (during the implementation of the WDC method). In both
232 cases, the optimal number of nearest locations and the exponent were selected based on 1) the time
233 series (TS) evaluation using the Pearson correlation coefficient between observed and estimated
234 WSTS and 2) the probability distribution (PD) evaluation by calculating the coefficient of
235 determination (R^2) between observed and estimated WSQ derived from the WSTS. The R^2 is
236 presented in equation S4 of the supporting material. The results of the models were presented for
237 each evaluation metric (TS and PD) separately.

238 The second spatial interpolation method implemented in the study was the Random Forest for
239 Spatial Interpolation model [RFSI: 44]. The model uses nearby observations and their distance from
240 the target location as covariates with a random forest regression model to interpolate at unsampled
241 locations. The general formula of the model is [44]:

$$\hat{w}_t(s_0) = f(x_1(s_0), \dots, x_m(s_0), w_t(s_1), d_1, \dots, w_t(s_k), d_k) \quad (7)$$

243 where: $x_{i=1:m}(s_0)$ are covariates available at the target location s_0 , $f(\cdot)$ is a regression function
 244 linking the covariates and the WS values at the unsampled location. A comparative analysis carried
 245 out by Sekulić, et al. [44] revealed that in real-world conditions, the RFSI model outperformed Space-
 246 time regression kriging, and the approach can scale and perform better than another spatial
 247 interpolation method based on the random forest model [45]. Furthermore, as RFSI does not require
 248 semi-variogram modelling, it is easier to implement than kriging methods with less restrictive
 249 assumptions (e.g., stationarity and linearity). In the original RFSI model, the authors used the
 250 random forest model to learn the regression function. Due to its efficiency and scalability for large
 251 datasets, the LGBM implementation of the gradient boosting algorithm was used in place of the
 252 random forest model, and the approach was renamed LGBMSI for this study. The tuned LGBMSI
 253 parameters were the same parameters presented in Table 1 of the present paper. These parameters
 254 were also tuned using a random search with 1000 iterations. As done for the QR model, the
 255 available covariates were ranked using the MRMR algorithm. The number of covariates to include
 256 in the model was treated as a parameter to be tuned during random search. Two versions of
 257 LGBMSI were tested: The version presented in equation 7 (it will be referred to as simply LGBMSI
 258 in the following sections) and a version which uses as additional covariate the WS values from the
 259 nearest ERA5 grid point to the unsampled location ($w_t(ERA5_{s_0})$). The LGBMSI model with the
 260 ERA5 covariates will be referred to as LGBMSI-ERA5 in the following sections of the paper and
 261 is presented in equation 8:

$$\hat{w}_t(s_0) = f\left(x_1(s_0), \dots, x_m(s_0), w_t(s_1), d_1, \dots, w_t(s_k), d_k, w_t(ERA5_{s_0})\right) \quad (8)$$

263 **3.4. Global Wind Atlas mean bias correction**

264 The GWA version 3 (<https://globalwindatlas.info/>) feeds the output from a mesoscale atmospheric
265 model into a microscale model to downscale the ERA5 wind data. The resulting wind data has a
266 spatial resolution of 250m and accounts for the effect of the local topography and surface
267 roughness. Several studies used the GWA to bias-correct reanalysis WS data [46-48]. The
268 procedure involves applying a scaling factor to the reanalysis WS data to ensure that their mean
269 matches the mean WS from GWA. The scaling factor is computed as the ratio between the mean
270 WS from GWA and the reanalysis during the overlapping period of both datasets. The mean WS
271 from GWA and ERA5 at 10 m estimated for the period between 2008 and 2017 were used to
272 calculate the scaling factor. Nearest neighbour interpolation was used to interpolate the GWA data
273 at locations of interest. The bias-corrected ERA5 using GWA will be referred to as GWA-ERA5
274 in the remainder of the paper.

275 **3.5. Validation**

276 The model validation strategy adopted in this study is aligned with the modelling procedure's
277 primary task, which consisted of predicting WSTS at unsampled locations. During the models
278 tuning, random k-fold cross-validation across the training locations was implemented to estimate
279 the model's performance for prediction at (pseudo) unsampled locations. In 5-fold cross-validation,
280 the training locations are randomly split into five groups. Training is carried out with the data of 4
281 groups, and the model is evaluated on the remaining group. This procedure was repeated five times,
282 using each group once as the validation set. The final evaluation of the selected model was
283 performed on a group of locations (test samples) held back and comprising approximately 30% (97
284 locations) of the available locations (303) for the entire study.

285 The estimated WSTS at locations of the test samples were evaluated according to the following
286 criteria:

287 1. Time series evaluation: The Pearson correlation (PC), mean absolute error (MAE) and root-
288 mean-squared error (RMSE) were calculated between observed and estimated WSTS. The PC,
289 MAE and RMSE are presented in equations S1, S2, and S3 of the supporting material,
290 respectively.

291 2. Probability distribution evaluation: Two approaches were used to evaluate the probability
292 distribution of the estimate WSTS. First, quantiles with non-exceedance probabilities between
293 10% and 90% and a spacing of 10% were calculated from the WSTS using equation S8 in the
294 supporting material. The R^2 , MAE and RMSE were used to compare the observed and
295 estimated WSQ. Lastly, the Overlap percentage [OP: 49] was used to assess the overlap between
296 estimated and observed empirical probability distribution function (PDF). The OP is presented
297 in equation S6 of the supporting material. For a review of criteria used for the selection of PD
298 for WS data the reader is referred to [50].

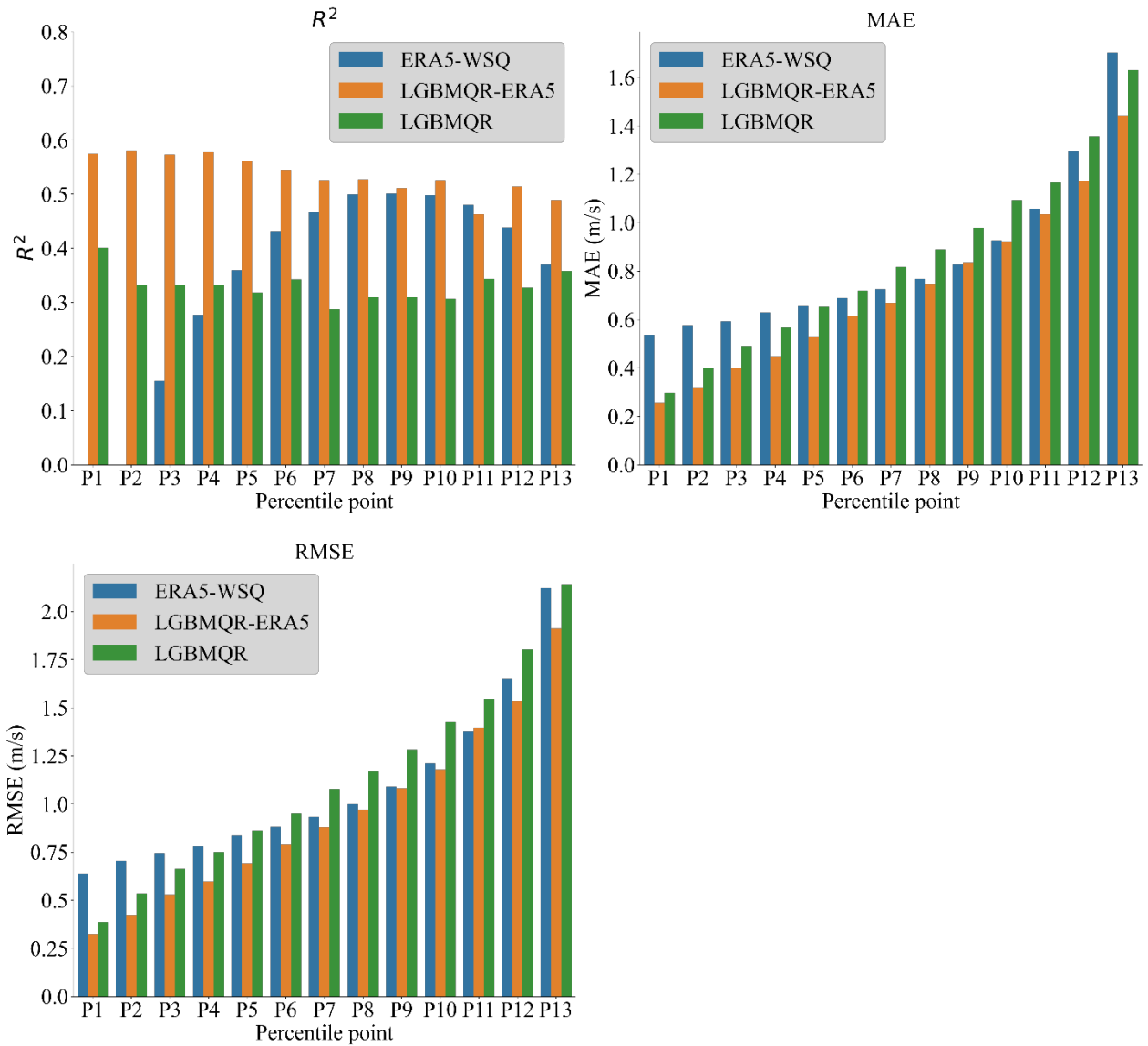
299 3. Interannual variability (IAV) evaluation: The robust coefficient of variation [RCov: 51] of
300 annual median WS was calculated to assess IAV. RCov serves as a robust and resistant measure
301 of variability analogue to the coefficient of variation, which lacks robustness and resistance to
302 outliers. The MAE and mean error (ME) between observed and estimated RCov were used to
303 evaluate the performance of the models in reproducing the observed IAV. The RCov and the
304 ME are presented in equations S7 and S5 of the supporting material, respectively.

