Environmental Research

Improved daily estimates of relative humidity at high resolution across Germany: a Random Forest approach --Manuscript Draft--

Manuscript Number:	ER-23-1226			
Article Type:	Research paper			
Section/Category:	Environmental Health & Risk Assessment			
Keywords:	relative humidity; spatiotemporal modeling; machine learning; external validation; exposure assessment; environmental epidemiology			
Corresponding Author:	Nikolaos Nikolaou Helmholtz Center Munich German Research Center for Environmental Health Neuherberg, Bavaria GERMANY			
First Author:	Nikolaos Nikolaou			
Order of Authors:	Nikolaos Nikolaou			
	Laurens M. Bouwer			
	Marco Dallavalle			
	Mahyar Valizadeh			
	Massimo Stafoggia			
	Annette Peters			
	Kathrin Wolf			
	Alexandra Schneider			
Abstract:	The lack of readily available methods for estimating high-resolution near-surface relative humidity (RH) and the incapability of weather stations to fully capture the spatiotemporal variability can lead to exposure misclassification in studies of environmental epidemiology. We therefore aimed to predict German-wide 1 × 1 km daily mean RH during 2000-2021. RH observations, modelled air temperature, precipitation and wind speed as well as remote sensing information on topographic elevation, vegetation, and the true color band composite were incorporated in a Random Forest (RF) model, in addition to date for capturing the temporal variations of the response-explanatory variables relationship. The model achieved high accuracy (R2 = 0.80) and low errors (Root Mean Square Error (RMSE) of 5.42%), calculated via ten-fold cross-validation. A comparison of our RH predictions with measurements from a dense monitoring network in the city of Augsburg, South Germany confirmed the good performance (R2 ≥ 0.84, RMSE ≤ 5.91%). The model displayed high Germanwide RH (21y-average of 79.05%) and high spatial variability across the country, exceeding 15% on yearly averages. Our findings indicate that the proposed RF model is suitable for estimating RH for a whole country in high-resolution and provide a reliable RH dataset for epidemiological analyses and other environmental research purposes.			
Suggested Reviewers:	Itai Kloog Professor, Ben-Gurion University of the Negev ikloog@bgu.ac.il Experience in the field of modeling meteorological parameters Liuhua Shi Assistant Professor, Emory University liuhua.shi@emory.edu Experience in the field of modeling meteorological parameters Timo Lanki Professor, University of Eastern Finland timo.lanki@thl.fi Experience in the field of environmental epidemiology			

Rémy Slama Senior Investigator, Inserm Public Health Thematic Institute Remy.Slama@univ-grenoble-alpes.fr Experience in the field of environmental epidemiology
lan Hough Postdoctoral researcher, University Grenoble Alpes ian.hough@univ-grenoble-alpes.fr Experience in the field of environmental epidemiology
Bin Zhou Researcher, University of Augsburg bin.zhou@med.uni-augsburg.de Experience in the field of modeling meteorological parameters

HELMHOLTZ

Helmholtz Zentrum München Deutsches Forschungszentrum für Gesundheit und Umwelt (GmbH), Postfach 11 29, 85758 Neuherberg

Submission of manuscript: "Improved daily estimates of relative humidity at high resolution across Germany: a Random Forest approach" Nikolaos Nikolaou PhD Student

+49 89 3187 49700 +49 176 377 488 68 nikolaos.nikolaou@helmholtz-muenchen.de

February 08, 2023

Dear Editor(s),

On behalf of all authors, I hereby submit the manuscript entitled "Improved daily estimates of relative humidity at high resolution across Germany: a Random Forest approach" for consideration of publication in *Environmental Research* as an original research article.

Many disciplines of environmental science stand in need of highly-resolved spatiotemporal datasets of meteorological parameters such as of relative humidity. For instance, a main challenge in the field of environmental epidemiology is the occurring error in exposure assessment of participants of cohort studies, as the commonly used weather station observations are incapable to fully capture the spatiotemporal variability of most meteorological exposures, including relative humidity. However, there is a lack of readily available methods for modeling relative humidity in high-resolution country-wide. To contribute to the existing knowledge and tackle this issue for subsequent health-environment analyses in Germany, we aimed to extend and improve the country's spatiotemporal coverage of relative humidity data by using the random forest methodology and go beyond the conventional interpolation of meteorological observations, also incorporating several spatiotemporal predictors from multiple sources, e.g., remote sensing.

Thus, we generated a high resolution relative humidity dataset in the complex geo-climate terrain of Germany and the random forest model achieved high predictive accuracy and low errors. Particular strengths of our work is the applied cross-validation scheme which allowed us to simulate the prediction step of our model as well as the conducted external validation of our model predictions with data from an independent and dense network of weather stations. Most importantly, we introduced a reliable, straightforward and also generalizable approach for estimating relative humidity in different spatial settings, which would be of high interest for researchers in various fields of environmental science.

We believe this manuscript is appropriate for publication in *Environmental Research* due to its research focus and novelty. In addition, this work fits greatly in *Environmental Research* in continuation of our recently published paper on the

Helmholtz Zentrum München Deutsches Forschungszentrum für Gesundheit und Umwelt (GmbH), Ingolstädter Landstr. 1, 85764 Neuherberg, Telefon +49 89 3187 0, Fax +49 89 3187 3322, info@helmholtz-munich.de | Geschäftsführung: Prof. Dr. med. Dr. h.c. Matthias H. Tschöp, Kerstin Günther, Daniela Sommer (kom.) | Aufsichtsratsvorsitzende: Prof. Dr. Veronika von Messling | Registergericht: Amtsgericht München HRB 6466 | USt-IdNr. DE 129521671 | Bankverbindung: Münchner Bank eG, Konto-Nr. 2 158 620, BLZ 701 900 00, IBAN DE0470190000002158620, BIC GENODEF1M01

HELMHOLTZ MUNICI)

German-wide and highly-resolved air temperature models (Nikolaou et al., 2022), with the estimated mean air temperature to serve as a very important predictor for our relative humidity model as well. Best of our knowledge, this is the first highly-resolved relative humidity dataset of high performance, also externally validated and generated for more than two decades in a country level. This manuscript may be of particular interest to the journal's audience as it addresses a very important challenge in the field of environmental epidemiology.

