

1 **A merged continental planetary boundary layer height**
2 **dataset based on high-resolution radiosonde measurements,**
3 **ERA5 reanalysis, and GLDAS**

4
5 Jianping Guo^{a★}, Jian Zhang^{b*}, Jia Shao^{d★}, Tianmeng Chen^a, Kaixu Bai^c, Yuping
6 Sun^a, Ning Li^a, Jingyan Wu^a, Rui Li^e, Jian Li^a, Qiyun Guo^f, Jason B. Cohen^g, Panmao
7 Zhai^a, Xiaofeng Xu^h, Fei Hu^{i*}

8
9
10 ^aState Key Laboratory of Severe Weather, Chinese Academy of Meteorological
11 Sciences, Beijing 100081, China

12 ^bHubei Subsurface Multi-scale Imaging Key Laboratory, Institute of Geophysics and
13 Geomatics, China University of Geosciences, Wuhan 430074, China

14 ^cKey Laboratory of Geographic Information Science (Ministry of Education), School of
15 Geographic Sciences, East China Normal University, Shanghai 200241, China

16 ^dCollege of Informatics, Huazhong Agricultural University, Wuhan 430070, China

17 ^eMinistry of Education Key Laboratory for Earth System Modeling, Department of
18 Earth System Science, Tsinghua University, Beijing 100084, China

19 ^fMeteorological Observation Center, China Meteorological Administration, Beijing
20 100081, China

21 ^gSchool of Environment and Spatial Informatics, China University of Mining and
22 Technology, Xuzhou, China

23 ^hChina Meteorological Administration, Beijing 100081, China

24 ⁱState Key Laboratory of Atmospheric Boundary Layer Physics and
25 Atmospheric Chemistry, Institute of Atmospheric Physics, Beijing 100029, China

26
27 ★Both authors Jianping Guo and Jia Shao contributed equally and should be considered
28 co-first authors.

29
30
31 Correspondence to:

32 Dr. Jian Zhang (Email: zhangjian@cug.edu.cn) and Dr./Prof. Fei Hu (Email:
33 hufei@mail.iap.ac.cn)

ABSTRACT

35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63

The planetary boundary layer (PBL) is the lowermost part of the troposphere that governs the exchange of momentum, mass and heat between surface and atmosphere. To date the radiosonde measurements have been extensively used to estimate PBLH; suffering from low spatial coverage and temporal resolution, the radiosonde data is incapable of providing the diurnal description of PBLH across the globe. To fill this data gap, this paper aims to produce a temporally continuous PBLH dataset during the course of a day over the global land by applying the machine learning algorithms to integrate high-resolution radiosonde measurements, ERA5 reanalysis, and the Global Land Data Assimilation System (GLDAS) product. This dataset covers the period from 2011 to 2021 with a temporal resolution of 3-hour and a horizontal resolution of $0.25^{\circ} \times 0.25^{\circ}$. The radiosonde dataset contained around 180 million profiles over 370 stations across the globe. The machine learning model was established by taking 18 parameters derived from ERA5 reanalysis and GLDAS as input variables while the PBLH biases between radiosonde observations and ERA5 reanalysis were used as the learning targets. The input variables were presumably representative regarding the land properties, near-surface meteorological conditions, terrain elevations, lower tropospheric stabilities, and solar cycles. Once a state-of-the-art model had been trained, the model was then used to predict the PBLH bias at other grids across the globe with parameters acquired or derived from ERA5 and GLDAS. Eventually, the merged PBLH can be taken as the sum of the predicted PBLH bias and the PBLH retrieved from ERA5 reanalysis. Overall, this merged high-resolution PBLH dataset was globally consistent with the PBLH retrieved from radiosonde observations both in magnitude and spatiotemporal variation, with a mean bias of as low as -0.9 m. The dataset and related codes are publicly available at <https://doi.org/10.5281/zenodo.6498004> (Guo et al., 2022), which are of significance for a multitude of scientific research and applications, including air quality, convection initiation, climate and climate change, just to name a few.

64 **1. Introduction**

65 Planetary boundary layer (PBL), the lowermost part of the troposphere where the
66 turbulence and convection mainly occur, is of significance in modulating the exchange
67 of momentum, heat, moisture, and mass between the surface and the free atmosphere
68 over a range of scales (Stull 1988; Cooper and Eichinger, 1994; Edson et al., 2013).
69 The turbulence in the PBL is largely generated mechanically, which is owing to both
70 wind shear and friction, and is generated convectively, which is owing to buoyancy and
71 surface heating (Degrazia et al., 2020). Within the PBL, vertical turbulent mixing of air
72 masses is rapid and constant, on the order of 30 minutes or less (Wallace and Hobbs,
73 2006). Therefore, the reliable parameterization of the PBL is crucial for the accurate
74 representations of vertical diffusion, cloud formation/development, and pollutant
75 deposition in numerical weather prediction (NWP), climate, air quality and coupled
76 atmosphere–hydrosphere–biosphere models (Seibert, 2000; Hu et al., 2010; Baklanov
77 et al., 2011). It has been well recognized that the variation of PBL height (PBLH)
78 significantly impacts the near-surface air quality (Petäjä et al., 2016; Wang and Wang,
79 2016; Lou et al., 2019; Li et al., 2021) and climate system as well (Esau and
80 Zilitinkevich, 2010; Davy and Esau, 2016).

81 The development of PBL is subject to the changes of the energy balance near the
82 ground surface, largely through the linkages between soil moisture and sensible heat
83 flux, latent heat flux and net radiation (Dirmeyer et al., 2014; Xu et al., 2021). In
84 particular, the sensible heat flux is closely associated with the variation in
85 evapotranspiration, land type, and cloud cover. Also, the daytime convective PBL is
86 modulated by cloud radiative effects, particularly in the early afternoon (Guo et al.,
87 2016; Zhang et al., 2018; Davis et al., 2020). Furthermore, the aerosol radiative effect
88 (due to both aerosol scattering and absorption) indirectly affects the evolution of PBL
89 by changing the atmospheric heating rate and the solar radiation reaching the surface
90 (Wang et al., 2013; Li et al., 2017; Yang et al., 2016). Besides, the entrainment of air
91 from above the PBL can also significantly drive the evolution of PBL (Hu et al., 2010).

92 To date, a variety of methods have been applied on vertical profiles of aerosol
93 properties, water vapor, temperature, refractivity, and wind to estimate PBLH (e.g.,
94 Holzworth 1964; Seibert 2000; Lammert and Bösenberg 2006; McGrath-Spangler and
95 Denning 2012; Chan and Wood 2013; Su et al., 2018; Liu et al., 2019; Ding et al., 2021).
96 The estimate varies considerably with data sources, algorithms, and data vertical
97 resolutions (Seibert et al., 2000; Seidel et al., 2010). For instance, PBLH determined by
98 the minimum vertical gradient relative humidity is about 1 km larger than that from the
99 parcel method, even though the latter algorithm is generally thought to be one of the
100 most reliable methods for the estimation of the convective boundary layer (CBL) height
101 (Hennemuth and Lammert, 2006; Seidel et al., 2010). In addition, different data sources,
102 such as ceilometer Lidar, COSMIC GPS RO satellite, radiosonde, and the fifth
103 generation ECMWF (European Centre for Medium-Range Weather Forecasts)
104 atmospheric reanalysis (ERA5) reanalysis dataset can reach quite different estimates of
105 PBLH (Saha et al., 2022). Recently, as suggested by Teixeira et al. (2021), the PBLH
106 should be ideally estimated using direct observations of vertical profiles of turbulent
107 quantities, which is due in large part to the turbulent nature of PBL. But only a few
108 places have such observations. A wide range of complex physical and chemical
109 processes involved in the PBL further make PBLH estimates quite elusive and tricky
110 (Seidel et al., 2010; Teixeira et al., 2021).

111 Among the instruments, radiosonde is the most accepted instrument for deriving
112 the PBLH for both CBL and stable boundary layer (SBL), due to its unprecedented
113 capability of providing in situ observations of the thermodynamic and dynamic states
114 of the PBL (Seidel et al., 2010; de Arruda Moreira et al., 2018, Guo et al., 2019). In
115 addition, the bulk Richardson number method has been proved to be the most suitable
116 PBLH algorithm for application to a large radiosonde dataset (Seidel et al., 2012). The
117 dataset with a full vertical resolution (5–8 m) has previously been used to study PBLHs
118 over China and near-globe (Guo et al., 2016; 2021). The limitation of this dataset is its
119 poor coverage over the ocean and some continental areas without high-resolution
120 radiosonde observations.

121 By contrast, reanalysis datasets, such as ERA5 reanalysis and the Modern-Era
122 Retrospective-analysis for Research and Applications version 2 product (MERRA-2),
123 have a unique advantage in spatial-temporal coverage. Our recent study (Guo et al.,
124 2021) suggests that ERA5 is the most promising reanalysis data source in terms of
125 characterizing the evolution of PBLH, with an underestimation of daytime PBLH at
126 around 130 m, when compared to high-resolution radiosonde. Nevertheless, the
127 underestimation of PBLH in ERA5 reanalysis can be as high as 500 m in the afternoon
128 when the PBL is fully developed. This underestimation could be attributed to, but not
129 limited to, the gradient of terrain elevation and the lower tropospheric stability.
130 Particularly, a higher terrain gradient or a more unstable troposphere generally lead to
131 a lower PBLH in ERA5 reanalysis.

132 Rather, by exploiting both the advantages of in situ atmospheric measurements
133 from radiosonde and the high-resolution model products from ERA5 reanalysis, it is
134 quite desirable to generate a new PBLH dataset by seamlessly blending these versatile
135 products. The biases between PBLHs retrieved from the ERA5 and radiosonde could
136 be represented by the land properties, near-surface meteorological conditions, among
137 others, and further be minimized or optimized via a machine learning model. The
138 Global Land Data Assimilation System (GLDAS) incorporates satellite- and ground-
139 based observations and produces a global, high-resolution product regarding land states
140 and fluxes (Rodell et al., 2004). To this end, the present analyses used the radiosonde
141 dataset that contained around 180 million profiles over 370 stations across the world,
142 in combination with the ERA5 reanalysis and GLDAS data. A long-term merged PBLH
143 dataset covering the period 2011 to 2021 were generated, which could have crucial
144 implications for the development and evaluation of weather and climate, environmental
145 meteorology, and boundary layer parameterization. The rest of the paper is organized
146 as follows. Section 2 describes the fundamental data sets and the PBLH methodology
147 we use in this study, Sections 3 and 4 report on the machine learning algorithm used to
148 generate the merged PBLH dataset, also revealed are the data quality, and Section 5

149 represents the climatological merged continental PBLH, and Section 6 ends with a brief
150 summary and conclusion.