305 **4. Results**

306 **4.1. Quantile regression models**

307 A thousand random combinations of the LGBM hyperparameters (Table 1) were tested with
308 LGBMQR and LGBMQR-ERA5 models. Table S2 in the supporting material shows the best
309 parameters found using a random search, including the number of selected covariates. Figure 2
310 illustrates the R^2 , MAE, and RMSE between estimated and observed WSQ from the test samples.
311 For reference, the same metrics between ERA5-WSQ and observed WSQ are also presented in
312 Figure 2. Figure 3 shows boxplots of the metrics calculated over the different percentile points (P1
313 – P13) at each test sample. The Wilcoxon signed-rank test was used to test the statistical
314 significance of these metrics between pairs of models (the test P-values are shown in Table S3 of
315 the supporting material). The P-values associated with LGBMQR-ERA5 are all less than 0.05, and
316 the P-values between LGBMQR and ERA5-WSQ are more significant than 0.05. LGBMQR and
317 ERA5-WSQ had significantly lower median R^2 and higher median MAE and RMSE than
318 LGBMQR-ERA5. LGBMQR underperformed compared to ERA5-WSQ, but the difference
319 between the methods was not significant according to the Wilcoxon signed-rank test. LGBMQR
320 outperformed ERA5-WSQ for WSQ with low exceedance probabilities (P1, P2, P3) while ERA5-
321 WSQ were more accurate in the middle and upper tail of the distributions. It is evident from these
322 results that the inclusion of the ERA5-WSQ improves the QR model performance; thus, WSQs
323 from LGBMQR-ERA5 were used in subsequent analyses of the study.

324

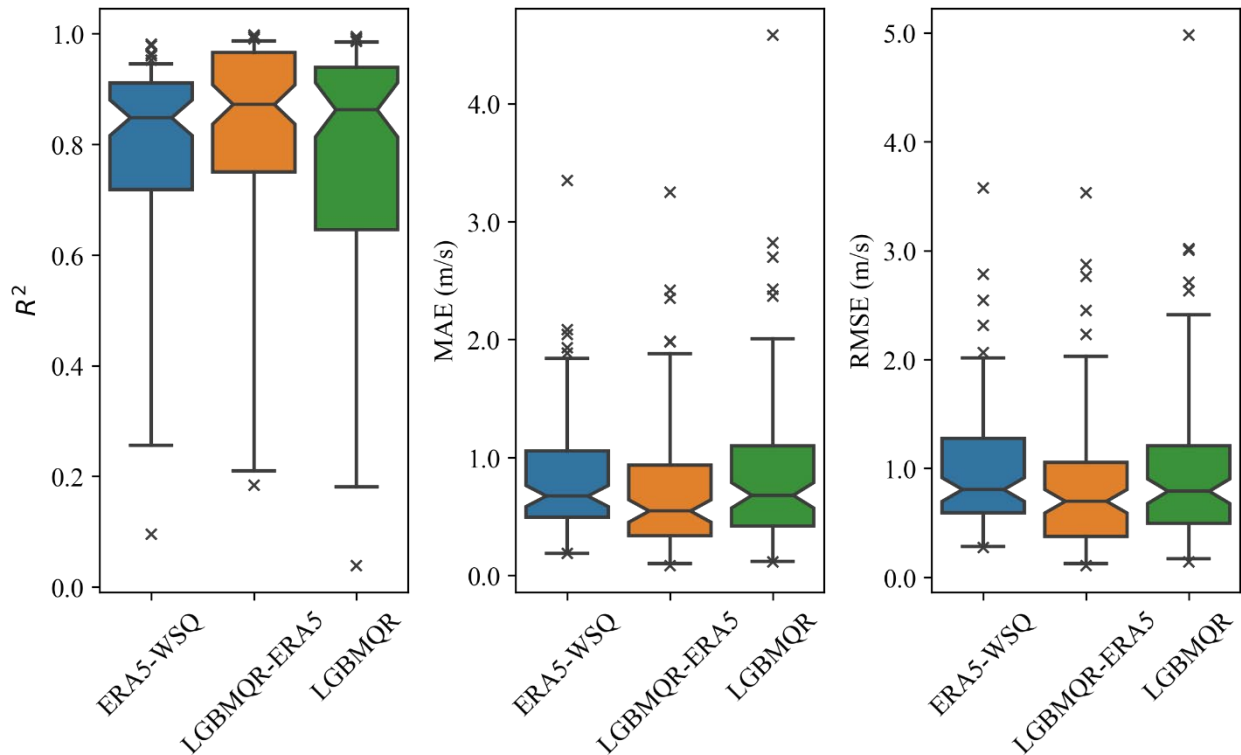


325

326 Figure 2: Results of the R^2 , RMSE and MAE between estimated and predicted WSQ at various

327 percentile points (P1-P13). The metrics were calculated across the test samples for each

328 percentile point.



329

330 Figure 3: Result of the R^2 , MAE and RMSE between observed and estimated WSQ. The metrics
331 were calculated across the percentile points (P1-P13) at each location in the test samples.

332 4.2. Inverse distance weighting parameters

333 Table 2 shows the optimal parameters (p and k) for IDW based on the TS and PD evaluation. The
334 selection of the best parameters was performed with the training set. The optimal k and p was
335 contingent upon the evaluation criteria. For the interpolation of WSNEP, the optimal number of
336 nearest neighbours (optimal $k = 1$) based on the PD evaluation is equated to the nearest neighbour
337 interpolation. Generally, it was observed smaller values of k were optimal for the PD criteria. The
338 results of the different evaluation criteria will be presented and discussed separately in the
339 following section.