This work is original. All of the authors have read and approved the paper and it has not been published previously nor is it being considered by any other peer-reviewed journal. Our study does not include human subjects. We have no conflicts of interest to disclose.

Please address all correspondence concerning this manuscript to: <u>nikolaos.nikolaou@helmholtz-muenchen.de</u> Thank you a lot for your consideration.

Yours Sincerely, Nikolaos Nikolaou Institute of Epidemiology, Helmholtz Zentrum München - German Research Center for Environmental Health, Ingolstädter Landstr. 1, D-85764 Neuherberg, Germany

References

Nikolaou, N.; Dallavalle, M.; Stafoggia, M.; Bouwer, L. M.; Peters, A.; Chen, K.; Wolf K.; Schneider, A., 2022. Highresolution spatiotemporal modeling of daily near-surface air temperature in Germany over the period 2000–2020. Environ. Res. 115062. <u>https://doi.org/10.1016/j.envres.2022.115062</u>

Highlights

- We estimated daily 1 × 1 km near-surface relative humidity (RH) in Germany, 2000-2021
- The random forest model achieved good performance ($R^2 = 0.80$, RMSE = 5.42 %)
- Validation with external data confirmed the model's high accuracy and low errors
- We propose an RH modeling process generalizable to other study domains / countries
- We provide reliable and highly-resolved RH data for epidemiological studies

1 Improved daily estimates of relative humidity at high resolution across Germany: a Random

- 2 Forest approach
- 3 Nikolaos Nikolaou^{1,2*}, Laurens M. Bouwer³, Marco Dallavalle^{1,2}, Mahyar Valizadeh¹, Massimo Stafoggia⁴,
- 4 Annette Peters^{1,2}, Kathrin Wolf^{1§}, Alexandra Schneider^{1§}
- ⁵ ¹Institute of Epidemiology, Helmholtz Zentrum München, German Research Center for Environmental
- 6 Health, Neuherberg, Germany
- 7 ²Institute for Medical Information Processing, Biometry, and Epidemiology, Pettenkofer School of Public
- 8 Health, LMU Munich, Munich, Germany
- 9 ³Climate Service Center Germany (GERICS), Helmholtz-Zentrum Hereon, Hamburg, Germany
- ⁴Department of Epidemiology, Lazio Regional Health Service ASL Roma 1, Rome, Italy
- 11 *Corresponding author
- 12 Address: Ingolstädter Landstr. 1, D-85764 Neuherberg, Germany
- 13 Phone: +49 176 377 488 68
- 14 E-mail: <u>nikolaos.nikolaou@helmholtz-muenchen.de</u>
- 15 [§]Shared last authorship

16 Abstract

17 The lack of readily available methods for estimating high-resolution near-surface relative humidity (RH) 18 and the incapability of weather stations to fully capture the spatiotemporal variability can lead to 19 exposure misclassification in studies of environmental epidemiology. We therefore aimed to predict 20 German-wide 1 × 1 km daily mean RH during 2000-2021. RH observations, modelled air temperature, 21 precipitation and wind speed as well as remote sensing information on topographic elevation, vegetation, 22 and the true color band composite were incorporated in a Random Forest (RF) model, in addition to date 23 for capturing the temporal variations of the response-explanatory variables relationship. The model 24 achieved high accuracy ($R^2 = 0.80$) and low errors (Root Mean Square Error (RMSE) of 5.42 %), calculated 25 via ten-fold cross-validation. A comparison of our RH predictions with measurements from a dense 26 monitoring network in the city of Augsburg, South Germany confirmed the good performance ($R^2 \ge 0.84$, 27 RMSE ≤ 5.91 %). The model displayed high German-wide RH (21y-average of 79.05 %) and high spatial 28 variability across the country, exceeding 15 % on yearly averages. Our findings indicate that the proposed 29 RF model is suitable for estimating RH for a whole country in high-resolution and provide a reliable RH 30 dataset for epidemiological analyses and other environmental research purposes.

31 Keywords: relative humidity; spatiotemporal modeling; machine learning; external validation; exposure
 32 assessment; environmental epidemiology

33 Acknowledgements

This work was supported by the Helmholtz Climate Initiative (HI-CAM) project, which is funded by the Helmholtz Association's Initiative and Networking Fund, the Helmholtz Information & Data Science Academy (HIDA), financially supported by the HIDA Trainee Network program, and by the Digital Earth project, supported by the Helmholtz Association's Initiative and Networking Fund (funding code ZT-0025).

38 The authors are responsible for the content of this publication.

39 1 Introduction

Relative humidity (RH) refers to the water vapor content of air and quantifies how far the atmosphere is from its saturation point. RH is a key parameter for many fields such as agriculture (Zhang et al., 2015), hydrology (Forootan, 2019) and climatology (Sherwood et al., 2010) as it contributes among others to the soil moisture, the hydrological cycle and the weather and climate conditions. Thus, RH plays an important role in plant and animal life (Xiong et al., 2017) as well as in human comfort and well-being (Davis et al., 2016; Yang et al., 2018).

46 RH has mostly been used as a confounder or effect modifier in studies focusing on air temperature (T_{air}) 47 (Armstrong, 2006; Zeng et al., 2017), or as part of an index, e.g., apparent temperature (Analitis et al., 48 2008). Nevertheless, there is evidence that RH is potentially an independent risk factor for mortality (Ou 49 et al., 2014) and morbidity (Luo et al., 2020). In epidemiology, RH data are usually retrieved from weather 50 monitors. But their locations are irregularly distributed over space, usually in rural areas, and their number 51 is limited. Hence, weather stations are inadequate to fully represent the spatiotemporal RH variations in 52 complex geo-climatic urban and rural landscapes, and by using their observations, error is introduced in 53 the exposure assessment of study participants leading to estimates biased towards the null (Zeger et al., 54 2000). Climate reanalysis data could be an alternative source for environmental health research (Mistry 55 et al., 2022), but the resolution is usually coarser than 9 km and the data fail to capture the city-level 56 exposure variability effectively. We therefore suggest to extend the methods and datasets in order to 57 improve the predictions of RH exposure for people participating in epidemiological studies, such as 58 prospective cohorts with data on the residential addresses of the participants.