151 **2. Data sources and conventional PBLH determination method**

152 2.1 High-resolution radiosonde measurements

153 As described in Guo et al. (2021) and Zhang et al. (2022), a high-resolution
154 radiosonde dataset gained from several organizations was adopted, spanning the years
155 from 2011 to 2021. The organizations include the China Meteorological Administration
156 (CMA), the National Oceanic and Atmospheric Administration (NOAA), the Global
157 Climate Observing System (GCOS) Reference Upper-Air Network (GRUAN), the
158 Centre for Environmental Data Analysis of the United Kingdom (CEDA), University
159 of Wyoming, and German Deutscher Wetterdienst. The detailed information on the
160 provided data is listed in Table 1. In total, over 185 million radiosonde profiles were
161 collected to determine PBLH, 95% of which were released at regular synoptic times of
162 0000 UTC and 1200 UTC, and the rest of which were irregularly launched at other
163 times during the intensive observational periods. Note that those soundings with the
164 lowest burst height lower than 10 km above ground level (a.g.l) were eliminated. In
165 addition, all the original soundings were evenly interpolated to the profiles with a
166 vertical resolution of 10 m by cubic spline interpolation.

167 The spatial distribution of sample numbers over each radiosonde station at four
168 different synoptic times (0000 UTC, 0600 UTC, 1200 UTC, 1800 UTC) is presented in
169 Fig. 1. It is noticeable that the radiosonde stations over Europe, the U.S., China, and
170 Australia have an unprecedented rich geographic coverage. Furthermore, the
171 radiosonde measurements over China and the U.S. have a fair temporal continuity at
172 0000 UTC and 1200 UTC, with a total sample size reaching up to as large as 3000 for
173 each station. In comparison, the stations are poorly distributed over regions or countries
174 such as southern America, the Pacific islands, Russia, the Middle East, India, and Africa.

175 2.2 ERA5 and GLDAS

176 ERA5 is the latest version of ECMWF reanalysis, benefiting from a decade of
177 developments in model physics, core dynamics, and data assimilation (Hersbach et al.,
178 2020). The PBLH product is resolved by the ERA5 reanalysis on a 1440×721
179 longitude/latitude grid, with a spatial resolution of 0.25°×0.25° and a temporal
180 resolution of 1 hour, which is realistically simulated by the bulk Richardson number
181 method. In addition, the parameters, such as the lower tropospheric stability (LTS), the
182 standard deviation of digital elevation model (SDDDEM), 10-m surface wind speed, 2-
183 m air temperature, and 2-m pressure, are either computed or directly extracted from
184 ERA5 reanalysis. LTS is defined as the difference in potential temperature between 700
185 and 1000 hPa (Guo et al., 2016). As a result, a total of six parameters were obtained
186 from ERA5 reanalysis.

187 The land property parameters were taken from GLDAS, which include downward
188 short-wave radiation (DSWR), downward long-wave radiation (DLWR), surface heat
189 net flux (SHF), surface latent heat net flux (LHF), evapotranspiration, transpiration, soil
190 moistures in 0–10 cm, 10–40 cm, 40–100 cm, and 100–200 cm, and total precipitation
191 amount. Totally, 11 parameters were extracted from the GLDAS product. GLDAS has
192 a temporal resolution of 3 hours and the same spatial resolution as that of ERA5
193 reanalysis. However, GLDAS has no data over Antarctica. It should be noted that there
194 exists a 0.125° lag between the start latitude and longitude of GLDAS and those of
195 ERA5 and therefore, the latitude and longitude of GLDAS were minus 0.125° have to
196 be used to match with ERA5 reanalysis.

197 According to the methods proposed by Guo et al. (2021), the collocation
198 procedures between the grid products from ERA5 and GLDAS and station-based
199 radiosonde observations were mainly implemented as follows. (1) The grid should
200 contain the radiosonde station. (2) The UTC time (hour) of grid product and radiosonde
201 stay the same.

202 2.3 PBLH determination by using bulk Richardson number method

203 The bulk Richardson number (Ri) is widely used for the climatological study of
204 PBLH from radiosonde measurements thanks to its applicability and reliability for all
205 atmospheric conditions (Anderson 2009; Seidel *et al.*, 2012). Ri, a good indicator of
206 turbulence and thermodynamic stability, is calculated as the ratio of turbulence due to
207 buoyancy to that due to mechanical shear, which is formulated as

$$208 \quad \text{Ri}(z) = \frac{\left(\frac{g}{\theta_{vs}}\right)(\theta_{vz}-\theta_{vs})z_{AG}}{(u_z-u_s)^2+(v_z-v_s)^2+(bu_*^2)} \quad (1)$$

209 where g is the gravitational acceleration, z_{AG} the AGL, θ_v the virtual potential
210 temperature, u_* the surface friction velocity, u and v the horizontal wind component,
211 and b the constant which is usually set to zero since friction velocity is much weaker
212 compared with the horizontal wind (Seidel *et al.*, 2012). The subscripts of z and s
213 denote the parameters at z height above ground and the ground level, respectively.

214 The critical value of $\text{Ri}(z)$ can be used to identify a statically stable layer atop the
215 PBL (Seibert *et al.*, 2000), and it is commonly taken as 0.25. Meanwhile, PBLH
216 estimates were found varying little by differing the input of critical values ($\text{Ri} =$
217 $0.2; 0.25; 0.3$) (Guo *et al.*, 2016). Therefore, the PBLH here is identified as the
218 interpolated height where $\text{Ri}(z)$ profile crosses the critical value of 0.25. The
219 determined PBLH was set invalid in the following two scenarios: (1) $\text{Ri}(z)$ in Eq. (1)
220 exceeds 0.25, where z is the second level of radiosonde measurement; (2) the estimated
221 PBLH is extremely high (for instance, 10 km), and it could mistake free-tropospheric
222 features.

223 3. Methodology

224 As shown in Fig. 2, there exist discernable biases between PBLH retrieved from
225 radiosonde (hereinafter referred to as PBLH_{RS}) and PBLH determined from ERA5
226 reanalysis (hereinafter referred to as $\text{PBLH}_{\text{ERA5}}$). The match procedures between
227 PBLH_{RS} and $\text{PBLH}_{\text{ERA5}}$ follow Guo *et al.* (2021). Noticeably, the PBLH bias (PBLH_{RS}
228 minus $\text{PBLH}_{\text{ERA5}}$) is less dependent on years, with a mean bias of 95.7 m, indicative of

229 a possible systematic PBLH underestimation of the ERA5 reanalysis. By contrast, the
230 underestimation is around 137 m during the daytime (Guo et al., 2021), which is
231 systematically larger than that during all days obtained in the present study. However,
232 the bias is found varying with seasons and local solar times (LST). More precisely, the
233 mean bias varies from 150 m in the March–April–May (MAM) to 64 m in the
234 September–October–November (SON), and from 309 m at 1700 LST to 1.8 m at 0000
235 LST. Moreover, the standard deviation of bias greatly changes from 64 m at 0100 LST
236 to 807 m at 1700 LST. The large uncertainty raised by $PBLH_{ERA5}$ during the daytime
237 motivated this study to establish a new PBLH dataset that would be more consistent
238 with observations.

239 Previous studies indicate that the bias could be physically attributed to the variables
240 such as SDDM and LTS (Guo et al., 2021). However, the potential correlations with
241 other variables, including DLWR, DSWR, SHF, LHF, evapotranspiration, transpiration,
242 total precipitation rate (TPR), soil moistures (SMs), as well as wind speed, pressure,
243 and air temperature at the near surface, have yet to be systematically investigated.
244 Figure 3 shows that the bias is positively correlated with SHF, transpiration, LTS, and
245 2-m near-surface temperature, with a correlation coefficient ranging from 0.39 to 0.9
246 based on 10 evenly split bins. However, these parameters could be independent. For
247 instance, evapotranspiration is determined by surface features which include plant
248 physiology, land cover, and soil moisture, and it is the most important non-radiative
249 process transmitting latent heat from the surface to the atmosphere (Cuxart and Boone,
250 2020). In addition, soil moisture probably contributes to decreases in the surface
251 sensible flux locally (Basha and Ratnam, 2009). We further perform correlation
252 analyses between the aforementioned variables and PBLH biases between radiosonde
253 and ERA5 reanalysis, and the statistical results are shown in Table 2

254 It is found that the PBLH bias is highly associated with the variations in land
255 properties, near-surface meteorological conditions, terrain elevations, LTS, and solar
256 cycles. Consequently, it is possible to predict the PBLH bias based on these potential
257 influential variables. Once the spatially resolved bias is available, a bias corrected

258 PBLH dataset, namely, a merged PBLH product (denoted as $PBLH_{merged}$ hereafter), can
259 be acquired by perturbing $PBLH_{ERA5}$ with the addition of predicted bias. This process
260 can be formulated as

$$261 \quad PBLH_{merged} = PBLH_{bias} + PBLH_{ERA5} \quad (2)$$

262 where $PBLH_{bias}$ denotes the PBLH bias to be predicted. Under this philosophy, here
263 we established a data-driven $PBLH_{bias}$ prediction model, with abovementioned factors
264 used as the potential input variables while the PBLH bias over radiosonde sites as the
265 learning target. Considering the possible dependence on magnitude of $PBLH_{ERA5}$ and
266 its corresponding LST, these two factors were also used as covariates in predicting
267 PBLH bias.

268 After testing with several machine learning models, such as the ridge regression,
269 the decision tree regressor, the support vector regressor, the multilayer perceptron
270 regression, and random forest (RF), we find the latter method gives the most proper and
271 robust prediction. Therefore, a RF regressor is established to give a prediction of
272 $PBLH_{bias}$, and it can be described as

$$273 \quad PBLH_{bias} = RF(DSWR, DLWR, LHF, SHF, EP, TP, SM10, SM40, SM100, \\ 274 \quad SM200, TPR, PBLHE, LTS, SDDEM, NSP, NST, NSWS, LST) \quad (3)$$

275 where the abbreviation RF represents the random forest regressor, and the other
276 acronyms and abbreviations are listed in Table 2. In the RF model, the hyper-parameters
277 of the maximum depth of the tree and the random state of the bootstrapping of the
278 samples are compiled to 20 and 5 in this analysis, respectively. The dataset that contains
279 the input array and the learning target is randomly divided into two parts, with 70% for
280 training and 30% for validation. All the data from 2011–2021 were included in the
281 model training stage. The following statistical metrics, including the mean squared error
282 (MSE), root mean square error (RMSE), arithmetic mean, and arithmetic mean of the
283 absolute difference, are applied to evaluate the performance of the prediction model.

