

# **Spatially-explicit correction of simulated urban air temperatures using crowd-sourced data**

Oscar Brousse,<sup>a</sup> Charles Simpson,<sup>a</sup> Owain Kenway,<sup>b</sup> Alberto Martilli,<sup>c</sup> E. Scott Krayenhoff,<sup>d</sup>  
Andrea Zonato,<sup>e</sup> and Clare Heaviside,<sup>a</sup>

<sup>a</sup>*Institute of Environmental Design and Engineering, University College London*

<sup>b</sup>*Centre for Advanced Research Computing, University College London*

<sup>c</sup>*Center for Energy, Environment and Technology (CIEMAT)*

<sup>d</sup>*School of Environmental Sciences, University of Guelph*

<sup>e</sup>*Department of Civil, Environmental and Mechanical Engineering, University of Trento*

*Corresponding author: O. Brousse, o.brousse@ucl.ac.uk*

11 ABSTRACT: Urban climate model evaluation often remains limited by a lack of trusted urban  
12 weather observations. The increasing density of personal weather stations (PWS) make them  
13 a potential rich source of data for urban climate studies that address the lack of representative  
14 urban weather observations. In our study, we demonstrate that PWS data not only improve urban  
15 climate models' evaluation, but can also serve for bias-correcting their output prior to any urban  
16 climate impact studies. After simulating near-surface air temperatures over London and south-  
17 east England during the hot summer of 2018 with the Weather Research Forecast (WRF) model  
18 and its Building Effect Parameterization with the Building Energy Model (BEP-BEM) activated,  
19 we evaluated the modelled temperatures against 402 urban PWS and showcased a heterogeneous  
20 spatial distribution of the model's cool bias that was not captured using official weather stations  
21 only. This finding indicated a need for spatially-explicit urban bias corrections of air temperatures,  
22 which we performed using an innovative method using machine learning to predict the models'  
23 biases in each urban grid cell. Our technique is the first to consider that urban temperatures are  
24 heterogeneously accurate in space and that this accuracy is not linearly correlated to the urban  
25 fraction. Our results showed that the bias-correction was beneficial to bias-correct daily-minimum,  
26 -mean, and -maximum temperatures in the cities. We recommend that urban climate modellers  
27 further investigate the use of PWS for model evaluation and derive a framework for bias-correction  
28 of urban climate simulations that can serve urban climate impact studies.

29 SIGNIFICANCE STATEMENT: Urban climate simulations are subject to spatially heteroge-  
30 neous biases in urban air temperatures. Common validation methods using official weather stations  
31 do not suffice for detecting these biases. Using a dense set of personal weather stations in London  
32 we detect these biases before proposing an innovative way for correcting them with machine learn-  
33 ing techniques. We argue that any urban climate impact study should use such technique if possible  
34 and that urban climate scientists should continue investigating paths to improve our methods.

## 35 1. Introduction

36 Although decades following the 1960s have seen an increase in the body of literature on urban  
37 climates (Oke et al. 2017), the scales of applicability and the transferability of their outcomes are  
38 often limited. This can partially be attributed to the lack of observations representative of the  
39 variety of existing urban climates in cities. To address this limitation, two major solutions were  
40 proposed over the past 20 years: firstly, the development of urban surface energy balance coupled  
41 to regional climate models (e.g., Masson (2000), Martilli et al. (2002), Wouters et al. (2016)),  
42 and secondly, the increased interest towards crowd-sourced and low-cost weather sensors (e.g.,  
43 Muller et al. (2015), Chapman et al. (2017), Fenner et al. (2017), Meier et al. (2017)). After  
44 proper validation and parameterization, urban climate models (UCMs) offer an unprecedented  
45 opportunity to represent the impact of cities on a wide variety of weather variables at very high  
46 spatial and temporal resolutions. This has been further supported by the recent development of  
47 global standardized land use land cover datasets designed for urban climate studies that permit  
48 their parameterization in cities formerly deprived of these data (see the World Urban Dataset and  
49 Access Portal Tool (WUDAPT) project; Ching et al. (2018), Demuzere et al. (2022)). Likewise,  
50 after proper filtering and quality control (Napoly et al. 2018; Fenner et al. 2021), crowd-sourced  
51 personal weather sensors (PWS) permit the extension of sensing networks into urban environments  
52 that were formerly not studied despite the fact that PWS often do not meet the standards imposed  
53 by official meteorological offices for implementation of weather stations. Several studies have  
54 demonstrated their range of applications since then (e.g., Fenner et al. (2019), Venter et al. (2020),  
55 Potgieter et al. (2021), Benjamin et al. (2021), Varentsov et al. (2021), Venter et al. (2021), Brousse  
56 et al. (2022)).

57 One of the major limitations induced by the lack of official weather stations in cities is that  
58 quantifying existing uncertainties as a function of urban climate archetype is not feasible. This  
59 means that certain urban environments are poorly evaluated and hence modelled, assuming that  
60 UCMs will perform similarly under all constraints imposed by the variety of urban environments  
61 that compose a city. In face of this challenge, crowd-sourced PWS could improve the evaluation  
62 of UCMs, as Hammerberg et al. (2018) demonstrated over Vienna. But the potential of PWS may  
63 even be greater, particularly when used jointly with or in parallel to UCMs. In fact, a recent study  
64 by Sgoff et al. (2022) improved the weather forecasting of the Icosahedral Nonhydrostatic Model  
65 (ICON; Zängl et al. (2015)) at a horizontal resolution of 2 km over Germany by assimilating the  
66 data provided by PWS for air temperature and relative humidity at 2 m height. Although data  
67 assimilation occurs at runtime, PWS could also be used to bias-correct urban climate simulations  
68 as a post-processing step. Oleson et al. (2018) already noted the need for a global dataset of  
69 urban weather observations to properly bias-correct simulated urban climates. We indeed expect  
70 urban climate simulations to have systematic biases that can be induced for a variety of reasons,  
71 such as: urban canopy parameters (Demuzere et al. 2017; Hammerberg et al. 2018; Zonato et al.  
72 2020); complexity of urban climate models (Grimmond et al. 2011; Loridan and Grimmond 2012;  
73 Lipson et al. 2021); time at which the simulation is initialised (Bassett et al. 2020); choice of initial  
74 and boundary conditions for lateral and vertical forcing (Brisson et al. 2015); or choice of model  
75 parameterizations – such as the two evaluated in this work (see Methods). Hence, UCM will always  
76 present a certain degree of uncertainty that has to be allowed for prior to performing urban climate  
77 impact studies that use climatic variables derived from modelled simulations to estimate the impact  
78 of the urban climate on other things (e.g. mortality, biodiversity, etc.). Using PWS could thus be  
79 beneficial for obtaining realistic urban weather data of present and future urban climates that can  
80 be used to perform urban climate impact studies and guide decision-making.

81 In this study, we propose to leverage the increasingly dense network of PWS over south-east  
82 England since 2015 (Brousse et al. 2022) to evaluate and bias-correct urban climate simulations  
83 that were run for the hot summer of 2018 – the hottest summer on average in the UK. Common  
84 practices in bias-correction include adding the mean bias to the modelled variable distribution or  
85 applying a separate correction to each quantile of the distribution (Maraun and Widmann 2018).  
86 Model biases are usually measured at official weather stations at rural sites, thereby assuming



87 that the urban heat island phenomenon is accurately represented by the UCM (e.g., Lauwaet et al.  
88 (2015), or Oleson et al. (2018)). Some studies however tried considering the urban effect by linearly  
89 transforming the bias-correction coefficient via an urbanization ratio calculated at each grid cell,  
90 like in Wouters et al. (2017) over Belgium. Assuming that urban climate simulations biases cannot  
91 be linearly related to the urban fraction only, we decided to test whether urban in-situ observations  
92 can be used to perform an urban-specific bias-correction of air temperatures driven by machine  
93 learning.

94 We chose to use machine learning regressors to correct the air temperature biases because machine  
95 learning allows us to perform spatially explicit bias-corrections that are directly derived from the  
96 observed biases at all PWS locations and that are related to a set of spatially explicit covariates.  
97 Machine learning regressors of ranging complexities allow for the statistical discretisation of a  
98 single relationship between the covariates and the variety of biases. To our knowledge, such  
99 a technique has never been proposed as a viable approach for bias-correction of urban climate  
100 simulations, probably because of the lack of observations in urban areas. We hereby hypothesize  
101 that such an innovative bias-correction method would be beneficial for urban heat impact studies  
102 by improving the UCM outputs on which they rely. Such innovations are needed to better assess  
103 the heat burden in cities (Nazarian et al. 2022).

104 To respond to these issues through the scope of urban near-surface temperatures, we: i) evaluated  
105 the ability of the complex three-dimensional UCM embedded in WRF – the Building Effect  
106 Parameterization coupled with its Building Energy Model (BEP-BEM) – to accurately represent  
107 the urban impact on air temperatures under two boundary layer schemes for the summer of 2018 in  
108 south-east England using official weather stations and PWS separately to show their added value for  
109 detecting spatially heterogeneous urban temperature biases; ii) used machine learning regressions  
110 to predict the models' daily air temperature biases in the urban environment and bias-correct the  
111 two simulations suggested in part i – which allowed us to determine an optimal time-step at which  
112 the bias-correction should be performed to optimize the outputs.; and iii) compared the two bias-  
113 corrected products against the predicted daily air temperatures using only PWS measurements to  
114 investigate how realistic the bias-corrected products are. In parallel, to illustrate the benefit gained  
115 from the bias-correction for impact studies, we showcase how the bias-correction leads to different  
116 population weighted temperatures in the Greater London area. We also estimated the amount of

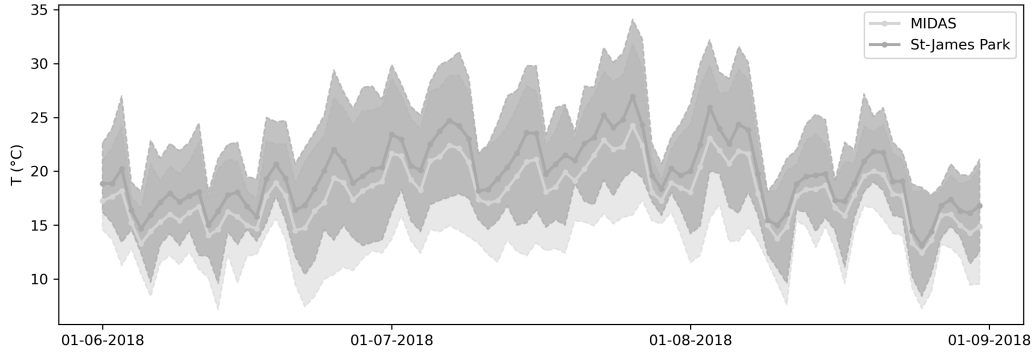


FIG. 1. Diurnal ranges of temperatures observed by the Met Office MIDAS automatic weather stations. The urban St-James' Park station in central London (dark grey) is always hotter than the average temperature of all MIDAS stations in south-east England (light grey) for daily average, minimum and maximum temperatures. The thick lines represent the daily average temperature and the shading represent the spread between daily maxima and minima.

PWS that are necessary to achieve optimal machine learning regressors performance and tested the added value of official weather stations for bias-correction.

It is important to consider that our study does not try to estimate how a bias-corrected modelled product is better compared to a predicted product from observations for urban climate impact studies. We hereby simply try to demonstrate that any urban climate impact work that is based on urban climate modelling should pursue a spatially explicit bias-correction specific to urban areas.

## 2. Methods

### *a. Model setup and region of interest*

We focused our study on the south-eastern parts of England, centred over the metropolis of London, host to approximately 9 million inhabitants. We chose to model the impact of urbanization on 2 m air temperature in London during the summer of 2018, since it was the hottest summer on average in the UK (McCarthy et al. 2019). During the the British Isles heatwaves, maximum daily temperatures often surpassed 30 °C (Figure 2) with a maximum of 34.4 °C measured at London's Heathrow airport on the 26<sup>th</sup> of July. This former record has yet been broken in 2019 and 2022.

136 To model the impact of the urban areas of London and south-east England on local meteorology,  
137 we used the Weather Research Forecast (WRF) regional climate model version 4.3 and activated  
138 the embedded Building Effect Parameterization (BEP; Martilli et al. (2002)) urban climate model  
139 with its partner Building Energy Model (BEM; Salamanca et al. (2010); Salamanca and Martilli  
140 (2010)) – hereafter referred to as BEP-BEM. We ran the model at a horizontal resolution of 1 x  
141 1 km following a two-way nesting strategy where the outer domain is forced by ERA5 6-hourly  
142 data at 25 km with 199 by 199 grid points and the two intermediate domains are run at horizontal  
143 resolutions of 9 and 3 kilometres with 252 by 241 and 210 by 180 grid points, respectively (Figure 2,  
144 upper panel). Initial land surface conditions were provided by the default MODIS 5-arc-second  
145 land use dataset provided by the WRF community while sea surface temperatures were updated  
146 6-hourly out of ERA-5. No lake models were activated, hence meaning that inland fresh water  
147 bodies are given the MODIS Water land cover class and are not updated on 6-hourly time steps as  
148 sea-surface temperatures. We ran the model in parallel over 200 CPUs using restarts every four  
149 days of simulation. We started the simulations on the 25<sup>th</sup> of May 2018 and end them on the 31<sup>st</sup>  
150 of August 2018, considering the first 7 days of simulation as spin-up time.

151 All domains used the same physical and dynamical parameterizations which we obtained out of  
152 preliminary testing done over the two hottest days of the summer 2018 – 26<sup>th</sup> and 27<sup>th</sup> of July 2018  
153 (see Appendix A). We thereby used the WRF Single-moment 3-class microphysics scheme (Hong  
154 et al. 2004), the Dudhia shortwave and RRTM longwave schemes (Dudhia 1989; Mlawer et al.  
155 1997), and the revised MM5 surface layer scheme (Jiménez et al. 2012). In the first domain, the  
156 Kain-Fritsch convection scheme was activated (Kain 2004) and then turned off in the second and  
157 third domains, which were at convection-permitting scales. We set the model top at 50 hPa with an  
158 additional 5000 m damping layer and subdivided the atmosphere into 56 vertical layers. We used  
159 the Noah-MP land surface scheme (Niu et al. 2011; Yang et al. 2011) in its default parameterization  
160 over 4 soil layers.

163 Urban canopy parameters required by the WRF BEP-BEM model were provided via the newly  
164 standardized WUDAPT-TO-WRF (W2W) python package developed by Demuzere et al. (2021),  
165 following the Fortran version used by Brousse et al. (2016). This allowed the transfer of spatially-  
166 explicit morphological urban canopy parameters suitable for urban climate simulations via Local  
167 Climate Zones (LCZ) maps covering the inner domain (Figure 2, lower panel). We use the

161

162

Table 1. Thermal and radiative parameters per LCZ based on Stewart et al. (2014). Road parameters are considering a mixture of asphalted and concrete road pavements and grass.

Heat capacity			Thermal conductivity			Albedo			Emissivity			
[J.m <sup>-3</sup> .K <sup>-1</sup> ]			[J.m <sup>-1</sup> .s <sup>-1</sup> .K <sup>-1</sup> ]									
Roof	Wall	Road	Roof	Wall	Road	Roof	Wall	Road	Roof	Wall	Road	
LCZ 1	1.80E+06	1.80E+06	1.75E+06	1.25	1.09	0.77	0.13	0.25	0.15	0.91	0.90	0.95
LCZ 2	1.80E+06	2.67E+06	1.65E+06	1.25	1.50	0.73	0.18	0.20	0.16	0.91	0.90	0.95
LCZ 3	1.44E+06	2.05E+06	1.63E+06	1.00	1.25	0.69	0.15	0.20	0.18	0.91	0.90	0.95
LCZ 4	1.80E+06	2.00E+06	1.54E+06	1.25	1.45	0.60	0.13	0.20	0.20	0.91	0.90	0.95
LCZ 5	1.80E+06	2.00E+06	1.50E+06	1.25	1.45	0.62	0.13	0.25	0.20	0.91	0.90	0.95
LCZ 6	1.44E+06	2.05E+06	1.47E+06	1.00	1.25	0.60	0.13	0.25	0.21	0.91	0.90	0.95
LCZ 7	2.00E+06	7.20E+05	1.38E+06	2.00	0.50	0.51	0.15	0.20	0.24	0.28	0.90	0.92
LCZ 8	1.80E+06	1.80E+06	1.80E+06	1.25	1.25	0.80	0.18	0.25	0.17	0.91	0.90	0.95
LCZ 9	1.44E+06	2.56E+06	1.37E+06	1.00	1.00	0.55	0.13	0.25	0.23	0.91	0.90	0.95
LCZ 10	2.00E+06	1.69E+06	1.49E+06	2.00	1.33	0.61	0.10	0.20	0.21	0.91	0.90	0.95

European LCZ map by Demuzere et al. (2019). Thermal and radiative parameters are also directly derived from the LCZ classification and follow those used by Stewart et al. (2014), who used these parameters for the city of Basel, Switzerland. Each parameter for roofs, walls and roads is related to each modal LCZ of the 1 km grid cell via the `URBPARM_LCZ.TBL` (see Table 1). We decided to keep the roughness length for momentum and the lower boundary for temperatures of roofs, walls, and roads identical across each LCZ. We fixed the roughness length at  $1.00\text{E-}4$  m for walls and at 0.01 m for roofs and roads, respectively. This does not mean that the effective roughness length at the bulk level does not differ between urban morphologies. Although materials composing them are considered identical in the drag they impose on the flow, their density and height will matter. Urban canyons with buildings above 25 m and another with buildings below 5 m will effectively have a different roughness length. For the boundary temperatures, we set it at 299 K for the roofs and the walls, respectively, and at 293 K for the road. We chose to deactivate the air conditioning in our simulation because air conditioning systems are not common in residential areas across London and surrounding cities, which compose the major part of the land use land cover.