59 There is a clear methodological gap in RH modeling, especially for high spatiotemporally-resolved RH 60 predictions and for timespans up to multiple years. Li et al. (2014) mapped RH every 3 hours at 1 km by 61 using a two-step interpolation procedure of re-analysis data based on a partial thin-plate spline (TPS) and

62 simple kriging (Root Mean Square Error (RMSE) = 11.06 %). The traditional interpolation techniques have 63 limited efficiency when mapping meteorological exposures in spatially highly heterogeneous areas, and 64 are characterized by neighboring effects on exposures predictions, without being capable of capturing 65 small-scale and intra-city variations. Li and Zha (2018) used a Random Forest (RF) model and satellite data, 66 to estimate RH during the summer of 2009 ($R^2 = 0.70$, RMSE = 7.4 %). Spatiotemporal predictors which 67 could explain a large amount of the remaining RH variance, e.g., Tair, were not included. Longer periods and more predictors need to be tested to capture the full annual and inter-annual RH variability. For China, 68 69 the RF model had better results than TPS and kriging, but improvements are needed for better RH 70 variability representation, higher prediction accuracy and further temporal extension to the annual level.

Remote sensing data are progressively used in environmental exposures modeling (Rosenfeld et al., 2017; Yao et al., 2022) being publicly available in high spatiotemporal resolution. There is also a growing body of machine learning (ML) methods applied in the field (Jin et al., 2022; Silibello et al., 2021; Stafoggia et al., 2019). The RF algorithm, which consists of a multitude of decision trees, copes greatly with big data, with potentially correlated predictors and their non-linearity, and with overfitting. In addition, RF can be applied without the burden of complex hyperparameters tuning and is also robust against outliers.

The specific objectives of this study were (a) to estimate highly spatiotemporal resolved RH for Germany based on T_{air} and other observation, remote sensing and modelled data by using a RF model, (b) to evaluate the model's performance and (c) to produce a reliable German-wide RH dataset for subsequent epidemiological analyses and various research purposes. Thereby, we aimed to extend the current literature and provide a generalizable method for other countries to produce highly resolved RH datasets.

82 2 Materials and Methods

83 2.1 Study domain

Germany extends in an area of 357,021 km², having a strongly diverse landscape and a high elevation range (-3.54 to 2,962 m). In the south-eastern regions, the climate is classified as warm summer humid continental, while in north-western regions it is characterized as temperate oceanic (Beck et al., 2018b). We divided Germany's land area into 366,536 grid cells of 1 × 1 km resolution, following the European INSPIRE (Infrastructure for Spatial Information in the European Community) standard for gridded datasets and using the Lambert Azimuthal Equal-Area (LAEA) projection, EPSG: 3035 (©GeoBasis-DE/BKG (2021)).

90 2.2 Input data

Large amounts of input data were incorporated in the RF modeling process. We used meteorological
 observations, remote sensing and spatiotemporally resolved modelled data, all retrieved from 2000 to
 2021 across the study area.

94 2.2.1 RH data

95 We used daily mean RH observations (DWD, 2022) from 406 weather stations of the German 96 Meteorological Service (DWD) (**Figure S1**). The RH data has been quality controlled by the DWD and all 97 the needed information such as station location as well as relocations was included in their metadata files.

98 2.2.2 T_{air} data

In our previous work (Nikolaou et al., 2022), we estimated daily mean T_{air} in high-resolution (1 × 1 km) across Germany using a regression-based method incorporating two linear mixed models. In brief, we predicted T_{air} by calibrating the strong relationship between the weather stations' T_{air} observations and the satellite-based land surface temperature (LST) also adjusting for various spatial predictors. We also applied a TPS interpolation in T_{air} data in order to achieve a full German-wide coverage. Extensive validation showed high performance ($R^2 \ge 0.96$) and low errors (RMSE ≤ 1.41 °C).

105 2.2.3 Elevation data

We downloaded elevation data at 30-arc-second spatial resolution, provided by the U.S. Geological Survey's Earth Resources Observation Systems (EROS) Data Center (Gesch et al., 1999). We aggregated these data to a 1 × 1 km grid, including the land borders and the shorelines in the North and Baltic Seas to match our intended spatial resolution (**Figure S2**).

110 2.2.4 Greenness data

The normalized difference vegetation Index (NDVI) is a proxy of vegetation greenness on the Earth surface, quantifying the vegetation cover and quality over space. We retrieved NDVI data of 1 × 1 km from the TERRA MODIS product MOD13A3v006 (Didan, 2015). These are monthly data, which is sufficient, as vegetation does not change considerably during a month.

115 2.2.5 True color band composite data

The visible red, green and blue light bands demonstrate how we see Earth's surface from space. We retrieved the daily true color band composite, i.e. the surface spectral reflectance for the red (band 1), blue (band 3) and green (band 4) bands at 500 m spatial resolution from the TERRA MODIS product MOD09GAv006, corrected for atmospheric conditions (Vermote, 2015). We aggregated the data to a 1 × 1 km grid, to suit the output's spatial resolution.

121 2.2.6 Precipitation data

We used daily precipitation data of 1 × 1 km developed by the REGNIE (Regionalisierte Niederschlagshöhen) method which are publicly available from the DWD Climate Data Center (DWD, 2022). REGNIE is based on the interpolated DWD weather station precipitation measurements, using a combination of multiple linear regressions and Inverse Distance Weighting (IDW), with orographic conditions considered (Rauthe et al., 2013).

127 2.2.7 Wind speed data

We retrieved daily mean wind speed (DWD, 2022) of the same 406 weather stations as for the RH data (Figure S1). We interpolated this dataset to 1 × 1 km spatial resolution using TPS, since studies have suggested that TPS outperformed other interpolation methods such as kriging or IDW for mapping climate variables (Wu et al., 2013; Wu et al., 2015).

132 2.3 Modeling

133 We trained the RF model, trying to evaluate its efficiency in reproducing the observed RH values measured 134 by the weather stations, i.e. the ground-based truth. We used 500 trees and 8 randomly sampled variables 135 as candidates at every split (num.trees = 500, mtry = 8), training the model for each year separately to 136 capture annual variations. The daily observed mean RH (%) at the DWD stations was the response variable. 137 The predictors were our previously modelled daily mean T_{air} (Celsius), the daily red, green and blue bands 138 (dimensionless), the daily mean precipitation height (mm) and the daily mean wind speed ($m \cdot s^{-2}$) as well 139 as elevation (meters) and monthly NDVI (dimensionless). We also included the day of the year (1 to 140 365 366) in order to capture daily variations in the response-predictor variables relationship.