284 4 Validation

285 Table 3a presents the prediction accuracy on the training and testing sets. Overall,
286 the RMSE and arithmetic mean on the training subset are 243 and -0.2 , respectively.
287 In comparison, these two metrics are 370 and -2.8 on the testing subset, implying the
288 presence of slight overfitting. To demonstrate the merit of $PBLH_{merged}$, we further
289 compare the PBLH bias before and after merging. As illustrated in Fig.4a, the mean
290 bias between $PBLH_{RS}$ and $PBLH_{merged}$ is -0.9 m, which is smaller than the bias
291 between $PBLH_{RS}$ and $PBLH_{ERA5}$. In addition, the mean of absolute bias decreases from
292 260 m ($PBLH_{RS}$ minus $PBLH_{ERA5}$) to 168 m ($PBLH_{RS}$ minus $PBLH_{merged}$), and the
293 standard derivation declines from 472 m to 241 m, as listed in Table 3b. Moreover, the
294 correlation coefficient between $PBLH_{RS}$ and $PBLH_{ERA5}$ is 0.59, and it increases to 0.92
295 between $PBLH_{RS}$ and $PBLH_{merged}$. More importantly, the bias between $PBLH_{RS}$ and
296 $PBLH_{merged}$ during the daytime is dramatically decreased to 20 m, compared to the bias
297 between $PBLH_{RS}$ and $PBLH_{ERA5}$ (300 m). These metrics clearly demonstrate a better
298 accuracy of $PBLH_{merged}$ than $PBLH_{ERA5}$, indicative of the merit of correcting modeling
299 biases in $PBLH_{ERA5}$.

300 Furthermore, the overview of PBLH bias ($PBLH_{RS}$ minus $PBLH_{merged}$) in terms of
301 spatial variation, and the seasonal variations over the four regions of interest are
302 presented in Fig. 5. As compared to the finding in Guo et al. (2021), the bias
303 dramatically decreases to dozens of meters for all the stations (Fig. 5d), many of which
304 slightly overestimate PBLH. More specifically, the PBLH over East Asia is
305 overestimated by around 6 m (Fig.5f), whereas it is underestimated by around 1 m over
306 Northern America (Fig. 5a). Based on the bias with near-global coverage, we could
307 infer that the merged model gives a more realistic PBLH estimate.

308 Intensive radiosonde observation is conducted across China in boreal summer
309 season at 0600 UTC (1400 Beijing Time) when the PBL is fully developed (Zhang et
310 al., 2018). In addition to the overall near-global spatial distribution, a deeper
311 investigation of $PBLH_{merged}$ across China at 0600 UTC is presented in Fig. 6. The spatial

312 distribution of $PBLH_{merged}$ exhibits a pronounced “Northwest High Southeast Low”
313 spatial pattern (Fig. 6a), which generally agrees with Zhang et al. (2018). The
314 correlation coefficient between $PBLH_{merged}$ and $PBLH_{RS}$ is as high as 0.99, indicating
315 their extreme consistencies in terms of spatial variations. The annual variations in
316 $PBLH_{merged}$, $PBLH_{RS}$, and $PBLH_{ERA5}$ follow a similar trend, achieving a maximum in
317 2013 and a minimum in 2019 (Fig. 6b). The variations in $PBLH_{merged}$ and $PBLH_{RS}$ are
318 rather close to each other. However, $PBLH_{ERA5}$ creates a different temporal variation,
319 and it is systematically underestimated, compared to $PBLH_{RS}$.

320 As a good case in point for the comparison of fine structures, we show the diurnal
321 variation of $PBLH_{merged}$ and $PBLH_{RS}$ at 0600 UTC over three stations in Fig. 7. Three
322 sites, including one in northwestern China where the highest PBLH is usually obtained,
323 one in northern China where the most intensive observations can be found, and one in
324 southern China where the lowest PBLH can be detected. The diurnal variations of
325 $PBLH_{merged}$ and $PBLH_{RS}$ are strongly correlated with the lowest correlation of 0.88
326 (Fig.7d). From Figs. 5-7, we can observe that the spatial-temporal variations of
327 $PBLH_{merged}$ and $PBLH_{RS}$ are in good agreement.

328 **5 Merged continental planetary boundary layer height**

329 The climatological mean of $PBLH_{merged}$ in four seasons at 0000 and 1200 UTCs
330 during the years from 2011 to 2021 is illustrated in Fig. 8, and the $PBLH_{RS}$ at the same
331 UTC and in the same season are overlaid as filled circles. At all UTCs and in all seasons
332 the $PBLH_{merged}$ is considerably high during the daytime and reaches a maximum of
333 around 2 km, especially in the afternoon, as compared to the nighttime. In addition,
334 $PBLH_{merged}$ experiences a noticeable seasonal variation. For instance, over Australia,
335 the $PBLH_{ERA5}$ in SON and December–January–February (DJF) seasons is about 400 m
336 larger than those of the other two seasons (Fig.8a–d), and vice versa in the Northern
337 Hemisphere. Moreover, we can observe that $PBLH_{merged}$ has a clear latitude- and
338 elevation-dependent. It decreases from approximately 2 km at low and middle latitudes

339 to around 0.8 km at high latitudes during the daytime. At similar latitudes, the
340 $PBLH_{merged}$ over terrain with a high elevation could be substantially larger than that
341 with a low elevation. For example, in DJF season and at 0000 UTC the $PBLH_{ERA5}$ over
342 the Andes Mountain is about 0.4 km higher than that over the surrounding flat region
343 (Fig. 8d). In a short conclusion, the spatial-temporal variability of the $PBLH_{merged}$ is
344 inevitably associated with local times, seasons, latitudes, terrain elevations, and
345 hemispheres.

346 In general, $PBLH_{merged}$ is remarkably consistent with $PBLH_{RS}$ in terms of seasonal
347 variation and diurnal cycle, especially at 0000 UTC and 1200 UTC when the radiosonde
348 measurement is comparatively sufficient. These findings suggest that the $PBLH_{merged}$
349 could adequately resolve the climatological variation of PBLH.

350 The difference in $PBLH_{merged}$ and $PBLH_{ERA5}$ during the years 2011–2021 at four
351 typical times is further illustrated in Fig. 9. Compared to $PBLH_{ERA5}$, the $PBLH_{merged}$ is
352 overall overestimated, with a mean overestimation of approximately 90 m. The
353 overestimation appears very close to the difference in $PBLH_{RS}$ and $PBLH_{ERA5}$. The
354 overestimation over North America at 0000 UTC, over East Asia and South Asia at
355 1200 UTC, and over Africa at 1800 UTC can be as high as 500 m. However, PBLH
356 over some areas, such as the Middle East at 0600 UTC and the Western United States
357 at 1800 UTC, is slightly underestimated by around 200 m.

358 **6 Conclusions and summary**

359 The general underestimation of PBLH by reanalysis dataset, especially during the
360 daytime, motivates the present analysis to generate a merged long-term high-resolution
361 seamless continental PBLH dataset (i.e., $PBLH_{merged}$) by integrating multi-modal data
362 products, which includes 185 million high-resolution radiosondes from the years 2011
363 to 2021, ERA5 reanalysis, and GLDAS product. The $PBLH_{merged}$ generated in this study
364 has a horizontal resolution of $0.25^\circ \times 0.25^\circ$ and a temporal resolution of 3 hours,
365 identical to $PBLH_{ERA5}$, but with much higher data accuracy.

366 Compared to the $PBLH_{RS}$, the $PBLH_{merged}$ is overestimated by around -0.9 m, which
367 is considerably smaller than the bias between $PBLH_{RS}$ and $PBLH_{ERA5}$ (95.7 m). During
368 the daytime, the mean and the standard derivation of bias are remarkably decreased
369 from 300 m and 600 m ($PBLH_{RS}$ minus $PBLH_{ERA5}$) to 20 m and 300 m ($PBLH_{RS}$ minus
370 $PBLH_{merged}$), respectively. In addition, the climatological variation of the merged PBLH
371 dataset is highly correlated with $PBLH_{RS}$, both in magnitude and spatial-temporal
372 variation. Moreover, the climatological mean of continental $PBLH_{merged}$ is around 90 m
373 higher than that of $PBLH_{ERA5}$, which is quantitatively consistent with the comparison
374 result of $PBLH_{RS}$ and $PBLH_{ERA5}$. Overall, the merged dataset closely agrees with the
375 radiosonde-derived PBLH in terms of magnitude and spatial-temporal variation.

376 In conclusion, the $PBLH_{merged}$ dataset is outstanding in terms of both spatiotemporal
377 coverage and good accuracy. This dataset could be of importance for advancing our
378 understanding of the PBL processes involved in air quality prediction, weather forecast,
379 and climate projection under global warming. In the future, with more dataset available
380 over the ocean, the global seamless PBLH dataset is warranted, and this needs more
381 field campaigns to be deployed over the open ocean or islands in the ocean in which
382 more intensive radiosonde balloons are launched. Besides, it is imperative to improve
383 the observational capability of satellite-based instruments in characterizing the
384 temperature and humidity profiles in the PBL, which no doubt helps fill the gaps of
385 atmospheric sounding over the ocean.

386 **Author contributions**

387 JG and FH conceptualized this study. JG and JZ carried out the dataset production with
388 comments from other co-authors. JG, JZ and JS drafted the first manuscript, and JS,
389 KB, and RL further revised it. JS established the model and its optimization. All authors
390 contributed to the discussion of result interpretation and helped finalized the submission.

391
392

393 **Competing interests**

394 The contact author has declared that neither they nor their co-authors have any
395 competing interests.

396

397 **Financial support**

398 This study is jointly supported by the Natural Science Foundation of China under grants
399 U2142209 and 62101203, the Hubei Provincial Natural Science Foundation of China
400 under grant KZ22Z3021, the Fundamental Research Funds for the Central Universities,
401 China University of Geosciences (Wuhan) under grant 162301192698, the
402 Fundamental Research Funds for the Central Universities, Huazhong Agricultural
403 University under grant 2662021XXQD002, the Chinese Academy of Sciences under
404 grant GXDA20040502, and Chinese Academy of Meteorological Sciences under grant
405 2021KJ029.

406 **Data availability**

407 The merged PBLH dataset and the related codes can be accessed at
408 <https://doi.org/10.5281/zenodo.6498004> (Guo et al., 2022).