In this study, two potential planetary boundary layers (PBL) schemes are compared in terms of performance and need of bias correction: the commonly used Bougeault-Lacarrère scheme (BouLac; Bougeault and Lacarrere (1989)) for urban simulations that use BEP-BEM, and the recently coupled YSU scheme to BEP-BEM (Hong et al. 2006; Hong and Kim 2008; Hendricks et al. 2020). Although we found that the latter performed better over the two hottest days of summer 2018 (see Appendix A), we decided to keep a simulation with BouLac as YSU has only been applied over Dallas (Wang and Hu 2021) whereas BouLac has been used in multiple studies already (e.g., Salamanca et al. (2011), Salamanca et al. (2012), Gutiérrez et al. (2015), Tewari et al. (2017), Mughal et al. (2019)). The Mellor-Yamada-Janjic (MYJ; Janjić (1994), Janić (2001)) scheme, also available for BEP-BEM simulations, is disregarded in this study since this PBL scheme is especially used for mountainous terrain (Zonato et al. 2022), and we are modelling the relatively flat terrain of south-east England.

#### *b. Model evaluation prior to bias correction*

We evaluated the model's performances against 35 official weather stations' measurements of air temperature at 2 m obtained from the UK Met Office MIDAS network (Sunter (2021), UKMO

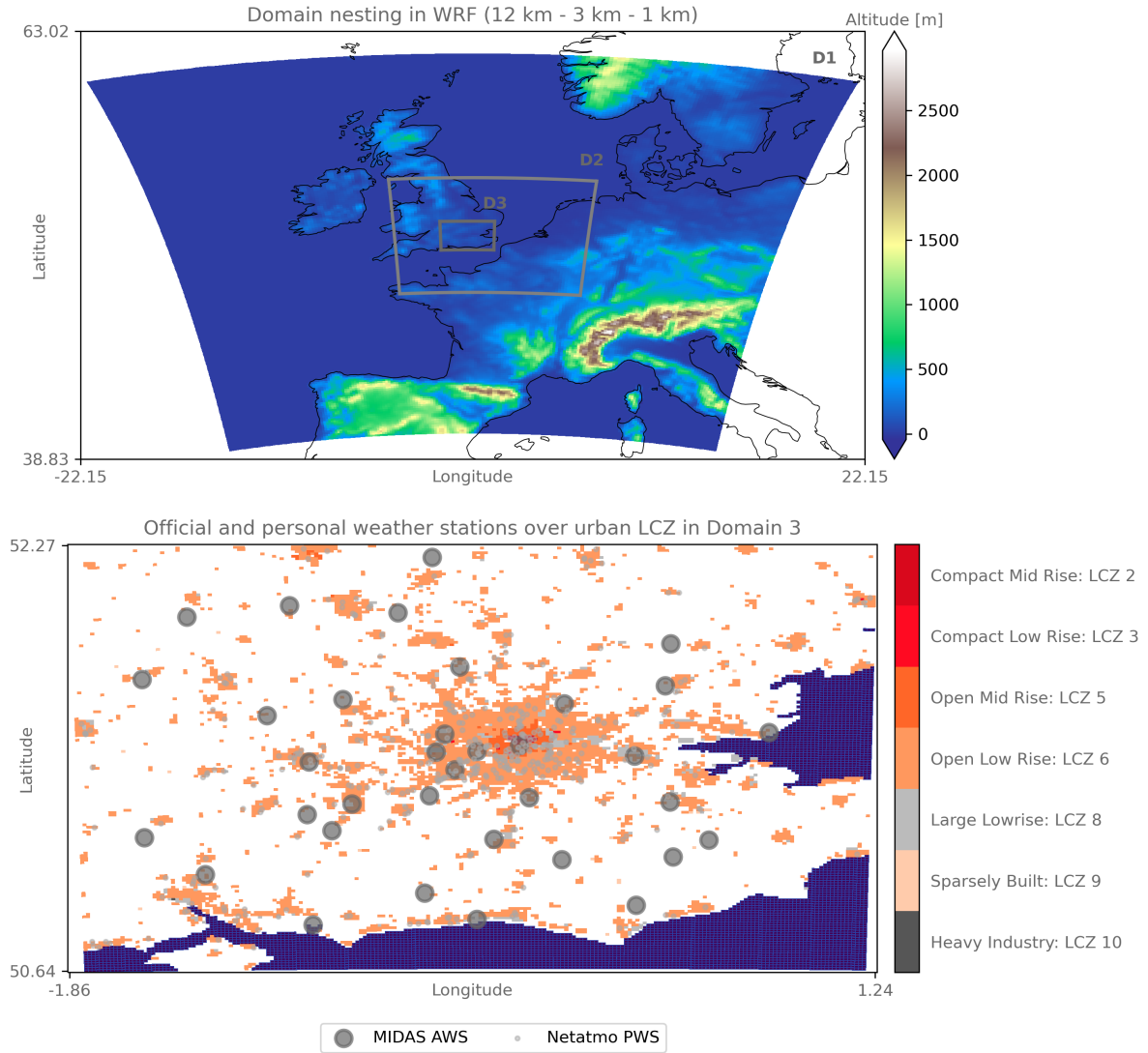


FIG. 2. Domain nesting (upper) and urban land cover in the inner domain (lower). The WRF nesting strategy consists of three nested domains at 12 km (D1), 3 km (D2) and 1 km (D3) horizontal resolution. The altitude is plotted to highlight the flat terrain of south-east England covered in D3. In the lower panel, the resulting urban landcover in D3 after using the WUDAPT-TO-WRF python tool is presented in the form of Local Climate Zones (LCZ). The MIDAS official automatic weather stations (AWS) and the Netatmo personal weather stations (PWS) used for the evaluation of the model and the subsequent bias-correction using PWS only are overlaid in grey. The sea is shown in blue in the lower panel while coastlines are drawn in black in the upper panel.

(2021); Figure 1, lower panel). To address the issue of lack of official observations amongst the urban environment, we used Netatmo PWS to complement the model evaluation (Figure 1, lower panel). The Netatmo PWS measurements were obtained through the Netatmo App developer API and were collected for all PWS contained within the inner most domain of WRF and that were running over the 2015 to 2020 period (more information can be found in (Brousse et al. 2022)). Prior to the evaluation, unrealistic PWS measurements were filtered out using the Crowd-QC v1.0 R package from Grassmann et al. (2018). This statistical quality check and filtering method is based on the assumption that the whole set of PWS should be regarded as a reference to individual stations specificities. Through four main obligatory quality-checks – potentially complemented by three optional – erroneous data are removed. Details of this filtering method can be found in other publications like Napoly et al. (2018) and or Brousse et al. (2022) who used the same dataset over London. For the summer 2018, the filtering reduced the dataset from 935 potential PWS to 909 potential stations over the whole domain. Such filtering has already been applied over several studies, including a large scale study by Venter et al. (2021) over a European city, and has recently been ameliorated into the *CrowdQC+* package (Fenner et al. 2021). The purpose of this study is not to test the effect of PWS quality check on the model evaluation and bias correction.

After quality-checking the PWS we also added an additional filtering where we removed PWS that did not have sufficient temporal data coverage and that were not located in an urban pixel according to WRF. Only PWS that have less than 4 hours per day without data and that are located in urban pixels with an urban fraction greater than 0 are retained – where the WRF land-use land-cover at 1 km horizontal resolution refers to an LCZ. This ensures that we do not include measurements that are not representative of the daily variations in air temperatures or built-up environments. Additionally, the prior filtering performed using the *CrowdQC* package also ensures that measurements that are not representative of outdoor thermal variations (e.g., indoor sensors) or that are resulting from defective sensors are taken out. Overall, the filtering step is necessary to ensure that our model outputs are evaluated against measurements of sufficient quality and that the subsequent bias-correction is deprived of unnecessary noise in the data that could lower its performance. This resulted in a sample of 402 PWS usable for model evaluation and bias correction. Out of these, 354 were located in WRF grids classified as LCZ 6, 30 in LCZ 5, 8 in LCZ 2, 6 in LCZ 8, 3 in LCZ 9 and 1 in LCZ 3.

Each model simulation was evaluated using a set of common statistical indicators: the root mean squared error (RMSE), the mean absolute error (MAE), the mean bias error (MB), Spearman's coefficient of correlation ( $r$ ) and the square of Pearson's coefficient of correlation ( $r^2$ ). These metrics are obtained using the Python scikit-learn and scipy's stats packages from Pedregosa et al. (2011) and Virtanen et al. (2020).

### *c. Bias correction using personal Netatmo weather stations*

In our study, we propose an innovative method to bias-correct urban temperatures at a horizontal scale of 1 km by using machine learning regression. The advantage of using machine learning regression compared to more common bias-correction strategies (e.g., the definition of a single bias coefficient) is that we are able to relate our model output biases out of spatially varying and explicit sets of parameters. In our case, we make the assumption that the spatial variation in the bias of the model is dependent only upon the spatial morphological inputs to the UCM. These include the urban fraction, the surface height, the average building height, the building surface to plan area fraction ( $\lambda_b$ ), the plan area fraction ( $\lambda_p$ ) and the frontal area fraction ( $\lambda_f$ ). Using this set of predictive covariates, we train our regressors to predict the bias in the modelled air temperature at 2 m ( $T_2$ ) based on observed biases at urban PWS locations. This way, we are able to bias-correct the modelled temperatures in each urban pixel based on the predicted bias ( $T_2 - \text{bias}_{pred}$ ). Our bias-correction does not make use of official MIDAS weather stations as their use is considered detrimental to the bias correction following an analysis on sample size and sensor types given in Appendix B.

We chose to bias-correct the simulated daily minimum, maximum and average  $T_2$  ( $T_{2min}$ ,  $T_{2max}$ , and  $T_{2mean}$ ) using filtered PWS observations in London and south-east England. Daily temporal scale is considered optimal as it combines a higher spatial density of measurements compared to hourly data and a lower computational requirement; it is also a commonly used temporal scale for urban heat impact studies. Daily minimum and maximum air temperatures at 2 m are defined following the Met Office Had-UK definition: minimum temperature observed from 9AM of the previous day  $d-1$  to 9AM of the  $d$  day, and maximum temperature observed from 9AM of the  $d$  day to 9AM of the next day  $d+1$  (Hollis et al. 2019).



TABLE 2. Hyperparameter tuning used by each regressors

Model	Parameters Dictionary
Linear	'normalize': False
Ridge	'alpha': 1, 'normalize': True, 'random_state': 42, 'solver': 'lsqr', 'tol': 0.01
Lasso	'alpha': 1, 'normalize': False, 'random_state': 42, 'selection': 'random', 'tol': 1e-10
Random Forest	'max_features': 'sqrt', 'min_samples_leaf': 11, 'min_samples_split': 2, 'n_estimators': 400, 'random_state': 42
Gradient Boosting	'learning_rate': 0.2, 'max_depth': 3, 'max_features': 'sqrt', 'min_samples_leaf': 10, 'min_samples_split': 22, 'n_estimators': 200, 'random_state': 42, 'subsample': 0.2

We test the ability of 6 different regressors of increasing complexity available in the Python scikit-learn packages (Pedregosa et al. 2011) to predict the model bias based on WRF spatial urban canopy parameters only. These regressors are: dummy regression (which simply returns the mean bias), linear regression, Ridge regression, Lasso regression, Random Forest regression, and Gradient Boosting regression. Each of the different regressors, except the dummy regression, offers a set of parameters that can be fine-tuned to increase each regressor's performance. Hence, prior to running the daily bias-correction we use a 5 K-fold cross-validation using the *Grid Search CV* package from scikit-learn in Python to evaluate the impact of hyperparameter tuning on the regressors' performances based on RMSE, MAE and  $r^2$ . The cross-validation is done over the summertime average daily mean temperature bias from the YSU run only, for computational reasons. We retain RMSE as the refitting score to better capture the spatial spread and extremes of T2. The resulting parameterizations are given in Table2. We chose to keep the same hyperparameter tuning for all bias correction and predictions to ease comparability between the outcomes.

Once the hyperparameter tuning is done and prior to performing the final bias-correction, we test if the bias-correction is beneficial for palliating to the models' bias and if it also benefits from training the regressors at the daily time-step or if a training using the time-mean bias is sufficient. To perform this evaluation using the same metrics as in the model evaluation, we bootstrap each regressors 25 times per day, randomly sampling 80 % of the PWS locations that had data available on that day as training and keeping the remaining 20 % as testing – for both the daily-minimum, -maximum and -average, and their respective summer time-mean average. We then first average all bootstrapped T2\_BC at the testing PWS sites before performing a subsequent averaging to obtain an average T2\_BC at the daily time step representative of all randomly selected testing PWS sites. These are evaluated against the daily average of all observed temperature at the PWS sites – for

285 daily minimum, maximum and average. In short, we are measuring how well do the two different  
286 types of bias correction perform under all regressors for capturing the daily variation (n=92 days)  
287 of temperature on average.

288 After this final step, we bias-correct both the BouLac and the YSU runs using 100 % of the  
289 measured biases and related covariates at PWS locations to compare the spatial outcomes. We  
290 also predict T2 out of PWS' observed T2 with the same set of covariates used to predict the model  
291 bias to illustrate how divergent each bias-corrected model outputs are to a simplified predicted T2  
292 that is not a derivative of any model constraint. Because more refined and complex techniques  
293 exist to predict air temperature from PWS and very high-resolution earth observations (e.g., Venter  
294 et al. (2020), Venter et al. (2021)), we do not evaluate these predicted temperatures which should  
295 simply be considered as an illustration of how bias-corrected products are similar or divergent to  
296 observational data.

297 Lastly, to illustrate the potential benefit of modelled air temperature bias-correction prior to  
298 urban heat impact studies, we calculate the average population weighted temperatures – based on  
299 the United Kingdom census data from 2011 – in Greater London before and after the bias-correction.

### 300 **3. Results**

#### 301 *a. WRF simulation evaluation*

302 When we evaluate the two model simulations against MIDAS official weather stations only, they  
303 perform similarly, demonstrating a systematic negative bias of  $\sim 0.55$  °C on average (Table 3). The  
304 average correlation with the automatic weather stations following the squared Pearson's  $r^2$  is of  
305 0.77 for BouLac and 0.79 for YSU, while using Spearman's  $r$  it is of 0.86 and 0.88, respectively. A  
306 slight decreased performance is found in urban pixels for YSU, with an average MAE of 1.83 °C  
307 and a negative MB of 0.79 °C compared to BouLac's 1.82 °C for MAR and -0.56 °C for MB.  
308 In general, the bias is more important at night, and, in non-urban stations, performances are  
309 similar. Hence, looking only at the models' performances using standard in-situ observations  
310 doesn't provide information on which model represents the urban climate more accurately.

315 On the other hand, comparison with PWS observations identifies differences in performance in  
316 urban areas between the models, as shown by the performance metrics plotted in Figure 3 and C1.  
317 The BouLac simulation has a stronger cool bias of  $-1.46$  °C  $\pm$  0.6 °C on average in the urban area,

TABLE 3. Average of all performance metrics calculated at each MIDAS official weather stations for hourly air temperature at 2 m for the summer period (1<sup>st</sup> June 2018 to the 31<sup>st</sup> of August 2018). Urban stations are stations located in a pixel classified as an urban LCZ in WRF and rural stations are located in other natural land-use land-cover.

	BouLac					YSU				
	RMSE	MAE	MB	r <sup>2</sup>	r	RMSE	MAE	MB	r <sup>2</sup>	r
<b>All</b>	2.33	1.82	-0.56	0.77	0.86	2.31	1.83	-0.57	0.79	0.88
<b>Urban</b>	2.42	1.88	-0.73	0.76	0.86	2.42	1.92	-0.93	0.77	0.87
<b>Rural</b>	2.32	1.81	-0.53	0.78	0.86	2.28	1.81	-0.50	0.80	0.88

compared to YSU’s MB of  $-0.97\text{ }^{\circ}\text{C} \pm 0.81\text{ }^{\circ}\text{C}$ . RMSE and MAE are similar, with values of  $2.79\text{ }^{\circ}\text{C} \pm 0.36\text{ }^{\circ}\text{C}$  and  $2.19\text{ }^{\circ}\text{C} \pm 0.31\text{ }^{\circ}\text{C}$  for BouLac and  $2.65\text{ }^{\circ}\text{C} \pm 0.40\text{ }^{\circ}\text{C}$  and  $2.14\text{ }^{\circ}\text{C} \pm 0.34\text{ }^{\circ}\text{C}$  for YSU. These metrics are consistent with the MIDAS observations, highlighting a systematic cool bias of the model and a coefficient of determination ( $r^2$ ) of 80 %. Importantly, the variability in the model’s performance is more greater in the YSU run – reflected by greater standard deviations of performance metrics – and, in the BouLac simulation, the metrics are more heterogeneously distributed amongst the urban area. Indeed, when we look at the YSU simulation, we can see that the model has a smaller MB in suburban areas and a greater MB in the city centre. Yet, in parallel, the correlation with the PWS is lower in the suburban areas and higher in the centre of the city. This could mean that YSU accurately represents the urban temperatures on average due to compensating effects, which we do not intend to evaluate in this study. Nevertheless, this shows how PWS are beneficial for capturing the spatial heterogeneity of each model’s performance and therefore supports the use of spatially-varying bias-correction.

### *b. Bias correction of urban climate simulations*

Over our domain of study covering south-east England during the Summer 2018, both models are subject to a cold negative bias of  $\sim -0.5\text{ }^{\circ}\text{C}$  on average according to official stations and of  $\sim -1.0\text{ }^{\circ}\text{C}$  to  $\sim -1.5\text{ }^{\circ}\text{C}$  according to PWS. But as demonstrated above, the bias of the models against PWS observations has substantial spatial variation and so the bias correction for urban heat impact studies should be spatially explicit.

We find that each machine learning regressors give similar performance(Figure 4; values numerically given in Tables C1 and C2 ). All bias-corrections were however beneficial compared to the

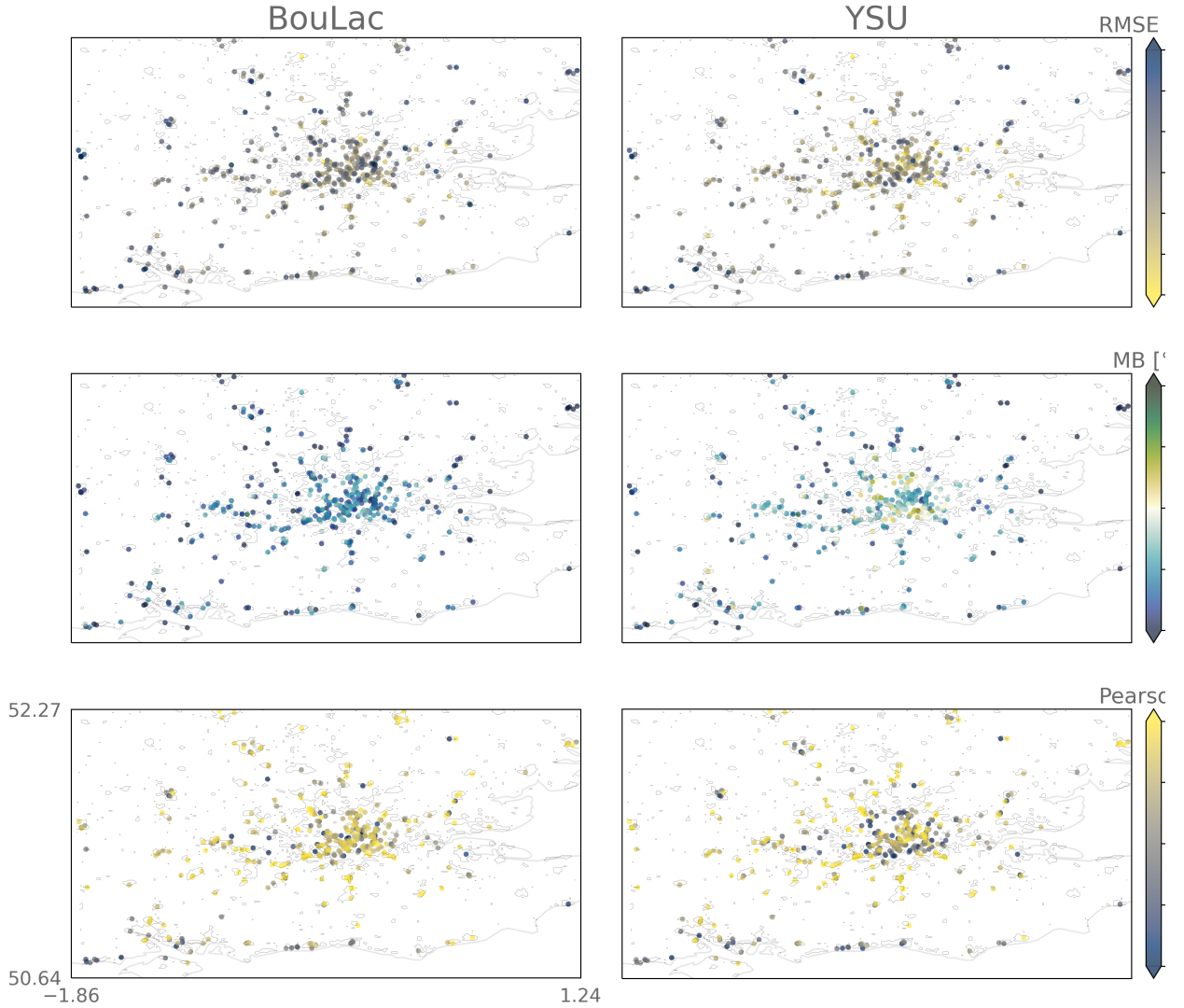


FIG. 3. Performance metrics calculated at location of each citizen personal weather station (PWS) for the two model simulations using different planetary boundary layer schemes (YSU and BouLac). The metrics are calculated over the whole summer 2018 with hourly outputs of near surface air temperature at 2 m. Root mean square error (RMSE) and mean bias (MB) are given in degrees Celsius ( $^{\circ}\text{C}$ ). Coefficients of correlation measured with the squared Pearson's  $r$  are also provided. Mean absolute error (MAE) and Spearman's  $r$  are given in Figure C1 to increase clarity.

original outputs from the WRF model, reducing RMSE, MAE and MB by  $0.29^{\circ}\text{C}$ ,  $0.32^{\circ}\text{C}$  and  $1.02^{\circ}\text{C}$  on average. The bias-correction was most efficient for daily-minimum temperatures and less for daily-maximum temperatures, where RMSE was not diminished – if not slightly increased

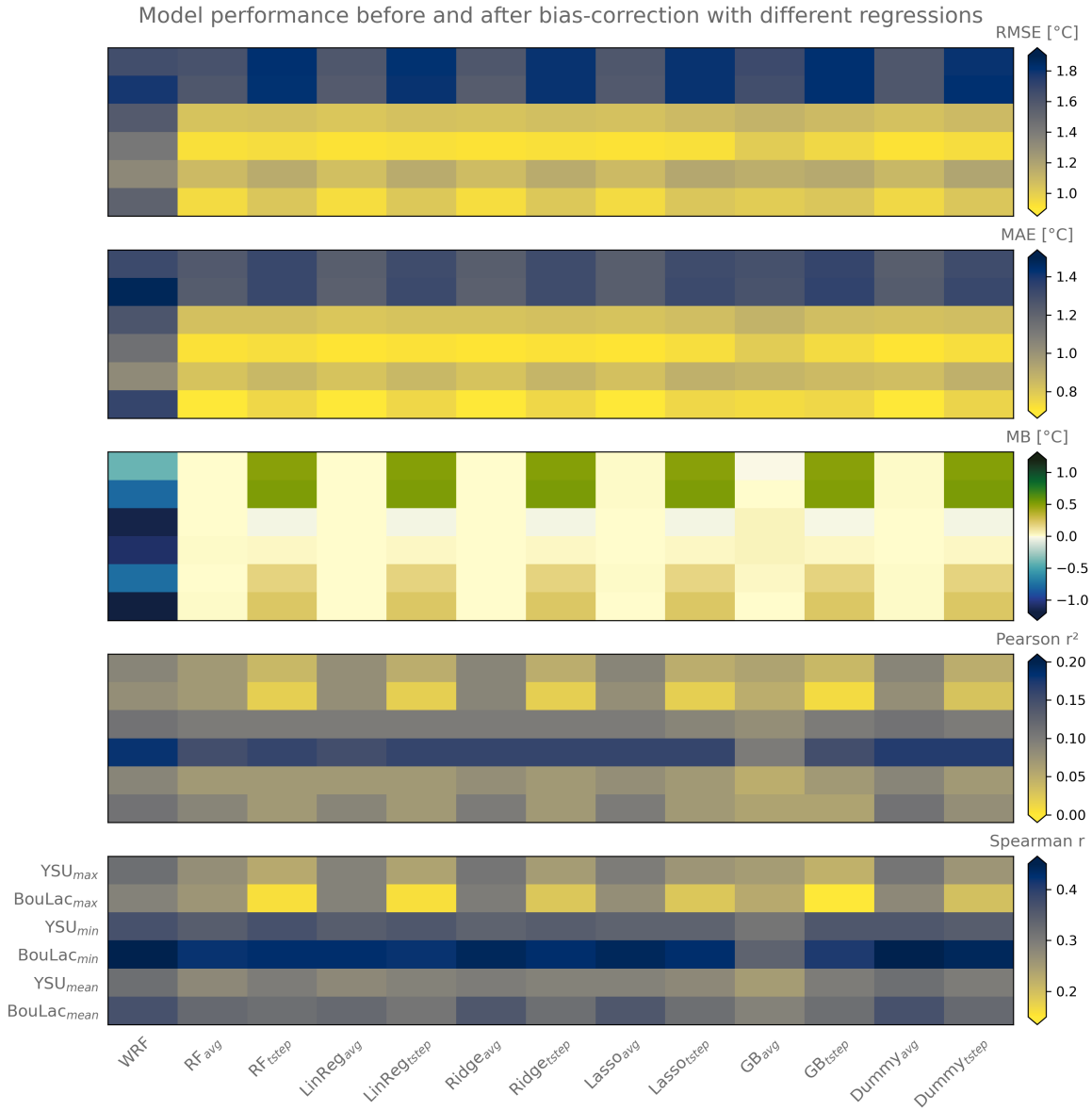


FIG. 4. Performance metrics for the model prior to the bias-correction (WRF) and all the different regressions (random forest: RF; linear regression: LinReg; Ridge regression: Ridge; Lasso regression: Lasso; gradient boosting: GB; and dummy regression: Dummy). The different regressions are assigned a suffix: “avg” for regressions that were trained on the summer time-mean average of daily-minimum, -mean or -maximum temperatures, and “tstep” for those that were trained with the temperatures at each daily time-step.

(by 0.05 °C for YSU daily-maximum temperatures for example) – by the time-step bias-correction. Interestingly, the spatial correlation between the bias-corrected and the observed temperatures are low, with values ranging from around 0.02 to 0.2 for the squared Pearson's  $r$  and from around 0.15 to 0.45 for Spearman's  $r$ . This can be expected as machine learning algorithms have difficulties representing a time-varying variable with static spatial elements only (Georganos et al. 2021; Venter et al. 2021). Unexpectedly, we find that the training at the daily time-step does not outperform the training at the summer time-mean in terms of spatial correlation with the heat distribution across London. Nonetheless, if we take the average daily-minimum, -mean and -maximum temperatures of all PWS and compare it to the modelled temperatures, we find that the time-step bias-correction is closer to the observations (Figures C2 to C4). Lastly, we find that greater model performance is achieved with a minimum of ~24 % (96 PWS) of the whole sample of PWS and that official weather stations are detrimental to the regressors performance (see Appendix B).

Comparing the spatial differences of the bias-corrected products related to the complexities of each regressors, we find that although each regressor is performing similarly on average, important disparities are found between the outputs. For example, when looking at the average bias-correction imposed to daily-minimum temperatures after training the regressors at each time-step, the Lasso and the Ridge regressors impose a flat bias-correction, similar to the dummy regression, while the random forest and gradient boosting regressors' degrees of freedom result in a spatially diverse bias-correction (Figure 5 and Figures C5 and C6). Besides, the linear regression imposes an average bias-correction spatially-correlated to the modal LCZ. In general, the signal is consistent across each regressors, apart from the Lasso and the dummy regression, where, for YSU, central London requires a stronger bias-correction by 1 °C to 2 °C compared to the suburban areas where the bias-correction is around 0.5 °C ; for BouLac, the central bias-correction is lower than YSU. We find that these spatial tendencies are also found for daily-maximum and daily-average temperatures, defending our hypothesis of a systematic bias correlated to spatially explicit input parameters. The spatial differences in bias-correction are however less important for daily-maximum temperatures, which is the time at which the urban heat island is also expected to be the lowest.

# Modelled temperatures and respective bias-corrections with multiple regressors

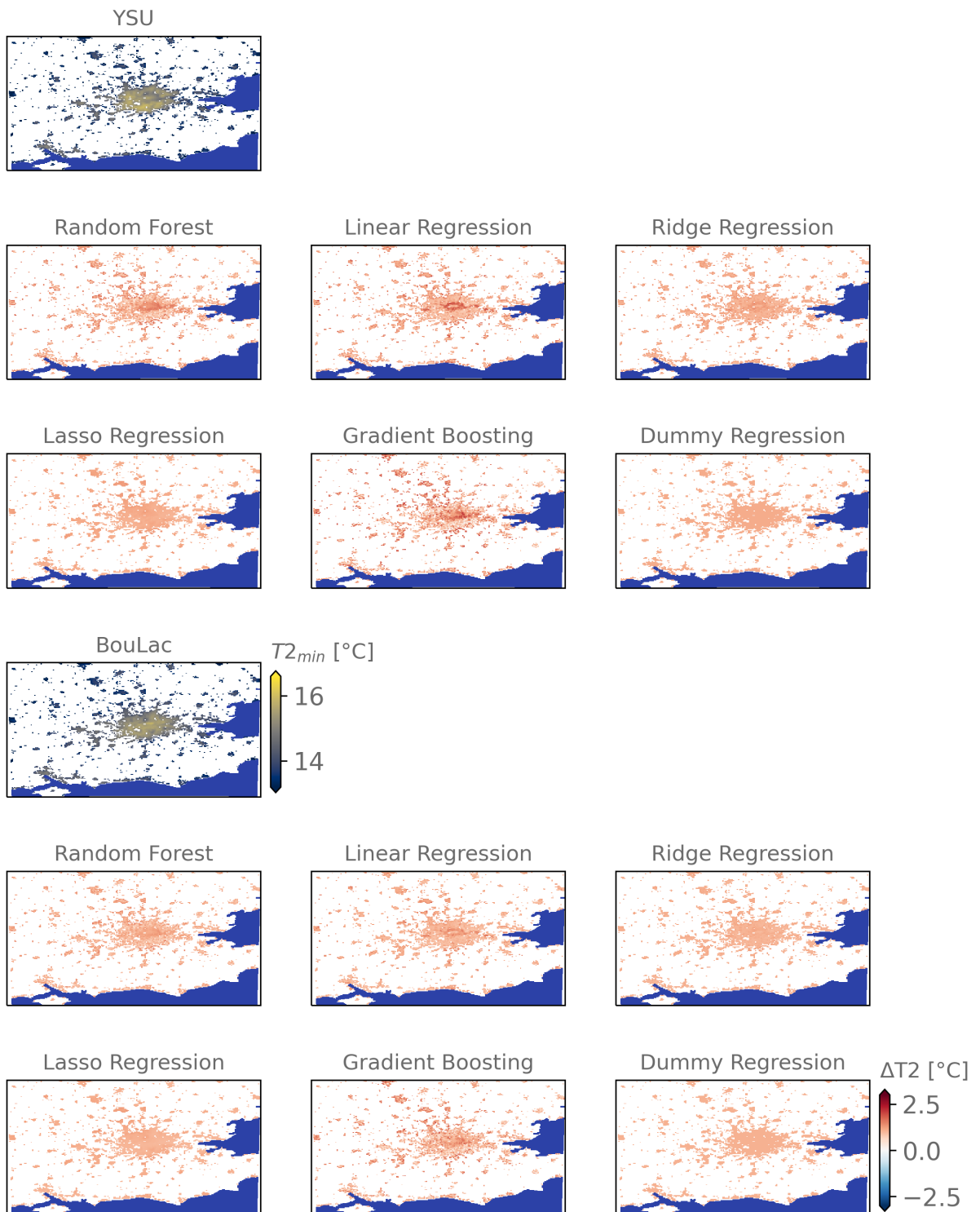


FIG. 5. All regressions propose different bias-corrections ( $\Delta T_2$ ) of the average modelled absolute daily minimum urban temperature ( $T_{2_{min}}$ ). Differences of bias-correction are observed between the runs with different planetary boundary layer schemes (Bougeault-Lacarrère – BouLac, and Yonsei University – YSU). The centre of London is subject to a stronger bias-correction. Rural lands are masked in grey and the seas in blue. Bias corrections of daily mean and maximum temperatures are given in Figures C5 and C6

Finally, we find that the bias-corrected BouLac simulation corresponds spatially to predicted temperatures using PWS more than YSU – something we find equally across all regressors (Figure 6 and Figures C7 to C11). As an example, when comparing the average bias-corrected products using the time-step trained random forest regressor we can see that YSU urban heat is more homogeneously distributed than BouLac’s or the predicted temperatures from PWS only. BouLac’s bias-corrected product shows stronger urban heat in central London compared to suburban areas, coherent with the predicted temperatures. Nonetheless, BouLac’s suburban areas are hotter by 0.5 °C to 1.0 °C than the predicted ones with PWS only. This remains less pronounced than in YSU. Lastly, we can see that both bias-corrected products show similar trends when compared to the PWS-only predicted temperatures with hotter suburban areas and cooler secondary cities as well as coastlines. Again, this does not show which product between the PWS-only predicted temperatures and the bias-corrected products is better since we do not evaluate this here.

These results show that bias-correction of modelled air temperature change their spatio-temporal distributions. When focusing on the potential impact bias-correction may have in estimated urban heat impact on urban health, we find that using the random forest regression trained at each daily time-step leads to an increased average population weighted temperature by 0.77 °C in the YSU case, and of 1.24 °C in the BouLac case. Raw model outputs are thereby lowering the impact of heat on the urban population.

#### 4. Discussion

In this study, we argue that the joint use of crowd-sourced personal weather stations (PWS) and urban climate models (UCMs) can add value to urban climate research and in particular to urban climate impact research. This is supported by two major outcomes of our case-study focused over London during the summer 2018. First, we showed that evaluation of urban climate simulations using PWS enables the detection of spatially-varying systematic biases in urban areas related to the



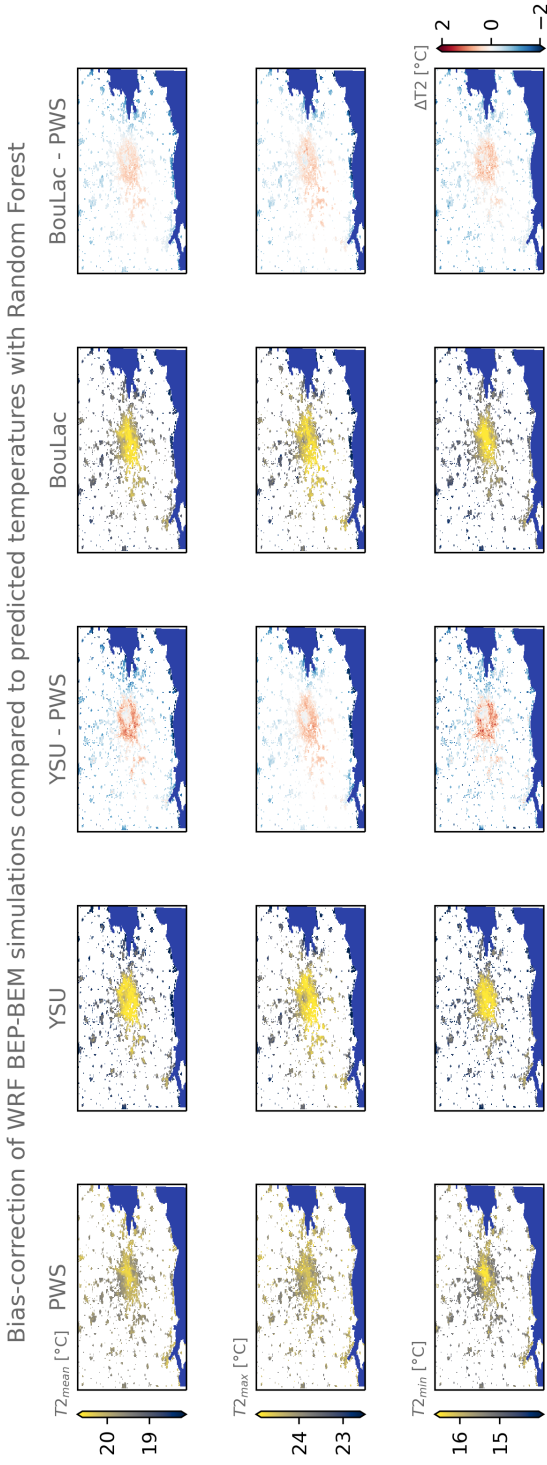


Fig. 6. The random forest regressor leads to different bias-corrections of the two WRF simulations parameterized with different turbulence schemes – the Yonsei University (YSU) and the Bougeault-Lacarrère (BouLac) – and with the BEP-BEM urban canopy model activated. This holds for average daily mean, minimum and maximum temperatures ( $T2_{mean}$ ,  $T2_{min}$  and  $T2_{max}$ ) after the daily time-step bias-correction. Compared to the predicted temperatures using the personal weather stations data only (PWS), the bias-corrected products are hotter in the suburban areas of the Greater London and cooler in the rural areas. The difference is more pronounced in YSU (see YSU – PWS). Greyed areas represent natural areas where the bias-correction is not performed and the sea is shown in dark blue. The same figures for the other regressors are given in Figures C7 to C11

415 UCMs’ parameterization, which are not detectable using only official weather stations. Second,  
416 we demonstrated that PWS, combined with detailed morphological data derived from LCZ maps,  
417 can be used to derive a spatially varying bias-correction via commonly used machine-learning  
418 regressors. This latter point has major implications for urban climate impact research – and  
419 especially future urban climate impact studies – as we hereby propose the first bias-correction  
420 technique that considers the existence of a non-linear spatially heterogeneous bias in modelled  
421 urban climates.