141 2.3.1 Model performance

142 Ten-fold cross-validation (CV) was used to assess the model performance by randomly dividing the set of 143 the DWD weather monitors to a training and a testing set (90:10) ten times. Each time, the model was re-144 fitted using the training set and then the RH was predicted in the respective testing set. Our aim was to 145 estimate a full time series of RH in locations without weather stations and therefore in grid cells where 146 the RF model was not previously trained, and consequently to simulate the prediction step of the 147 modeling procedure. Regressing the observed mean RH vs. the predicted mean RH by the RF model's testing set, we calculated the corresponding R² and RMSE each of the ten times and then we took their 148 149 average to represent each year's CV-R² and CV-RMSE.

In the prediction step, we applied the RF model to all grid cells and days combinations without available
 RH measurements of DWD weather stations in order to obtain a complete RH dataset for entire Germany.

152 2.3.2 Validation with external data

153 An additional validation was conducted by comparing our daily mean RH predictions with measurements 154 of an independent dense monitoring network during 2015-2019. The network included RH measurements 155 of 4 minutes temporal resolution from 82 HOBO-Logger devices (ONSET, Type Pro v2), which were located 156 in the city of Augsburg and in two adjacent counties (Augsburg county and Aichach-Friedberg) (Figure S3). 157 Detailed information for the monitoring network and the measurements' quality assurance can be found 158 in the corresponding paper (Beck et al., 2018a). For our comparison, we aggregated the 4-min RH values to daily means and then 7-day averages. We generated the corresponding R² and RMSE as derived from 159 160 linearly regressing the predicted RH from the model against the observed RH from the HOBO-Logger 161 monitors.

The majority of the HOBO-Logger stations were located in the city center of Augsburg or close to it, where no DWD measurements were available in the training step of the RF model (closest stations were approx. 10 and 18 km apart from the city center, see **Figure S3**). Thus, we investigated the performance of the model in an area without prior information but of great epidemiological interest since highly populated implicating that more people are exposed here.

167 2.4 Descriptive analyses and case study

Descriptive statistics [mean, standard deviation (SD), minimum (min), first quartile (Q1), median, third quartile (Q3) and maximum (max)] were calculated from our German-wide RH predictions and from the DWD observations. We also investigated the spatiotemporal RH patterns over the last 2 decades, overall and by season. To demonstrate the improvement in our exposure assessment, we compared the spatial distributions of the daily mean RH predictions from the RF model and the daily mean RH measurements from the DWD stations in an urban location for the two last decades. The city of Regensburg covers an area of 80.76 km² with about 150,000 inhabitants, and, as one of the study sites of the German National Cohort (NAKO) (German National Cohort (GNC) Consortium, 2014), has also an epidemiological research interest. We performed our analysis in R, v. 4.2.2 (R Core Team, 2022). The RF model was developed with the R

178 package "ranger" (Wright and Ziegler, 2017).

179 3 Results

Figure S4 shows the Spearman correlation coefficients for the models' variables. Briefly, RH was found to be highly and positively associated with the true color band composite ($r \approx 0.5$) while there was a strong negative correlation with T_{air} ($r \approx -0.5$). In Figure S5, we demonstrate the variable importance plot findings. Date played a very important role. We also observed that T_{air} and the blue band were the most important spatiotemporal predictors of the RF model for estimating RH. They were followed by precipitation, wind speed and elevation and then NDVI, green and red band. The order of the predictors was slightly different through the years, but there were main trends as described.

187 **3.1 Model performance**

The model achieved high accuracy [22-year average $R^2 = 0.80$ (range: 0.75 - 0.86)] and small errors [22year average RMSE = 5.42 % (range: 4.72 % - 6.60 %), **Table 1**]. We observed an increase of the model performance (increase of R^2 and decrease of RMSE), together with an increase of the total number of available weather station data over the years. Autumn months (September-November) had the lowest RMSE = 4.97 % (range: 4.11 % - 6.16 %) while spring months (March-May) had the highest RMSE = 5.76 % (range: 4.99 % - 6.78 %) (**Figure 1**). We also observed that predictions belonging to the lower 10 % of the dataset gave higher errors (RMSE = 8.36 %, range: 7.30 % - 9.71 %) compared to the predictions of the
upper 10 % of the dataset (RMSE = 5.74 %, range: 4.81 % - 7.14 %) (Figure 1).

196 **3.2** Validation with external data

We found a strong correspondence between our RH model predictions and the external HOBO-Logger network measurements with a 5-year average R² of 0.84 (range: 0.80 - 0.87) and a 5-year average RMSE of 5.91 % (range: 5.40 % - 6.90 %) for the daily average RH exposure (**Table 2**). For the 7-day average RH exposure, as expected, the accuracy was even higher ($R^2 = 0.86$, range: 0.83 - 0.91) and the errors lower (RMSE = 4.82 %, range: 4.17 % - 5.92 %). Density scatterplots confirmed the good correlation (**Figure S6**).

202 3.3 Case study - Regensburg

In Figure 2, we display the average spatial RH patterns for the region of Regensburg for the period 2000-204 2021. The city area showed up to 4 % lower RH values than the surrounding rather rural county area. 205 However, the variability of the daily values which will be also considered in subsequent epidemiological 206 analysis is much larger than the 22-year average - e.g., up to 14 % (randomly selected example day in 207 Figure S7). Yet, the rural region was characterized by variations even in neighbouring tiles. The average 208 RH exposure in Regensburg measured by the available DWD weather station of the region was far below 209 the Q1 of the RH predictions of the RF model for the region (Figure 3).

210 3.4 Spatiotemporal RH patterns

Table 3 shows descriptive statistics of measured and modelled RH across Germany for 2000-2021.
Germany was characterized by high RH values with Q1 of both DWD stations' and model's RH distribution
to be 71 % and 71.98 %, respectively. The observed and predicted 22-year average RH derived by the DWD
stations and the RF model were 79.05 % (SD = 12.38 %) and 79.05 % (SD = 10.44 %), respectively.