340

341 Table 2: Optimal parameters of the IDW for WSTS and WSNEP interpolation

Interpolated variable	Evaluation criteria	Optimal k	Optimal p	Abbreviation used for the model herein
WSTS	PD	6	0.3	IDW-PD
	TS	11	1.7	IDW-TS
WSNEP	PD	1	-	WDC-PD
	TS	9	1	WDC-TS

342

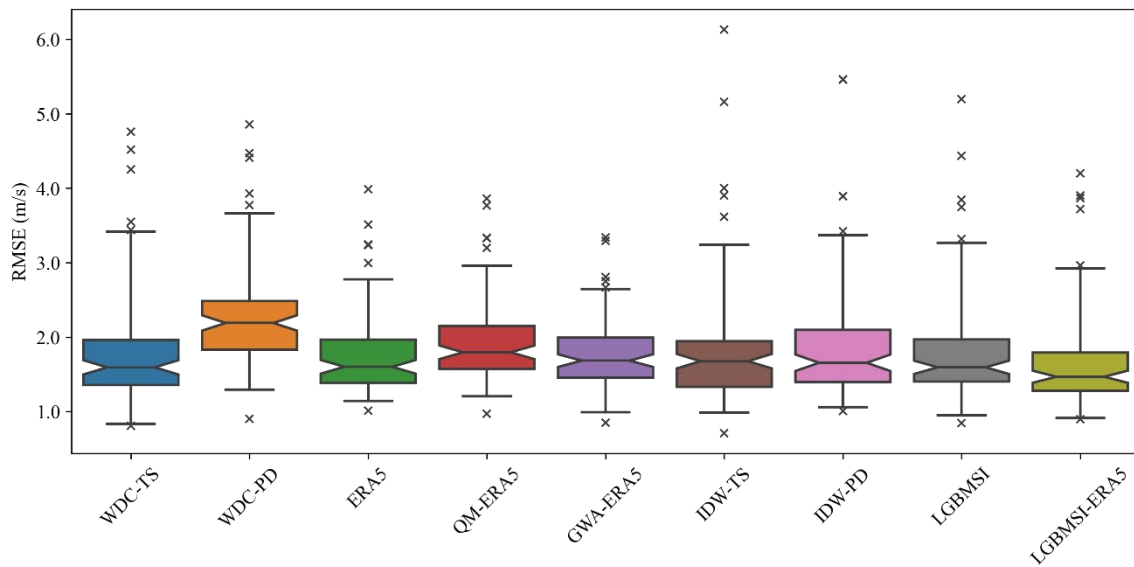
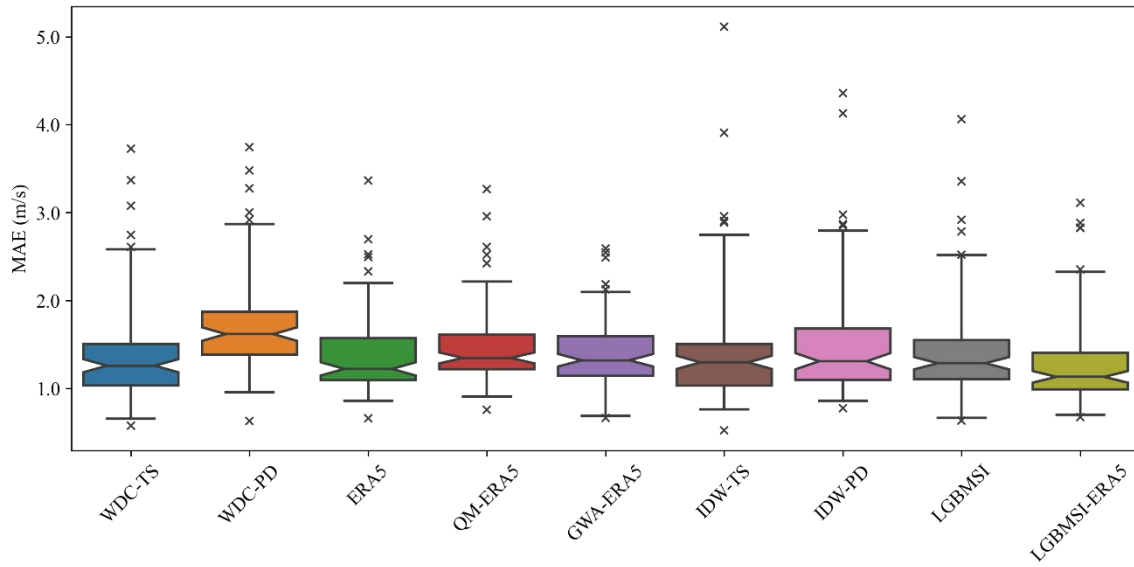
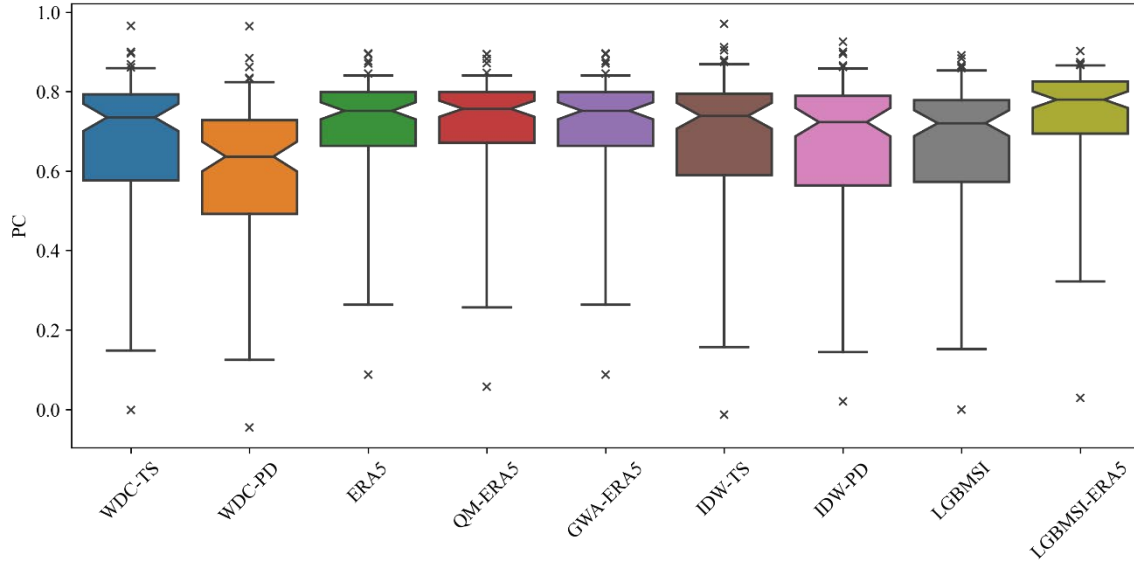
343 4.3. Time series evaluation

344 Figure 4 shows a boxplot of the PC, MAE and RMSE between observed and estimated WSTS from
 345 the test samples, while the median values of the metrics are given in Table 3. LGBMSI-ERA5 had
 346 the highest median PC alongside the lowest median MAE and RMSE. In contrast, WDC-PD had
 347 the lowest median PC and the highest median MAE and RMSE. WDC-TS performed better than
 348 WDC-PD, with performances comparable to the IDW model. The ERA5 WSTS showed a
 349 relatively high median PC and methods directly exploiting this dataset (GWA-ERA5, LGBMSI-
 350 ERA5, QM-ERA5) maintained a higher median PC with less variability in the distribution of the
 351 metric in comparison to methods solely using observations from nearby locations (WDC-PD,
 352 WDC-TS, LGBMSI, IDW-PD and IDW-TS). Despite QM-ERA5 showing a relatively high median
 353 PC, it also had a high median MAE and RMSE. Table S4 in the supporting material gives the P-
 354 value of the Wilcoxon signed-rank test between pairs of models for the different evaluation metrics.
 355 From the results of the Wilcoxon test, it was found that LGBMSI-ERA5 was the only method with
 356 an MAE and RMSE statistically inferior to IDW-TS. For the time series evaluation criteria,
 357 LGBMSI-ERA5 was the best-performing method, WDC-PD was the least effective method and
 358 most other models had performances comparable (in a statistical sense) to IDW-TS.