409 ERA5 data is publicly accessible at
410 <https://cds.climate.copernicus.eu#!/search?text=ERA5&type=dataset> (ECMWF,
411 2019). NASA GLDAS can be accessed at:
412 [https://disc.gsfc.nasa.gov/datasets/GLDAS_NOAH025_3H_2.1/summary?keywords=
413 GLDAS \(NASA, 2021\).](https://disc.gsfc.nasa.gov/datasets/GLDAS_NOAH025_3H_2.1/summary?keywords=)

414 **References**

- 415 Anderson, P. S: Measurement of Prandtl number as a function of Richardson number
416 avoiding self-correlation, *Boundary Layer Meteorol*, 131, 345–362,
417 <https://doi.org/10.1007/s10546-009-9376-4>, 2009.
- 418 Baklanov, A. A., Grisogono, B., Bornstein, R., Mahrt, L., Zilitinkevich, S. S., Taylor,
419 P., Larsen, S.E., Rotach, M.W. and Fernando, H. J. S.: The nature, theory, and
420 modeling of atmospheric planetary boundary layers, *Bull Am Meteorol*
421 *Soc*, 92(2), 123–128, <https://doi.org/10.1175/2010BAMS2797.1>, 2011
- 422 Basha, G., and Ratnam, M. V.: Identification of atmospheric boundary layer height over
423 a tropical station using high-resolution radiosonde refractivity profiles:
424 Comparison with GPS radio occultation measurements, *J. Geophys. Res.*
425 *Atmos.*, 114(D16), <https://doi.org/10.1029/2008JD011692>, 2009.
- 426 Chan, K. M., and Wood, R.: The seasonal cycle of planetary boundary layer depth
427 determined using COSMIC radio occultation data, *J. Geophys. Res. Atmos.*,
428 118, 12 422–12 434, <https://doi.org/10.1002/2013JD020147>, 2013.
- 429 Cooper, D. I. and Eichinger, W. E.: Structure of the atmosphere in an urban planetary
430 boundary layer from lidar and radiosonde observations, *J. Geophys. Res.*
431 *Atmos.*, 99(D11), 22937–22948, <https://doi.org/10.1029/94JD01944>, 1994.
- 432 Cuxart, J., and Boone A. A.: Evapotranspiration over Land from a Boundary-Layer
433 Meteorology Perspective. *Boundary Layer Meteorol.*, 177, 427–459,
434 <https://doi.org/10.1007/s10546-020-00550-9>, 2020.
- 435 Davis, E. V., Rajeev, K., and Mishra, M.K.: Effect of clouds on the diurnal evolution
436 of the atmospheric boundary-layer height over a tropical coastal
437 station, *Boundary Layer Meteorol.*, 175(1), 135–152,
438 <https://doi.org/10.1007/s10546-019-00497-6>, 2020.
- 439 Davy, R., and Esau, I.: Differences in the efficacy of climate forcings explained by
440 variations in atmospheric boundary layer depth, *Nat. Commun.*, 7(1), 11690.
441 <https://doi.org/10.1038/ncomms11690>, 2016.

442 de Arruda Moreira, G., Guerrero-Rascado, J. L., Bravo-Aranda, J. A., Benavent-Oltra,
443 J. A., Ortiz-Amezcuca, P., Róman, R., Bedoya-Velásquez, A. E., Landulfo, E.
444 and Alados-Arboledas, L.: Study of the planetary boundary layer by microwave
445 radiometer, elastic lidar and Doppler lidar estimations in Southern Iberian
446 Peninsula, *Atmos Res.*, 213, 185–195,
447 <https://doi.org/10.1016/j.atmosres.2018.06.007>, 2018.

448 Degrazia, G. A., D. Anfossi, J. C. Carvalho, C. Mangia, T. Tirabassi and Campos Velho,
449 H. F.: Turbulence parameterisation for PBL dispersion models in all stability
450 conditions, *Atmos. Environ.*, 34(21), 3575–3583,
451 [https://doi.org/10.1016/S1352-2310\(00\)00116-3](https://doi.org/10.1016/S1352-2310(00)00116-3), 2000.

452 Ding, F., Iredell, L., Theobald, M., Wei, J., and Meyer, D.: PBL height from AIRS,
453 GPS RO, and MERRA-2 products in NASA GES DISC and their 10-year
454 seasonal mean intercomparison, *Earth Space Sci.*, 8,
455 e2021EA001859, <https://doi.org/10.1029/2021EA001859>, 2021.

456 Dirmeyer, P. A., Wang, Z., Mbulu, M. J. and Norton, H. E.: Intensified land surface
457 control on boundary layer growth in a changing climate, *Geophys. Res.
458 Lett.*, 41(4), 1290–1294, <https://doi.org/10.1002/2013GL058826>, 2014.

459 ECMWF.: ERA5 reanalysis [data set], Retrieved from
460 <https://cds.climate.copernicus.eu/#!/search?text=ERA5&type=dataset>, 2019.

461 Edson, J. B., Jampana, V., Weller, R. A., Bigorre, S. P., Plueddemann, A. J., Fairall, C.
462 W., Miller, S. D., Mahrt, L., Vickers, D., and Hersbach, H.: On the Exchange of
463 Momentum over the Open Ocean, *J Phys Oceanogr.*, 43(8), 1589–1610,
464 <https://doi.org/10.1175/JPO-D-12-0173.1>, 2013.

465 Esau, I., and Zilitinkevich, S.: On the role of the planetary boundary layer depth in the
466 climate system. *Adv. Sci. Res.*, 4, 63, <https://doi.org/10.5194/asr-4-63-2010>,
467 2010.

468 Guo, J., Li, Y., Cohen, J. B., Li, J., Chen, D., Xu, H., Liu, L., Yin, J., Hu, K., and Zhai.
469 P.: Shift in the temporal trend of boundary layer height in China using long-

470 term (1979–2016) radiosonde data, *Geophys. Res. Lett.*, 46, 6080–6089,
471 <https://doi.org/10.1029/2019GL082666>, 2019.

472 Guo, J., Miao, Y., Zhang, Y., Liu, H., Li, Z., Zhang, W., He, J., Lou, M., Yan, Y., Bian,
473 L., and Zhai, P.: The climatology of planetary boundary layer height in China
474 derived from radiosonde and reanalysis data, *Atmos. Chem. Phys.*, 16, 13309–
475 13319, <https://doi.org/10.5194/acp-16-13309-2016>, 2016.

476 Guo, J., Zhang, J., Yang, K., Liao, H., Zhang, S., Huang, K., Lv, Y., Shao, J., Yu, T.,
477 Tong, B., Li, J., Su, T., Yim, S. H. L., Stoffelen, A., Zhai, P., and Xu, X.:
478 Investigation of near-global daytime boundary layer height using high-
479 resolution radiosondes: first results and comparison with ERA5, MERRA-2,
480 JRA-55, and NCEP-2 reanalyses, *Atmos. Chem. Phys.*, 21, 17079–17097,
481 <https://doi.org/10.5194/acp-21-17079-2021>, 2021.

482 Guo, J., Zhang, J., Shao, J.: A Harmonized Global Continental High-resolution
483 Planetary Boundary Layer Height Dataset Covering 2017-2021 [data set],
484 <https://zenodo.org/record/6498004>, 2022.

485 Hennemuth, B., and Lammert, A.: Determination of the atmospheric boundary layer
486 height from radiosonde and lidar backscatter, *Boundary Layer Meteorol.*,
487 120(1), 181–200, <https://doi.org/10.1007/s10546-005-9035-3>, 2006.

488 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J.,
489 Nicolas, J., Peubey, C., Radu, R., Schepers, D. and Simmons, A.: The ERA5
490 global reanalysis, *Q. J. R. Meteorol. Soc.*, 146(730), 1999–2049,
491 <https://doi.org/10.1002/qj.3803>, 2020.

492 Holzworth, G. C.: Estimates of mean maximum mixing depths in the contiguous United
493 States, *Mon. Wea. Rev.*, 92, 235–242, [https://doi.org/10.1175/1520-0493\(1964\)092,0235:EOMMMD.2.3.CO;2](https://doi.org/10.1175/1520-0493(1964)092,0235:EOMMMD.2.3.CO;2), 1964.

495 Hu, X. M., Nielsen-Gammon, J. W. and Zhang, F.: Evaluation of three planetary
496 boundary layer schemes in the WRF model, *J Appl Meteorol Climatol.*, 49(9),
497 1831–1844, <https://doi.org/10.1175/2010JAMC2432.1>, 2010.

498 Lammert, A., and Bösenberg, J.: Determination of the con- vective boundary-layer
499 height with laser remote sensing, *Bound.-Layer Meteor.*, 119, 159–170,
500 <https://doi.org/10.1007/s10546-005-9020-x>, 2006.

501 Li, Q., Zhang, H., Cai, X. et al.: The impacts of the atmospheric boundary layer on
502 regional haze in North China, *npj Clim Atmos Sci.*, 4(1), 1–10.
503 <https://doi.org/10.1038/s41612-021-00165-y>, 2021.

504 Li, Z., Guo, J., Ding, A., Liao, H., Liu, J., Sun, Y., Wang, T., Xue, H., Zhang, H. and
505 Zhu, B.: Aerosol and boundary-layer interactions and impact on air quality. *Natl.*
506 *Sci. Rev.*, 4(6), 810–833, <https://doi.org/10.1093/nsr/nwx117>, 2017.

507 Liu, B., Y. Ma, J. Guo, W. Gong, Y. Zhang, F. Mao, J. Li, X. Guo, and Shi, Y.:
508 Boundary layer heights as derived from ground-based radar wind profiler in
509 Beijing, *IEEE Trans. Geosci. Remote Sens.* 57(10), 8095–8104,
510 <https://doi.org/10.1109/TGRS.2019.2918301>, 2019.

511 Lou, M., J. Guo, L. Wang, H. Xu, D. Chen, Y. Miao, Y. Lv, Y. Li, X. Guo, S. Ma, and
512 Li, J.: On the relationship between aerosol and boundary layer height in summer
513 in China under different thermodynamic conditions. *Earth Space Sci.*, 6(5),
514 887–901, <https://doi.org/10.1029/2019EA000620>, 2019.

515 McGrath-Spangler, E. L., and Denning, A. S.: Estimates of North American
516 summertime planetary boundary layer depths derived from space-borne lidar. *J.*
517 *Geophys. Res.*, 117, D15101, <https://doi.org/10.1029/2012JD017615>, 2012.

518 Min, M., Bai, C., Guo, J., Sun, F., Liu, C., Wang, F., Xu, H., Tang, S., Li, B., Di, D.
519 and Dong, L.: Estimating summertime precipitation from Himawari-8 and
520 global forecast system based on machine learning, *IEEE Trans Geosci Remote*
521 *Sens.*, 57(5), 2557–2570, <https://doi.org/10.1109/TGRS.2018.2874950>, 2018.

522 NASA.: Global Land Data Assimilation System [data set], Retrieved from
523 https://disc.gsfc.nasa.gov/datasets/GLDAS_CLSM025_DA1_D_2.2/summary
524 [?keywords=GLDAS](https://disc.gsfc.nasa.gov/datasets/GLDAS_CLSM025_DA1_D_2.2/summary?keywords=GLDAS), 2021.

525 Petäjä, T., Järvi, L., Kerminen, VM. et al.: Enhanced air pollution via aerosol-boundary
526 layer feedback in China, *Sci. Rep.*, 6, 18998. <https://doi.org/10.1038/srep18998>,
527 2016.

528 Rodell, M., Houser, P. R., Jambor, U. E. A., et al.: The global land data assimilation
529 system. *Bull. Am. Meteorol. Soc.*, 85(3), 381–394,
530 <https://doi.org/10.1175/BAMS-85-3-381>, 2004.

531 Saha, S., Sharma, S., Kumar, K.N., Kumar, P., Lal, S. and Kamat, D.: Investigation of
532 atmospheric boundary layer characteristics using ceilometer lidar, COSMIC
533 GPS RO satellite, radiosonde and ERA-5 reanalysis dataset over Western Indian
534 region, *Atmos Res.*, 268, 105999,
535 <https://doi.org/10.1016/j.atmosres.2021.105999>, 2022.

536 Seibert, P., Beyrich, F., Gryning, S.-E., Joffre, S., Rasmussen, A., and Tercier,
537 P.: Review and intercomparison of operational methods for the determination
538 of the mixing height, *Atmos. Environ.*, 34, 1001–1027,
539 [https://doi.org/10.1016/S1352-2310\(99\)00349-0](https://doi.org/10.1016/S1352-2310(99)00349-0), 2000.

540 Seidel, D. J., Ao, C. O., and Li, K.: Estimating climatological planetary boundary layer
541 heights from radiosonde observations: Comparison of methods and uncertainty
542 analysis, *J. Geophys. Res. Atmos.*, 115(D16).
543 <https://doi.org/10.1029/2009JD013680>, 2010.

544 Seidel, D. J., Zhang, Y., Beljaars, A., Golaz, J. C., Jacobson, A.R. and Medeiros, B.:
545 2012. Climatology of the planetary boundary layer over the continental United
546 States and Europe, *J. Geophys. Res. Atmos.*, 117(D17),
547 <https://doi.org/10.1029/2012JD018143>, 2012.

548 Stull, R. B.: *An Introduction to Boundary Layer Meteorology*. Kluwer Academic, 666
549 pp, 1988.

550 Su, T., Li, Z., and Kahn, R.: Relationships between the planetary boundary layer height
551 and surface pollutants derived from lidar observations over China: regional
552 pattern and influencing factors, *Atmos. Chem. Phys.*, 18, 15921–15935,
553 <https://doi.org/10.5194/acp-18-15921-2018>, 2018.

554 Teixeira, J., Piepmeier, J. R., Nehrir, A. R., Ao, C. O., Chen, S. S., Clayson, C. A.,
555 Fridlind, A. M., Lebsock, M., Mc-Carty, W., Salmun, H., Santanello, J. A.,
556 Turner, D. D., Wang, Z., and Zeng, X.: Toward a global planetary boundary
557 layer observing system: the NASA PBL incubation study team report, NASA
558 PBL Incubation Study Team, 134 pp., available at:
559 [https://science.nasa.gov/science-red/s3fs-](https://science.nasa.gov/science-red/s3fs-public/atoms/files/NASAPBLIncubationFinalReport.pdf)
560 [public/atoms/files/NASAPBLIncubationFinalReport.pdf](https://science.nasa.gov/science-red/s3fs-public/atoms/files/NASAPBLIncubationFinalReport.pdf), last access: 28 April
561 2022.

562 Wallace, J. M. and Hobbs, P. V: Atmospheric Science: An Introductory Survey,
563 Academic Press, Burlington, MA., 2006.

564 Wang, X. and Wang, K.: Homogenized variability of radiosonde-derived atmospheric
565 boundary layer height over the global land surface from 1973 to 2014, *J.*
566 *Clim.*, 29(19), 6893–6908, <https://doi.org/10.1175/JCLI-D-15-0766.1>, 2016.

567 Wang, Y., A. Khalizov, M. Levy, and Zhang, R.: New Directions: Light Absorbing
568 Aerosols and Their Atmospheric Impacts, *Atmos. Environ.*, 81, 713–715,
569 <https://doi.org/10.1016/j.atmosenv.2013.09.034>, 2013.

570 Xu, Z., Chen, H., Guo, J., and Zhang, W.: Contrasting effect of soil moisture on the
571 daytime boundary layer under different thermodynamic conditions in summer
572 over China, *Geophys Res. Lett.*, 48, e2020GL090989. [https://doi.](https://doi.org/10.1029/2020GL090989)
573 [org/10.1029/2020GL090989](https://doi.org/10.1029/2020GL090989), 2021.

574 Yang, X., Zhao, C., Guo, J., and Wang, Y.: Intensification of aerosol pollution
575 associated with its feedback with surface solar radiation and winds in Beijing,
576 *J. Geophys. Res. Atmos.*, 121, 4093–4099,
577 <https://doi.org/10.1002/2015JD024645>, 2016.

578 Zhang, J., Guo, J. P., Zhang, S. D., and Shao, J.: Inertia-gravity wave energy and
579 instability drive turbulence: evidence from a near-global high-resolution
580 radiosonde dataset, *Clim. Dyn.*, 1–14, [https://doi.org/10.1007/s00382-021-](https://doi.org/10.1007/s00382-021-06075-2)
581 [06075-2](https://doi.org/10.1007/s00382-021-06075-2), 2022.

582 Zhang, W., Guo, J., Miao, Y., Liu, H., Song, Y., Fang, Z., He, J., Lou, M., Yan, Y., Li,
583 Y., and Zhai, P.: On the summertime planetary boundary layer with different
584 thermodynamic stability in China: A radiosonde perspective, *J. Clim.*, 31(4),
585 1451–1465, <https://doi.org/10.1175/JCLI-D-17-0231.1>, 2018.

586

587

588

589

590

591

592

593

594

595

596

597

598

599

600

601

602 **Table 1.** Basic information of data used in the present study, including data source, the
 603 number of stations, vertical resolution, and the years with data curation.

Data source	Number of station	Vertical resolution	Years
CMA	120	5–8 m	2011–2021
NOAA	89	5 m	2011–2021
GRUAN	8	5 m	2011–2021
CEDA	12	10 m	2011–2021
University of Wyoming	125	5–10 m	2017–2021
German Deutscher Wetterdienst	14	10 m	2011–2021

604

605

606

607

608

609

610

611

612

613

614

615

616 **Table 2.** Summary of input parameters of machine learning algorithms, and the
617 corresponding statistical metrics for their correlation analyses between with PBLH bias
618 between radiosonde and ERA5 reanalysis, including correlation coefficient and
619 confidence level.

Parameters	Acronyms	Data sources	Correlation coefficient	Confidence level
Downward shortwave radiation	DSWR	GLDAS	0.14	100%
Downward longwave radiation	DLWR	GLDAS	0.02	100%
Latent heat flux	LHF	GLDAS	0.14	100%
Sensible heat flux	SHF	GLDAS	0.10	100%
Evapotranspiration	EP	GLDAS	0.14	100%
Transpiration	TP	GLDAS	-0.02	100%
Soil moisture 0-10cm	SM10	GLDAS	-0.04	100%
Soil moisture 10-40cm	SM40	GLDAS	-0.03	100%
Soil moisture 40-100cm	SM100	GLDAS	-0.02	100%
Soil moisture 100-200cm	SM200	GLDAS	-0.03	100%
Total precipitation rate	TPR	GLDAS	-0.02	100%
Boundary layer height	PBLH _{ERA5}	ERA5	-0.10	100%
Lower tropospheric stability	LTS	ERA5	0.10	100%
Standard deviation of orography height	SDDEM	ERA5	0.06	100%
Near-surface pressure	NSP	ERA5	-0.11	100%
Near-surface temperature	NST	ERA5	0.05	100%
Near-surface wind speed	NSWS	ERA5	-0.08	100%
Local solar time	LST	-	0.17	100%

620

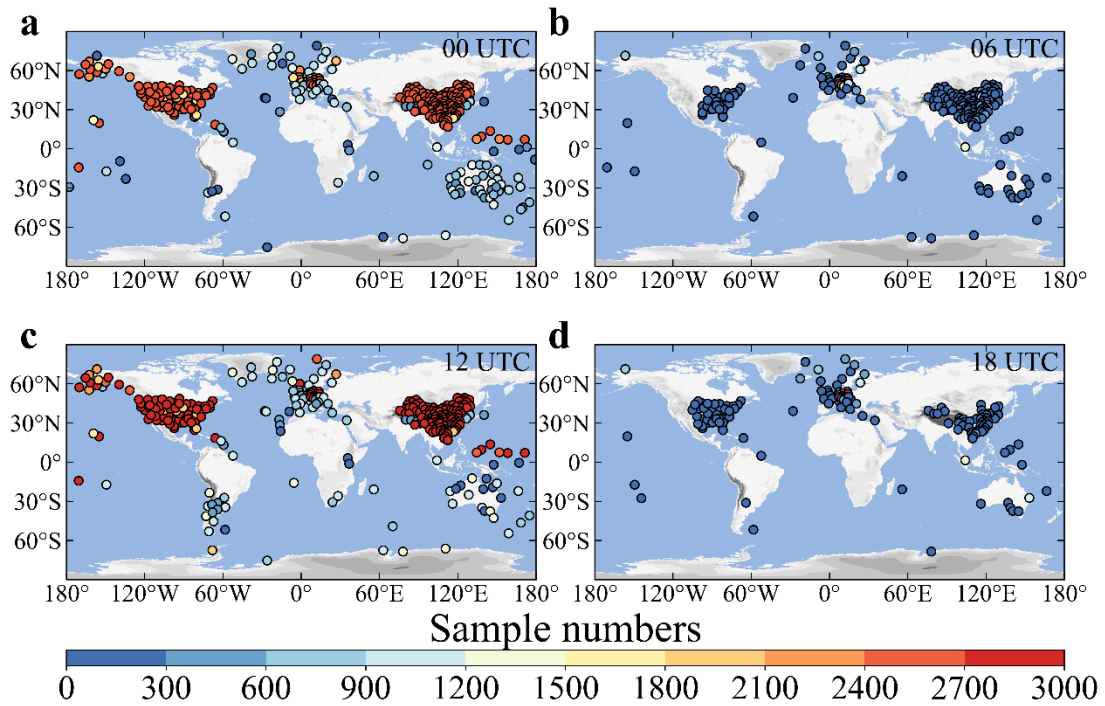
621

622 **Table 3.** Basic information on evaluation indices. MSE, mean squared error; RMSE,
 623 root mean square error; ABSmean, mean of the absolute bias; STD, standard derivation;
 624 RMS, root mean square.

(a) evaluation indices of the training set and test set				
	MSE	RMSE	Mean	ABSmean
Train set	59176	243	-0.2	152
Predict set	136971	370	-2.8	204

(b) evaluation indices of PBLH bias				
	Mean	ABSmean	STD	RMS
$PBLH_{RS} - PBLH_{ERA5}$	95.7	260	472	481
$PBLH_{RS} - PBLH_{merged}$	-0.9	168	241	287

625
 626
 627
 628
 629
 630
 631
 632
 633
 634



635

636 **Figure 1.** Spatial distribution of sample number (color circles) for each radiosonde

637 station at 0000 (a), 0600 (b), 1200 (c), and 1800 UTC from the years 2011 to 2021.

638 Stations with less than 10 samples are not indicated.

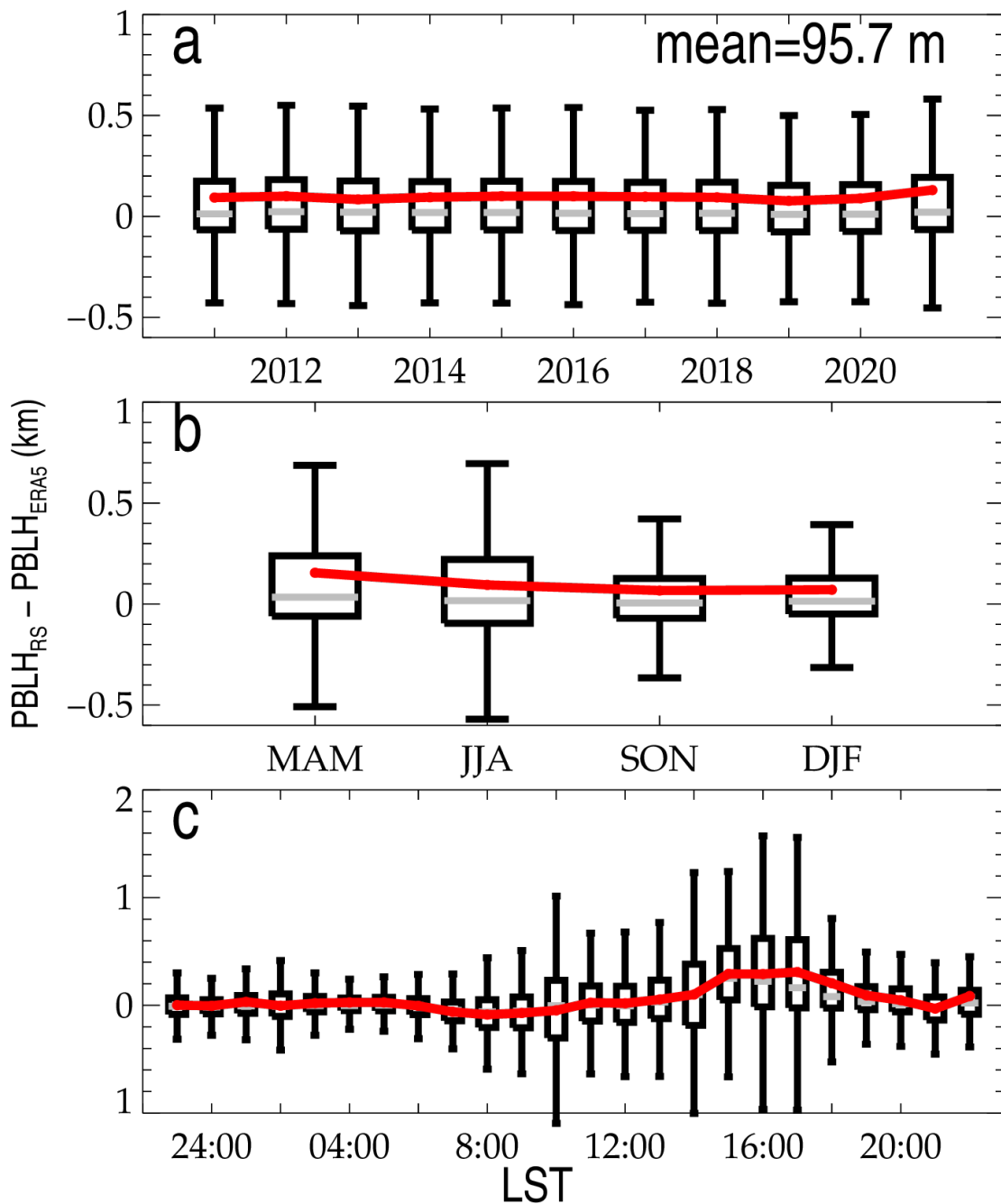
639

640

641

642

643



645

646

Figure 2. Evolution of the difference between $PBLH_{ERA5}$ and $PBLH_{RS}$ at various time

647

scales: different years (a), different seasons (b), and at different local times (c). MAM,

648

March–April–May; JJA, June–July– August; SON, September–October–November;

649

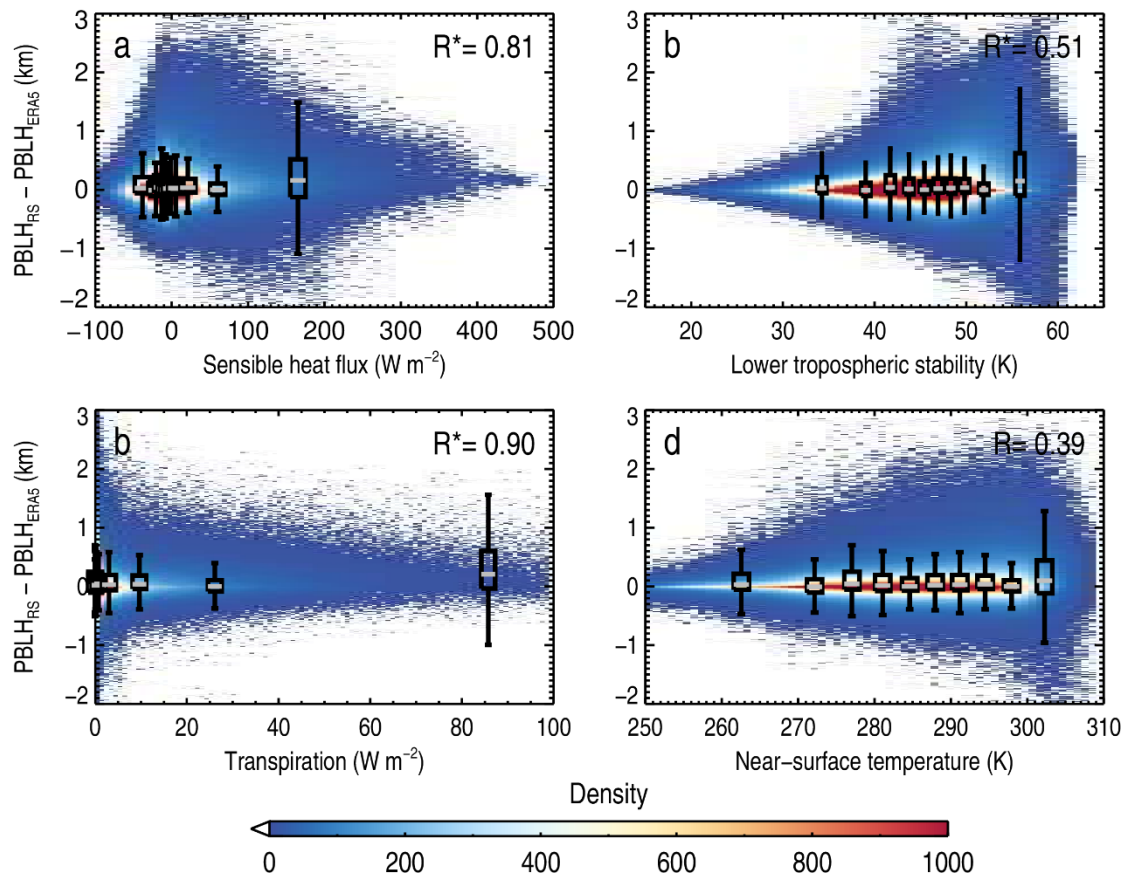
DJF, December–January–February. The mean bias is labelled in the upper right corner

650

of panel (a). Note that the southern hemisphere DJF (JJA) is combined with northern

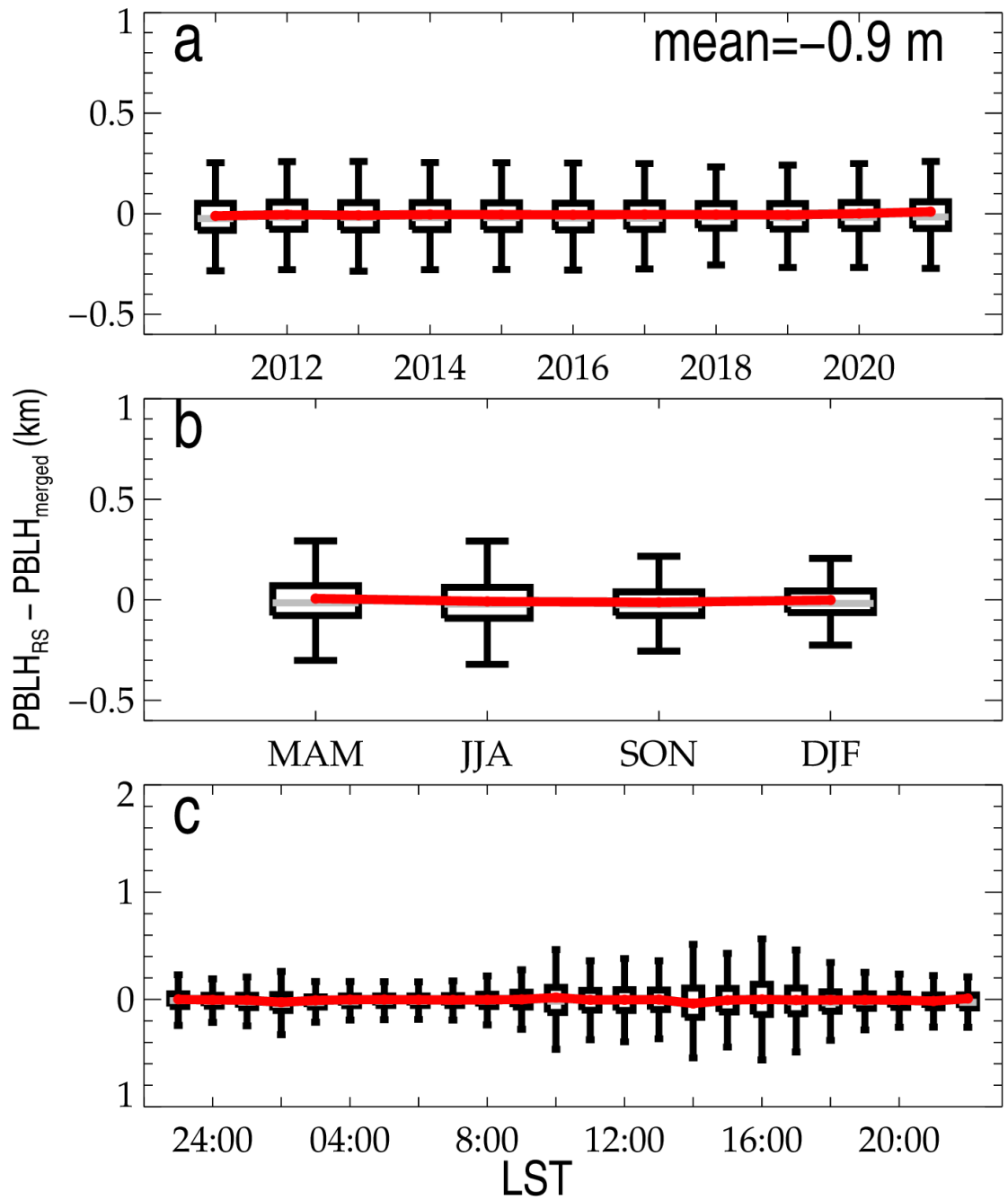
651

hemisphere JJA (DJF).