215 Figure 4 displays the 22-year averaged predicted RH output map of Germany (plot 1) which indicates 216 spatial RH patterns, including urbanization, mountainous regions, rivers, forests and coastlines. 217 Metropolitan areas such as those of Berlin, Hamburg and Munich and the extended and other dense urban 218 cores (e.g., from Karlsruhe to Frankfurt) had much lower RH values compared to the neighbouring rural 219 settings. In Figure S8, we zoomed in the Augsburg region, which consists of the city center and two 220 adjacent counties, to give an example of the high spatial difference between a city center and its 221 neighbouring but less urbanized areas. Additionally, dense mountainous regions such as the Alps and Harz, 222 coastlines as the North Sea coast and rivers as Elbe in a large part of it, had the highest RH values country-223 wide (Figure 4). Also, the temporal RH variability in Germany is presented for 2001-2021, by exhibiting 224 the differences between the predicted RH yearly averages and the 21-year average (Figure 4, plot 2). We 225 excluded the year 2000 because the model predictions are only available from late February of that year 226 due to the missing T_{air} values until then. There were some fluctuations over the years but without 227 indication of an increasing or decreasing trend. The most humid years were 2001 (81.31 %), 2014 (81.25 228 %) and 2013 (80.99 %) while the most arid were 2003 (75.38 %), 2020 (75.58 %) and 2018 (75.60 %), which 229 are known hot and dry years from the recent climatological record.

230 Mapping the 22-year average RH by season (Figure S9) identified winter and fall as the most humid

231 seasons. High spatial RH variability was also observed within each season.

232 4 Discussion

In this paper, we introduced an approach for spatial and temporal modeling of RH using RF, a popular ML
 method for prediction tasks. The approach goes beyond the conventional interpolation of meteorological
 observations, and uses several other data sources. We produced a reliable spatiotemporally-resolved RH
 dataset at 1 × 1 km spatial resolution across Germany for the period 2000-2021. The RF model achieved
 good performance with high predictive accuracy and low errors, validated with both internal data using

238 cross-validation ($R^2 = 0.80$, RMSE = 5.42 %), and with independent observational data ($0.84 \le R^2 \le 0.86$, 239 4.82 % \leq RMSE \leq 5.91 %). A case study for the city of Regensburg shows that our dataset is capable of 240 capturing the full range of spatial variability of RH compared to the standard use of meteorological 241 observations. These DWD station observations could not represent the high RH values of the peripheral 242 areas in Regensburg, but also not the very low RH values of the city center. This clearly demonstrates the 243 added value of our approach and how the use of additional data sources supplementing the conventional 244 use of meteorological observations improved the RH prediction. It is especially important to capture the 245 RH spatial variability for assessing differences in human's individual exposure in epidemiological studies. 246 We also presented an analysis of the spatiotemporal RH patterns in Germany during 2000-2021.

247 The RH-health relevance has not been clarified adequately (Bind et al., 2014). RH adverse effects on 248 human health could be partially explained by its interplay with the excessive heat stress and the body 249 dehydration, as described in Davis et al. (2016). During extended and excessive heat events such as 250 heatwaves, the human body struggles against heat-driven physiological responses and a key mechanism 251 for its temperature regulation is evaporation. However, when RH is high and therefore air contains a lot 252 of moisture, it is difficult for the sweat to be relieved and thus cooling becomes insufficient. Hence, the 253 body core temperature increases while this increase is associated with a variety of detrimental health 254 effects (Schneider et al., 2017). Additionally, low RH can affect the human skin sensitivity to mechanical 255 stress (Engebretsen et al., 2016). RH is also associated with the transition of vector-borne diseases e.g., 256 from mosquitos and ticks (Davis et al., 2016) as well as with the development and stability of 257 microorganisms in aerosols, facilitating airborne diseases (Božič et al., 2021).

258 So far there is a literature gap in the investigation of the RH exposure's direct effects on human health 259 and the accompanying underlying mechanisms. Further and more detailed research is needed. Hence, it 260 is critically important for epidemiologists to have access to high-resolution and reliable RH datasets.

261 Most epidemiological studies retrieve the participants' exposure information, in this case RH, from 262 available meteorological stations that do not capture the full variability of RH, especially at the city scale. 263 In the Regensburg area, an epidemiological study would usually assign RH measurements from the station 264 most closely located to each participant's residential address, but fails to account for the spatial variability 265 of RH that is actually occurring. Therefore, some measurement error would be introduced and the 266 variability would be lost. Focusing on the city area, participants who live there would be assigned with a 267 higher RH value than their actual one. At the same time, those living outside the city center would be 268 assigned with RH values that are too low. This clearly demonstrates the urgent need for high 269 spatiotemporal RH datasets for health studies for less biased exposure estimates.

270 Compared to other studies that use interpolation techniques such as TPS or kriging, our RF model is 271 capable of reducing errors by half. Li et al. (2014) introduced a two-step procedure to map RH every 3 272 hours at 1 km resolution over China during 1958-2010. They fitted a partial TPS interpolation to reanalysis 273 data, location and elevation as predictors, to estimate a trend surface, and then a simple kriging was 274 applied to the residuals for trend surface correction. They reported a RMSE of 11.06 % whereas our model 275 showed a RMSE of 5.42 %. More recently, Li and Zha (2018) also used an RF model, combining station and 276 satellite data, to estimate RH during the hot summer of 2009 over China. Elevation and vegetation were 277 found to be the most important predictors for RH. Comparing our model with their work, it seems that 278 our additional inclusion of Tair, date information, precipitation and wind speed data in the modeling 279 process, significantly improved the model's performance. Li and Zha (2018) reported a $R^2 = 0.70$ and RMSE = 7.4 %, whereas our model could improve the R^2 to 0.80 and lower the errors to RMSE = 5.42 %. In 280 281 addition, our RF model allowed us to model RH for entire years and not only for one season.