359 Table 3: Median PC, MAE and RMSE between observed and estimated WSTS

Model	PC	MAE (m/s)	RMSE (m/s)
WDC-TS	0.73	1.26	1.59
WDC-PD	0.64	1.62	2.19
ERA5	0.75	1.22	1.60
QM-ERA5	0.76	1.34	1.80
GWA-ERA5	0.75	1.32	1.68
IDW-TS	0.74	1.29	1.68
IDW-PD	0.72	1.31	1.66
LGBMSI	0.72	1.28	1.59
LGBMSI-ERA5	0.78	1.13	1.47

360



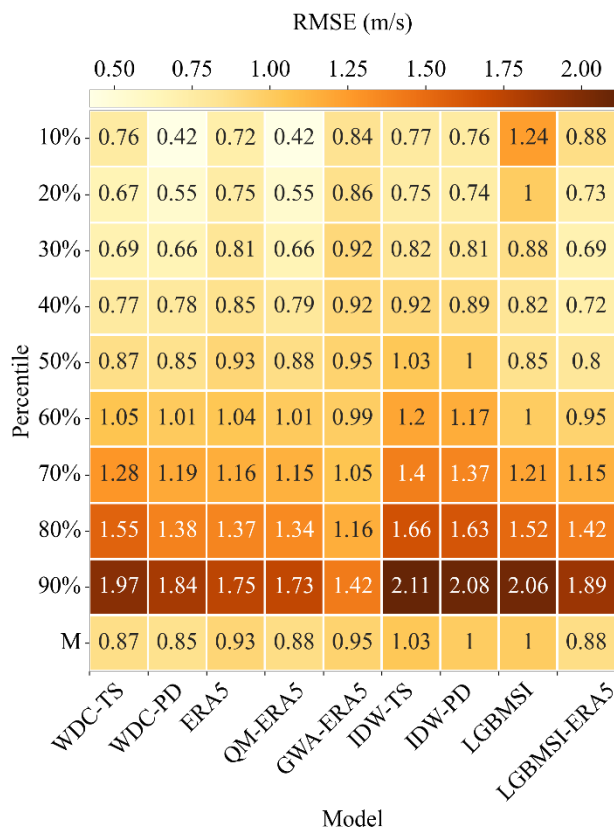
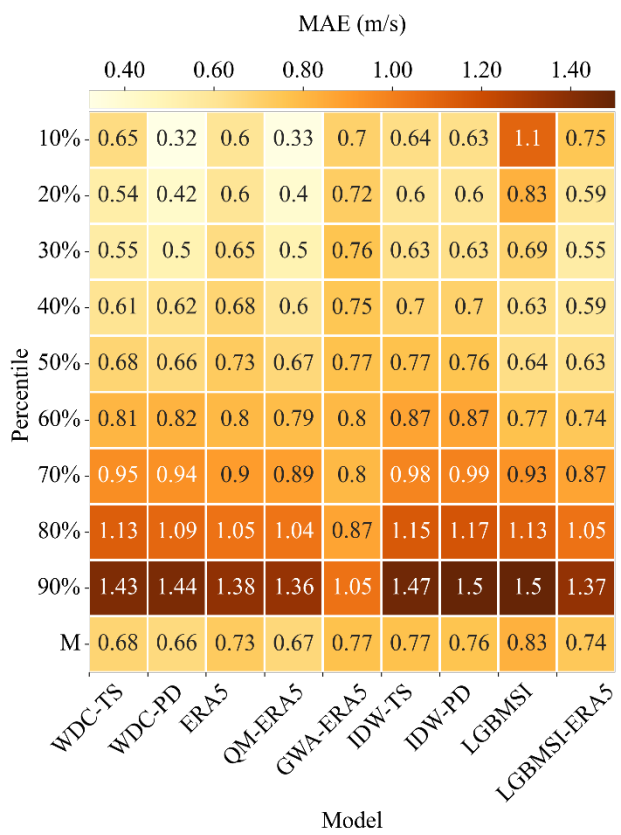
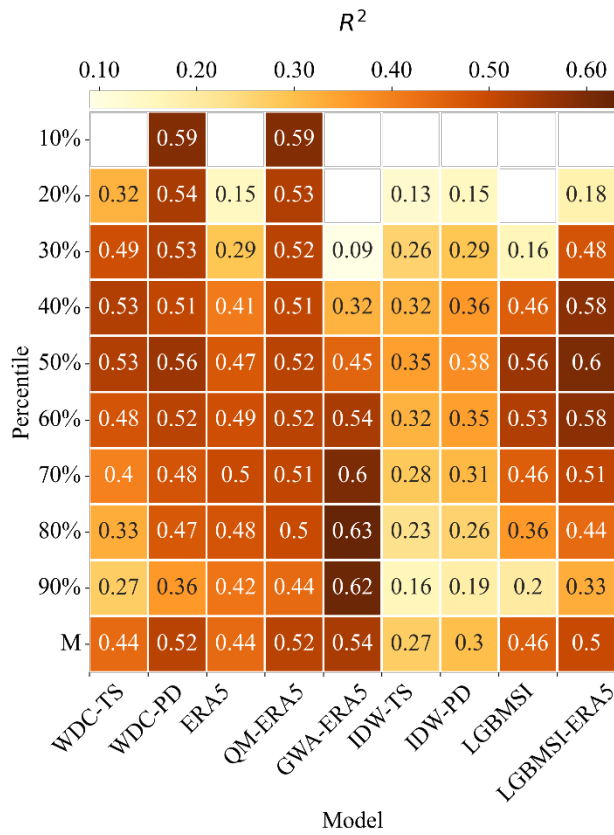
362 Figure 4: Result of the PC, MAE and RMSE between observed and estimated WSTS.

363 **4.4. Probability distribution evaluation**

364 Figure 5 shows matrices detailing the R^2 , MAE and RMSE calculated between the estimated and
365 observed WSQ across various percentile points. The last row of these matrices (labelled M)
366 presents the median value (calculated over the different percentile points). QM-ERA5 and WDC-
367 PD were the top-performing methods overall, mainly due to their relatively strong performance in
368 estimating WSQ in the lower and middle tail of the distribution. Both LGBMSI-ERA5 and
369 LGBMSI performed relatively well in the middle of the distribution but were less effective in
370 estimating WSQ in the lower tail. GWA-ERA5 was the best method for estimating WSQ in the
371 upper tail of the distribution, yet it performed poorly for low exceedance probabilities WSQ. The
372 IDW methods demonstrated an overall lack of effectiveness in estimating WSQ across the
373 distribution.

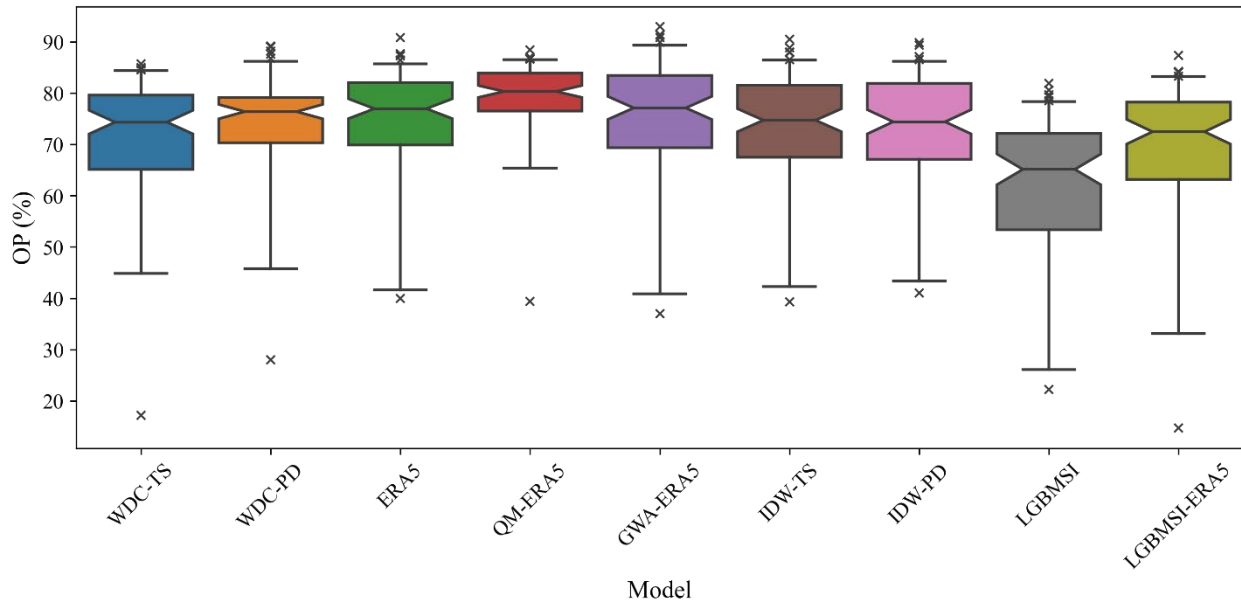
374 The OP metric measured the overlap between the empirical PDF computed from the estimated and
375 observed WSTS. Figure 6 presents boxplots illustrating the distribution of the OP metric. QM-
376 ERA5 had the highest median OP at 80%, followed by ERA5 at 77%, GWA-ERA5 at 77% and
377 WDC-PD at 76%. Also, QM-ERA5 and WDC-PD displayed less spread in the distribution of the
378 metric compared to ERA5 and GWA-ERA5. LGBMSI and LGBMSI-ERA5 had the lowest median
379 OP values at 65% and 72% respectively. The statistical significance of the results was tested with
380 the Wilcoxon signed-ranked test between pairs of models (Table S5 of the supporting material).
381 The P-values associated with QM-ERA5, LGBMSI, LGBMSI-ERA5 and WDC-TS were always
382 small (ex.: less than 0.05). The differences between IDW, ERA5, GWA-ERA5 and WDC-PD were
383 not statistically significant (P-values greater than 0.05) for the OP metric. Overall, QM-ERA5 was
384 the top performer for the OP metric, followed by (listed in no particular order) IDW, ERA5, GWA-

385 ERA5 and WDC-PD. WDC-TS performed slightly better than LGBMSI-ERA5, while LGBMSI
386 was the least effective method.



388 Figure 5: Result of the R^2 , MAE and RMSE between observed and estimated WSQ. The last row
 389 of the matrices gives the median of the metric calculated across the different percentile points.
 390 Values of R^2 less than 0 were omitted from the matrices.

391



392

393 Figure 6: Boxplots of OP metrics calculated between observed and estimated empirical PDF.

394

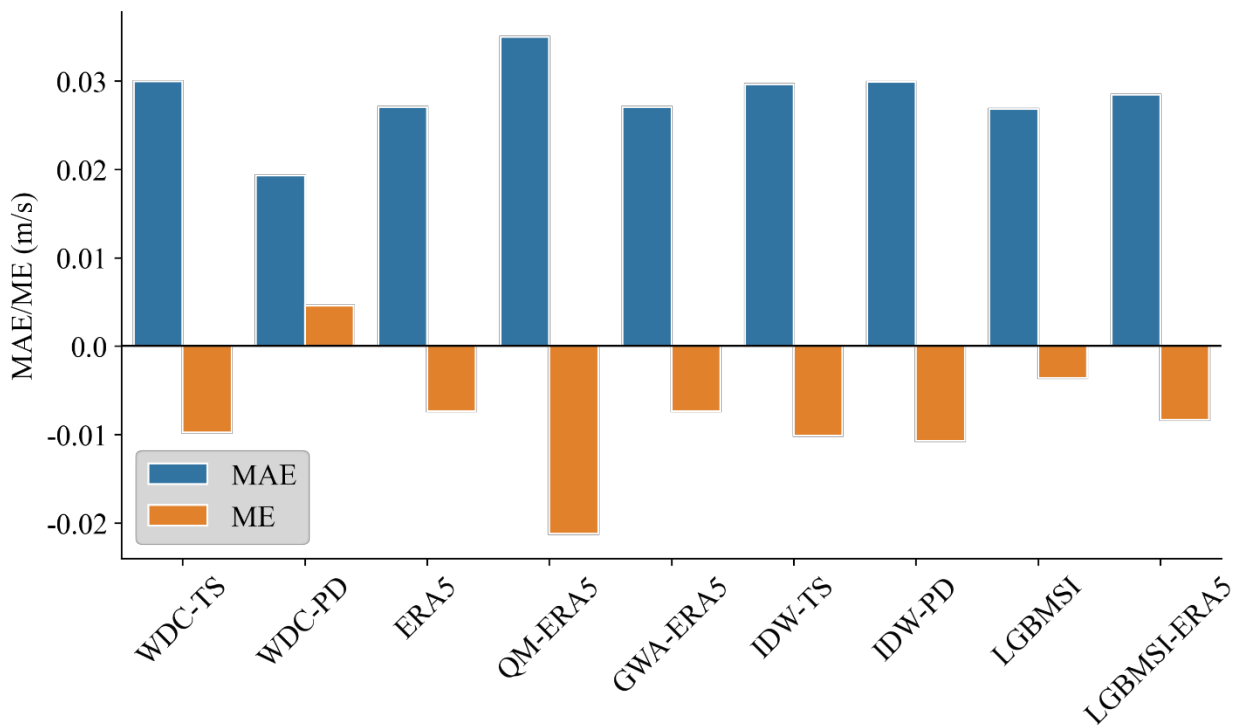
395 4.5. Interannual variability evaluation

396 The IAV assesses the fluctuation of wind speed across multiple years. Studies have indicated that
 397 wind speed exhibits IAV in many parts of the world [52-54]. The IAV has been linked to
 398 atmospheric teleconnections [54-56] such as the El Niño-Southern Oscillation and the North
 399 Atlantic Oscillation. Accurately assessing the IAV of wind resources is essential for providing
 400 adequate information for the long term planning of wind energy projects [57]. Some attempts have

401 been made to develop teleconnection-based long term forecasting models for wind speed that use
402 low frequency atmospheric circulation patterns as covariates [58].

403 Figure 7 presents a bar plot representing the MAE and ME between observed and estimated RCov
404 of median annual WS.

405 WDC-PD gave the smaller MAE at 2%, while the other methods gave a slightly higher MAE at
406 3%. Notably, WDC-PD was the only method that overestimated, on average, the IAV (positive
407 ME). The other methods showed, on average, an underestimation of the IAV (negative ME). There
408 was no substantial difference in the performance among the various methods based on the IAV.



409
410 Figure 7: Result of MAE and ME between observed and estimated RCov of median annual WS.

411

412 5. Discussion

413 The study results indicate that no single method excelled according to all evaluation criteria,
414 suggesting potential for improvement through combining specific methods. For instance, it was
415 found that WSQ derived from the GWA-ERA5 time series was the most accurate in the upper tail
416 of the distribution. Conversely, in the lower tail, WSQs from GWA-ERA5 were inaccurate
417 compared to QM-ERA5 and WDC-PD. Based on these outcomes, future studies are recommended
418 to explore using the mean WS from the GWA dataset as covariates of the QR model to potentially
419 improve the estimation of the conditional WSQ in the upper tail of the distribution, thus enhancing
420 the performance of QM-ERA5 and WDC-PD.

421 LGBMSI-ERA5 was the top performer based on the time series evaluation. In the case of the
422 evaluation based on the PD, QM-ERA5 was the top performer. Generally, more complex methods
423 yielded superior performances compared to the baseline model (IDW), suggesting some benefits
424 in implementing complex methods in part due to their ability to integrate various WS covariates.
425 The ERA5 dataset was a valuable covariate. For instance, ERA5 WSTs are well correlated with
426 ground measurements, and this correlation could be improved significantly (in a statistical sense)
427 by using the dataset as a covariate with LGBMSI. Also, ERA-WSQ significantly improved (in a
428 statistical sense) the performance of the QR model. It should be noted that other covariates used as
429 input of the QR models demonstrated a higher ability to predict WSQ in the distribution's lower
430 tail than ERA5-WSQ, which seemed less accurate in the lower tail.

431 QM-ERA5 improved the performance of ERA5 in most cases. The approach is relatively easy to
432 implement and relies on a reasonable estimation of the WSD at unsampled locations. One reason
433 that could explain the improved performance of QM-ERA5 is its higher accuracy in the lower tail
434 of the distribution compared to ERA5 wind data. It was also revealed that the WDC method was

435 competitive. However, the approach is sensitive to the evaluation criteria used to select the optimal
436 parameters of the IDW for interpolating the WSNEP. Different evaluation criteria lead to different
437 optimal parameters, which leads, in turn, to different performances during evaluation. For instance,
438 WDC-PD performed relatively well based on the evaluation of PD, while it performed poorly based
439 on the TS evaluation. In contrast, WDC-TS performed relatively well based on the TS evaluation
440 and was less effective than WDC-PD based on the evaluation of the PD. In future studies, it is
441 recommended that different methods to interpolate the WSNEP are explored to improve the
442 performance of the WDC method. For instance, a more complex interpolation method, such as
443 RFSI, could be applied to interpolate the WSNEP.

444 In this study, LGBM with the pinball lost function was used as the QR model (LGBMQR). Other
445 quantile regression models could be viable alternatives, such as quantile regression forests [59] and
446 quantile regression neural networks [60]. LGBMQR was adopted because it is efficient during
447 training, and in general, gradient-boosting models have demonstrated superior performance on
448 tabular data [61]. In upcoming research, a comparative analysis can be performed to evaluate the
449 performance of different QR models for conditional WSQ mapping.

450 For practical reasons, the analysis in the present study was carried out at the World Meteorological
451 Organization (WMO) recommended wind speed measurement height of 10 m. Modern wind
452 turbines operate at hub heights of 100 m and beyond. It would be ideal to assess the wind resource
453 directly at these hub heights. However, there is lack of extensive wind speed time series data at
454 these heights and even when available, accessing such data from private wind farm operators can
455 pose challenges. To account for this disparity, vertical wind profile equations such as the
456 logarithmic and power law are employed to extrapolate the estimated wind speed from 10 m to the
457 hub height [15, 21]. This procedure inevitably introduces additional uncertainty to the estimated

458 wind resource. Future research should be conducted to evaluate and quantify this layer of
459 uncertainty more comprehensively.

460 **6. Conclusions and future research**

461 This study conducted a comprehensive evaluation of various approaches for the prediction of wind
462 speed time series at unsampled locations. It was found that no single method consistently
463 outperformed the other methods according to all evaluation criteria. However, complex methods
464 that include various covariates were more effective than the baseline method. Mainly, two
465 approaches (QM-ERA5 and LGBMSI-ERA5) applied to bias-correct ERA5 wind speed data
466 seemed promising and showed improved results compared to the most common ERA5 bias
467 correction method (GWA-ERA5). It should be noted that both methods are more complex and
468 computationally demanding than GWA-ERA5. However, LGBMSI-ERA5 significantly improved
469 the accuracy of the ERA5 data when evaluating the time series correlations, while QM-ERA5
470 significantly improved the overlap percentage between the observed and estimated empirical PDF.
471 In future studies, it is recommended that the performance of LGBMSI-ERA5 and QM-ERA5 be
472 explored further in different regions with different wind regimes. Another promising research route
473 is the potential to combine different approaches to produce a more accurate model across multiple
474 evaluation criteria.

475 Also, with the QR model, there is a potential to account for the non-stationarity of the WSD by
476 using related covariates. For instance, Ouarda and Charron [54] found that the North-Atlantic
477 Oscillation and the Pacific North American indices of atmospheric circulation were good predictors
478 of the IAV of WS in the province of Québec, Canada. In future studies, these climate indices can
479 be used as covariates with a QR model in the region to map conditional WSQ that accounts for the

480 resource's IAV. This analysis could lead to a better evaluation of the wind resources at unsampled
481 locations, thus reducing the risk associated with future projects.

482 The comprehensive evaluation provided in the present study aims to assist practitioners in choosing
483 the most suitable methodologies for their specific projects. Furthermore, it is anticipated that this
484 research will inspire future studies to systematically evaluate various approaches for predicting
485 wind speed time series at unsampled locations. This will foster in the long run a better
486 understanding of the strengths and limitations of these approaches and encourage their refinement
487 and the development of more robust techniques for the prediction of wind speed time series at
488 unsampled locations.

489 **Acknowledgments**

490 The authors thank the Natural Sciences and Engineering Research Council of Canada (NSERC)
491 and the Canada Research Chair Program for funding this research. The authors would like to extend
492 their gratitude to the Editor, and two anonymous reviewers for their comments and suggestions,
493 which significantly improved the quality of the paper.

494 **Data availability statement**

495 The data used in the study are available from public source. The measure wind speed data were
496 acquired from Environment and climate change Canada
497 (https://collaboration.cmc.ec.gc.ca/cmc/climate/Get_More_Data_Plus_de_donnees/), the digital
498 elevation model from the Japan Aerospace Exploration Agency
499 (https://www.eorc.jaxa.jp/ALOS/en/dataset/aw3d30/aw3d30_e.htm), the ERA5 10 m wind
500 components from ECMWF (<https://doi.org/10.24381/cds.adbb2d47>), the land use map from

501 Natural Resources Canada (<https://doi.org/10.3390/rs9111098>), the GWA mean wind speed from
502 the Technical University of Denmark (<https://globalwindatlas.info/en>).

503 **References**

- 504 [1] A. Cherp, V. Vinichenko, J. Tosun, J.A. Gordon, J. Jewell. National growth dynamics of wind
505 and solar power compared to the growth required for global climate targets. *Nature Energy*. 6
506 (2021) 742-54. [10.1038/s41560-021-00863-0](https://doi.org/10.1038/s41560-021-00863-0)
- 507 [2] R. Wiser, J. Rand, J. Seel, P. Beiter, E. Baker, E. Lantz, P. Gilman. Expert elicitation survey
508 predicts 37% to 49% declines in wind energy costs by 2050. *Nature Energy*. 6 (2021) 555-65.
509 [10.1038/s41560-021-00810-z](https://doi.org/10.1038/s41560-021-00810-z)
- 510 [3] G. Ren, J. Liu, J. Wan, Y. Guo, D. Yu. Overview of wind power intermittency: Impacts,
511 measurements, and mitigation solutions. *Applied Energy*. 204 (2017) 47-65.
512 <https://doi.org/10.1016/j.apenergy.2017.06.098>
- 513 [4] J.C.Y. Lee, M.J. Fields. An overview of wind-energy-production prediction bias, losses, and
514 uncertainties. *Wind Energ Sci*. 6 (2021) 311-65. [10.5194/wes-6-311-2021](https://doi.org/10.5194/wes-6-311-2021)
- 515 [5] A. Lopez, T. Mai, E. Lantz, D. Harrison-Atlas, T. Williams, G. Maclaurin. Land use and
516 turbine technology influences on wind potential in the United States. *Energy*. 223 (2021) 120044.
517 <https://doi.org/10.1016/j.energy.2021.120044>
- 518 [6] D. Niermann, M. Borsche, Andrea K. Kaiser-Weiss, F. Kaspar. Evaluating renewable-energy-
519 relevant parameters of COSMO-REA6 by comparison with satellite data, station observations
520 and other reanalyses. *Meteorologische Zeitschrift*. 28 (2019) 347-60. [10.1127/metz/2019/0945](https://doi.org/10.1127/metz/2019/0945)
- 521 [7] H. Hersbach, B. Bell, P. Berrisford, S. Hirahara, A. Horányi, J. Muñoz-Sabater, et al. The
522 ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological Society*. 146 (2020)
523 1999-2049. <https://doi.org/10.1002/qj.3803>
- 524 [8] A. Molod, L. Takacs, M. Suarez, J. Bacmeister. Development of the GEOS-5 atmospheric
525 general circulation model: evolution from MERRA to MERRA2. *Geosci Model Dev*. 8 (2015)
526 1339-56. [10.5194/gmd-8-1339-2015](https://doi.org/10.5194/gmd-8-1339-2015)
- 527 [9] G. Gualtieri. Analysing the uncertainties of reanalysis data used for wind resource
528 assessment: A critical review. *Renewable and Sustainable Energy Reviews*. 167 (2022) 112741.
529 <https://doi.org/10.1016/j.rser.2022.112741>
- 530 [10] K. Gruber, P. Regner, S. Wehrle, M. Zeyringer, J. Schmidt. Towards global validation of
531 wind power simulations: A multi-country assessment of wind power simulation from MERRA-2
532 and ERA-5 reanalyses bias-corrected with the global wind atlas. *Energy*. 238 (2022) 121520.
533 <https://doi.org/10.1016/j.energy.2021.121520>
- 534 [11] R.K. Samal. Assessment of wind energy potential using reanalysis data: A comparison with
535 mast measurements. *Journal of Cleaner Production*. 313 (2021) 127933.
536 <https://doi.org/10.1016/j.jclepro.2021.127933>
- 537 [12] A. Ayik, N. Ijumba, C. Kabiri, P. Goffin. Preliminary wind resource assessment in South
538 Sudan using reanalysis data and statistical methods. *Renewable and Sustainable Energy Reviews*.
539 138 (2021) 110621. <https://doi.org/10.1016/j.rser.2020.110621>
- 540 [13] J. Ramon, L. Lledó, V. Torralba, A. Soret, F.J. Doblas-Reyes. What global reanalysis best
541 represents near-surface winds? *Quarterly Journal of the Royal Meteorological Society*. 145
542 (2019) 3236-51. <https://doi.org/10.1002/qj.3616>

543 [14] G. Gualtieri. Reliability of ERA5 Reanalysis Data for Wind Resource Assessment: A
544 Comparison against Tall Towers. *Energies*2021.

545 [15] K. Gruber, C. Klöckl, P. Regner, J. Baumgartner, J. Schmidt. Assessing the Global Wind
546 Atlas and local measurements for bias correction of wind power generation simulated from
547 MERRA-2 in Brazil. *Energy*. 189 (2019) 116212. <https://doi.org/10.1016/j.energy.2019.116212>

548 [16] A.-J. Collados-Lara, L. Baena-Ruiz, D. Pulido-Velazquez, E. Pardo-Igúzquiza. Data-driven
549 mapping of hourly wind speed and its potential energy resources: A sensitivity analysis.
550 *Renewable Energy*. 199 (2022) 87-102. <https://doi.org/10.1016/j.renene.2022.08.109>

551 [17] M. Cellura, G. Cirrincione, A. Marvuglia, A. Miraoui. Wind speed spatial estimation for
552 energy planning in Sicily: A neural kriging application. *Renewable Energy*. 33 (2008) 1251-66.
553 <https://doi.org/10.1016/j.renene.2007.08.013>

554 [18] F. González-Longatt, H. Medina, J. Serrano González. Spatial interpolation and orographic
555 correction to estimate wind energy resource in Venezuela. *Renewable and Sustainable Energy*
556 *Reviews*. 48 (2015) 1-16. <https://doi.org/10.1016/j.rser.2015.03.042>

557 [19] F. Veronesi, S. Grassi, M. Raubal. Statistical learning approach for wind resource
558 assessment. *Renewable and Sustainable Energy Reviews*. 56 (2016) 836-50.
559 <https://doi.org/10.1016/j.rser.2015.11.099>

560 [20] C. Jung. High Spatial Resolution Simulation of Annual Wind Energy Yield Using Near-
561 Surface Wind Speed Time Series. *Energies*. 9 (2016) 344. doi:10.3390/en9050344

562 [21] C. Jung, D. Schindler. Integration of small-scale surface properties in a new high resolution
563 global wind speed model. *Energy Conversion and Management*. 210 (2020) 112733.
564 <https://doi.org/10.1016/j.enconman.2020.112733>

565 [22] F. Houndekindo, T.B.M.J. Ouarda. A non-parametric approach for wind speed distribution
566 mapping. *Energy Conversion and Management*. 296 (2023) 117672.
567 <https://doi.org/10.1016/j.enconman.2023.117672>

568 [23] A. Castellarin, G. Botter, D.A. Hughes, S. Liu, T.B.M.J. Ouarda, J. Parajka, et al. Prediction
569 of flow duration curves in ungauged basins. in: G. Blöschl, H. Savenije, M. Sivapalan, A.
570 Viglione, T. Wagener, (Eds.), *Runoff Prediction in Ungauged Basins: Synthesis across Processes,*
571 *Places and Scales*. Cambridge University Press, Cambridge, 2013. pp. 135-62.

572 [24] A.I. Requena, T.B.M.J. Ouarda, F. Chebana. Flood Frequency Analysis at Ungauged Sites
573 Based on Regionally Estimated Streamflows. *Journal of Hydrometeorology*. 18 (2017) 2521-39.
574 10.1175/JHM-D-16-0143.1

575 [25] T.B.M.J. Ouarda, C. Charron, A. St-Hilaire. Regional estimation of river water temperature
576 at ungauged locations. *Journal of Hydrology X*. (2022). 10.1016/j.hydroa.2022.100133

577 [26] A.J. Cannon, S.R. Sobie, T.Q. Murdock. Bias Correction of GCM Precipitation by Quantile
578 Mapping: How Well Do Methods Preserve Changes in Quantiles and Extremes? *Journal of*
579 *Climate*. 28 (2015) 6938-59. <https://doi.org/10.1175/JCLI-D-14-00754.1>

580 [27] M.A. Ben Alaya, F. Chebana, T.B.M.J. Ouarda. Multisite and multivariable statistical
581 downscaling using a Gaussian copula quantile regression model. *Climate Dynamics*. 47 (2016)
582 1383-97. 10.1007/s00382-015-2908-3

583 [28] T. Tadono, H. Ishida, F. Oda, S. Naito, K. Minakawa, H. Iwamoto. Precise Global DEM
584 Generation by ALOS PRISM. *ISPRS Annals of Photogrammetry, Remote Sensing and Spatial*
585 *Information Sciences*. II4 (2014) 71-6. 10.5194/isprsannals-II-4-71-2014

586 [29] R. Latifovic, D. Pouliot, I. Olthof. Circa 2010 Land Cover of Canada: Local Optimization
587 Methodology and Product Development. *Remote Sensing*. 9 (2017) 1098.

- 588 [30] Y. He, H. Li. Probability density forecasting of wind power using quantile regression neural
589 network and kernel density estimation. *Energy Conversion and Management*. 164 (2018) 374-84.
590 <https://doi.org/10.1016/j.enconman.2018.03.010>
- 591 [31] D. Ouali, F. Chebana, T.B.M.J. Ouarda. Quantile Regression in Regional Frequency
592 Analysis: A Better Exploitation of the Available Information. *Journal of Hydrometeorology*. 17
593 (2016). 10.1175/JHM-D-15-0187.1
- 594 [32] G. Ke, Q. Meng, T. Finley, T. Wang, W. Chen, W. Ma, et al. LightGBM: A Highly Efficient
595 Gradient Boosting Decision Tree. *Neural Information Processing Systems*2017.
- 596 [33] J. Fan, X. Ma, L. Wu, F. Zhang, X. Yu, W. Zeng. Light Gradient Boosting Machine: An
597 efficient soft computing model for estimating daily reference evapotranspiration with local and
598 external meteorological data. *Agricultural Water Management*. 225 (2019) 105758.
599 <https://doi.org/10.1016/j.agwat.2019.105758>
- 600 [34] J. Park, J. Moon, S. Jung, E. Hwang. Multistep-Ahead Solar Radiation Forecasting Scheme
601 Based on the Light Gradient Boosting Machine: A Case Study of Jeju Island. *Remote Sensing*. 12
602 (2020) 2271.
- 603 [35] E. Genov, C.D. Cauwer, G.V. Kriekinge, T. Coosemans, M. Messagie. Forecasting
604 flexibility of charging of electric vehicles: Tree and cluster-based methods. *Applied Energy*. 353
605 (2024) 121969. <https://doi.org/10.1016/j.apenergy.2023.121969>
- 606 [36] M. Feurer, F. Hutter. Hyperparameter Optimization. in: F. Hutter, L. Kotthoff, J.
607 Vanschoren, (Eds.), *Automated Machine Learning: Methods, Systems, Challenges*. Springer
608 International Publishing, Cham, 2019. pp. 3-33.
- 609 [37] C. Ding, P. Hanchuan. Minimum redundancy feature selection from microarray gene
610 expression data. *J Bioinform Comput Biol*. 3 (2005) 185-205. 10.1142/s0219720005001004
- 611 [38] F. Houndekindo, T.B.M.J. Ouarda. Comparative study of feature selection methods for wind
612 speed estimation at ungauged locations. *Energy Conversion and Management*. 291 (2023)
613 117324. <https://doi.org/10.1016/j.enconman.2023.117324>
- 614 [39] H.A. Mombeni, B. Mansouri, M. Akhoond. Asymmetric kernels for boundary modification
615 in distribution function estimation. *REVSTAT-Statistical Journal*. 19 (2021) 463–84–84.
- 616 [40] P. Geurts, D. Ernst, L. Wehenkel. Extremely randomized trees. *Machine Learning*. 63 (2006)
617 3-42. 10.1007/s10994-006-6226-1
- 618 [41] A.W. Wood, E.P. Maurer, A. Kumar, D.P. Lettenmaier. Long-range experimental
619 hydrologic forecasting for the eastern United States. *Journal of Geophysical Research:*
620 *Atmospheres*. 107 (2002) ACL 6-1-ACL 6-15. <https://doi.org/10.1029/2001JD000659>
- 621 [42] C. Shu, T.B.M.J. Ouarda. Improved methods for daily streamflow estimates at ungauged
622 sites. *Water Resources Research*. 48 (2012). <https://doi.org/10.1029/2011WR011501>
- 623 [43] C. Jung, D. Schindler. Introducing a new wind speed complementarity model. *Energy*. 265
624 (2023) 126284. <https://doi.org/10.1016/j.energy.2022.126284>
- 625 [44] A. Sekulić, M. Kilibarda, G.B.M. Heuvelink, M. Nikolić, B. Bajat. Random Forest Spatial
626 Interpolation. *Remote Sensing*2020.
- 627 [45] T. Hengl, M. Nussbaum, M.N. Wright, G.B.M. Heuvelink, B. Gräler. Random forest as a
628 generic framework for predictive modeling of spatial and spatio-temporal variables. *PeerJ*. 6
629 (2018) e5518. 10.7717/peerj.5518
- 630 [46] J.P. Murcia, M.J. Koivisto, G. Luzia, B.T. Olsen, A.N. Hahmann, P.E. Sørensen, M. Als.
631 Validation of European-scale simulated wind speed and wind generation time series. *Applied*
632 *Energy*. 305 (2022) 117794. <https://doi.org/10.1016/j.apenergy.2021.117794>

633 [47] S.C. de Aquino Ferreira, F.L. Cyrino Oliveira, P.M. Maçaira. Validation of the
634 representativeness of wind speed time series obtained from reanalysis data for Brazilian territory.
635 Energy. 258 (2022) 124746. <https://doi.org/10.1016/j.energy.2022.124746>
636 [48] G. Luzia, M.J. Koivisto, A.N. Hahmann. Validating EURO-CORDEX climate simulations
637 for modelling European wind power generation. Renewable Energy. 217 (2023) 118989.
638 <https://doi.org/10.1016/j.renene.2023.118989>
639 [49] S.E. Perkins, A.J. Pitman, N.J. Holbrook, J. McAneney. Evaluation of the AR4 Climate
640 Models' Simulated Daily Maximum Temperature, Minimum Temperature, and Precipitation over
641 Australia Using Probability Density Functions. Journal of Climate. 20 (2007) 4356-76.
642 <https://doi.org/10.1175/JCLI4253.1>
643 [50] T.B.M.J. Ouarda, C. Charron, F. Chebana. Review of criteria for the selection of probability
644 distributions for wind speed data and introduction of the moment and L-moment ratio diagram
645 methods, with a case study. Energy Conversion and Management. 124 (2016) 247-65.
646 <https://doi.org/10.1016/j.enconman.2016.07.012>
647 [51] S. Watson. Quantifying the Variability of Wind Energy. Advances in Energy Systems2019.
648 pp. 355-68.
649 [52] C. Jung, D. Taubert, D. Schindler. The temporal variability of global wind energy – Long-
650 term trends and inter-annual variability. Energy Conversion and Management. 188 (2019) 462-
651 72. <https://doi.org/10.1016/j.enconman.2019.03.072>
652 [53] P.E. Bett, H.E. Thornton, R.T. Clark. Using the Twentieth Century Reanalysis to assess
653 climate variability for the European wind industry. Theoretical and Applied Climatology. 127
654 (2017) 61-80. 10.1007/s00704-015-1591-y
655 [54] T.B.M.J. Ouarda, C. Charron. Non-stationary statistical modelling of wind speed: A case
656 study in eastern Canada. Energy Conversion and Management. 236 (2021) 114028.
657 <https://doi.org/10.1016/j.enconman.2021.114028>
658 [55] F. Zhou, Z. Zhao, C. Azorin-Molina, X. Jia, G. Zhang, D. Chen, et al. Teleconnections
659 between large-scale oceanic-atmospheric patterns and interannual surface wind speed variability
660 across China: Regional and seasonal patterns. Science of The Total Environment. 838 (2022)
661 156023. <https://doi.org/10.1016/j.scitotenv.2022.156023>
662 [56] M.S. Naizghi, T.B.M.J. Ouarda. Teleconnections and analysis of long-term wind speed
663 variability in the UAE. International Journal of Climatology. 37 (2017) 230-48.
664 <https://doi.org/10.1002/joc.4700>
665 [57] S.C. Pryor, T.J. Shepherd, R.J. Barthelmie. Interannual variability of wind climates and wind
666 turbine annual energy production. Wind Energ Sci. 3 (2018) 651-65. 10.5194/wes-3-651-2018
667 [58] H. Woldesellasse, P.R. Marpu, T.B.M.J. Ouarda. Long-term forecasting of wind speed in the
668 UAE using nonlinear canonical correlation analysis (NLCCA). Arabian Journal of Geosciences.
669 13 (2020) 962. 10.1007/s12517-020-05981-9
670 [59] N. Meinshausen, G. Ridgeway. Quantile regression forests. Journal of machine learning
671 research. 7 (2006).
672 [60] A.J. Cannon. Quantile regression neural networks: Implementation in R and application to
673 precipitation downscaling. Computers & Geosciences. 37 (2011) 1277-84.
674 <https://doi.org/10.1016/j.cageo.2010.07.005>
675 [61] L. Grinsztajn, E. Oyallon, G. Varoquaux. Why do tree-based models still outperform deep
676 learning on typical tabular data? Advances in Neural Information Processing Systems. 35 (2022)
677 507-20.

678