422 Of course, using PWS for evaluating UCM simulations should always cautiously be considered  
423 because of the lower accuracy of PWS and the potential uncertainties related to user-driven mistakes  
424 in the set-up of their PWS (e.g., indoor sensors instead of outdoor, poor shading conditions, height  
425 of the sensor, etc.). However, reliable tools have now been developed since the first use of PWS for  
426 model evaluation by Hammerberg et al. (2018) to filter dubious measurements out (e.g., *CrowdQC*  
427 from Napoly et al. (2018) or *CrowdQC+* by Fenner et al. (2021)), thus making PWS observations  
428 increasingly reliable. This does not resolve the question of the representativity of measurements,  
429 i.e., “how is one PWS measurement representative of the simulated urban pixel?” Yet, the increasing  
430 density of PWS in the urban environments begins to alleviate this uncertainty – despite a recognised  
431 unequal distribution of PWS amongst a variety of environmental, socio-economic and demographic  
432 indicators (Brousse et al. 2023). For example, Venter et al. (2020) found that a density of one  
433 PWS per square kilometre is optimal for predicting seasonal air temperature in Oslo. Dense PWS  
434 networks hence permit the detection of systematic biases that would otherwise pass undetected.  
435 Therefore, to support the development of PWS as a source of urban weather observations for model  
436 evaluation, urban climate scientists should identify an optimal density of PWS for UCM evaluation,  
437 to define which cities are in need of urban weather observations, and to start instigating common  
438 frameworks and standards.

439 We consider our study innovative and supportive of future advances in the field because it is the  
440 first bias-correction technique in urban environments which considers that UCMs’ simulated UHI is  
441 spatially heterogeneous in its accuracy and that the UHI is not solely linearly correlated to the urban  
442 fraction. Aided by the expanding fields of crowd-sourcing weather observations through PWS,  
443 machine learning, and potentially deep learning, we infer that our work should serve as the basis of  
444 future research that would try, but not restricted to, improving the bias-correction of urban climate

models using PWS. For instance, we did not find any machine learning regressor to be more efficient at predicting the model bias. This could be explained by the rather restricted set of covariates we used for training the regressors as well as the coarse horizontal resolution of 1 km at which the covariates were aggregated to be consistent with the model's spatial resolution. Higher spatial resolutions and more specific satellite earth observations could be used to improve regressors' performance, following up on the work by Venter et al. (2021), for example. When modelling the near-surface UHI, which is not a model bias, their regressor achieved similar performances as ours, with an RMSE of 1.05 °C and a Pearson's  $r^2$  of 0.23. Although the common use of model's input parameters and earth observations as covariates could be beneficial, a particular attention should be given to the choice of earth observations since these should not be decorrelated to the model's physics and dynamics as the purpose would remain the bias-correction.

Independent of the set of covariates used in this study we found that the regressors performances greatly improved when trained over a certain number of PWS (more than ~90) before plateauing. Because of this, future research should try to investigate how machine learning regressors could benefit from unfiltered PWS data and other PWS data sources. Interestingly, we found that official sources of data like MIDAS were detrimental to the regressors, potentially because official weather stations tend to be placed in open fields or parks without surrounding built-up areas to increase measurement accuracies. This would explain why our regressors tended to further increase the systematic cool bias when using only MIDAS stations for training as parks are typically cooler at night and on average than more urbanised areas where PWS are located. In addition, we found that training regressors at the daily time-step did not outperform a training with the summer time-mean average. Regressors could therefore gain in performance by adding a temporal component to the covariates. Following up on this idea, the recent work by (Zumwald et al. 2021) tried predicting the near-surface air temperature in Zurich for the 30<sup>th</sup> of June 2019 out of ~650 Netatmo PWS' measurements during the preceding week. Their set of covariates consisted of spatial earth observations as well as 35 meteorological predictors that were all derived from one official automatic weather stations. The latter predictors helped training the model to recognise how the temperature measured at each PWS location was related to the meteorological variables measured at the automatic weather stations. Their predictions at hourly time-steps achieved reasonable performances with RMSEs around 1.70 °C. Bias-correction of UCM simulations could

475 hence be improved by incorporating temporally explicit meteorological observations from official  
476 weather stations. Notwithstanding, this would require extensive investigation on the area down  
477 to which each official station is representative for training the regressors. More geographically  
478 oriented machine learning regressors, like the geographical random forests (Georganos et al. 2021),  
479 could also help integrate these spatial heterogeneities for an improved bias-correction.

480 In general, we support the use of PWS observations for bias-correction of urban climate simula-  
481 tions. As shown in this case study, model outputs prior to any bias-correction could lead to under-  
482 or over-estimation of urban heat impact on public health. We indeed find that for the summer 2018  
483 in London, average population weighted temperatures were higher after bias-correcting the model  
484 outputs, suggesting higher urban heat related mortality during this period. This simple example  
485 shows that bias-correction of urban climate simulations could have important implications for  
486 calculating the exposure of urban citizen to heat or estimating the urban heat-related mortality.  
487 Although preferring bias-corrected model outputs to predicted urban air temperatures from earth  
488 observations for present-day urban heat impact studies is not covered in this study – and must be  
489 further explored – we still argue that bias-correction should be done prior to any urban heat impact  
490 studies that imply using climate model outputs. This argument is especially valid for future climate  
491 projections at urban scale and we encourage future research to investigate how to transfer present  
492 urban bias-correction coefficients to simulated future urban climates. Doing so, bias-corrected  
493 simulations could help targeting areas where heat mitigation or adaptation strategies could be more  
494 beneficial as their efficiency is dependent on their location and scales of implementation (Yang and  
495 Bou-Zeid 2019; Broadbent et al. 2022). We also suggest that our methods could be extended to  
496 other fields of urban climatology and urban air quality. Several devices already offer the possibility  
497 to obtain information on air quality, precipitation or wind speed, to name a few (De Vos et al. 2020).  
498 Hence bias-correction of regional climate models’ outputs using crowd-sourced data should not be  
499 restricted only to air temperatures.

## 500 **5. Conclusions**

501 We demonstrate that the higher density of personal weather stations (PWS) measurements of  
502 temperatures in cities like London is beneficial for urban climate model evaluation. We then show  
503 that PWS could be helpful for bias-correcting modelled temperatures using a set of machine learning

504 statistical regressors. We did not observe tangible differences in performance of the regressors  
505 to predict the bias at various locations. A minimum of ~24 % of the total sample size of PWS  
506 (96 stations of the 402 used in this study) was required to efficiently train our regressors; official  
507 weather sources like MIDAS were detrimental to the urban bias-correction, probably because of  
508 site specificities. Our work has important implications for urban climate impact studies that would  
509 make use of urban climate model outputs.

510 *Acknowledgments.* We personally thank Stefanos Georganos for his help and his comments on  
511 machine learning classifiers and regressors. We also thank Daniel Fenner and Fred Meier for their  
512 valuable insights concerning data acquisition, filtering and treatment of crowd-sourced citizen  
513 weather stations. Lastly, we are grateful to Matthias Demuzere and other committed members of  
514 the WUDAPT project for providing the European LCZ map and the python W2W tools. CH is  
515 supported by a NERC fellowship (NE/R01440X/1) and acknowledges funding for the HEROIC  
516 project (216035/Z/19/Z) from the Wellcome Trust, which funds OB and CS.

517 OB designed the study and led the conception of the manuscript with the support of CH and  
518 CS. OB was responsible for the WRF modelling, the model evaluation and the bias-correction.  
519 CS provided support in the python coding and in the statistical analysis for the bias-correction.  
520 OK was responsible for technical support of the installation of WRF on the University College  
521 London's "Kathleen" and "Myriad" super-computers. AZ and AM offered guidance in the set-up  
522 of the WRF model v4.3 and urban heat modelling expertise with SK. All authors contributed to  
523 the writing of the manuscript.

524 The authors declare no conflicts of interest.

525 *Data availability statement.* The simulations done in this research were performed using the WRF  
526 model v4.3 (<https://github.com/wrf-model/WRF.git>). The scripts and WRF namelists used  
527 in this study are accessible at [https://github.com/oscarbrousse/JAMC\\_BiasCorrection\\_](https://github.com/oscarbrousse/JAMC_BiasCorrection_PWS/)  
528 [PWS/](https://github.com/oscarbrousse/JAMC_BiasCorrection_PWS/). The related outputs presented in this research available upon reasonable request addressed  
529 to the corresponding author.

## APPENDIX A

### Model sensitivity testing over the two hottest days of Summer 2018

Prior to running the 3-months simulation, we tested the model's sensitivity to a set of parameterization to assess which model is the best performing model for the 3-months simulation. We perform the sensitivity in a progressive way; parameters are kept if beneficial, removed if detrimental. We chose to run the simulations over the two hottest days of the summer 2018 with one additional day as spin-up time – from the 25<sup>th</sup> to the 27<sup>th</sup> of July 2018 – to see how the model is capable of accurately representing an extreme condition in terms of air temperature at 2 m – tested against official MIDAS automatic weather stations and personal Netatmo PWS. The model was also tested for relative humidity and wind speed at 10 m at MIDAS locations where records were available. All wind-speed measurements are converted from knots to  $\text{m}\cdot\text{s}^{-1}$ .

We start from Heaviside et al. (2015) model's parameterization, who simulated the impact of urbanization on the local climate in the West Midlands in England, but supplement the CORINE land-use land-cover by the Local Climate Zones classification instead since Brousse et al. (2016) compared both products and proved the added value of LCZ over Madrid. We chose the work by Heaviside et al. (2015) as a starting point since it also uses the BEP urban climate model, coupled to the WRF model and is one of the only WRF simulations done over England.

From there, our simulations tested: i) the use of YSU, recently coupled to the BEP-BEM model (Hendricks et al. 2020), instead of Bougeault-Lacarrere; ii) the use of the more complex land surface scheme Noah-MP in its default parameterization instead of the default Noah land surface model; iii) the forcing by ERA5 reanalysis data at 25 km horizontal resolution instead of ERA-Interim; iv) the reduction of soil moisture by 50 % and its increase by 200 %, following suggestions provided by Martilli et al. (2021). We chose not to test the impact of urban canopy parameters in this case to keep our simulations standardized and universally coherent through the LCZ scheme. Their simulation used the same micro-, clouds, convection and radiation physics than ours.

We found that all steps taken from the original parameterization by Heaviside et al. (2015) were beneficial to the model's performance. Through an intermediate simulation where we tested again the BouLac turbulence scheme after step iii, we found that YSU was still performing better.

### Sensitivity of machine learning regressors to data quality and quantity

Before running our bias-correction and our bootstrapping we needed to evaluate the degradation in performance of all the regressors in relation to the quantity of data available for training. This way, we could ascertain that the chosen amount of 80 % for running the bootstrapping procedure was not detrimental to the regressors' performances. Additionally, despite the fact that official weather data coming from MIDAS is usually coming from open fields like airports or parks, we still chose to test how our model performs if only this data was available for bias-correction; thereby ensuring that the use of the dense network of PWS is useful for bias-correction. To test this we trained all the regressors over both WRF boundary layer conditions to bias-correct the summertime average daily mean, minimum and maximum temperatures. This means that we are testing the ability of the regressors to predict the bias at certain PWS locations to correct the modelled temperature. In this case, we evaluate the bias-corrected temperatures against the observed temperatures. We chose not to run over daily time steps as this would be too computationally expensive.

We followed a bootstrapping procedure, where 20 % of the PWS temperature data were randomly selected and kept for testing the regressors performance. Random samples with increasing ratios of the remaining 80 % of PWS temperature data and covariates were used to train the regressors 25 times. We ensured that the randomly sampled 20 % and ratios are kept constant between regressors. We first started with 1 % of the remaining 80 % and increased the ratio by steps of 1 % until 10 % of the remaining 80 %. Steps of 10 % were then used until reaching 90 % of the remaining 80 %. We chose to use these steps as we expect our regressors performance to rapidly increase with a low amount of data before plateauing with a greater amount of data. Then, to test the added value of urban PWS density and data we trained the same regressors over the modelled bias at the 10 urban MIDAS stations locations and evaluated the bias correction against the 20 % of the PWS data kept for evaluation at each bootstrapping step. As a comparison, we also evaluated the WRF output prior to bias correction against the same 20 % of PWS temperature data at each bootstrapping step to demonstrate the added value of bias correction using a certain amount of PWS.

We found that all regressors benefited from a greater amount of PWS data which reduced the root mean squared error (RMSE), the mean absolute error (MAE) and the mean bias (MB) on average and also reduced the variability of performances between each bootstrap sample. Only

588 gradient boosting showed a slightly deteriorated performance by having more than 30 % of the 80 %  
589 PWS data used for training (96 PWS) – probably due to overfitting. Below 40 PWS, all models  
590 performed poorly. We also showed that training the regressors over official MIDAS data only led  
591 to a poor bias correction for both summertime average daily minimum and mean temperatures.  
592 For the maximum, no clear benefit was demonstrable, which was also the case with PWS and  
593 which could be explained by the lower UHII during hot hours of the day, as discussed in the  
594 manuscript. We argue that this general outcome is explicable by the standard location of MIDAS  
595 weather stations – typically located in open parks or fields – which would explain why the bias  
596 correction for minimum temperatures further increases the cool bias already existing in WRF. This  
597 supports the use of PWS for bias correction of urban temperatures for two reasons: first, the need  
598 for a sufficiently dense network of weather stations in urban environments; second, the necessity  
599 of weather stations located in typical built-up environments to accurately represent the effect of  
600 built-up surfaces on the local climate.



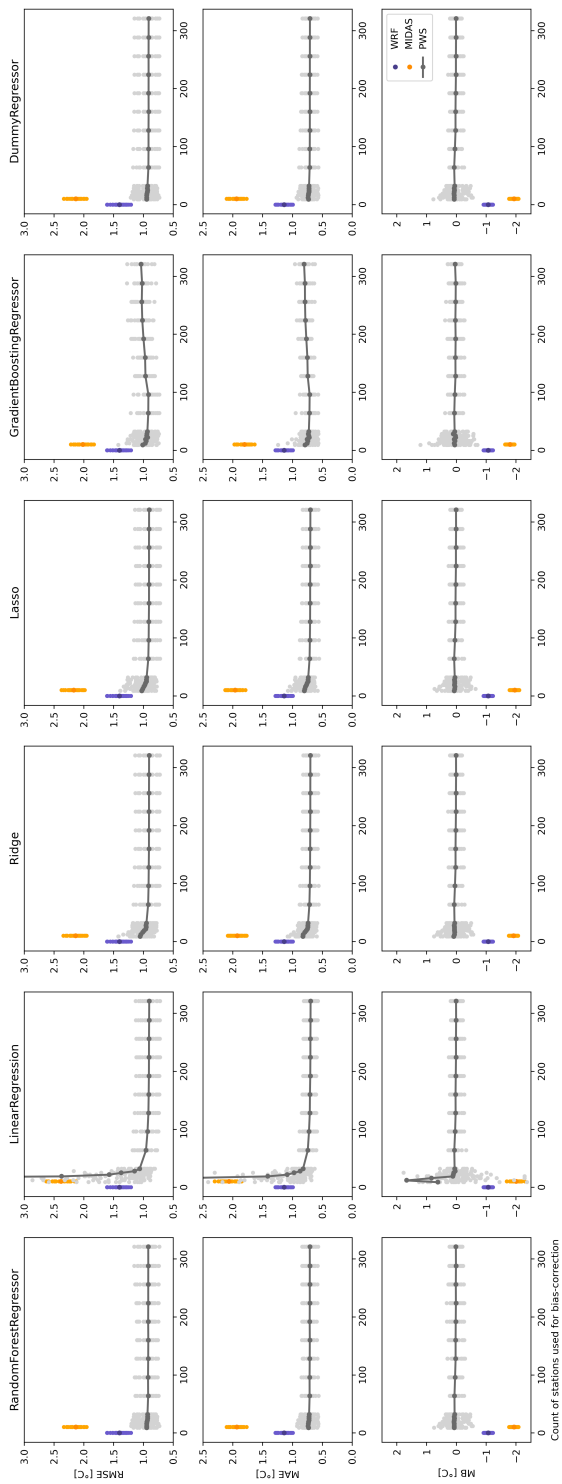


Fig. B1. Regressors performance for bias correction of the summer average daily-minimum air temperature depending on the amount of weather stations' used for training. The performance is evaluated with the mean absolute error (MAE; in °C), root mean squared error (RMSE; °C) and mean bias (MB; °C). Blue dots represent the WRF model performance prior to bias-correction, in orange are the performance of the WRF model after bias-correction using MIDAS official weather stations, and in grey are the performance of the WRF model after bias-correction using subsets of the available Netatmo personal weather stations. Small lighter dots are representative of performances measured at each bootstrapping steps (n=25) and large darker dots are the average of all bootstraps. Here the WRF model was run with the Bougeault-Lacarrère boundary layer scheme (BouLac).

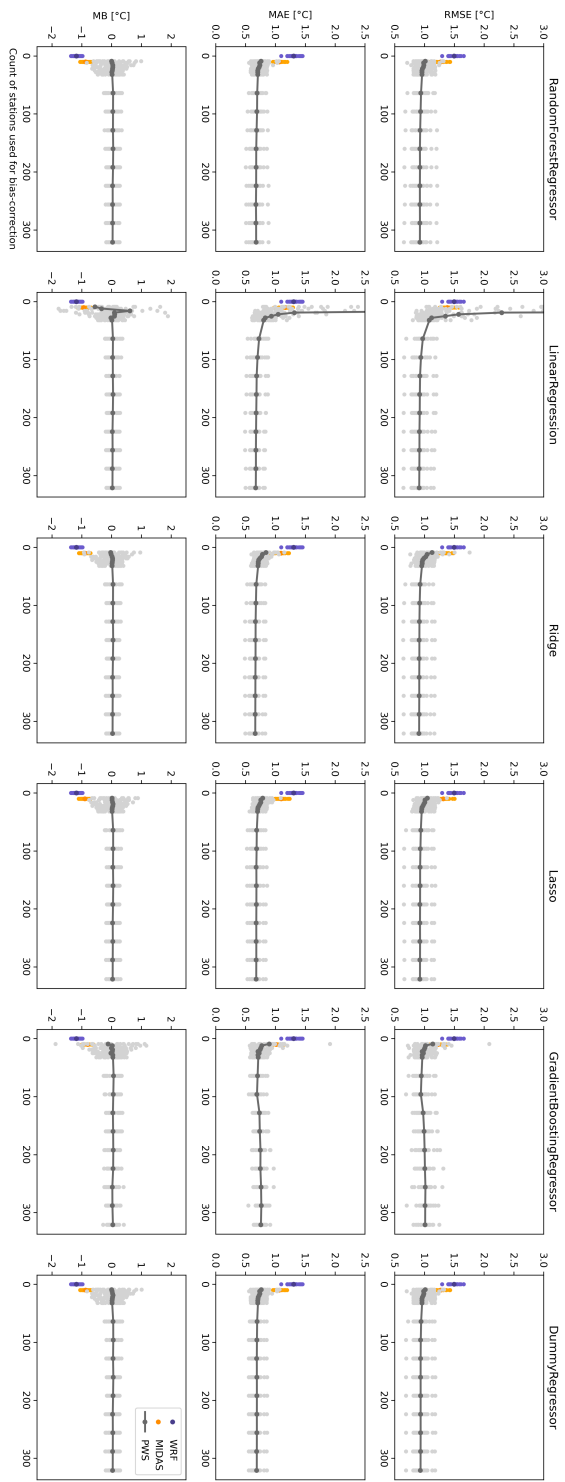


Fig. B2. Same as Fig. B1 but for summer average daily-mean temperatures.

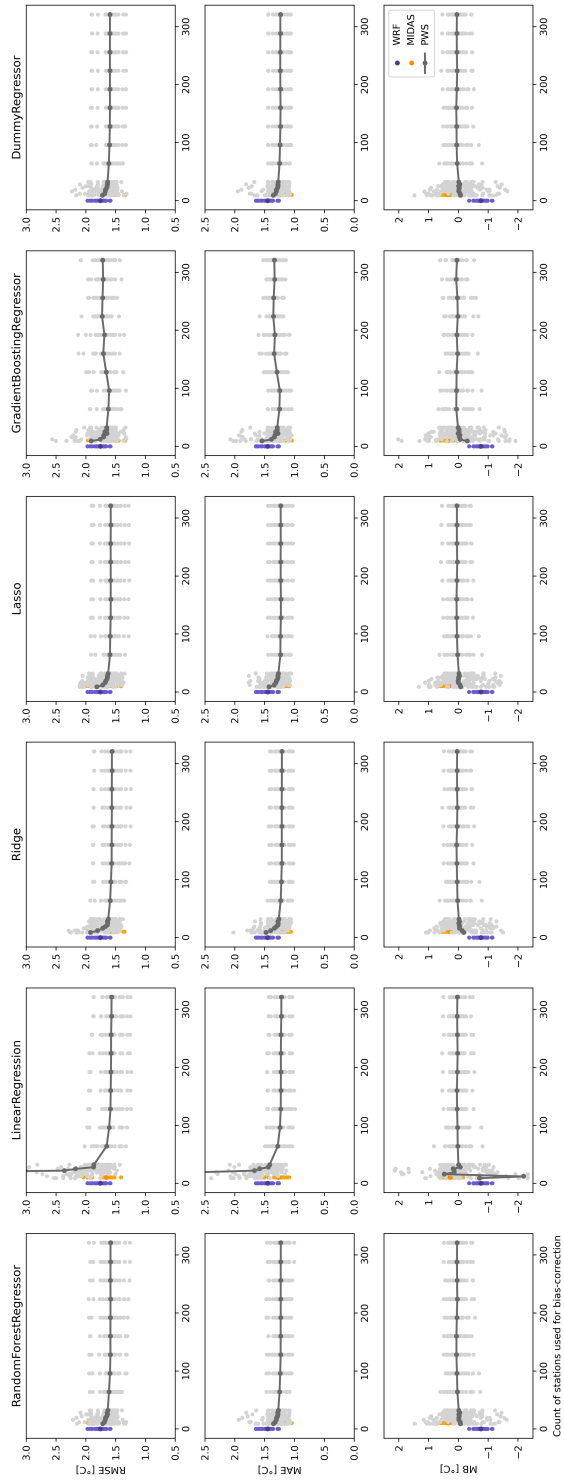


FIG. B3. Same as Fig. B1 but for summer average daily-maximum temperatures.

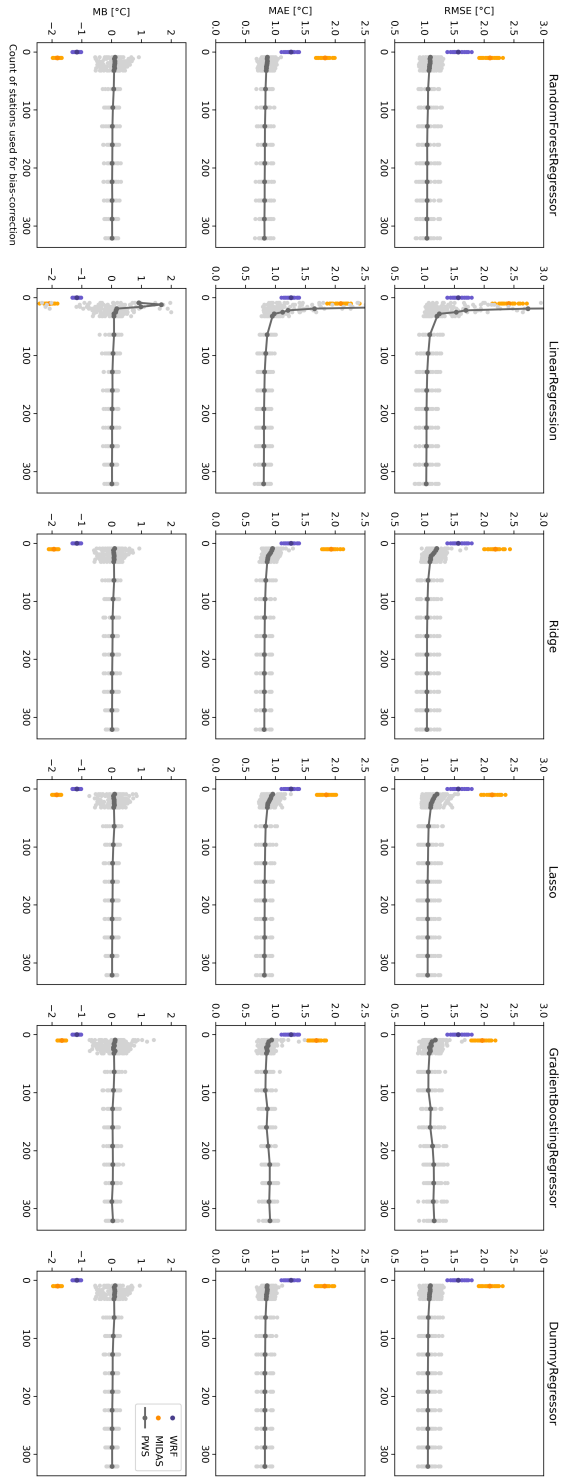


Fig. B4. Same as Fig. B1 but WRF model used the YSU planetary boundary layer scheme.

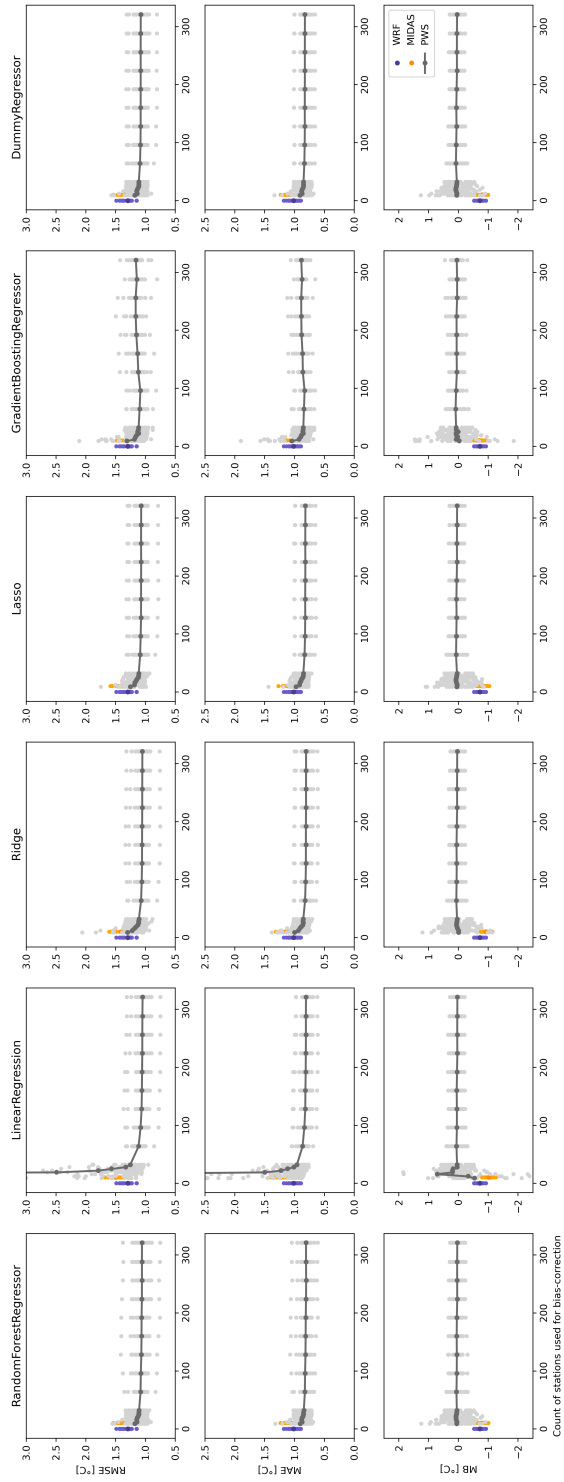


FIG. B5. Same as Fig. B4 but for summer average daily-mean temperatures.

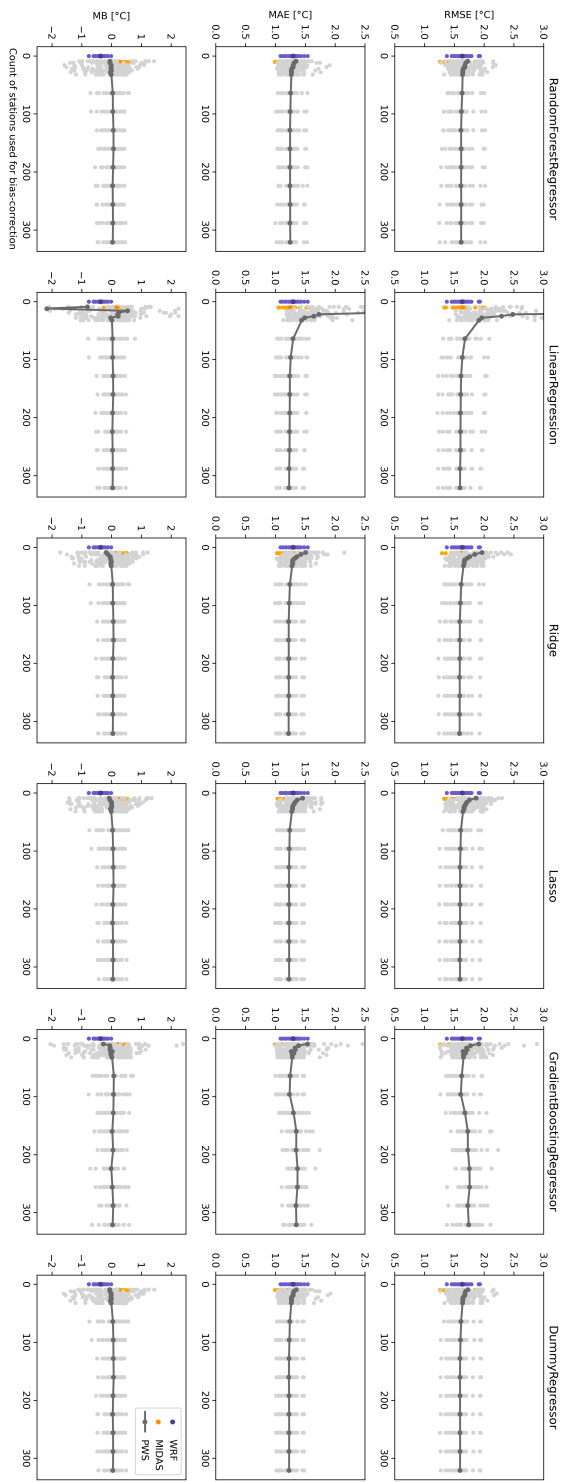


Fig. B6. Same as Fig. B4 but for summer average daily-mean temperatures.

APPENDIX C

Additional Figures and Tables

This section presents all the figures that are not given in the main text.

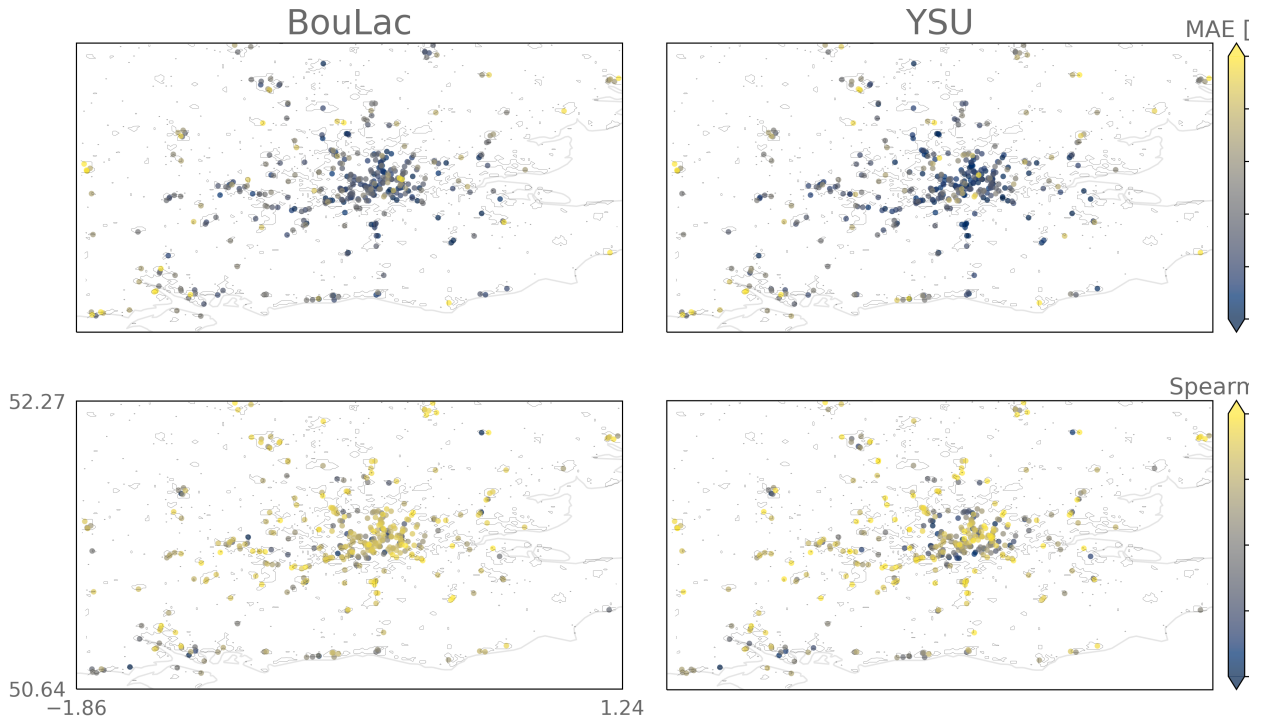


FIG. C1. Same as figure 3, but for MAE and Spearman's r.

TABLE C1. Performance metrics used in Figure 4 for the model using Boulac prior to the bias-correction (WRF) and all the different regressors (random forest: RF; linear regression: LR; Ridge regression: RD; Lasso regression: LA; gradient boosting: GB; and dummy regression: DU). The different regressions are assigned a suffix: “avg” for regressions that were trained on the summer time-mean average of daily-minimum, -mean or -maximum temperatures, and “1step” for those that were trained with the temperatures at each daily time-step.

<b>Boulac</b>													
	<b>WRF</b>	<b>RF<sub>avg</sub></b>	<b>RF<sub>1step</sub></b>	<b>LR<sub>avg</sub></b>	<b>LR<sub>1step</sub></b>	<b>RD<sub>avg</sub></b>	<b>RD<sub>1step</sub></b>	<b>LA<sub>avg</sub></b>	<b>LA<sub>1step</sub></b>	<b>GB<sub>avg</sub></b>	<b>GB<sub>1step</sub></b>	<b>DU<sub>avg</sub></b>	<b>DU<sub>1step</sub></b>
<b>MEAN</b>													
<b>RMSE</b>	1.54	0.95	1.04	0.94	1.03	0.94	1.03	0.95	1.04	1.01	1.04	0.96	1.04
<b>MAE</b>	1.34	0.69	0.75	0.69	0.75	0.68	0.75	0.69	0.75	0.74	0.75	0.7	0.76
<b>MB</b>	-1.2	0.01	0.23	0	0.23	0	0.23	0	0.23	0	0.23	0.01	0.23
<b>Pearson r<sup>2</sup></b>	0.11	0.09	0.07	0.09	0.07	0.1	0.07	0.1	0.07	0.06	0.06	0.11	0.08
<b>Spearman r</b>	0.37	0.33	0.32	0.33	0.31	0.36	0.32	0.36	0.32	0.29	0.32	0.37	0.33 0.88
<b>MIN</b>													
<b>RMSE</b>	1.42	0.93	0.94	0.92	0.93	0.92	0.93	0.92	0.93	1.01	0.96	0.92	0.94
<b>MAE</b>	1.15	0.72	0.73	0.71	0.72	0.71	0.72	0.71	0.73	0.79	0.74	0.71	0.73
<b>MB</b>	-1.08	0.01	0.02	0	0.02	0	0.02	0	0.02	0.04	0.02	0	0.02
<b>Pearson r<sup>2</sup></b>	0.18	0.15	0.16	0.15	0.16	0.16	0.16	0.16	0.16	0.1	0.15	0.17	0.17
<b>Spearman r</b>	0.46	0.42	0.43	0.43	0.42	0.44	0.43	0.44	0.43	0.34	0.41	0.46	0.44
<b>MAX</b>													
<b>RMSE</b>	1.78	1.6	1.81	1.58	1.8	1.57	1.8	1.59	1.8	1.65	1.82	1.6	1.82
<b>MAE</b>	1.48	1.24	1.33	1.22	1.32	1.22	1.31	1.23	1.32	1.28	1.35	1.24	1.33
<b>MB</b>	-0.79	0	0.52	0	0.52	0	0.53	0.01	0.52	0	0.51	0.01	0.53
<b>Spearman r</b>	0.08	0.07	0.02	0.08	0.02	0.09	0.02	0.08	0.02	0.05	0.01	0.08	0.03
<b>Spearman r</b>	0.29	0.26	0.16	0.29	0.16	0.3	0.19	0.27	0.19	0.23	0.14	0.28	0.2



TABLE C2. Performance metrics used in Figure 4 for the model using YSU prior to the bias-correction (WRF) and all the different regressors (random forest: RF; linear regression: LR; Ridge regression: RD; Lasso regression: LA; gradient boosting: GB; and dummy regression: DU). The different regressions are assigned a suffix: “avg” for regressions that were trained on the summer time-mean average of daily-minimum, -mean or -maximum temperatures, and “tstep” for those that were trained with the temperatures at each daily time-step.

YSU														
	WRF	RF <sub>avg</sub>	RF <sub>tstep</sub>	LR <sub>avg</sub>	LR <sub>tstep</sub>	RD <sub>avg</sub>	RD <sub>tstep</sub>	LA <sub>avg</sub>	LA <sub>tstep</sub>	GB <sub>avg</sub>	GB <sub>tstep</sub>	DU <sub>avg</sub>	DU <sub>tstep</sub>	
MEAN														
RMSE	1.33	1.09	1.16	1.07	1.16	1.08	1.16	1.09	1.18	1.15	1.17	1.1	1.19	
MAE	1.04	0.82	0.86	0.82	0.86	0.82	0.87	0.83	0.89	0.87	0.85	0.84	0.89	
MB	-0.76	0	0.17	0	0.17	0	0.17	0.01	0.16	0.02	0.17	0.01	0.17	
Pearson r <sup>2</sup>	0.09	0.07	0.07	0.07	0.07	0.08	0.07	0.08	0.07	0.05	0.07	0.09	0.07	
Spearman r	0.32	0.28	0.3	0.28	0.29	0.3	0.29	0.29	0.28	0.25	0.3	0.32	0.3	
MIN														
RMSE	1.58	1.05	1.06	1.04	1.06	1.05	1.07	1.06	1.09	1.12	1.09	1.06	1.09	
MAE	1.27	0.83	0.83	0.81	0.82	0.82	0.83	0.82	0.84	0.88	0.84	0.83	0.84	
MB	-1.17	0	-0.03	0	-0.03	0	-0.03	0	-0.03	0.04	-0.02	0	-0.03	
Pearson r <sup>2</sup>	0.11	0.1	0.1	0.1	0.1	0.1	0.1	0.1	0.09	0.08	0.1	0.11	0.1	
Spearman r	0.37	0.35	0.37	0.35	0.36	0.34	0.35	0.34	0.34	0.31	0.36	0.36	0.35	
MAX														
RMSE	1.65	1.63	1.82	1.6	1.81	1.6	1.8	1.6	1.8	1.67	1.82	1.6	1.8	
MAE	1.32	1.25	1.33	1.23	1.31	1.23	1.31	1.23	1.31	1.29	1.34	1.23	1.31	
MB	-0.41	0	0.49	0	0.5	0	0.5	0.01	0.49	-0.01	0.49	0.01	0.5	
Pearson r <sup>2</sup>	0.09	0.07	0.04	0.08	0.05	0.09	0.05	0.09	0.05	0.06	0.04	0.09	0.05	
Spearman r	0.32	0.27	0.23	0.29	0.24	0.31	0.25	0.3	0.26	0.25	0.22	0.31	0.26	

Average model's bias correction of daily min temperature after 25 bootstrap

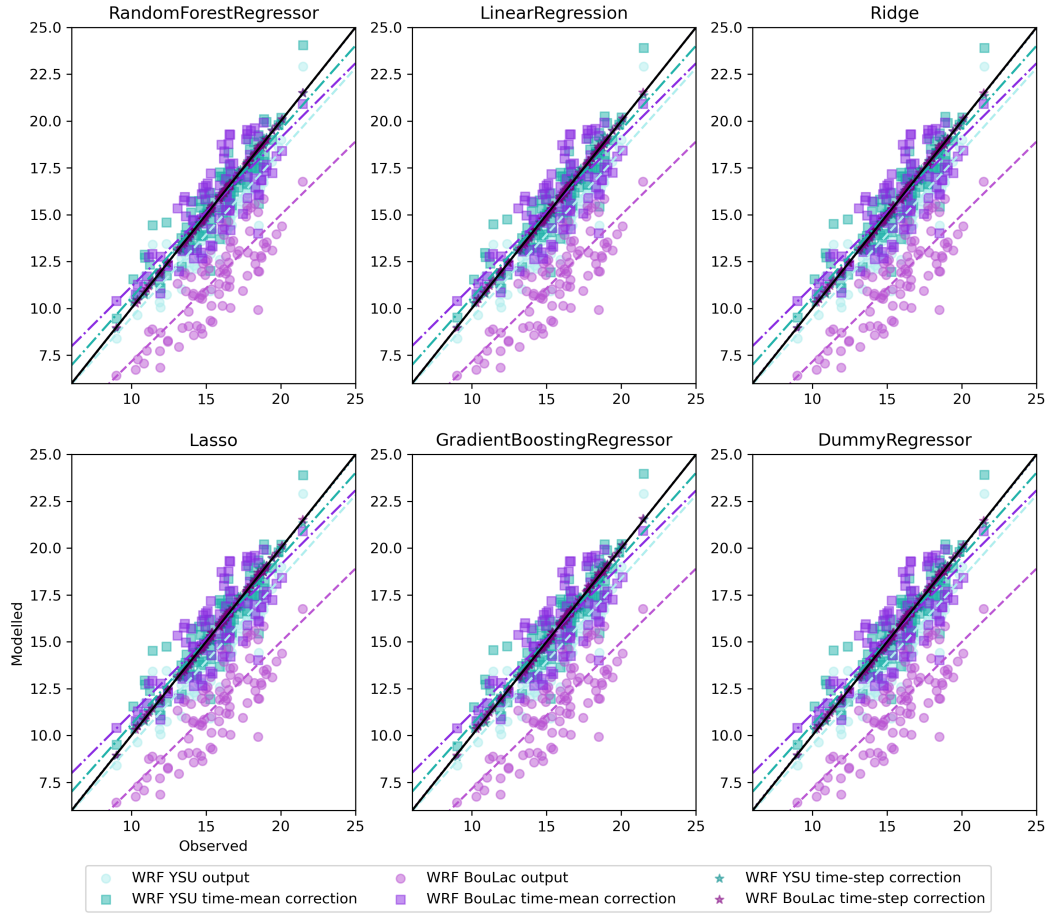


FIG. C2. Average modelled daily minimum air temperature at 2 m against observed at citizens' personal weather stations locations show that all machine learning regressors perform a similar bias-correction on average. In blue, modelled temperatures at 2 m are from the model simulation that used the Yonsei University (YSU) planetary boundary layer scheme before the bias correction (circles), after the summer time-mean bias correction (squares) and after the daily time-step bias correction (stars). In purple, the same values are given for the simulation which used the Bougeault-Lacarrère (BouLac) scheme. Dashed lines represent the least squares polynomial fitted lines and the black full line represents the identity line.

Average model's bias correction of daily max temperature after 25 bootstrap

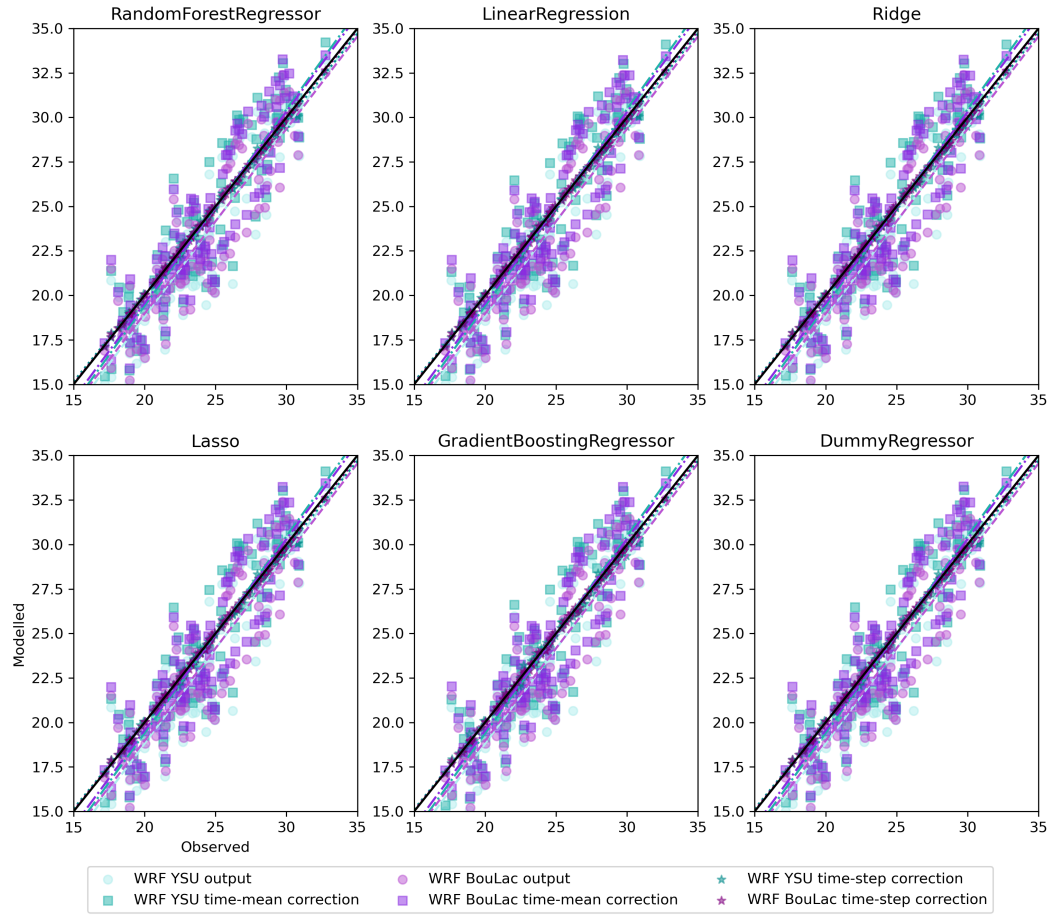


FIG. C3. Same as figure C2, but for daily maximum temperatures.

Average model's bias correction of daily mean temperature after 25 bootstrap

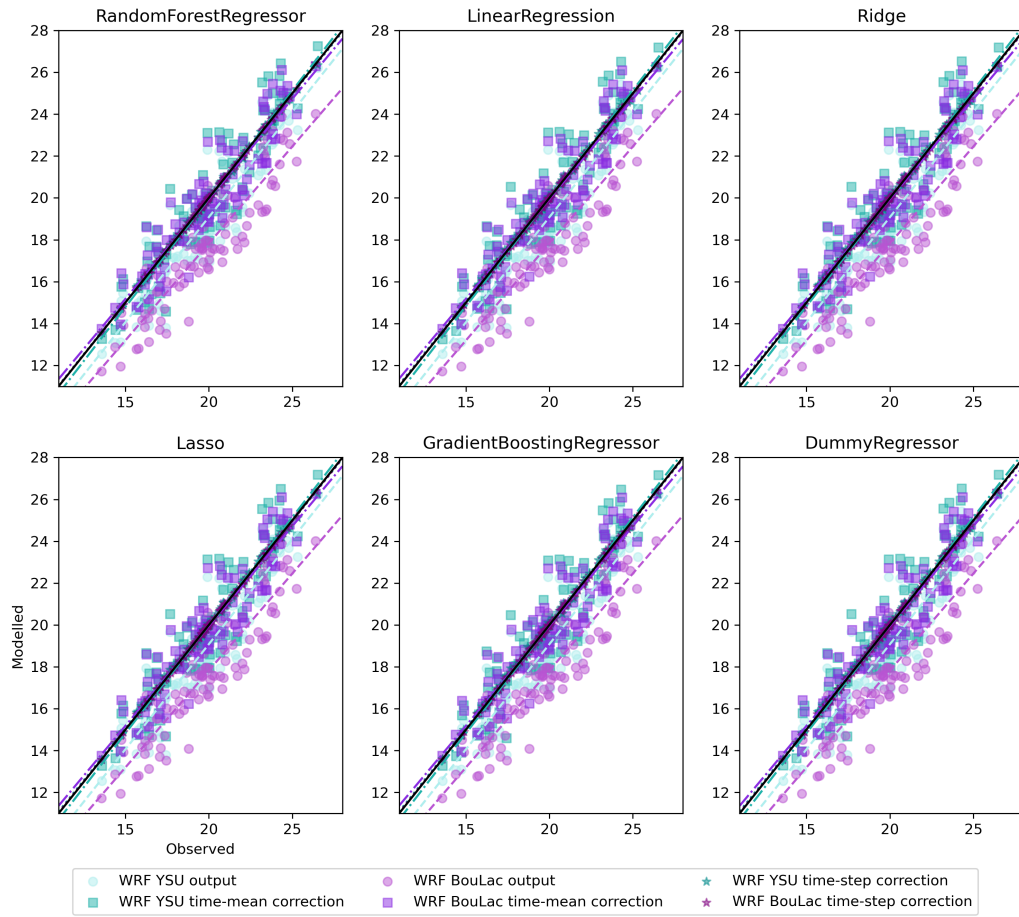


FIG. C4. Same as figure C2, but for daily mean temperatures.

# Modelled temperatures and respective bias-corrections with multiple regressors

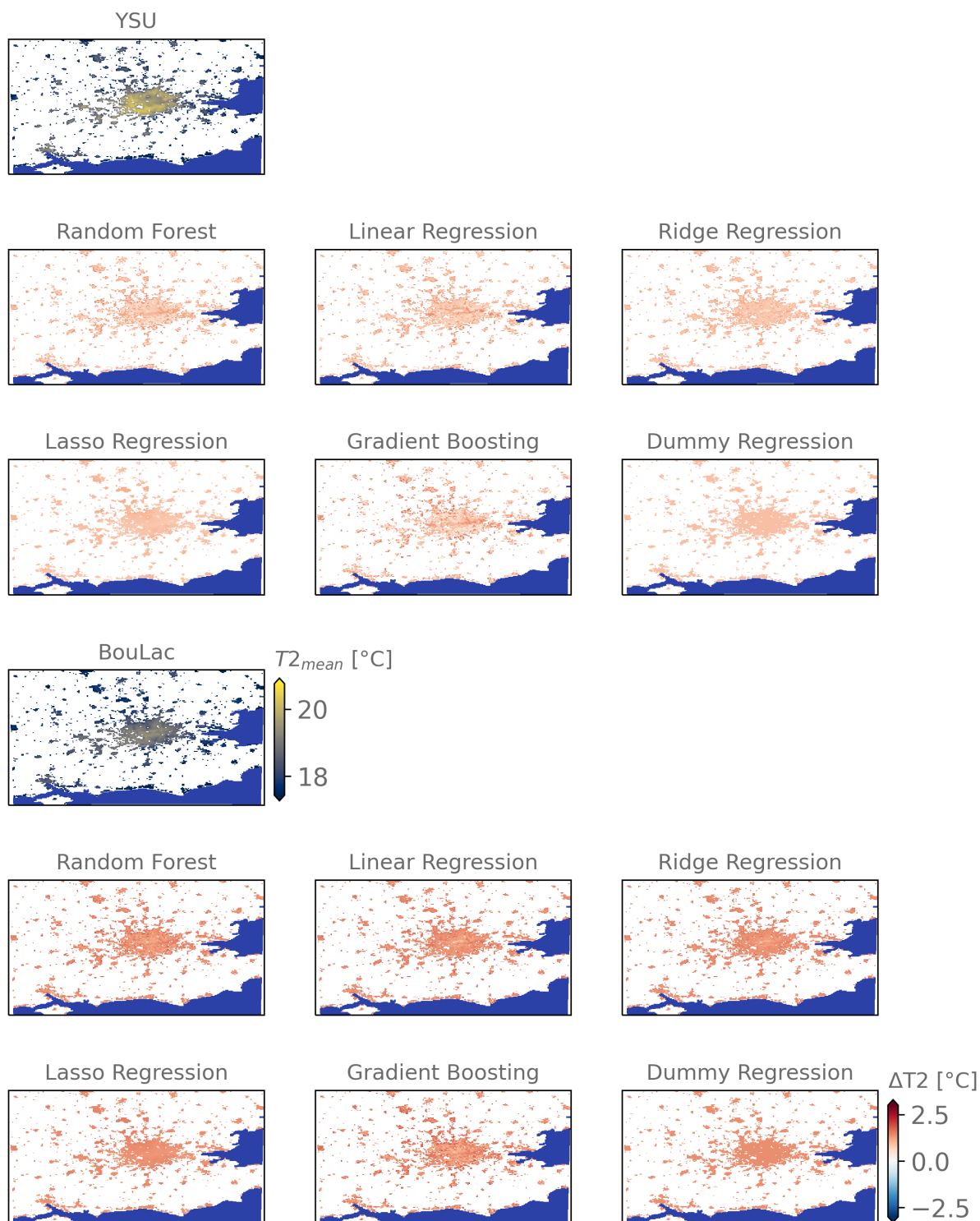


FIG. C5. Same as figure 5, but for daily mean temperatures.

# Modelled temperatures and respective bias-corrections with multiple regressors

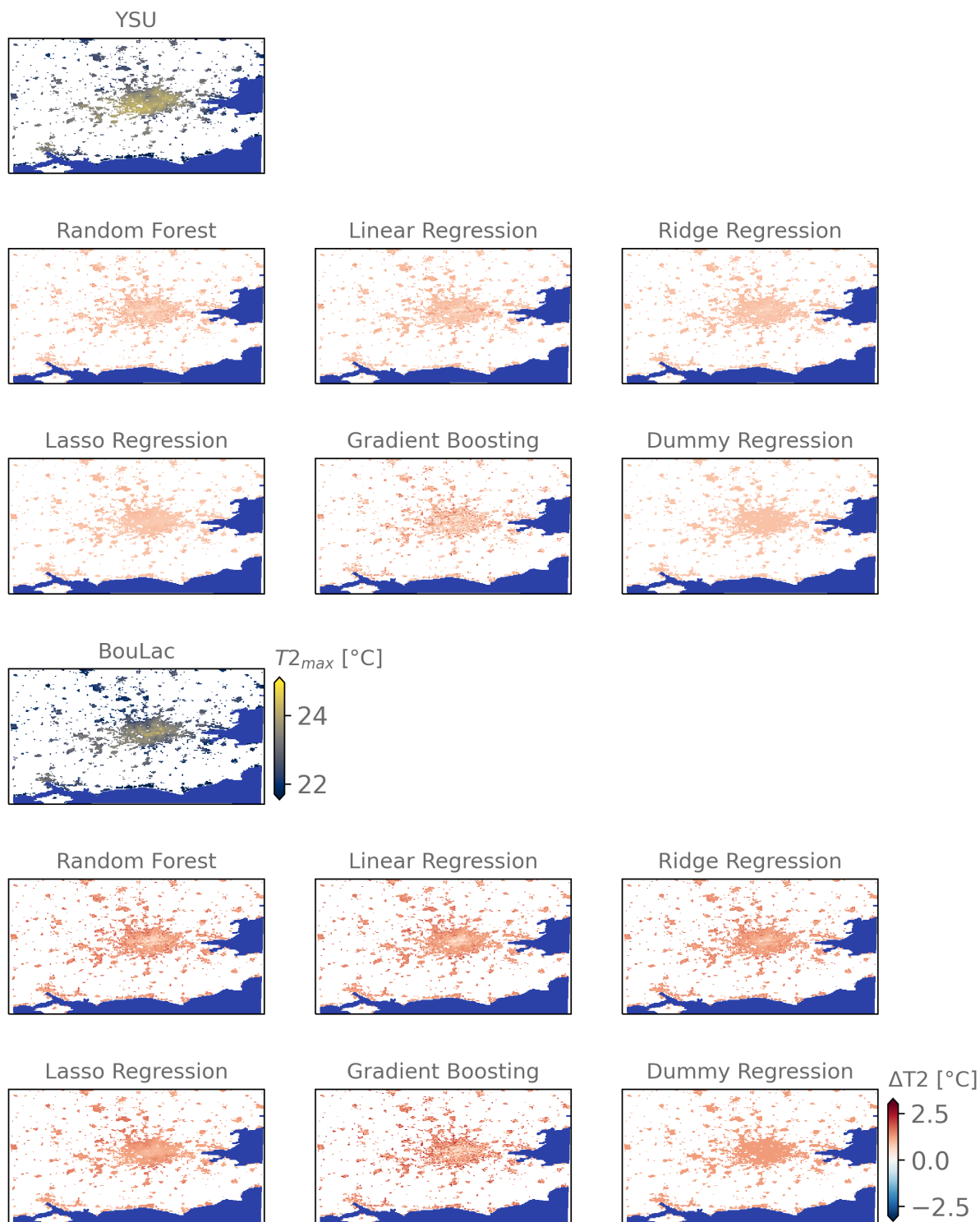


FIG. C6. Same as figure 5, but for daily maximum temperatures.

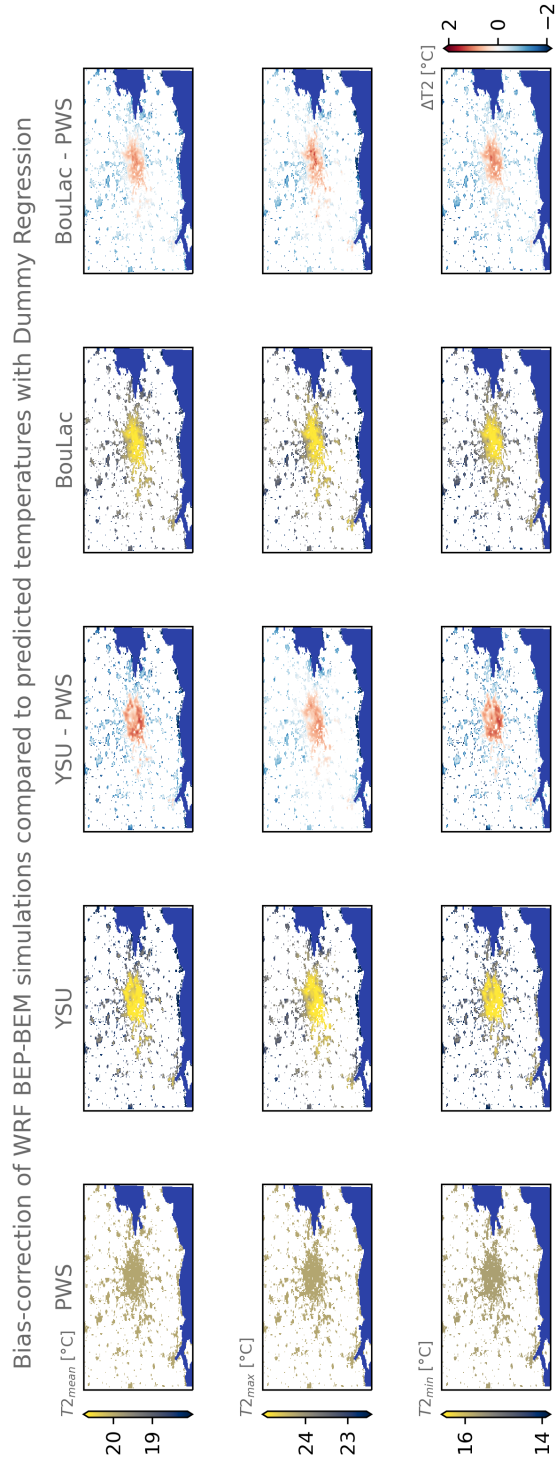


FIG. C7. Same as figure 6, but for dummy regression.

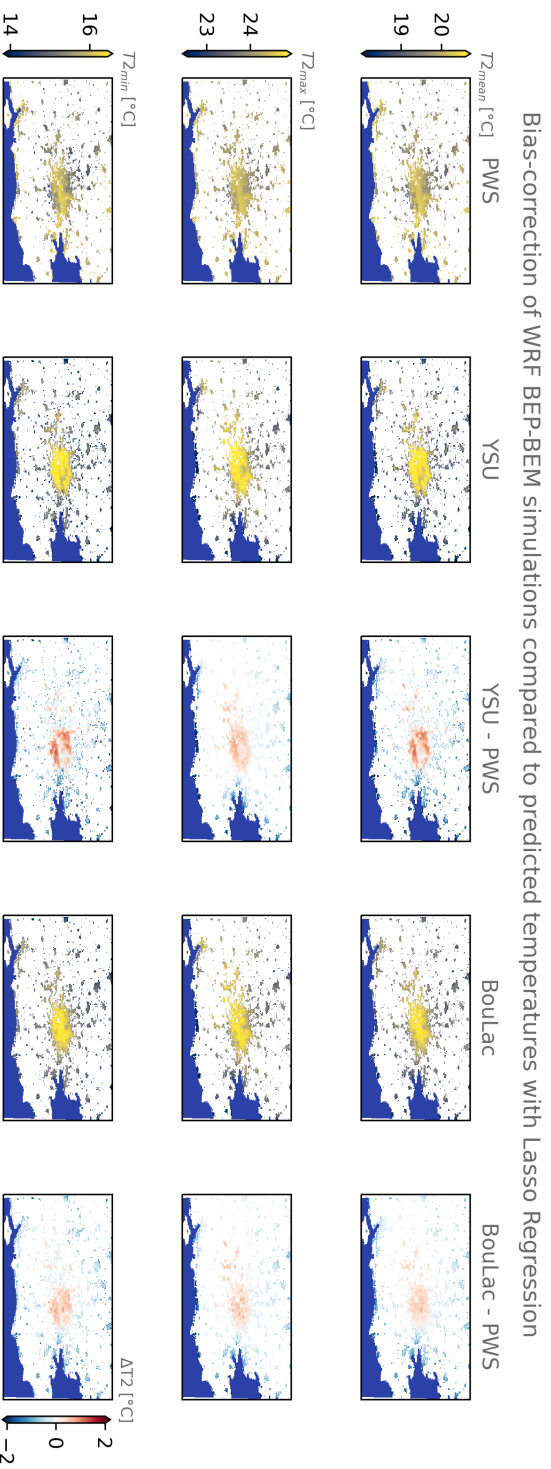


FIG. C8. Same as figure 6, but for Lasso regression.



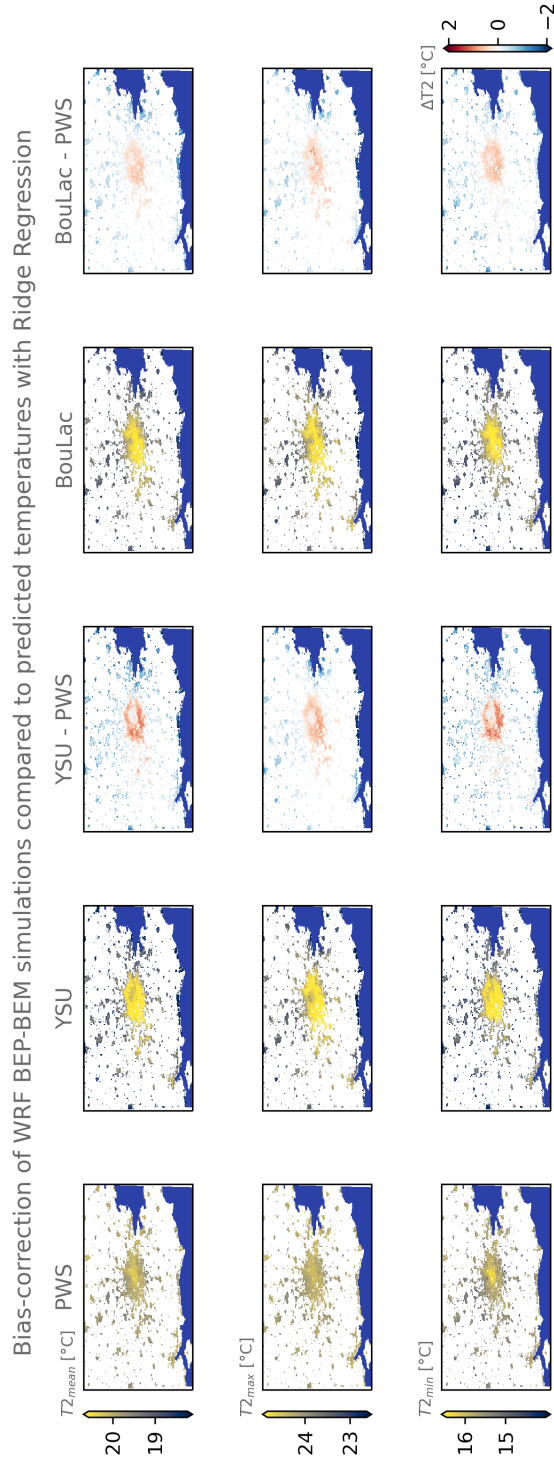


FIG. C9. Same as figure 6, but for Ridge regression.

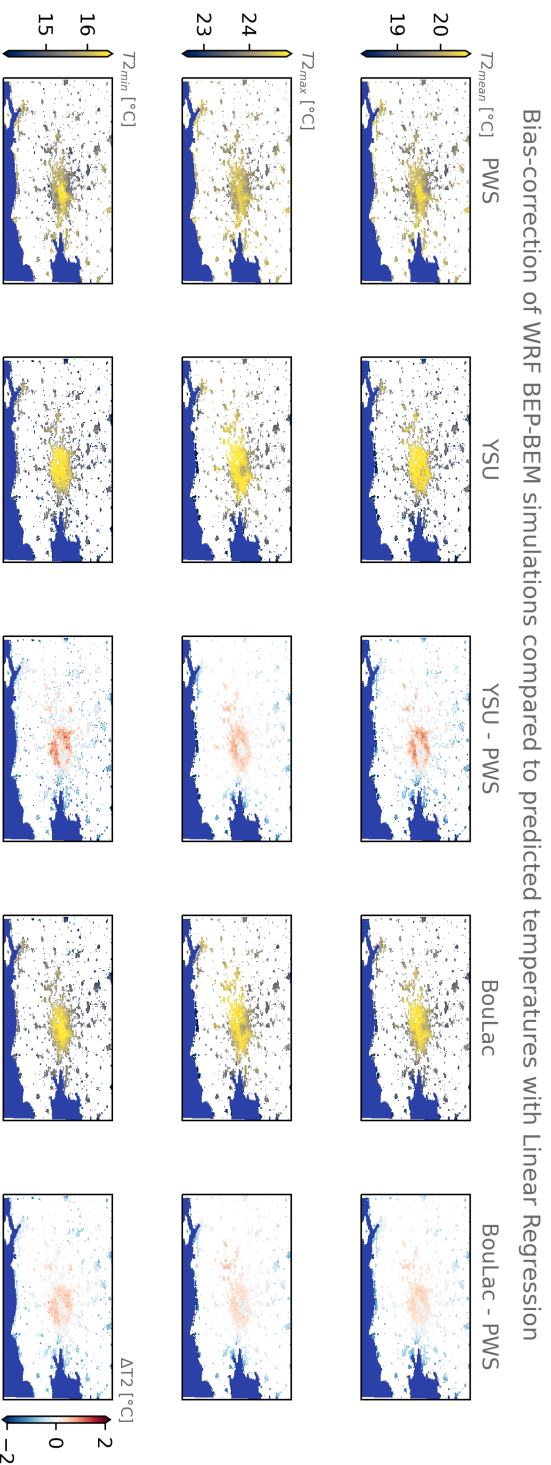


FIG. C10. Same as figure 6, but for linear regression.

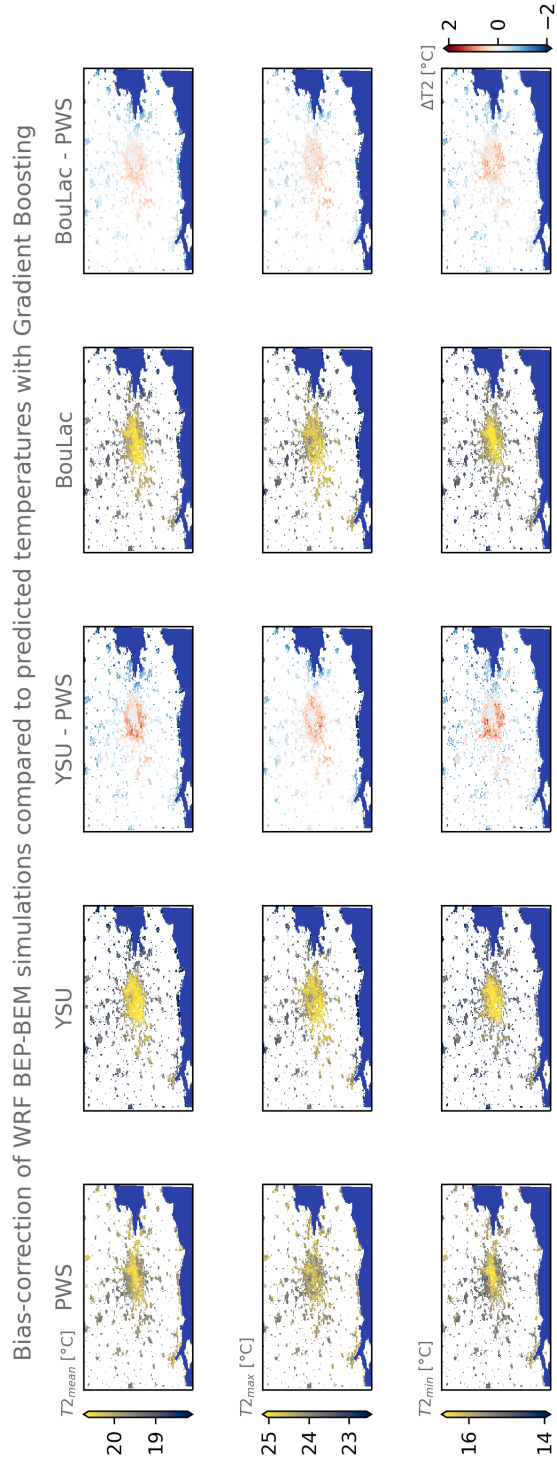


FIG. C11. Same as figure 6, but for gradient boosting regression.

## References

- Bassett, R., P. Young, G. Blair, F. Samreen, and W. Simm, 2020: A large ensemble approach to quantifying internal model variability within the wrf numerical model. *Journal of Geophysical Research: Atmospheres*, **125** (7), e2019JD031 286.
- Benjamin, K., Z. Luo, and X. Wang, 2021: Crowdsourcing urban air temperature data for estimating urban heat island and building heating/cooling load in london. *Energies*, **14** (16), 5208.
- Bougeault, P., and P. Lacarrere, 1989: Parameterization of orography-induced turbulence in a mesobeta-scale model. *Monthly weather review*, **117** (8), 1872–1890.
- Brisson, E., M. Demuzere, and N. Van Lipzig, 2015: Modelling strategies for performing convection-permitting climate simulations. *Meteorologische Zeitschrift*, **25** (2), 149–163.
- Broadbent, A. M., J. Declet-Barreto, E. S. Krayenhoff, S. L. Harlan, and M. Georgescu, 2022: Targeted implementation of cool roofs for equitable urban adaptation to extreme heat. *Science of the Total Environment*, **811**, 151 326.
- Brousse, O., A. Martilli, M. Foley, G. Mills, and B. Bechtel, 2016: Wudapt, an efficient land use producing data tool for mesoscale models? integration of urban lcz in wrf over madrid. *Urban Climate*, **17**, 116–134.
- Brousse, O., C. Simpson, A. Poorthuis, and C. Heaviside, 2023: Unequal distributions of crowd-sourced weather data in england and wales. Preprint available online, <https://doi.org/10.21203/rs.3.rs-2715073/v1>.
- Brousse, O., C. Simpson, N. Walker, D. Fenner, F. Meier, J. Taylor, and C. Heaviside, 2022: Evidence of horizontal urban heat advection in london using six years of data from a citizen weather station network. *Environmental Research Letters*, **17** (4), 044 041.
- Chapman, L., C. Bell, and S. Bell, 2017: Can the crowdsourcing data paradigm take atmospheric science to a new level? a case study of the urban heat island of london quantified using netatmo weather stations. *International Journal of Climatology*, **37** (9), 3597–3605.

Ching, J., and Coauthors, 2018: Wudapt: An urban weather, climate, and environmental modeling infrastructure for the anthropocene. *Bulletin of the American Meteorological Society*, **99** (9), 1907–1924.

De Vos, L., A. Droste, M. Zander, A. Overeem, H. Leijnse, B. Heusinkveld, G. Steeneveld, and R. Uijlenhoet, 2020: Hydrometeorological monitoring using opportunistic sensing networks in the amsterdam metropolitan area. *Bulletin of the American Meteorological Society*, **101** (2), E167–E185.

Demuzere, M., D. Argüeso, A. Zonato, and J. Kittner, 2021: W2w: A python package that injects wudapt’s local climate zone information in wrf (version v0.1.1). Retrieved online, <https://pypi.org/project/w2w/>.

Demuzere, M., B. Bechtel, A. Middel, and G. Mills, 2019: Mapping europe into local climate zones. *PloS one*, **14** (4), e0214474.

Demuzere, M., J. Kittner, A. Martilli, G. Mills, C. Moede, I. D. Stewart, J. van Vliet, and B. Bechtel, 2022: A global map of local climate zones to support earth system modelling and urban scale environmental science. *Earth System Science Data Discussions*, 1–57.

Demuzere, M., and Coauthors, 2017: Impact of urban canopy models and external parameters on the modelled urban energy balance in a tropical city. *Quarterly Journal of the Royal Meteorological Society*, **143** (704), 1581–1596.

Dudhia, J., 1989: Numerical study of convection observed during the winter monsoon experiment using a mesoscale two-dimensional model. *Journal of Atmospheric Sciences*, **46** (20), 3077–3107.

Fenner, D., B. Bechtel, M. Demuzere, J. Kittner, and F. Meier, 2021: Crowdqc+—a quality-control for crowdsourced air-temperature observations enabling world-wide urban climate applications. *Frontiers in Environmental Science*, 553.

Fenner, D., A. Holtmann, F. Meier, I. Langer, and D. Scherer, 2019: Contrasting changes of urban heat island intensity during hot weather episodes. *Environmental Research Letters*, **14** (12), 124013.

677 Fenner, D., F. Meier, B. Bechtel, M. Otto, and D. Scherer, 2017: Intra and inter ‘local climate  
678 zone’ variability of air temperature as observed by crowdsourced citizen weather stations in  
679 berlin, germany. *10.14279/depositonce-10378*.

680 Georganos, S., and Coauthors, 2021: Geographical random forests: a spatial extension of the  
681 random forest algorithm to address spatial heterogeneity in remote sensing and population  
682 modelling. *Geocarto International*, **36** (2), 121–136.

683 Grassmann, T., A. Napoly, F. Meier, and D. Fenner, 2018: Quality control for crowdsourced data  
684 from cws.

685 Grimmond, C. S. B., and Coauthors, 2011: Initial results from phase 2 of the international urban  
686 energy balance model comparison. *International Journal of Climatology*, **31** (2), 244–272.

687 Gutiérrez, E., J. E. González, A. Martilli, R. Bornstein, and M. Arend, 2015: Simulations of a  
688 heat-wave event in new york city using a multilayer urban parameterization. *Journal of Applied  
689 Meteorology and Climatology*, **54** (2), 283–301.

690 Hammerberg, K., O. Brousse, A. Martilli, and A. Mahdavi, 2018: Implications of employing  
691 detailed urban canopy parameters for mesoscale climate modelling: a comparison between  
692 wudapt and gis databases over vienna, austria. *International Journal of Climatology*, **38**, e1241–  
693 e1257.

694 Heaviside, C., X.-M. Cai, and S. Vardoulakis, 2015: The effects of horizontal advection on the  
695 urban heat island in birmingham and the west midlands, united kingdom during a heatwave.  
696 *Quarterly Journal of the Royal Meteorological Society*, **141** (689), 1429–1441.

697 Hendricks, E. A., J. C. Knierel, and Y. Wang, 2020: Addition of multilayer urban canopy models to  
698 a nonlocal planetary boundary layer parameterization and evaluation using ideal and real cases.  
699 *Journal of Applied Meteorology and Climatology*, **59** (8), 1369–1392.

700 Hollis, D., M. McCarthy, M. Kendon, T. Legg, and I. Simpson, 2019: Haduk-grid—a new uk  
701 dataset of gridded climate observations. *Geoscience Data Journal*, **6** (2), 151–159.

702 Hong, S.-Y., J. Dudhia, and S.-H. Chen, 2004: A revised approach to ice microphysical processes  
703 for the bulk parameterization of clouds and precipitation. *Monthly weather review*, **132** (1),  
704 103–120.

- 705 Hong, S.-Y., and S.-W. Kim, 2008: Stable boundary layer mixing in a vertical diffusion scheme.  
706 *18th Symposium on Boundary Layers and Turbulence B*, Vol. 16, 325.
- 707 Hong, S.-Y., Y. Noh, and J. Dudhia, 2006: A new vertical diffusion package with an explicit  
708 treatment of entrainment processes. *Monthly weather review*, **134** (9), 2318–2341.
- 709 Janjić, Z. I., 2001: Nonsingular implementation of the mellor-yamada level 2.5 scheme in the ncep  
710 meso model.
- 711 Janjić, Z. I., 1994: The step-mountain eta coordinate model: Further developments of the con-  
712 vection, viscous sublayer, and turbulence closure schemes. *Monthly weather review*, **122** (5),  
713 927–945.
- 714 Jiménez, P. A., J. Dudhia, J. F. González-Rouco, J. Navarro, J. P. Montávez, and E. García-  
715 Bustamante, 2012: A revised scheme for the wrf surface layer formulation. *Monthly weather*  
716 *review*, **140** (3), 898–918.
- 717 Kain, J. S., 2004: The kain–fritsch convective parameterization: an update. *Journal of applied*  
718 *meteorology*, **43** (1), 170–181.
- 719 Lauwaet, D., H. Hooyberghs, B. Maiheu, W. Lefebvre, G. Driesen, S. Van Looy, and K. De Ridder,  
720 2015: Detailed urban heat island projections for cities worldwide: dynamical downscaling  
721 cmip5 global climate models. *Climate*, **3** (2), 391–415.
- 722 Lipson, M., S. Grimmond, and M. Best, 2021: Urban-plumber model evaluation project: initial  
723 results. *EGU General Assembly Conference Abstracts*, EGU21–15 230.
- 724 Loridan, T., and C. Grimmond, 2012: Multi-site evaluation of an urban land-surface model: Intra-  
725 urban heterogeneity, seasonality and parameter complexity requirements. *Quarterly Journal of*  
726 *the Royal Meteorological Society*, **138** (665), 1094–1113.
- 727 Maraun, D., and M. Widmann, 2018: *Statistical downscaling and bias correction for climate*  
728 *research*. Cambridge University Press.
- 729 Martilli, A., A. Clappier, and M. W. Rotach, 2002: An urban surface exchange parameterisation  
730 for mesoscale models. *Boundary-layer meteorology*, **104** (2), 261–304.

731 Martilli, A., and Coauthors, 2021: Simulating the meteorology during persistent wintertime  
 732 thermal inversions over urban areas. the case of madrid. *Atmospheric Research*, **263**, 105 789.

733 Masson, V., 2000: A physically-based scheme for the urban energy budget in atmospheric models.  
 734 *Boundary-layer meteorology*, **94 (3)**, 357–397.

735 McCarthy, M., and Coauthors, 2019: Drivers of the uk summer heatwave of 2018. *Weather*, **74 (11)**,  
 736 390–396.

737 Meier, F., D. Fenner, T. Grassmann, M. Otto, and D. Scherer, 2017: Crowdsourcing air temperature  
 738 from citizen weather stations for urban climate research. *Urban Climate*, **19**, 170–191.

739 Mlawer, E. J., S. J. Taubman, P. D. Brown, M. J. Iacono, and S. A. Clough, 1997: Radiative transfer  
 740 for inhomogeneous atmospheres: Rrtm, a validated correlated-k model for the longwave. *Journal*  
 741 *of Geophysical Research: Atmospheres*, **102 (D14)**, 16 663–16 682.

742 Mughal, M. O., X.-X. Li, T. Yin, A. Martilli, O. Brousse, M. A. Dissegna, and L. K. Norford,  
 743 2019: High-resolution, multilayer modeling of singapore’s urban climate incorporating local  
 744 climate zones. *Journal of Geophysical Research: Atmospheres*, **124 (14)**, 7764–7785.

745 Muller, C., L. Chapman, S. Johnston, C. Kidd, S. Illingworth, G. Foody, A. Overeem, and R. Leigh,  
 746 2015: Crowdsourcing for climate and atmospheric sciences: current status and future potential.  
 747 *International Journal of Climatology*, **35 (11)**, 3185–3203.

748 Napoly, A., T. Grassmann, F. Meier, and D. Fenner, 2018: Development and application of a  
 749 statistically-based quality control for crowdsourced air temperature data. *Frontiers in Earth*  
 750 *Science*, **6**, 118.

751 Nazarian, N., and Coauthors, 2022: Integrated assessment of urban overheating impacts on human  
 752 life. *Earth’s Future*.

753 Niu, G.-Y., and Coauthors, 2011: The community noah land surface model with multiparameteri-  
 754 zation options (noah-mp): 1. model description and evaluation with local-scale measurements.  
 755 *Journal of Geophysical Research: Atmospheres*, **116 (D12)**.

756 Oke, T. R., G. Mills, A. Christen, and J. A. Voogt, 2017: *Urban climates*. Cambridge University  
 757 Press.



Oleson, K., G. Anderson, B. Jones, S. McGinnis, and B. Sanderson, 2018: Avoided climate impacts of urban and rural heat and cold waves over the us using large climate model ensembles for rcp8.5 and rcp4.5. *Climatic change*, **146** (3), 377–392.

Pedregosa, F., and Coauthors, 2011: Scikit-learn: Machine learning in python. *the Journal of machine Learning research*, **12**, 2825–2830.

Potgieter, J., N. Nazarian, M. J. Lipson, M. A. Hart, G. Ulpiani, W. Morrison, and K. Benjamin, 2021: Combining high-resolution land use data with crowdsourced air temperature to investigate intra-urban microclimate. *Frontiers in Environmental Science*, 385.

Salamanca, F., A. Krpo, A. Martilli, and A. Clappier, 2010: A new building energy model coupled with an urban canopy parameterization for urban climate simulations—part i. formulation, verification, and sensitivity analysis of the model. *Theoretical and applied climatology*, **99** (3), 331–344.

Salamanca, F., and A. Martilli, 2010: A new building energy model coupled with an urban canopy parameterization for urban climate simulations—part ii. validation with one dimension off-line simulations. *Theoretical and Applied Climatology*, **99** (3), 345–356.

Salamanca, F., A. Martilli, M. Tewari, and F. Chen, 2011: A study of the urban boundary layer using different urban parameterizations and high-resolution urban canopy parameters with wrf. *Journal of Applied Meteorology and Climatology*, **50** (5), 1107–1128.

Salamanca, F., A. Martilli, and C. Yagi e, 2012: A numerical study of the urban heat island over madrid during the desirex (2008) campaign with wrf and an evaluation of simple mitigation strategies. *International Journal of Climatology*, **32** (15), 2372–2386.

Sgoff, C., W. Acevedo, Z. Paschalidi, S. Ulbrich, E. Bauernschubert, T. Kratzsch, and R. Potthast, 2022: Assimilation of crowd-sourced surface observations over germany in a regional weather prediction system. *Quarterly Journal of the Royal Meteorological Society*.

Stewart, I. D., T. R. Oke, and E. S. Krayenhoff, 2014: Evaluation of the ‘local climate zone’ scheme using temperature observations and model simulations. *International journal of climatology*, **34** (4), 1062–1080.

785 Sunter, M., 2021: Midas data user guide for uk land observations, v20210705.

786 Tewari, M., F. Salamanca, A. Martilli, L. Treinish, and A. Mahalov, 2017: Impacts of projected  
 787 urban expansion and global warming on cooling energy demand over a semiarid region. *Atmo-  
 788 spheric Science Letters*, **18** (11), 419–426.

789 UKMO, 2021: Midas open: Uk hourly weather observation data, v202107. centre for envi-  
 790 ronmental data analysis, 08 september 2021. Data retrieved online, [https://doi.org/10.5285/  
 791 3bd7221d4844435dad2fa030f26ab5fd](https://doi.org/10.5285/3bd7221d4844435dad2fa030f26ab5fd).

792 Varentsov, M., D. Fenner, F. Meier, T. Samsonov, and M. Demuzere, 2021: Quantifying local  
 793 and mesoscale drivers of the urban heat island of moscow with reference and crowdsourced  
 794 observations. *Frontiers in Environmental Science*, 543.

795 Venter, Z. S., O. Brousse, I. Esau, and F. Meier, 2020: Hyperlocal mapping of urban air temperature  
 796 using remote sensing and crowdsourced weather data. *Remote Sensing of Environment*, **242**,  
 797 111 791.

798 Venter, Z. S., T. Chakraborty, and X. Lee, 2021: Crowdsourced air temperatures contrast satellite  
 799 measures of the urban heat island and its mechanisms. *Science Advances*, **7** (22), eabb9569.

800 Virtanen, P., and Coauthors, 2020: Scipy 1.0: fundamental algorithms for scientific computing in  
 801 python. *Nature methods*, **17** (3), 261–272.

802 Wang, J., and X.-M. Hu, 2021: Evaluating the performance of wrf urban schemes and pbl schemes  
 803 over dallas–fort worth during a dry summer and a wet summer. *Journal of Applied Meteorology  
 804 and Climatology*, **60** (6), 779–798.

805 Wouters, H., M. Demuzere, U. Blahak, K. Fortuniak, B. Maiheu, J. Camps, D. Tieleman, and  
 806 N. P. van Lipzig, 2016: The efficient urban canopy dependency parametrization (sury) v1. 0 for  
 807 atmospheric modelling: description and application with the cosmo-clm model for a belgian  
 808 summer. *Geoscientific Model Development*, **9** (9), 3027–3054.

809 Wouters, H., and Coauthors, 2017: Heat stress increase under climate change twice as large in  
 810 cities as in rural areas: A study for a densely populated midlatitude maritime region. *Geophysical  
 811 Research Letters*, **44** (17), 8997–9007.

- 812 Yang, J., and E. Bou-Zeid, 2019: Scale dependence of the benefits and efficiency of green and cool  
813 roofs. *Landscape and urban planning*, **185**, 127–140.
- 814 Yang, Z.-L., and Coauthors, 2011: The community noah land surface model with multiparam-  
815 eterization options (noah-mp): 2. evaluation over global river basins. *Journal of Geophysical*  
816 *Research: Atmospheres*, **116** (D12).
- 817 Zängl, G., D. Reinert, P. Rípodas, and M. Baldauf, 2015: The icon (icosahedral non-hydrostatic)  
818 modelling framework of dwd and mpi-m: Description of the non-hydrostatic dynamical core.  
819 *Quarterly Journal of the Royal Meteorological Society*, **141** (687), 563–579.
- 820 Zonato, A., A. Martilli, S. Di Sabatino, D. Zardi, and L. Giovannini, 2020: Evaluating the  
821 performance of a novel wudapt averaging technique to define urban morphology with mesoscale  
822 models. *Urban Climate*, **31**, 100 584.
- 823 Zonato, A., A. Martilli, P. A. Jimenez, J. Dudhia, D. Zardi, and L. Giovannini, 2022: A new  $k-\epsilon$   
824 turbulence parameterization for mesoscale meteorological models. *Monthly Weather Review*.
- 825 Zumwald, M., B. Knüsel, D. N. Bresch, and R. Knutti, 2021: Mapping urban temperature using  
826 crowd-sensing data and machine learning. *Urban Climate*, **35**, 100 739.