This study was also subject to limitations. The external validation set was not representative of the whole
Germany. The HOBO-Logger monitoring network was placed in Augsburg, South Germany. However, we

284 used the Augsburg's greater region which consists of a dense city center and two adjacent rural settings 285 and therefore the validation area was characterized by high spatial RH variability. Additionally, we were 286 already able to measure the model's predictive accuracy country-wide due to our monitor-based split in 287 the applied CV scheme (2.3.1 Model performance). The 1 × 1 km spatial resolution could be too coarse 288 for some studies, especially for local and small-scale analyses. However, as we demonstrated in the case 289 study of the city of Regensburg, the RF model of 1×1 km provided a valid representation of the RH 290 spatiotemporal variation at the city scale. For future analyses, we could consider downscaling methods 291 especially for cities (Hough et al., 2020).

292 5 Conclusion

We showed how observation, remote sensing and modelled data can be combined under a RF modeling process to reliably estimate RH in high temporal and spatial resolution across a country. Our product contributes substantially to reduce exposure errors for subsequent epidemiological studies, by better representing the spatiotemporal RH variability. We provide a reliable RH dataset for Germany and a wellfounded and generalizable approach for RH prediction for other study domains and countries.

298 References

- Analitis, A.; Katsouyanni, K.; Biggeri, A.; Baccini, M.; Forsberg, B.; Bisanti, L.; Kirchmayer, U.; Ballester, F.;
- 300 Cadum, E.; Goodman, P. G.; Hojs, A.; Sunyer, J.; Tiittanen, P.; Michelozzi, P., 2008. Effects of cold weather
- on mortality: results from 15 European cities within the PHEWE project. Am. J. Epidemiol. 168 (12), 1397-
- 302 1408. <u>https://doi.org/10.1093/aje/kwn266</u>
- Armstrong, B., 2006. Models for the relationship between ambient temperature and daily mortality.
 Epidemiology 624-631. <u>http://www.jstor.org/stable/20486290</u>
- 305 Božič, A.; Kanduč, M., 2021. Relative humidity in droplet and airborne transmission of disease. J. Biol.
- 306 Phys. 47 (1), 1-29. <u>https://doi.org/10.1007/s10867-020-09562-5</u>
- Beck, C.; Straub, A.; Breitner, S.; Cyrys, J.; Philipp, A.; Rathmann, J.; Schneider, A.; Wolf, K.; Jacobeit, J., 307 308 2018a. Air temperature characteristics of local climate zones in the Augsburg urban area (Bavaria, 309 southern Germany) under varying synoptic conditions. Urban Clim. 25, 152-166. 310 https://doi.org/10.1016/j.uclim.2018.04.007
- 311 Beck, H.E.; Zimmermann, N.E.; McVicar, T.R.; Vergopolan, N.; Berg, A.; Wood, E.F., 2018b. Present and
- 312 future Köppen-Geiger climate classification maps at 1-km resolution. Sci. Data 5, 1-12.
- 313 <u>https://doi.org/10.1038/sdata.2018.214</u>
- Bind, M. A.; Zanobetti, A.; Gasparrini, A.; Peters, A.; Coull, B.; Baccarelli, A.; Tarantini, L.; Koutrakis P.;
- 315 Vokonas P.; Schwartz, J., 2014. Effects of temperature and relative humidity on DNA methylation.
- 316 Epidemiology (Cambridge, Mass.) 25 (4), 561. <u>https://doi.org/10.1097%2FEDE.0000000000000120</u>
- 317 Davis, R.E.; McGregor, G.R.; Enfield, K.B., 2016. Humidity: A review and primer on atmospheric moisture
- 318 and human health. Environ. Res. 144, 106-116. <u>https://doi.org/10.1016/j.envres.2015.10.014</u>

- Didan, K., 2015. MOD13A3 MODIS/Terra vegetation Indices Monthly L3 Global 1km SIN Grid V006 [Data
 set]. NASA EOSDIS Land Processes DAAC. Accessed 2022-03-21 from
 https://doi.org/10.5067/MODIS/MOD13A3.006
- 322 DWD, 2022. DWD Climate Data Center (CDC): Historical daily station observations (temperature, pressure,
 323 precipitation, sunshine duration, etc.) for Germany, version v21.3, 2021.
- 324 DWD, 2022. DWD Climate Data Center (CDC): REGNIE grids of daily precipitation, last accessed:
 325 21.03.2022.
- 326 Engebretsen, K. A.; Johansen, J. D.; Kezic, S.; Linneberg, A.; Thyssen, J. P., 2016. The effect of
- 327 environmental humidity and temperature on skin barrier function and dermatitis. J. Eur. Acad. Dermatol.
- 328 Venereol. 30 (2), 223-249. <u>https://doi.org/10.1111/jdv.13301</u>
- Forootan, E., 2019. Analysis of trends of hydrologic and climatic variables. Soil Water Res. 14 (3), 163-171.
 10.17221/154/2018-SWR
- 331 German National Cohort (GNC) Consortium, 2014. The German National Cohort: aims, study design and
- 332 organization. Eur. J. Epidemiol. 29 (5), 371-382. <u>https://doi.org/10.1007%2Fs10654-014-9890-7</u>
- 333 Gesch, D.B.; Verdin, K.L.; Greenlee, S.K., 1999. New land surface digital elevation model covers the Earth.
- 334 EOS 80 (6), 69-70. <u>https://doi.org/10.1029/99EO00050</u>
- Hough, I.; Just, A.C.; Zhou, B.; Dorman, M.; Lepeule, J.; Kloog, I., 2020. A multi-resolution air temperature
- 336 model for France from MODIS and Landsat thermal data. Environ. Res. 183, 109244.
- 337 <u>https://doi.org/10.1016/j.envres.2020.109244</u>

- Jin, Z.; Ma, Y.; Chu, L.; Liu, Y.; Dubrow, R.; Chen, K., 2022. Predicting spatiotemporally-resolved mean air
 temperature over Sweden from satellite data using an ensemble model. Environ. Res. 204, 111960.
 https://doi.org/10.1016/j.envres.2021.111960
- Li, T.; Zheng, X.; Dai, Y.; Yang, C.; Chen, Z.; Zhang, S.; Wu, G.; Wang, Z.; Huang, C.; Shen, Y.; Liao R., 2014.
- 342 Mapping near-surface air temperature, pressure, relative humidity and wind speed over Mainland China
- 343 with high spatiotemporal resolution. Adv. Atmos. Sci. 31 (5), 1127-1135. <u>https://doi.org/10.1007/s00376-</u>