653

654 **Figure 3.** The joint distribution of the difference in PBLH_{RS} and PBLH_{ERA5} and the
 655 surface sensible heat flux (a), the lower tropospheric stability (b), transpiration (c), and
 656 the near-surface temperature (d). The box-and-whisker plots in 10 evenly intervals are
 657 overlaid in each panel, and the correlation coefficients are marked in the upper right
 658 corner of each panel, wherein the star superscripts indicate that the values are
 659 statistically significant ($p < 0.05$).



661

662 **Figure 4.** Similar to Figure 3, but for the difference between $PBLH_{RS}$ and $PBLH_{merged}$.

663

664

665

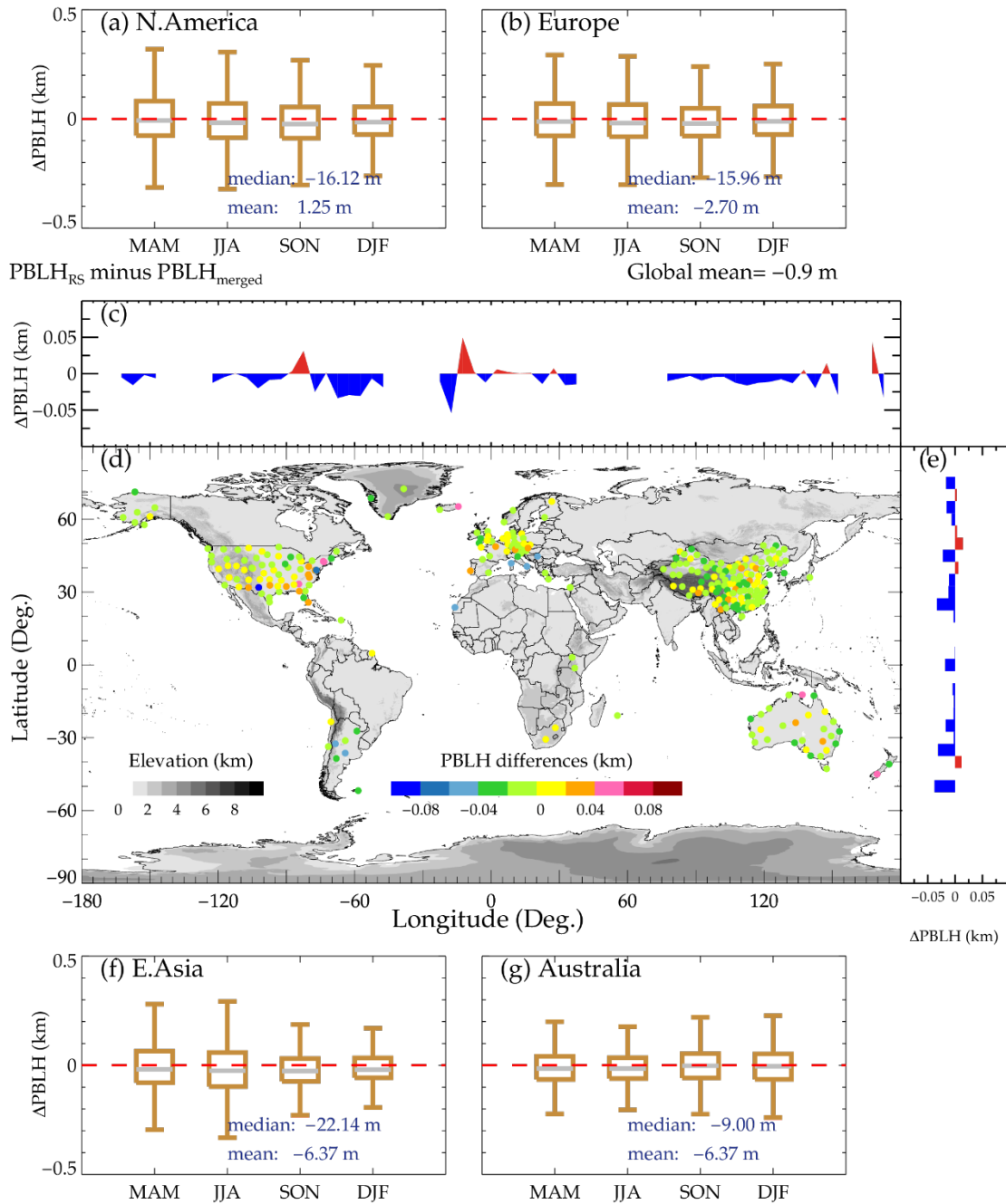
666

667

668

669

670



671

672 **Figure 5.** Spatial variations of PBLH differences between $PBLH_{RS}$ and $PBLH_{merged}$. (d)

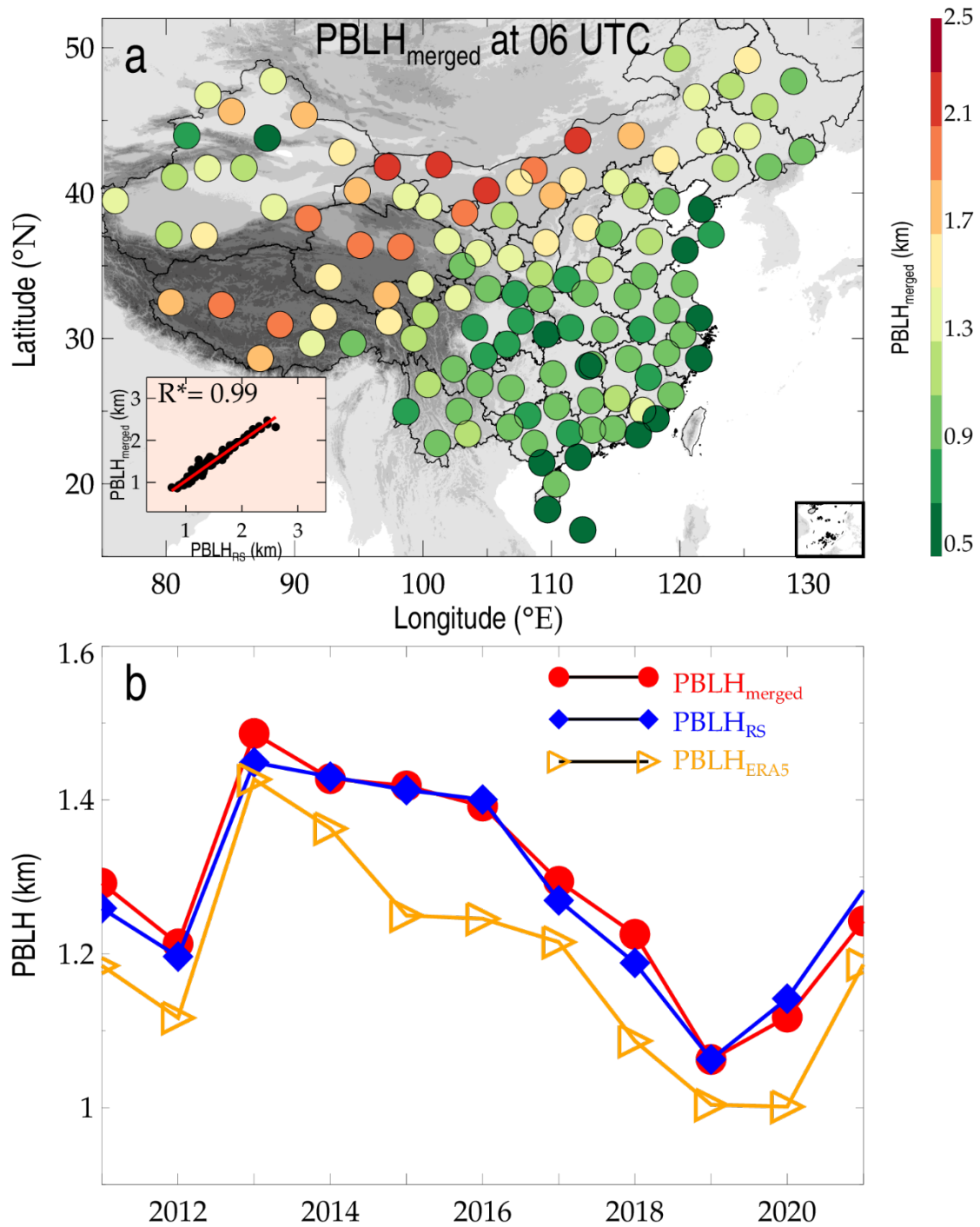
673 indicates the overall spatial distribution, and (c) and (d) illustrate its longitudinal and

674 latitudinal variations. (a), (b), (f), (g) represent the seasonal variations over the four

675 regions of interest, including North America, Europe, East Asia, and Australia. MAM,

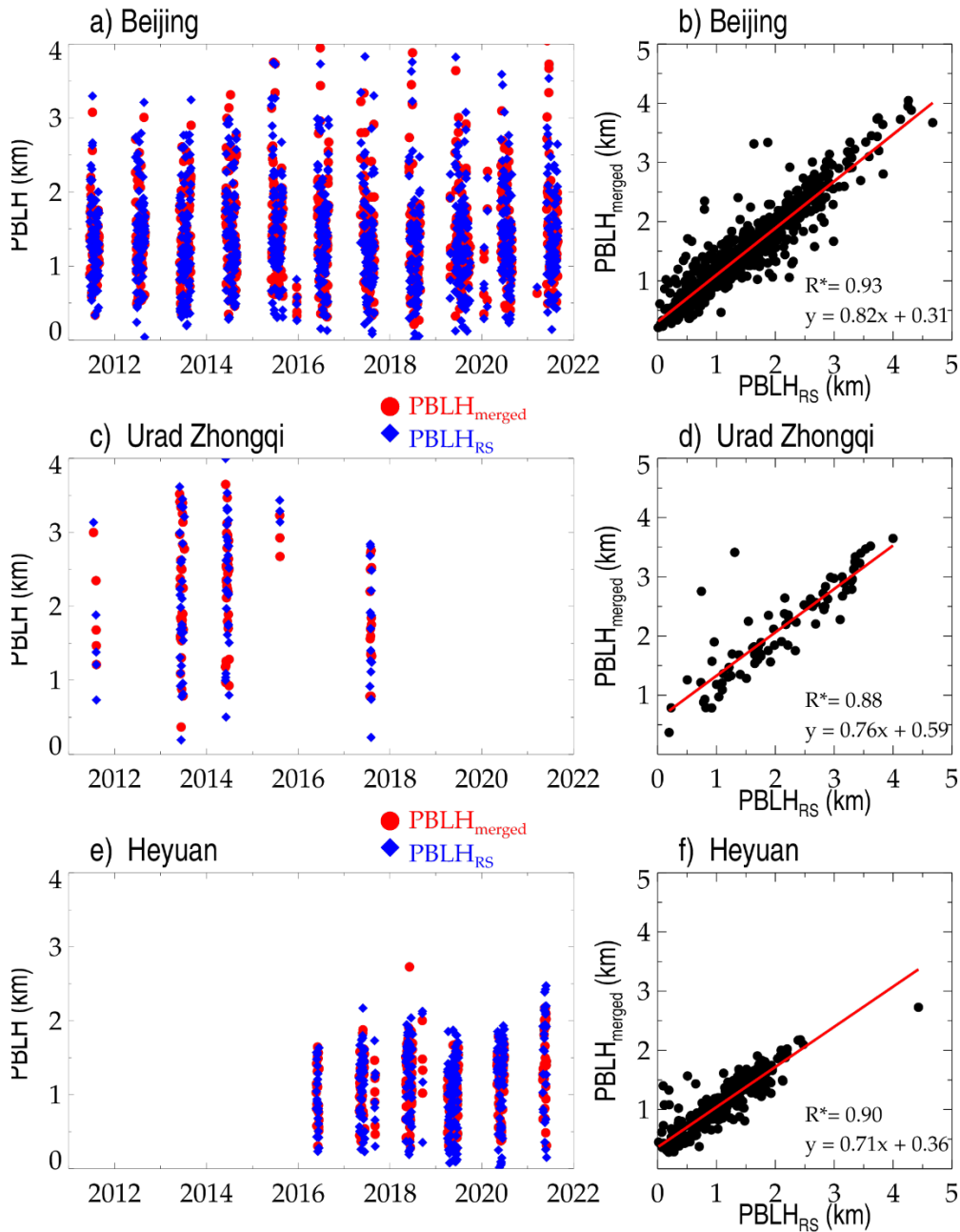
676 March–April–May; JJA, June–July–August; SON, September–October–November;

677 DJF, December–January–February.



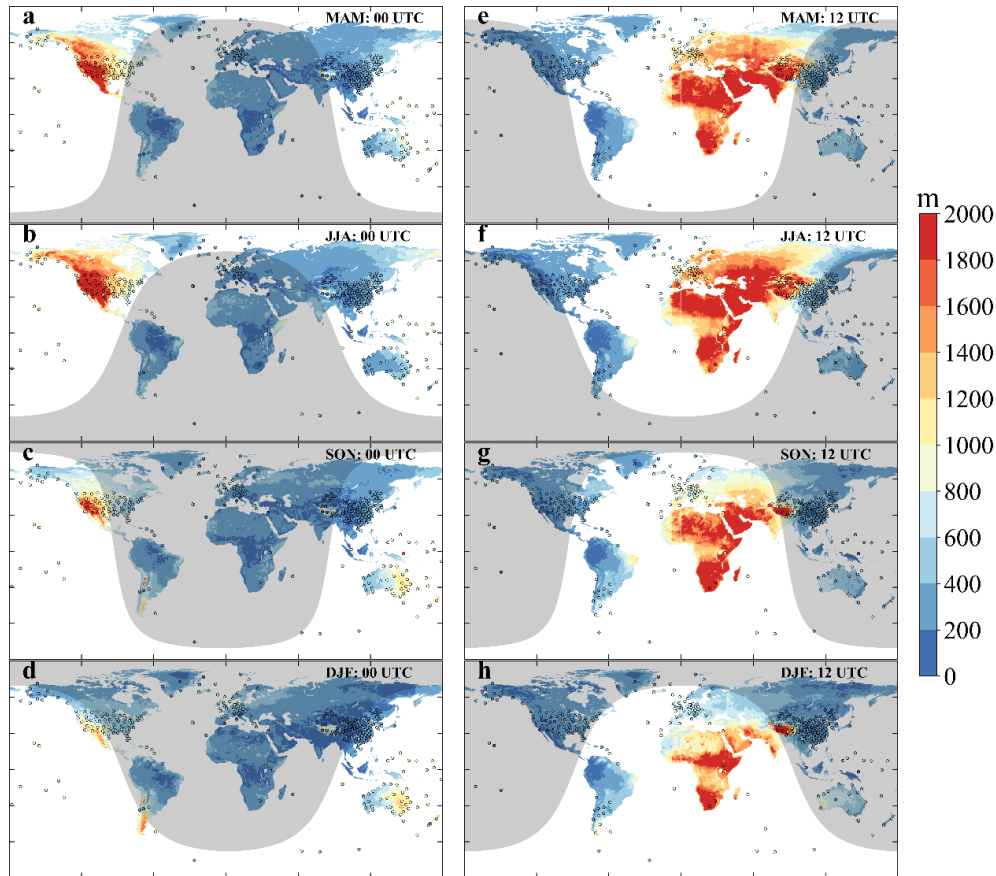
679

680 **Figure 6.** (a) Spatial distributions of the $PBLH_{merged}$ at 0600 UTC across China for the
 681 years 2011 to 2021. The scatter plot in the left bottom of the panel illustrates the
 682 statistical correlation between $PBLH_{merged}$ and $PBLH_{RS}$, where the star superscripts
 683 indicate that the values are statistically significant ($p < 0.05$). Also shown are the
 684 temporal evolution of annual average $PBLH_{merged}$, $PBLH_{RS}$, and $PBLH_{ERA5}$ during the
 685 period 2011 to 2021 (b).



687

688 **Figure 7.** Temporal variations of PBLH_{merged} (red) and PBLH_{RS} (blue) at Beijing
 689 (39.8°N, 116.47°E) (a), the Urad Zhongqi station (41.3°N, 108.3°E) (b) in the Nei
 690 Monggol Autonomous Region, and (c) the Heyuan (23.7°N, 114.7°E) station in the
 691 Guangdong province. (b), (d), and (f) demonstrate the joint-distributions of PBLH_{RS}
 692 and PBLH_{merged}, and correlation coefficients (R) and the fitted linear functions are given
 693 in the bottom right corner, where the star superscripts indicate that the values are
 694 statistically significant (p<0.05).



695

696 **Figure 8.** Spatial distribution of the PBLH at 0000 (a-d), and 1200 UTC (e-h) in four
 697 seasons over land produced by the merged algorithms proposed here (i-l). The colored
 698 solid circles indicate the PBLH retrieved from high-resolution radiosondes. The
 699 shadow zones show nighttime regions, depending on the solar zenith angle on 15 April
 700 2019 (MAM), 15 July 2019 (JJA), 15 October 2019 (SON), and 15 January 2019 (DJF).
 701 MAM, March–April–May; JJA, June–July–August; SON, September–October–
 702 November; DJF, December–January–February.

703

704

705

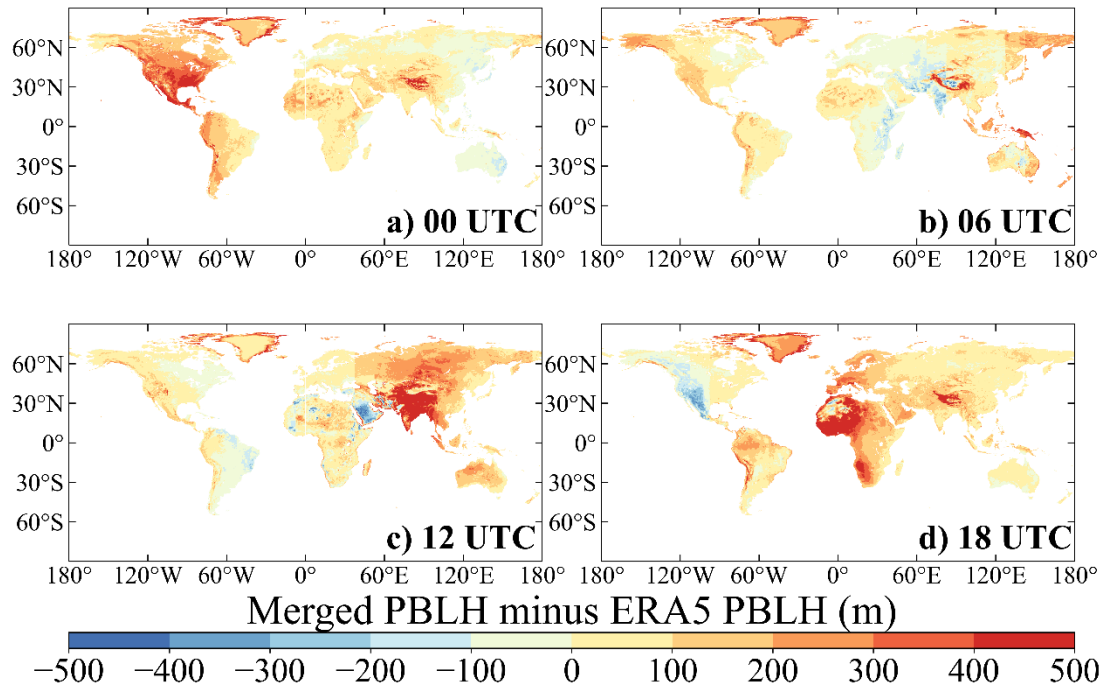
706

707

708

709

710



711

712

713

714

Figure 9. The spatial distributions of PBLH differences between the merged dataset and ERA5 reanalysis from the years 2011 to 2021 at 0000 (a), 0600 (b), 1200 (c), and 1800 UTC (d).