<u>344</u> <u>014-3190-8</u>

- Li, L.; Zha, Y., 2018. Mapping relative humidity, average and extreme temperature in hot summer over
- 346 China. Sci. Total Environ. 615, 875-881. <u>https://doi.org/10.1016/j.scitotenv.2017.10.022</u>
- 347 Luo, C.; Ma, Y.; Liu, Y.; Lv, Q.; Yin, F., 2020. The burden of childhood hand-foot-mouth disease morbidity
- 348 attributable to relative humidity: a multicity study in the Sichuan Basin, China. Sci. Rep. 10 (1), 1-10.
- 349 <u>https://doi.org/10.1038/s41598-020-76421-7</u>
- 350 Mistry, M.N.; Schneider, R.; Masselot, P.; Royé, D.; Armstrong, B.; Kyselý, J.; Orru, H.; Sera, F.; Tong, S.;
- Lavigne, É.; Urban, A.; Madureira, J.; García-León, D.; Ibarreta, D.; Ciscar, J.-C.; Feyen, L.; DeSchrijver, E.;
- 352 Coelho, M.S.Z.S.; Pascal, M.; Tobias, A.; Multi-Country Multi-City (MCC) Collaborative Research Network,
- 353 Guo, Y.; Vicedo-Cabrera, A.M.; Gasparrini, A., 2022. Comparison of weather station and climate reanalysis
- data for modelling temperature-related mortality. Sci. Rep. 12 (1), 1-14. <u>https://doi.org/10.1038/s41598-</u>
- 355 022-09049-4
- Nikolaou, N.; Dallavalle, M.; Stafoggia, M.; Bouwer, L. M.; Peters, A.; Chen, K.; Wolf K.; Schneider, A., 2022.
- 357 High-resolution spatiotemporal modeling of daily near-surface air temperature in Germany over the
- 358 period 2000–2020. Environ. Res. 115062. <u>https://doi.org/10.1016/j.envres.2022.115062</u>

Ou, C.Q.; Yang, J.; Ou, Q.Q.; Liu, H.Z.; Lin, G.Z.; Chen, P.Y.; Qian, J.; Guo Y.M., 2014. The impact of relative
humidity and atmospheric pressure on mortality in Guangzhou, China. BES 27 (12), 917-925.
https://doi.org/10.3967/bes2014.132

- 362 R Core Team, 2022. R: A language and environment for statistical computing. R Foundation for Statistical
- 363 Computing, Vienna, Austria. Retrieved from <u>https://www.R-project.org/</u>
- 364 Rauthe, M.; Steiner, H.; Riediger, U.; Mazurkiewicz, A.; Gratzki, A., 2013. A Central European precipitation
- 365 climatology Part I: Generation and validation of a high-resolution gridded daily data set (HYRAS).
- 366 Meteorologische Zeitschrift 22 (3), 235-256. <u>http://dx.doi.org/10.1127/0941-2948/2013/0436</u>
- 367 Rosenfeld, A.; Dorman, M.; Schwartz, J.; Novack, V.; Just, A.C.; Kloog, I., 2017. Estimating daily minimum,
- 368 maximum, and mean near surface air temperature using hybrid satellite models across Israel. Environ.
- 369 Res. 159, 297-312. <u>https://doi.org/10.1016/j.envres.2017.08.017</u>
- 370 Schneider, A.; Rückerl, R.; Breitner, S.; Wolf, K.; Peters, A., 2017. Thermal control, weather, and aging.
- 371 Curr. Environ. Health Rep. 4 (1), 21-29. <u>https://doi.org/10.1007/s40572-017-0129-0</u>
- 372 Sherwood, S.C.; Ingram, W.; Tsushima, Y.; Satoh, M.; Roberts, M.; Vidale, P.L.; O'Gorman, P.A., 2010.
- 373 Relative humidity changes in a warmer climate. J. Geophys. Res. Atmos. 115 (D9).
- 374 <u>https://doi.org/10.1029/2009JD012585</u>
- Silibello, C.; Carlino, G.; Stafoggia, M.; Gariazzo, C.; Finardi, S.; Pepe, N.; Radice, P.; Forastiere, F.; Viegi, G.,
 2021. Spatial-temporal prediction of ambient nitrogen dioxide and ozone levels over Italy using a Random
 Forest model for population exposure assessment. Air Qual. Atmos. Health 14 (6), 817-829.
- 378 <u>https://doi.org/10.1007/s11869-021-00981-4</u>

- Stafoggia, M.; Bellander, T.; Bucci, S.; Davoli, M.; De Hoogh, K.; De'Donato, F.; Gariazzo, C.; Lyapustin, A.;
 Michelozzi, P.; Renzi, M.; Scortichini, M.; Shtein, A.; Viegi, G.; Kloog, I.; Schwartz, J., 2019. Estimation of
 daily PM10 and PM2. 5 concentrations in Italy, 2013–2015, using a spatiotemporal land-use randomforect model. Env. Int. 124, 170, 170. https://doi.org/10.1016/j.envipt.2010.01.016
- 382 forest model. Env. Int. 124, 170-179. <u>https://doi.org/10.1016/j.envint.2019.01.016</u>
- Vermote, E.W., R., 2015. MOD09GA MODIS/Terra Surface Reflectance Daily L2G Global 1kmand 500m SIN
 Grid V006 [Data set]. NASA EOSDIS Land Processes DAAC. Accessed 2022-03-21 from
 https://doi.org/10.5067/MODIS/MOD09GA.006
- 386 Wright M.N.; Ziegler A., 2017. "ranger: A Fast Implementation of Random Forests for High Dimensional
- 387 Data in C++ and R." J. Stat. Softw. 77 (1), 1–17. <u>doi:10.18637/jss.v077.i01</u>.
- 388 Wu, W.; Tang, X.-P.; Yang, C.; Guo, N.-J.; Liu, H.-B.. 2013. Spatial estimation of monthly mean daily 389 sunshine across hours and solar radiation mainland China. RES 57, 546-553. https://doi.org/10.1016/j.renene.2013.02.027 390
- 391 Wu, W.; Xu, A.-D.; Liu, H.-B., 2015. High-resolution spatial databases of monthly climate variables (1961–
- 2010) over a complex terrain region in southwestern China. Theor. Appl. Climatol. 119 (1), 353-362.
- 393 <u>https://doi.org/10.1007/s00704-014-1123-1</u>
- Xiong, Y.; Meng, Q.-S.; Jie, G.; Tang, X.-F.; Zhang, H.-F., 2017. Effects of relative humidity on animal health
- 395 and welfare. J. Integr. Agric. 16 (8), 1653-1658. <u>https://doi.org/10.1016/S2095-3119(16)61532-0</u>
- 396 Yang, Y., You, E., Wu, J., Zhang, W., Jin, J., Zhou, M., Jiang, C., Huang, F., 2018. Effects of relative humidity
- on childhood hand, foot, and mouth disease reinfection in Hefei, China. Sci. Total Environ. 630, 820-826.
- 398 https://doi.org/10.1016/j.scitotenv.2018.02.262

- Yao, R., Wang, L., Huang, X., Cao, Q., Peng, Y., 2022. A method for improving the estimation of extreme
 air temperature by satellite. Sci. Total Environ. 155887. https://doi.org/10.1016/j.scitotenv.2022.155887
- 401 Zeger, S.L.; Thomas, D.; Dominici, F.; Samet, J.M.; Schwartz, J.; Dockery, D.; Cohen, A., 2000. Exposure
- 402 measurement error in time-series studies of air pollution: concepts and consequences. EHP 108 (5), 419-
- 403 426. <u>https://doi.org/10.1289/ehp.00108419</u>
- 404 Zeng, J.; Zhang, X.; Yang, J.; Bao, J.; Xiang, H.; Dear, K.; Liu, Q.; Lin, S.; Lawrence, W.R.; Lin, A.; Huang, C.,
- 405 2017. Humidity may modify the relationship between temperature and cardiovascular mortality in
- 406 Zhejiang Province, China. IJERPH 14 (11), 1383. <u>https://doi.org/10.3390/ijerph14111383</u>
- Zhang, P.; Zhang, J.; Chen, M., 2015. Economic impacts of climate change on Chinese agriculture: the
 importance of relative humidity and other climatic variables. Available at SSRN 2598810.
 <u>https://dx.doi.org/10.2139/ssrn.2598810</u>

410 Table legends

411 **Table 1**. Prediction accuracy for the RF model: 10-fold CV results for the daily mean RH predictions over

412 Germany during 2000-2021.

- 413 **Table 2**. Accuracy results from the validation with external data using the HOBO-Logger daily mean RH
- 414 observations and 7-day averages over the Augsburg region during 2015-2019.
- 415 **Table 3**. Observed and predicted mean RH (%) over Germany during 2000-2021.

416 Figure legends

- 417 Figure 1. Seasonal RMSE and RMSE to extremes for the model's RH predictions in Germany during 2000-
- 418 2021.
- 419 **Figure 2**. Spatial pattern of the averaged predicted RH in Regensburg during 2000-2021.
- 420 Figure 3. Distribution of predicted RH in the Regensburg region for 2000-2021 (histogram in blue and

421 corresponding boxplot above).

- 422 Figure 4. Spatiotemporal RH patterns in Germany during 2000-2021. Plot 1: Spatial patterns of the
- 423 predicted RH in Germany, averaged for 2000-2021. Plot 2: Difference between the predicted RH yearly
- 424 averages and the predicted RH 21-year average (2001-2021), German-wide.

Table 1. Prediction accuracy for the RF model: 10-fold CV results for thedaily mean RH predictions over Germany during 2000-2021.					
Year	R ²	RMSE	Sample size SE (number of cell-days)		
2000	0.75	6.07	100,699		
2001	0.75	5.84	121,225		
2002	0.75	5.99	123,946		
2003	0.79	6.60	123,364		
2004	0.75	6.02	126,604		
2005	0.78	5.58	134,386		
2006	0.80	5.66	135,600		
2007	0.81	5.22	139,482		
2008	0.80	5.40	140,135		
2009	0.79	5.41	142,295		
2010	0.83	5.06	142,629		
2011	0.84	5.30	141,781		
2012	0.82	5.12	141,820		
2013	0.81	5.13	140,928		
2014	0.82	4.88	142,641		
2015	0.82	5.36	142,908		
2016	0.81	5.05	139,491		
2017	0.80	5.00	143,206		
2018	0.84	5.39	143,026		
2019	0.83	5.20	140,866		
2020	0.86	5.32	116,670		
2021	0.83	4.72	116,544		
Overall	0.80	5.42	133,648		

Table 2 . Accuracy results from the validation with external data using the HOBO-Logger daily mean RHobservations and 7-day averages over the Augsburg region during 2015-2019.					
Year	R ²	RMSE	7-day average R ²	7-day average RMSE	
2015	0.85	5.45	0.89	4.17	
2016	0.80	5.59	0.83	4.41	
2017	0.83	5.40	0.84	4.39	
2018	0.87	6.21	0.91	5.22	
2019	0.83	6.90	0.85	5.92	
Overall	0.84	5.91	0.86	4.82	

Table 3. Observed and predicted mean RH (%) over Germany during 2000-2021.							
Source	Mean	SD	Min	Q1	Median	Q3	Max
DWD stations (n = 406)	79.05	12.38	3.00	71.00	81.00	88.75	100.00
RF model (n = 366,536 cells)	79.05	10.44	13.05	71.98	80.61	87.47	100.00



Seasonal RMSE (%) through the years 2000-2021

To extremes RMSE (%) through the years 2000-2021











Supplementary Material

Click here to access/download **Supplementary Material** Nikolaou et al_ rel humidity model_ supplementary.docx

Author contributions

Nikolaos Nikolaou: Conceptualization, Data curation, Methodology, Analysis, Visualization, Writing original draft, Writing - review & editing. Laurens M. Bouwer: Conceptualization, Data curation, Methodology, Writing - review & editing. Marco Dallavalle: Data curation, Writing - review & editing. Mahyar Valizadeh: Methodology, Writing - review & editing. Massimo Stafoggia: Methodology, Writing - review & editing. Annette Peters: Conceptualization, Writing - review & editing, Supervision. Kathrin Wolf: Conceptualization, Data curation, Methodology, Writing - review & editing, Supervision. Alexandra Schneider: Conceptualization, Methodology, Writing - review & editing, Supervision.

Declaration of interests

⊠The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

□The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: