

Verifying Operational Forecasts of Land-Sea Breeze and Boundary Layer

Mixing Processes

Ewan Short*

*School of Earth Sciences, and ARC Centre of Excellence for Climate Extremes, The University of
Melbourne, Melbourne, Victoria, Australia.*

**Corresponding author address:* School of Earth Sciences, The University of Melbourne, Melbourne, Victoria, Australia.

E-mail: `shorte1@student.unimelb.edu.au`

ABSTRACT

9 Forecasts issued by the Australian Bureau of Meteorology (BoM) are based
10 on automated post-processed model data that is edited by human forecasters.
11 Two types of edits are commonly made to the wind fields. These edits aim
12 to improve how the influence of boundary layer mixing and land-sea breeze
13 processes are represented in the forecast. In this study we compare the di-
14 urnally varying component of the BoM’s official edited wind forecast, with
15 that of station observations and unedited model datasets, to assess changes
16 to error and bias resulting from these edits. We consider coastal locations
17 across Australia over June, July and August 2018, aggregating data over three
18 spatial scales. The edited forecast generally only produces a lower mean ab-
19 solute error than model guidance at the coarsest spatial scale ($500 \times 150 \text{ km}^2$
20 to $2000 \times 150 \text{ km}^2$), but can achieve lower seasonal biases over all spatial
21 scales. However, the edited forecast only reduces errors or biases at partic-
22 ular times and locations, and rarely produces lower errors or biases than all
23 model guidance products simultaneously. To better understand the biases in
24 the diurnal wind cycles, we fit modified ellipses to the temporal hodographs
25 of seasonally averaged diurnal wind cycles. Biases in the official forecast di-
26 urnal cycle vary with location for multiple reasons, including biases in the
27 directions sea-breezes approach coastlines, amplitude and shape biases in the
28 hodographs, and disagreement as to whether sea-breeze or boundary layer
29 mixing processes contribute most to the diurnal cycle.

30 1. Introduction

31 Modern weather forecasts are typically produced by models in conjunction with human fore-
32 casters. Forecasters working for the Australian Bureau of Meteorology (BoM) construct a seven
33 day forecast by loading model data into a software package called the Graphical Forecast Edi-
34 tor (GFE), then editing this model data using tools within the GFE. Forecasters working for the
35 United States National Weather Service also use GFE and utilise a similar approach. Forecasters
36 can choose which model to base their forecast on, and refer to this as a choice of *model guidance*.
37 Edits are typically made to account for processes that are under-resolved at the resolutions of the
38 model guidance products, or to correct for perceived biases of the model guidance being used.
39 In Australia, the resulting gridded forecast datasets are provided to the public through the BoM's
40 online MetEye data browser (Bureau of Meteorology 2019), and are also translated into text and
41 icon forecasts algorithmically.

42 Forecasters, and the weather services that employ them, have good reasons for ensuring the
43 diurnally varying component of their wind forecasts are as accurate as possible. In addition to the
44 significant contribution diurnal wind cycles make to overall wind fields (e.g. Dai and Deser 1999),
45 diurnal wind cycles are important for the ventilation of pollution, with sea-breezes transporting
46 clean maritime air inland, where it helps flush polluted air out of the boundary layer (Miller et al.
47 2003; Physick and Abbs 1992). Furthermore, diurnal wind cycles affect the function of wind
48 turbines (Englberger and Dörnbrack 2018) and the design of wind farms (Abkar et al. 2016), as
49 daily patterns of boundary layer stability affect turbine wake turbulence, and the losses in wind
50 power that result.

51 Australian forecasters generally make two types of edits to the surface wind fields on a routine
52 daily basis. The first involves modifying the surface winds after sunrise at locations where the

53 forecaster believes the model guidance is providing a poor representation of boundary layer mixing
54 processes. Boundary layer mixing occurs as the land surface heats up, producing an unstable
55 boundary layer which transports momentum downward to the surface layer. Before this mixing
56 occurs, winds are typically both weaker and ageostrophically oriented due to surface friction (Lee
57 2018), and so mixing can affect both the speed and direction of the surface winds. Australian
58 forecasters perform these edits using a GFE tool which allows them to specify a region over which
59 to apply the edit, a height z and a percentage p , with the tool then calculating a weighted average
60 of the surface winds and winds at z , weighted by p .

61 The second type of edit involves changing the afternoon and evening surface winds around those
62 coastlines where the forecaster believes the model guidance is resolving the sea-breeze poorly.
63 Similarly to with boundary layer mixing, these edits are performed using a GFE tool that allows
64 forecasters to trace out the relevant coastline graphically, choose a wind speed and a time, with the
65 tool then smoothly blending in winds of the given speed perpendicular to the traced coastline at
66 the given time.

67 To our knowledge, no published work has assessed the diurnal component of human edited
68 forecasts, although some previous studies have assessed the performance of different operational
69 models at specific locations. Svensson et al. (2011) examined thirty different operational model
70 simulations, including models from most major forecasting centres utilising most commonly used
71 boundary layer parametrisation schemes, and compared their performance with a large eddy sim-
72 ulation (LES), and observations at Kansas, USA, during October 1999. They found that both the
73 models and LES failed to capture the roughly 6 kn ($1 \text{ kn} \approx 0.514 \text{ m s}^{-1}$) jump in wind speeds
74 shortly after sunrise, and underestimated morning low level turbulence and wind speeds.

75 Other studies have assessed near-surface wind forecasts, verifying the total wind speeds, not
76 just the diurnal component. Pinson and Hagedorn (2012) studied the 10 m wind speeds from the

77 European Centre for Medium Range Weather Forecasting (ECMWF) operational model ensemble
78 across western Europe over December, January, February 2008/09. They found that the worst
79 performing regions were coastal and mountainous areas, and attributed this to the small scale
80 processes, e.g. sea and mountain breezes, that are under-resolved by the ensemble’s coarse 50 km
81 spatial resolution.

82 The present study has two goals. First, to describe a method for comparing the diurnal cycles
83 of human edited wind forecasts to those of unedited model guidance forecasts, in order to assess
84 where and when human edits produce a reduction in error or bias. Second, to apply this methodol-
85 ogy across Australian coastal locations to assess both boundary layer mixing and land-sea breeze
86 forecaster edits. The remainder of this paper is organised as follows. Section 2 describes the
87 methodology, and datasets to which it is applied, section 3 provides results, and sections 4 and 5
88 provide a discussion and a conclusion, respectively.

89 **2. Data and Methods**

90 This study compares both human edited and unedited Australian Bureau of Meteorology (BoM)
91 wind forecasts with automatic weather station (AWS) data across Australia. The comparison is
92 performed by first isolating the diurnal perturbations of each dataset by subtracting 24-hour run-
93 ning means, then comparing these perturbations on an hour-by-hour basis.

94 *a. Data*

95 Four datasets are considered in this study; the human edited official BoM wind forecast data that
96 is issued to the public, observational data from automatic weather stations (AWS) across Australia,
97 unedited data from the ECMWF’s high resolution 10-day forecast model (HRES), and unedited
98 model data from the operational Australian Community Climate and Earth System Simulator (AC-

99 CESS), noting that HRES and ACCESS are two of the model guidance products most commonly
100 used by Australian forecasters for winds. We consider just the lead-day one forecasts of the official
101 forecast, HRES and ACCESS, for reasons discussed below.

102 This study primarily considers the austral winter months of June, July and August 2018. This
103 short time period was chosen to reduce the effect of changing seasonal and climatic conditions,
104 changing forecasting practice and staff, and of changes to the ACCESS and HRES models. Re-
105 sults for December, January and February 2017/18 are occasionally mentioned to strengthen con-
106 clusions or provide a seasonal contrast.

107 ACCESS is a nested model: in this study we consider the component covering the Australian
108 region from 65.0° south to 16.95° north, and 65.0° east to 184.57° east. This model runs at a 0.11°
109 (≈ 12 km) horizontal grid spacing, with a standard time-step of 5 minutes: occasionally a shorter
110 time step of 2.5 minutes is used to overcome numerical instabilities (Bureau of Meteorology 2016).
111 HRES runs at an ≈ 9 km horizontal grid spacing, with a 7.5 minute time-step (Modigliani and
112 Maass 2017).

113 Both ACCESS and HRES use parametrisation schemes to simulate sub-grid scale boundary
114 layer turbulence, and the resultant mixing. ACCESS uses the schemes of Lock et al. (2000) and
115 Louis (1979) for unstable and stable boundary layers respectively (Bureau of Meteorology 2010).
116 HRES uses similar schemes that the ECMWF develop in-house (European Center for Medium
117 Range Weather Forecasting 2018).

118 The Bureau's official forecast dataset is produced on a state by state basis at forecasting centres
119 located in most state capitals. To construct the official forecast dataset, forecasters make a choice
120 of model guidance in the GFE, which then interpolates or upscales the model data onto a standard
121 3 km spatial grid for Victoria and Tasmania, or a 6 km grid for the rest of the country. GFE displays
122 model data at hourly intervals by taking the model guidance output at each hour UTC, with the

exception of the HRES model data which is only provided to the BoM at 3 hourly intervals, and is therefore linearly interpolated to hourly intervals by the GFE. Forecasters then make edits to these 3 or 6 km hourly grids to produce the official forecast datasets.

We therefore compare the official forecast and model guidance datasets as they appear in the GFE, i.e. we compare the upscaled or interpolated datasets on the standardised 3 or 6 km, hourly grids. This both ensures a consistent comparison between model guidance products of different spatial resolutions, and an assessment of how the official forecast compares to the model guidance products as they actually appear to forecasters in the GFE. This is the standard approach the BoM takes when verifying any forecast variable.

These datasets are compared with observations from Australian automatic weather stations (AWS), which typically record wind speed and direction each minute. After basic quality control, 10 minute averages of speed and direction are taken at each station at each hour UTC, usually over the ten minutes leading up to each hour. To calculate verification results, each station is matched with the nearest 3 or 6 km grid-point in the datasets described above.

b. Assessing Diurnal Cycles

Forecasters edit model guidance wind data to account for under-resolved sea-breeze and boundary layer mixing processes. Instead of attempting to assess each type of edit individually, we study the overall diurnal signal by subtracting a twenty four hour centred running mean *background wind* from each zonal and meridional hourly wind data point, to create wind *perturbation* datasets.

To compare errors in the official forecast, ACCESS and HRES diurnal cycles we calculate the Euclidean distances between the official or model guidance perturbation vectors at each hour UTC, and the corresponding AWS perturbation vectors at each hour UTC, viewing the Euclidean distance as a measure of absolute error. For example, to assess whether the official forecast perturba-

tions, \mathbf{u}_O , or ACCESS perturbations, \mathbf{u}_A , produce lower absolute errors when compared with the observed AWS perturbations, \mathbf{u}_{AWS} , we calculate the *difference of absolute errors* (DAE),

$$\text{DAE}_{OA} = |\mathbf{u}_{AWS} - \mathbf{u}_A| - |\mathbf{u}_{AWS} - \mathbf{u}_O|. \quad (1)$$

The analogously defined quantities DAE_{OH} and DAE_{HA} provide a comparison of the official forecast and HRES perturbations, and of the HRES and ACCESS perturbations, respectively. We can then take means of the DAE on an hourly basis; i.e. average all the 00:00 UTC DAE values, all the 01:00 UTC values, and so forth, and denote such an average by $\overline{\text{DAE}}$.

Note that $\overline{\text{DAE}}$ compares just *one aspect* of the official forecast with model guidance: it does not, for instance, assess whether the variability of the official forecast is more realistic than that of model guidance. Thus, any statements about performance made throughout this paper refer solely to $\overline{\text{DAE}}$, or subsequently defined metrics, and no claim is being made that these are sufficient to completely characterise the accuracy, or value to the user, of how the diurnal wind cycle is represented in competing forecasts.

Sea-breeze and boundary layer mixing processes depend on the background atmospheric conditions in which they occur. By comparing wind perturbations rather than the overall wind fields we are not claiming these background conditions are irrelevant. However, when a forecaster makes an edit of a wind forecast to better resolve these processes, they are implicitly assuming that future background conditions will be close enough to the preceding 24 hour mean state, or to model predictions of the mean state, to justify making the edit. Thus, it makes sense to compare forecast perturbations to observed perturbations, as long as differences are interpreted as a consequence not only of how the forecaster or model resolves the diurnal cycle, but of how differences in the background state contribute to differences in the perturbations. To minimise the importance of background state differences, this study focuses exclusively on lead-day one forecasts.

Given the large degree of turbulence or random variability in both the AWS, official, and model datasets, care must be taken to ensure we do not pre-emptively conclude the official forecast has outperformed model guidance when $\overline{\text{DAE}} > 0$ purely by chance. The method for estimating confidence in $\overline{\text{DAE}}$ is based on a method proposed by Griffiths et al. (2017). Time series formed from the DAE values at a particular time, say 00:00 UTC, across the three month time period, are treated as an independent sample of a random variable E . The sampling distribution for each $\overline{\text{DAE}}$ can be modelled by a Student's t -distribution, and from this we calculate the probability that E is positive, denoted $\Pr(E > 0)$.

Although temporal autocorrelations of DAE, i.e. correlations between DAE values at a particular hour from one day to the next, are in practice small or non-existent, they are still accounted for by reducing the “effective” sample size to $n(1 - \rho_1) / (1 + \rho_1)$, where n is the actual sample size and ρ_1 is the lag-1 autocorrelation (Zwiers and von Storch 1995; Wilks 2011). In the language of statistical hypothesis testing, the null hypothesis that $E = 0$ would be rejected at significance level α if $\Pr(E > 0) > 1 - \frac{\alpha}{2}$ or $\Pr(E < 0) > 1 - \frac{\alpha}{2}$. However, in this study we prefer to simply state the value of $\Pr(E > 0)$, referring to this as a *confidence score*, and noting $\Pr(E < 0) = 1 - \Pr(E > 0)$. We say the official forecast outperforms model guidance with “high confidence” if $\Pr(E > 0) \geq 95\%$, or that model guidance outperforms the official forecast with “high confidence” if $\Pr(E > 0) \leq 5\%$, with high confidence implicit whenever it is not explicitly mentioned.

Following the “fuzzy verification” approach outlined by Ebert (2008), forecast and observational perturbation datasets are compared not only at individual stations, but are also averaged over two coarser spatial scales before being compared. The individual stations we consider are the 8 capital city *airport stations*, marked by stars in Fig. 1, as their high operational significance means that they are typically the most well maintained. An intermediate spatial scale is formed by averaging data over the 10 stations closest to each capital city airport station, with some flexibility allowed to

192 ensure stations are roughly parallel to the nearest coastline. These station groups are referred to as
 193 the *city station groups*. The coarsest spatial scale is formed by averaging over all stations within
 194 150 km of the nearest coastline, and grouping these by state. The Western Australian coastline is
 195 subdivided into three pieces, and stations along the Gulf of Carpentaria, north Queensland Penin-
 196 sula, and Tasmanian coastlines are neglected, in order to ensure each station group corresponds
 197 to an approximately linear segment of coastline to better resolve the land-sea breeze after spatial
 198 averaging (e.g. Vincent and Lane 2016). These eight station groups are referred to as the *coastal*
 199 *station groups*.

200 To compare errors in the perturbations over the two coarser spatial scales, we modify the defini-
 201 tion of DAE in equation (1) so that each perturbation dataset is first spatially averaged over either
 202 the city or coastal station groups. Confidence scores are calculated for the city and coastal station
 203 groups in the same way as for the individual airport stations, treating the spatially averaged data
 204 as a single time series. This provides a conservative way to deal with spatial correlation between
 205 the stations in each group (Griffiths et al. 2017).

206 To compare biases in the diurnal cycles of each dataset, we calculate the *difference of biases*
 207 (DB),

$$DB_{OA} = |\overline{u}_{AWS} - \overline{u}_O| - |\overline{u}_{AWS} - \overline{u}_A|, \quad (2)$$

208 with DB_{OH} and DB_{HA} defined analogously, where the over-bars denote temporal averages of the
 209 perturbations at a particular hour, over June, July and August 2018. These temporally averaged
 210 perturbations can be viewed as the climatological diurnal wind cycles over the three month study
 211 period for each dataset. Biases over the city and coastal station groups are calculated by taking the
 212 spatial average before the temporal average. Uncertainty in the DB is estimated through bootstrap-
 213 ping (Efron 1979). This is done by performing resampling with replacement on the underlying
 214 perturbation datasets, and calculating the DB multiple times using these resampled datasets. This

215 provides a distribution of DB values, which analogously to with DAE, we treat as a sample from
 216 a random variable B , and use this to estimate $\Pr(B > 0)$.

217 Another approach to forecast verification is to assess structural features of the phenomena being
 218 forecast rather than errors or biases of point predictions; this approach is particularly important
 219 at small spatiotemporal scales (e.g. Mass et al. 2002; Rife and Davis 2005). Gille et al. (2005)
 220 obtained summary statistics on the observed structure of mean diurnal wind cycles by using linear
 221 regression to calculate the coefficients u_i, v_i $i = 0, 1, 2$, for the fits

$$u = u_0 + u_1 \cos(\omega t) + u_2 \sin(\omega t), \quad (3)$$

$$v = v_0 + v_1 \sin(\omega t) + v_2 \cos(\omega t), \quad (4)$$

222 where ω is the angular frequency of the earth and t is the local solar time in seconds. These fits
 223 trace out ellipses in the x, y plane, and descriptive metrics like the eccentricity of the ellipse and
 224 the angle the semi-major axis makes with lines of latitude, can be calculated directly from the
 225 coefficients u_1, u_2, v_1 and v_2 . Gille et al. (2005) applied this fit to scatterometer data, which after
 226 temporal averaging resulted in just four zonal and meridional values per location, and as such the
 227 fit performed very well.

228 However, equations (3) and (4) do not provide a good fit for the hourly data considered here,
 229 primarily because they assume a twelve hour symmetry in the evolution of the diurnal cycle.
 230 In practice, asymmetries between daytime heating and nighttime cooling (e.g. Svensson et al.
 231 2011) result in surface wind perturbations accelerating rapidly just after sunrise, but remaining
 232 comparatively stagnant at night (e.g. Fig. 9). Thus, we instead fit the equations

$$u = u_0 + u_1 \cos(\alpha(\psi, t)) + u_2 \sin(\alpha(\psi, t)), \quad (5)$$

$$v = v_0 + v_1 \sin(\alpha(\psi, t)) + v_2 \cos(\alpha(\psi, t)), \quad (6)$$

to the climatological perturbations, with α the function from $[0, 24) \times [0, 2\pi) \rightarrow [0, 2\pi)$ given by

$$\alpha(\psi, t) \equiv \pi \left[\sin \left(\pi \frac{(t - \psi) \bmod 24}{24} - \frac{\pi}{2} \right) + 1 \right], \quad (7)$$

with t the time in units of hours UTC, and ψ providing the time when the wind perturbations vary least with time, noting that the same value of ψ is used for both the zonal and meridional perturbations. For each climatological diurnal wind cycle, we solve for the seven parameters u_0 , u_1 , u_2 , v_0 , v_1 , v_2 and ψ using non-linear regression.

3. Results

In this section, the methods described in section 2 are applied to Australian forecast and station data over the months of June, July and August 2018. First, differences in absolute errors (DAE) and differences in biases (DB) over this time period are assessed. Second, structural indices are compared to elucidate the physical reasons for biases.

a. Absolute Errors

Figure 2 provides the mean difference of absolute error $\overline{\text{DAE}}$ values and confidence scores defined in section 2 for the coastal station groups shown in Fig. 1, for $\overline{\text{DAE}}_{\text{OA}}$, $\overline{\text{DAE}}_{\text{OH}}$ and $\overline{\text{DAE}}_{\text{HA}}$, which represent the official forecast versus ACCESS, official forecast versus HRES, and HRES versus ACCESS comparisons, respectively. The results indicate that for the majority of station groups and hours, both the unedited ACCESS and HRES models outperform the official forecast. The lowest $\overline{\text{DAE}}$ values occur at the NT station group at 23:00 and 00:00 UTC for both $\overline{\text{DAE}}_{\text{OA}}$ and $\overline{\text{DAE}}_{\text{OH}}$. Although the official forecast outperforms at least one of ACCESS or HRES at multiple times and station groups, the only group and time where it outperforms both is 05:00 UTC over the South WA station group. HRES generally outperforms ACCESS from 10:00 - 14:00 UTC, with the South WA station group being the main exception.

254 Figures 3 and 4 provide case studies of the NT and South WA station groups, respectively. Figure
255 3 a) provides a time series of DAE for the NT station group at 23:00 UTC. The time series shows
256 significant temporal variability, with DAE frequently dropping below -2 kn. Figures 3 b) and c)
257 show hodographs of the winds and wind perturbations, respectively, at each hour UTC on the 3rd
258 of July, which provides an interesting example.

259 Figure 3 b) shows that the official wind forecast on this day was likely based on edited ACCESS
260 from 00:00 to 06:00 UTC, then edited HRES from 07:00 to 13:00 UTC, then unedited ACCESS
261 from 15:00 to 21:00 UTC. At 22:00 and 23:00 UTC, the official forecast winds acquire stronger
262 east-northeasterly components than the other datasets. Figure 5 a) shows the first ten values from
263 wind soundings at Darwin Airport at 12:00 UTC on July 3rd and 00:00 UTC on July 4th. In both
264 instances the winds are east-southeasterly, and so the rapidly changing wind perturbations at 22:00
265 UTC in the official forecast likely reflect a boundary layer mixing edit that has been applied either
266 too early, or has strengthened the southeasterly component of the winds too much. Similar issues
267 create low DAE values on the 8th of June and 9th and 10th of July.

268 Figure 4 a) provides a time series of DAE for the South WA station group at 05:00 UTC. As
269 with the NT station group there is significant temporal variability, with DAE frequently exceeding
270 1 kn. Figures 4 b) and c) provide hodographs of the winds and wind perturbations, respectively, on
271 the 9th of June, another interesting example. The perturbation hodograph shows both HRES and
272 ACCESS under-predicting the amplitude of the diurnal wind cycle on this day. Figure 5 shows
273 wind soundings at Perth Airport, the nearest station to provide wind soundings, between 12:00
274 UTC on the 8th June and 12:00 UTC on the 9th June. The 8th June 12:00 UTC sounding shows
275 surface northerlies of around 6 kn, becoming west to northwesterlies of over 20 kn 2.4 km above
276 the surface. However, the subsequent sounding at 00:00 UTC on the 9th of June shows that the

277 winds acquire a strong northerly component of 30 kn in the first 500 m of the atmosphere, with
278 the final sounding indicating a strong northwesterly wind at 725 m persisting until 12:00 UTC.

279 In Fig. 4 c), the official forecast perturbations from 04:00 to 07:00 UTC show stronger westerly
280 perturbations than either ACCESS or HRES, improving the amplitude of the official forecast's
281 diurnal wind cycle. However, the AWS perturbations are more northerly than those of the official
282 forecast, and so the official forecast winds have been strengthened in a slightly incorrect direction.
283 One explanation for this discrepancy is that the official forecast has been edited based on the June
284 8th 12:00 UTC sounding, with the winds above the surface changing direction in the subsequent 12
285 hours. A similar explanation can be given for the high DAE scores on the 3rd of August, although
286 in this case the official forecast slightly improves both the magnitude and direction of the 05:00
287 UTC wind perturbations.

288 Fig. 6 presents the $\overline{\text{DAE}}$ values and confidence scores for the airport stations, and city station
289 groups, for the official forecast versus HRES comparison, i.e. $\overline{\text{DAE}}_{\text{OH}}$. The results for the airport
290 stations are noisier than the results for the coastal station groups in Figs. 2 c) and d), although they
291 share some similarities. For instance, the official forecast outperforms HRES at 01:00 and 02:00
292 UTC at both the Darwin airport station and the NT coastal station group. There are four other
293 instances where the official forecast outperforms HRES with at least 90% confidence, although
294 this could simply be occurring by chance due repeated testing (Wilks 2011, p. 178).

295 For the city station groups, HRES outperforms the official forecast almost uniformly. The main
296 exception is the Darwin city station group, where the official forecast outperforms HRES at 02:00
297 UTC, and there is ambiguity as to whether the official forecast or HRES performs better at 01:00,
298 03:00 and 04:00 UTC, and from 15:00 to 22:00 UTC. The analogous $\overline{\text{DAE}}_{\text{OA}}$ official forecast
299 versus ACCESS comparisons (not shown) are similar, with the airport station results noisy, but
300 ACCESS outperforming the official forecast over the city station groups for the vast majority

301 of times and locations. Over the December, January, February 2017/18 season, HRES also out-
302 performs the official forecast almost uniformly over the city station groups, although the official
303 forecast versus ACCESS comparisons are more ambiguous.

304 Figure 7 provides the \overline{DAE} values and confidence scores for the airport stations, and city station
305 groups, for the HRES versus ACCESS comparison. As with Fig. 6, the results for the airport
306 stations are noisy, but more often than not show that HRES outperforms ACCESS. The results for
307 the city station groups show HRES usually outperforms ACCESS, the main exceptions being the
308 Darwin and Canberra city station groups. Results for the December, January, February 2017/18
309 season are again similar, but here HRES outperforms ACCESS over the city station groups almost
310 uniformly.

311 *b. Seasonal Biases*

312 Figure 8 provides the difference of biases (DB) and confidence scores defined in section 2, for the
313 coastal station groups for DB_{OA} , DB_{OH} and DB_{HA} , which represent the the official forecast versus
314 ACCESS, official forecast versus HRES, and HRES versus ACCESS comparisons, respectively.
315 At the NT station at 03:00 UTC, the official forecast outperforms both ACCESS and HRES with
316 confidence $\geq 93\%$. However, both ACCESS and HRES outperform the official forecast at 23:00
317 and 00:00 UTC, and from 05:00 to 11:00 UTC, consistent with the \overline{DAE} results of Fig. 2. Figure
318 9 a) shows that these biases are mostly a consequence of amplitude biases in the official forecast's
319 diurnal cycle.

320 At the South WA station group from 01:00 to 05:00 UTC, the official forecast outperforms
321 HRES with confidence scores of at least 88%. Figure 9 b) shows that HRES underestimates the
322 westerly perturbations at these times, with these perturbations likely associated with boundary
323 layer mixing processes, as discussed in section 3 a. Each of the official forecast, ACCESS and

HRES underestimate the amplitude of the diurnal cycle between 02:00 and 10:00 UTC, including both the westerly perturbations and the southerly sea-breeze perturbations.

At the NSW station group from 17:00 to 19:00 UTC, the official forecast outperforms both ACCESS and HRES with confidence scores of at least 95% and 75%, respectively. Figure 9 c) shows that these times correspond to “dimples” in the perturbation temporal hodographs that are present in all four datasets. The official forecast hodograph closely resembles that of ACCESS, except for this dimple, which has been exaggerated relative to ACCESS. Figure 9 c) also shows that although HRES exaggerates the amplitude of the easterly sea-breeze perturbations, it captures the narrower shape of the AWS hodograph better than the official forecast or ACCESS.

At the SA station group from 02:00 to 05:00 UTC and 09:00 to 12:00 UTC, the official forecast outperforms both ACCESS and HRES, although confidence scores do not exceed 88% and 65% respectively. Figure 9 d) shows that although the official forecast captures the amplitude of the perturbations from 01:00 to 05:00 UTC almost perfectly, its diurnal cycle is out of phase with that of the AWS during this period, explaining why the official forecast only slightly outperforms ACCESS in the results of Figures 8 a) and b).

For comparison, Fig. 10 presents the DB values and confidence scores for DB_{OH}, which represents the official forecast versus HRES comparison, for the airport stations and city station groups. Some regions exhibit consistent results across all three spatial scales, for example, the official forecast is less biased than HRES with at least 80% confidence at Sydney airport, the Sydney city station group, and the NSW coastal station group, from 14:00 to 18:00 UTC.

c. Ellipse Fits

The hodographs in Fig. 9 are roughly elliptical in shape, suggesting that descriptive quantities can be estimated by fitting equations (5) and (6) to the zonal and meridional climatological per-

347 turbations, as described in section 2. Figure 11 gives the R^2 values for the fits of the zonal and
348 meridional perturbations to equations (5) and (6), respectively. The fit performs best at the coastal
349 station group spatial scale, with R^2 generally above 95%.

350 Figure 12 provides four descriptive quantities based on the fits of equations (5) and (6) to the
351 averaged perturbations: these are maximum perturbation speed, eccentricity of the fitted ellipse,
352 angle the semi-major axis makes with lines of latitude, and the time at which the maximum pertur-
353 bation speed is achieved. Fig. 12 a) shows that at Brisbane airport the maximum AWS perturbation
354 is at least 1 kn greater than the official forecast, ACCESS and HRES, and Fig. 12 c) shows that the
355 orientation of the AWS fitted ellipse is at least 20 degrees anti-clockwise from the other datasets.
356 Figures 13 a) and b) show hodographs of the Brisbane airport climatological perturbations and
357 ellipse fits, respectively. Although the ellipse fits suppress some of the asymmetric details, they
358 capture the amplitudes and orientations of the real climatological diurnal cycles well. In this case
359 the results show that the average AWS sea-breeze approaches from the northeast, whereas the
360 official forecast, HRES and ACCESS sea-breezes approach more from the east-northeast.

361 To check whether this just represents a direction bias of the Brisbane Airport weather station,
362 Fig. 13 c) shows the climatological perturbations at the nearby Spitfire Channel station (see Fig. 1).
363 While the amplitude bias is smaller at Spitfire Channel than Brisbane Airport, the directional bias
364 is at least as high. A similar directional bias is evident at the nearby Inner Beacon station (not
365 shown), although the bias is smaller than at Spitfire Channel and Brisbane Airport. Similar biases
366 are also evident at these stations in analogous figures for December, January and February 2017/18
367 (not shown), with the semi-major axis of the official forecast's ellipse fit oriented 29° clockwise
368 from AWS's at Brisbane airport. Figure 1 shows there are two small islands to the east of Brisbane
369 airport; the more north-northeasterly orientation of the Brisbane Airport sea-breeze suggests these

islands may be redirecting winds between the east coast of Brisbane and the west coasts of these islands, and that this local effect is not being captured in the official forecast, ACCESS or HRES.

Another example is the Hobart Airport station. Figure 12 c) shows that the semi-major axis of the AWS ellipse fit is oriented 31, 35 and 62 degrees anti-clockwise from the semi-major axes of the HRES, official forecast and ACCESS ellipse fits, respectively. Figures 11 a) and b) show that the ellipse fit for the AWS perturbations at Hobart airport only achieve R^2 values of 59% and 68% for the u and v components, respectively, but figures 13 d) and e) show that the fit still captures orientations accurately, although it underestimates the maximum AWS perturbation. Figure 13 f) provides the climatological perturbations at the Hobart (city) station, which also show a large difference in orientation between ACCESS and AWS. Given the timing of the westerly perturbations in ACCESS, and the fact that the prevailing winds around Tasmania are westerly, these results suggest that ACCESS is exaggerating the boundary layer mixing processes involved in the diurnal cycle around Hobart. These biases are not present during December, January and February 2017/18, as strong south to southeasterly sea-breeze perturbations are now dominant in all four datasets, although the semi-major axis of ACCESS's ellipse fit is still oriented 14 degrees clockwise to that of AWS.

At the South WA station group (not shown) the semi-major axes of the ACCESS and official forecast ellipse fits are oriented at least 49 degrees anti-clockwise from those of the AWS and HRES ellipse fits, and the HRES perturbations peak between 1.2 and 2.5 hours after the other datasets. These differences occur because eccentricity values are low for this station group, and Figure 9 b) shows that the westerly perturbations associated with boundary layer mixing are weaker for HRES than the other datasets. A similar issue affects the VIC station group, explaining why the semi-major axes of the AWS ellipse fit is oriented at least 49 degrees anti-clockwise from those of the other datasets.

394 The Darwin Airport, Darwin Airport station group, and NT station group (not shown) provide
 395 further examples. Here the ellipse fits produce favourable R^2 values, although the fits slightly
 396 underestimate the AWS max perturbation speed at the Darwin Airport station due to this dataset's
 397 highly asymmetric hodograph. At all three spatial scales there are timing differences between
 398 the perturbation maximums of up to 8.2 hours. These timing differences occur because for some
 399 scales and datasets, the later north to northwesterly sea-breeze perturbations dominate the diurnal
 400 wind cycle, but for other scales and datasets the earlier easterly to southeasterly boundary layer
 401 mixing effects dominate.

402 **4. Synthesis**

403 For land-sea breeze and boundary layer mixing edits to reduce absolute errors in the subsequent
 404 days wind forecast, these edits should reduce the absolute errors in the diurnal component of
 405 the wind fields. However, Figs. 2 and 6 indicate that this is only possible when absolute error is
 406 considered at coarse spatial scales, as at individual airport stations results are noisy and ambiguous,
 407 and over the intermediate city station group scale HRES outperforms the official forecast almost
 408 uniformly.

409 Taking the effective resolutions of the models considered in this study to be approximately
 410 $7\Delta x$ (e.g. Skamarock 2004; Abdalla et al. 2013), where Δx is the horizontal grid spacing, we
 411 have effective resolutions of ≈ 84 km and ≈ 63 km for ACCESS and HRES respectively. From
 412 resolution considerations alone, one might expect that forecaster edits would be able to reduce
 413 errors at the individual airport station scale, and the intermediate city station group scale (see
 414 Fig. 1), as motion at these scales is unresolved or only partially resolved by ACCESS and HRES.

415 To further investigate the effect of spatial scale on error, consider first just the zonal components
 416 of the AWS and official forecast wind perturbations, denoted by u_{AWS} and u_{O} respectively. Con-

417 sidering just the values at a particular hour UTC, over the entire June, July, August time period,
 418 the mean square error $\text{mse}(u_{\text{AWS}}, u_{\text{O}}) = \overline{(u_{\text{AWS}} - u_{\text{O}})^2}$ can be decomposed $\text{mse}(u_{\text{AWS}}, u_{\text{O}}) =$

$$\underbrace{\text{var}(u_{\text{AWS}}) + \text{var}(u_{\text{O}}) - 2\text{cov}(u_{\text{AWS}}, u_{\text{O}})}_{\text{error variance}} + \underbrace{(\bar{u}_{\text{AWS}} - \bar{u}_{\text{O}})^2}_{\text{squared bias}} \quad (8)$$

419 where var, cov and the over-bar denote the sample variance, covariance and mean respectively.
 420 The first three terms are the variance of $u_{\text{AWS}} - u_{\text{O}}$, i.e. the error variance, and the last term is the
 421 square of the bias between u_{AWS} and u_{O} . Equation (8) can also be applied to the MSEs of HRES.
 422 Note that the mean square errors (MSEs) of the official forecast and HRES are closely related
 423 to $\overline{\text{DAE}}_{\text{OH}}$, which is the difference between the mean absolute errors of the official forecast and
 424 HRES; similarly, the squared bias components of the MSEs are closely related to DB_{OH} .

425 Figure 14 shows the terms of equation (8) for both the official forecast and HRES for Adelaide
 426 Airport, the Adelaide city station group, and the SA coastal station group. At all three scales the
 427 official forecast varies more than HRES, which is also the case at the other locations considered in
 428 this study. At Adelaide airport the variance of AWS is significantly larger than either the official
 429 forecast or HRES, but this additional variability is mostly uncorrelated to either dataset. This is
 430 unsurprising from representation considerations alone (e.g. Zaron and Egbert 2006), as the official
 431 forecast and HRES data represent averages over 6 km spatial grid-cells, whereas the AWS data
 432 represent point values. As a result, error variance terms are much larger than the squared bias
 433 terms, and of comparable magnitudes for both datasets. This is consistent with the comparatively
 434 noisy DAE results of Figs. 6 a) and b).

435 At the intermediate Adelaide city station group scale, the AWS variances are of similar magni-
 436 tudes to those of HRES, but smaller than those of the official forecast, with the official forecast's
 437 additional variability mostly uncorrelated to AWS. This results in larger error variance terms for
 438 the official forecast, consistent with HRESs almost complete outperformance of the official fore-

cast in Figs. 6 c) and d). Over the coarse SA coastal station group scale, variances in all three datasets are now small enough that the error variance terms no longer dwarf the bias terms. Although the error variance of the official forecast is still larger than that of HRES, HRES's zonal biases at 05:00 UTC are now sufficient to result in a larger MSE at this time, consistent with the DAE results of Fig. 2 c) and d).

Analogous points can be made for the other locations considered in this study, the main exception being Darwin airport, Darwin city station group, and the NT coastal station group, where zonal biases in HREF around 01:00 - 03:00 UTC are large enough to overcome the official forecast's larger error variance, producing the results of Fig. 6 and Figs. 2 c) and d). The results of Fig. 6 c) and d) are therefore generally a consequence of the official forecast being more variable than HRES, with this additional variability mostly random, in the sense of being uncorrelated with AWS. Similarly, the official forecast is generally more variable than ACCESS, explaining why the official forecast also struggles to outperform ACCESS at these scales, and ACCESS is generally more variable than HRES, explaining why HRES generally outperforms ACCESS in the DAE results of Fig. 7. In the coastal station group DAE results of Fig. 2, the random variability in each dataset is reduced, and biases are now large enough to actually affect errors in the diurnal component of the forecast.

These results show that switching model guidance products or performing edits can add more random noise to the diurnal component of the official forecast than what can be offset by reductions in bias, or improved correlations with AWS. Because the official forecast is built from multiple model datasets, most commonly HRES and ACCESS, blending datasets with different means will tend to produce greater variance than any of the component datasets. If the choice of model guidance is made primarily on which model best captures more slowly evolving synoptic scale features, then switching model guidance may add random variability to the diurnal component of

the official forecast. Furthermore, unless all forecasters follow identical thought processes when making edits, the edits will also add random variability. It is less clear why ACCESS shows greater random variability than HRES: one cause may be ACCESS's shorter time-step.

These results have implications for forecasting practice. Model guidance products are indeed biased in how they resolve diurnal wind cycles (e.g. Fig. 13), and there is therefore scope for forecaster edits to reduce these biases. However, editing model guidance generally fails to reduce error in the forecast diurnal cycle, even at scales finer than the effective resolutions of the models, as the cycle itself is mostly hidden by random variability. Averaging over large areas reduces this random variability, and so biases have a greater impact on forecast error, but even at large scales Fig. 2 shows model guidance still outperforms the official forecast more often than not.

Reducing the random variability of the official forecast, or the model guidance datasets that comprise it, will therefore improve the capacity of these types of edits to reduce error. One way to do this would be to move to an ensemble forecasting system, another would be to post process model guidance products, such as by averaging multiple time steps around the hour, before including them in GFE.

5. Conclusion

In this study we have presented methods for verifying the diurnal component of wind forecasts, with the intended application being the assessment of the edits Australian forecasters make to model guidance datasets in order to better resolve land-sea breeze and boundary layer mixing processes. We considered both errors and seasonal biases at each hour UTC, over three spatial scales, but the methods are immediately generalisable to other spatiotemporal scales.

When the methods are applied to Australian forecast data, the results indicate that the official edited forecast only produces lower absolute errors in the diurnal wind cycle when averaged over

coarse spatial scales of $500 \times 150 \text{ km}^2$ to $2000 \times 150 \text{ km}^2$: this scale corresponds to the aggregation of data within 150 km of the Australian coastline, subdivided into linear segments of coastline and by state (see Fig. 1). Even at these scales, reductions in error are isolated to particular locations and times of day, and the official forecast rarely has lower mean absolute error than both commonly used model guidance products simultaneously. This suggests that forecaster skill in improving diurnal wind processes lies more in making the choice of model guidance than in making edits.

By contrast, the official forecast can produce lower seasonal biases than model guidance at all three spatial scales, but again, it rarely produces lower biases than both standard model guidance products simultaneously. Reduced seasonal biases do not translate into reduced errors at the two smaller spatial scales because the diurnal cycle is mostly masked by the random variability in each dataset. Furthermore, because the official forecast exhibits much greater random variability than HRES, HRES almost uniformly outperforms the official forecast over the intermediate $50 \times 50 \text{ km}^2$ to $200 \times 200 \text{ km}^2$ spatial scale. The same is true for ACCESS, although to a slightly lesser extent, and also explains why HRES mostly outperforms ACCESS at this scale.

We also compare structural features of the diurnal wind cycles of each dataset by fitting modified ellipses to hodographs of seasonally averaged diurnal wind cycles, then deriving metrics from these ellipses. This approach reveals structural biases in the official forecast, including directional biases in the approach of the sea-breeze at Brisbane airport, eccentricity biases along the coast of NSW, and amplitude biases along the southwest coast of WA. It also reveals biases in model guidance datasets, such as ACCESS's overemphasis of boundary layer mixing processes around Hobart.

Future research could extend this study in multiple directions. One important question is whether the random variability in the official forecast, or the model guidance products that comprise it, could be reduced through ensemble forecasting or post-processing, as reducing random variability

would both decrease errors, and increase the value of land-sea breeze and boundary layer mixing edits. Another goal could be to identify precisely the spatiotemporal scales at which diurnal wind cycles can be identified against background noise, so as to better understand the scales at which land-sea breeze and boundary layer mixing edits can add value to a forecast.

Acknowledgments. Funding for this study was provided for Ewan Short by the Australian Research Council's Centre of Excellence for Climate Extremes (CE170100023). Datasets and software were generously provided by the Australian Bureau of Meteorology's Evidence Targeted Automation team. Thanks are due to Michael Foley, Deryn Griffiths, Nicholas Loveday, Ben Price and Alexei Hider for providing support at the Bureau of Meteorology's Melbourne and Darwin offices, and to Professors Craig Bishop and Todd Lane from the University of Melbourne, and Carly Kovacic from the United States' National Weather Service, for some helpful conversations.

References

Abdalla, S., L. Isaksen, P. A. E. M. Janssen, and N. Wedi, 2013: Effective spectral resolution of ECMWF atmospheric forecast models. 19–22, doi:10.21957/rue4o7ac, [Available online at <https://www.ecmwf.int/node/17358> - Accessed 07/12/2019].

Abkar, M., A. Sharifi, and F. Port-Agel, 2016: Wake flow in a wind farm during a diurnal cycle. *Journal of Turbulence*, **17** (4), 420–441, doi:10.1080/14685248.2015.1127379.

Bureau of Meteorology, 2010: Operational implementation of the ACCESS numerical weather prediction systems. Tech. Rep. NMOC Operations Bulletin No. 83, Bureau of Meteorology, Melbourne, Victoria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/apob83.pdf> - Accessed 25/04/2019].

531 Bureau of Meteorology, 2016: APS2 upgrade to the ACCESS-R numerical weather prediction sys-
 532 tem. Tech. Rep. BNOC Operations Bulletin No. 104, Bureau of Meteorology, Melbourne, Victo-
 533 ria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/apob107-external.pdf>
 534 - Accessed 25/04/2019].

535 Bureau of Meteorology, 2019: Meteye. Bureau of Meteorology, [Available online at [http://www.](http://www.bom.gov.au/australia/meteye/)
 536 [bom.gov.au/australia/meteye/](http://www.bom.gov.au/australia/meteye/)].

537 Dai, A., and C. Deser, 1999: Diurnal and semidiurnal variations in global surface wind
 538 and divergence fields. *Journal of Geophysical Research*, **104**, 31 109–31 125, doi:10.1029/
 539 1999JD900927.

540 Ebert, E. E., 2008: Fuzzy verification of high-resolution gridded forecasts: a review and proposed
 541 framework. *Meteor. Appl.*, **15** (1), 51–64, doi:10.1002/met.25.

542 Efron, B., 1979: Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7** (1),
 543 1–26, doi:10.1214/aos/1176344552.

544 Englberger, A., and A. Dörnbrack, 2018: Impact of the diurnal cycle of the atmospheric bound-
 545 ary layer on wind-turbine wakes: a numerical modelling study. *Boundary-Layer Meteorology*,
 546 **166** (3), 423–448, doi:10.1007/s10546-017-0309-3.

547 European Center for Medium Range Weather Forecasting, 2018: *Part IV: Physical processes*,
 548 223. No. 4, IFS Documentation, European Center for Medium Range Weather Forecasting,
 549 [Available online at <https://www.ecmwf.int/node/18714> - Accessed 25 April 2019].

550 Gille, S. T., S. G. Llewellyn Smith, and N. M. Statom, 2005: Global observations of the land
 551 breeze. *Geophysical Research Letters*, **32** (5), doi:10.1029/2004GL022139.

Griffiths, D., H. Jack, M. Foley, I. Ioannou, and M. Liu, 2017: Advice for automation of forecasts:
a framework. Tech. rep., Bureau of Meteorology, Melbourne, Victoria. [Available online at
<http://www.bom.gov.au/research/publications/researchreports/BRR-021.pdf>].

Lee, X., 2018: *Fundamentals of boundary-layer meteorology*. Springer atmospheric sciences,
Springer.

Lock, A. P., A. R. Brown, M. R. Bush, G. M. Martin, and R. N. B. Smith, 2000: A new bound-
ary layer mixing scheme. Part I: scheme description and single-column model tests. *Monthly
Weather Review*, **128** (9), 3187–3199, doi:10.1175/1520-0493(2000)128<3187:ANBLMS>2.0.
CO;2.

Louis, J.-F., 1979: A parametric model of vertical eddy fluxes in the atmosphere. *Boundary-Layer
Meteorology*, **17** (2), 187–202, doi:10.1007/BF00117978.

Mass, C. F., D. Ovens, K. Westrick, and B. A. Colle, 2002: Does increasing horizontal resolution
produce more skillful forecasts? *Bulletin of the American Meteorological Society*, **83** (3), 407–
430, doi:10.1175/1520-0477(2002)083<0407:DIHRPM>2.3.CO;2.

Miller, S. T. K., B. D. Keim, R. W. Talbot, and H. Mao, 2003: Sea breeze: Structure, forecasting,
and impacts. *Reviews of Geophysics*, **41** (3), doi:10.1029/2003RG000124.

Modigliani, U., and C. Maass, 2017: Detailed information of implementation of IFS cy-
cle 41r2. ECMWF, [Available online at [https://confluence.ecmwf.int/display/FCST/Detailed+](https://confluence.ecmwf.int/display/FCST/Detailed+information+of+implementation+of+IFS+cycle+41r2)
[information+of+implementation+of+IFS+cycle+41r2](https://confluence.ecmwf.int/display/FCST/Detailed+information+of+implementation+of+IFS+cycle+41r2) - Accessed 27/04/2019].

Physick, W. L., and D. J. Abbs, 1992: Flow and plume dispersion in a coastal valley. *Journal
of Applied Meteorology*, **31** (1), 64–73, doi:10.1175/1520-0450(1992)031<0064:FAPDIA>2.0.
CO;2.

574 Pinson, P., and R. Hagedorn, 2012: Verification of the ECMWF ensemble forecasts of wind speed
 575 against analyses and observations. *Meteor. Appl.*, **19** (4), 484–500, doi:10.1002/met.283.

576 Rife, D. L., and C. A. Davis, 2005: Verification of temporal variations in mesoscale numerical
 577 wind forecasts. *Monthly Weather Review*, **133** (11), 3368–3381, doi:10.1175/MWR3052.1.

578 Skamarock, W. C., 2004: Evaluating mesoscale NWP models using kinetic energy spectra.
 579 *Monthly Weather Review*, **132** (12), 3019–3032, doi:10.1175/MWR2830.1, URL [https://doi.org/](https://doi.org/10.1175/MWR2830.1)
 580 [10.1175/MWR2830.1](https://doi.org/10.1175/MWR2830.1), <https://doi.org/10.1175/MWR2830.1>.

581 Svensson, G., and Coauthors, 2011: Evaluation of the diurnal cycle in the atmospheric bound-
 582 ary layer over land as represented by a variety of single-column models: The second GABLS
 583 experiment. *Boundary-Layer Meteorology*, **140** (2), 177–206, doi:10.1007/s10546-011-9611-7.

584 Vincent, C. L., and T. P. Lane, 2016: Evolution of the diurnal precipitation cycle with the passage
 585 of a Madden-Julian Oscillation event through the Maritime Continent. *Monthly Weather Review*,
 586 **144** (5), 1983–2005, doi:10.1175/MWR-D-15-0326.1.

587 Wilks, D. S., 2011: *Statistical methods in the atmospheric sciences*. International geophysics
 588 series: v. 100, Elsevier.

589 Zaron, E. D., and G. D. Egbert, 2006: Estimating open-ocean barotropic tidal dissipation: The
 590 hawaiian ridge. *Journal of Physical Oceanography*, **36** (6), 1019–1035, doi:10.1175/JPO2878.
 591 1.

592 Zwiers, F. W., and H. von Storch, 1995: Taking serial correlation into account in tests of the mean.
 593 *Journal of Climate*, **8** (2), 336–351, doi:10.1175/1520-0442(1995)008<0336:TSCIAI>2.0.CO;2.

LIST OF FIGURES

Fig. 1.	Locations of the automatic weather stations considered in this study, where stars give the locations of the capital city <i>airport stations</i> . Stations are divided into the Darwin, Perth, Adelaide, Melbourne, Hobart, Canberra, Sydney and Brisbane <i>city station groups</i> , a) to h), respectively, and the <i>coastal station groups</i> , i). Height and depth shading intervals every 200 and 1000 m, respectively.	30
Fig. 2.	Heatmaps of mean difference of absolute error \overline{DAE} values, a), c), e), and confidence scores, b), d), f), for each coastal station group (see Fig. 1) and hour of the day, for the official forecast versus ACCESS, a) and b), official forecast versus HRES, c) and d), HRES versus ACCESS, e) and f). Positive \overline{DAE} values indicate that the former dataset in each pair is on average \overline{DAE} kn closer to observations than the latter dataset (see equation 1), where $1 \text{ kn} \approx 0.514 \text{ m s}^{-1}$. Confidence scores provide the probability the population or “true” value of \overline{DAE} is greater than zero (see section 2).	31
Fig. 3.	Time series, a), of the difference in absolute error DAE defined in equation (1) for the official forecast versus ACCESS, DAE_{OA} , and official forecast versus HRES, DAE_{OH} , for the NT coastal station group shown in Fig. 1 at 23:00 UTC. Also, temporal hodographs in hours UTC showing hourly changes in winds, b), and wind perturbations from a 24 hour running mean, c), at the NT coastal station group on the 3 rd of July 2018.	32
Fig. 4.	As in Fig. 3, but for, a), the South WA coastal station group at 05:00 UTC, and b) and c), the winds and wind perturbations, respectively, over the South WA coastal station group on the 9 th June 2018.	33
Fig. 5.	Vertical wind soundings at, a), Darwin Airport, and b), Perth Airport, with heights given in metres.	34
Fig. 6.	As in Fig. 2, but for the official versus HRES mean difference of absolute error \overline{DAE}_{OH} values, a) and c), and confidence scores, b) and d), for the airport stations, a) and b), and city station groups, c) and d).	35
Fig. 7.	As in Fig. 6, but for the HRES versus ACCESS mean difference in absolute error \overline{DAE}_{HA} values and confidence scores.	36
Fig. 8.	As in Fig. 2, but for the difference of biases (DB) values and confidence scores.	37
Fig. 9.	Temporal hodographs in hours UTC of wind perturbations spatially averaged over the, a), NT, b) South WA, c) NSW and d), SA coastal station groups (see Fig. 1), and temporally averaged over June, July and August 2018.	38
Fig. 10.	As in Fig. 6, but for the difference of biases (DB) values and confidence scores.	39
Fig. 11.	R^2 values as percentages for the fit of equation (5) to the zonal perturbations, a), c) and e), and equation (6) to the meridional perturbations, b), d) and f), for the airport stations, a) and b), city station groups, c) and d), and coastal station groups, e) and f), shown in Fig. 1.	40
Fig. 12.	Metrics derived from fitting ellipse equations (5) and (6) to wind perturbations at the Australian capital city airport stations, a) to d), and to wind perturbations spatially averaged over the city station groups and coastal station groups shown in Fig. 1, e) to h) and i) to l) respectively, with perturbations also temporally averaged over June, July and August 2018 in each case. Metrics given are the maximum perturbation speed, a), e) and i), eccentricity	

of fitted ellipse, b), f) and j), orientation semi-major axis makes with lines of latitude, c), g) and k), and time of maximum perturbation, d), h) and l). 41

Fig. 13. Temporal hodographs of wind perturbations at each hour UTC averaged over June, July and August 2018, at Brisbane and Hobart airports, a) and d), and the associated ellipse fits, b) and e). For comparison, c) and f) provide the hodographs of the averaged perturbations at the Spitfire Channel and Hobart city stations, respectively (see Fig. 1). 42

Fig. 14. Mean square error between the AWS and HRES zonal perturbations $\overline{(u_{\text{AWS}} - u_{\text{H}})^2}$, a), e), and i), decomposed into the error variance $\text{var}(u_{\text{AWS}} - u_{\text{H}})$ and squared bias $(\bar{u}_{\text{AWS}} - \bar{u}_{\text{H}})^2$ terms of equation (8). Also, the decomposed mean square error between the AWS and official forecast zonal perturbations, b), f) and j). Additionally, the HRES and AWS error variance term $\text{var}(u_{\text{AWS}} - u_{\text{H}})$ decomposed into the $\text{var}(u_{\text{AWS}})$, $\text{var}(u_{\text{H}})$ and $-2 \cdot \text{cov}(u_{\text{AWS}}, u_{\text{H}})$ terms, c), g) and k), and analogously for the official forecast and AWS error variance term $\text{var}(u_{\text{AWS}} - u_{\text{O}})$, d), h) and l). Decompositions given for Adelaide Airport, a) to d), the Adelaide city station group, e) to h), and the SA coastal station group, i) to l) (see Fig. 1.) . . . 43

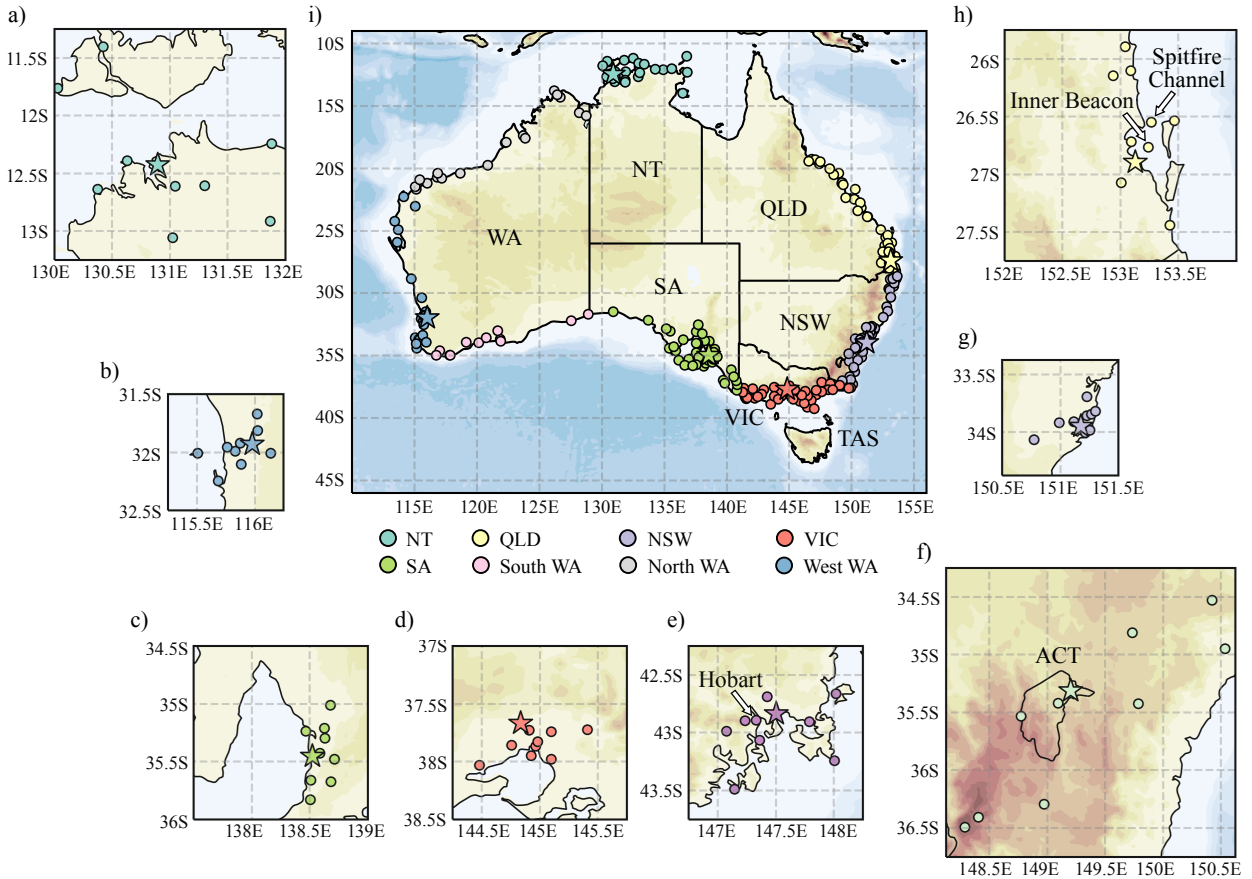


FIG. 1. Locations of the automatic weather stations considered in this study, where stars give the locations of the capital city *airport stations*. Stations are divided into the Darwin, Perth, Adelaide, Melbourne, Hobart, Canberra, Sydney and Brisbane *city station groups*, a) to h), respectively, and the *coastal station groups*, i). Height and depth shading intervals every 200 and 1000 m, respectively.

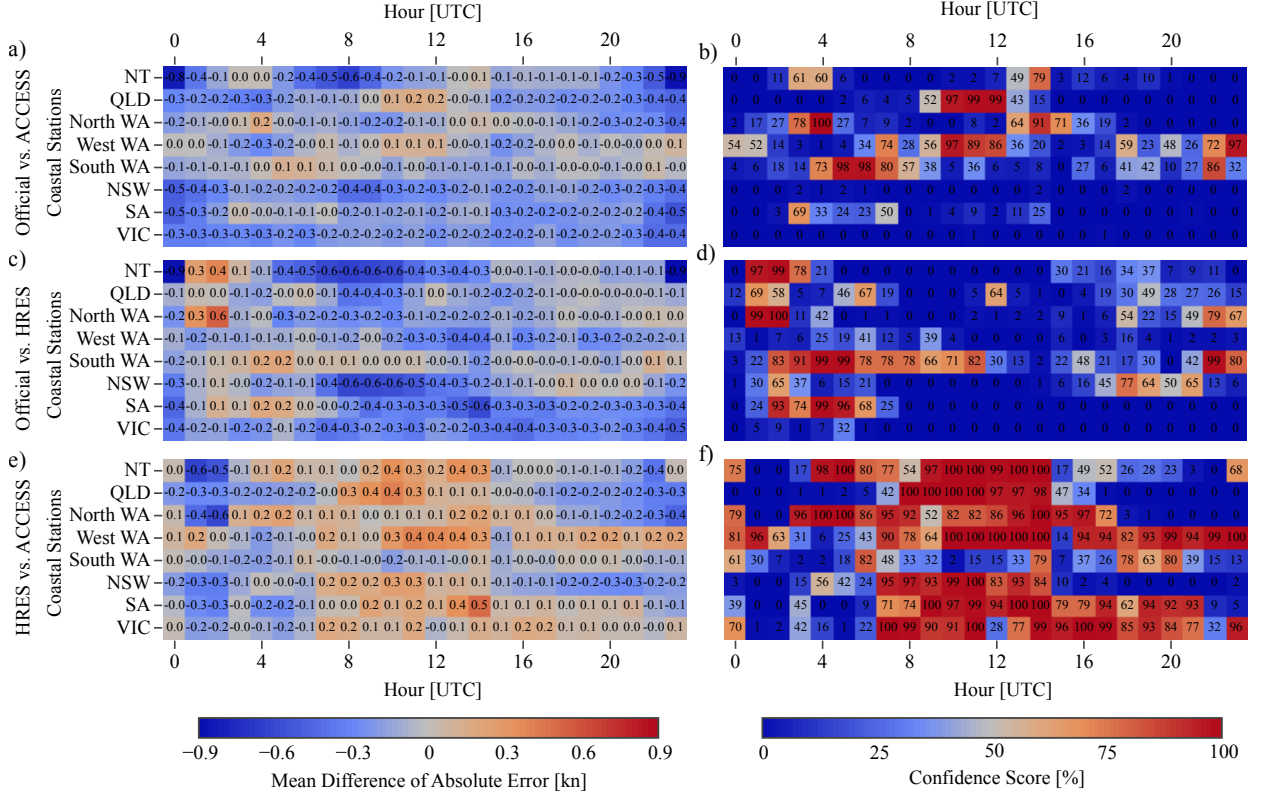


FIG. 2. Heatmaps of mean difference of absolute error $\overline{\text{DAE}}$ values, a), c), e), and confidence scores, b), d), f), for each coastal station group (see Fig. 1) and hour of the day, for the official forecast versus ACCESS, a) and b), official forecast versus HRES, c) and d), HRES versus ACCESS, e) and f). Positive $\overline{\text{DAE}}$ values indicate that the former dataset in each pair is on average $\overline{\text{DAE}}$ kn closer to observations than the latter dataset (see equation 1), where $1 \text{ kn} \approx 0.514 \text{ m s}^{-1}$. Confidence scores provide the probability the population or “true” value of $\overline{\text{DAE}}$ is greater than zero (see section 2).

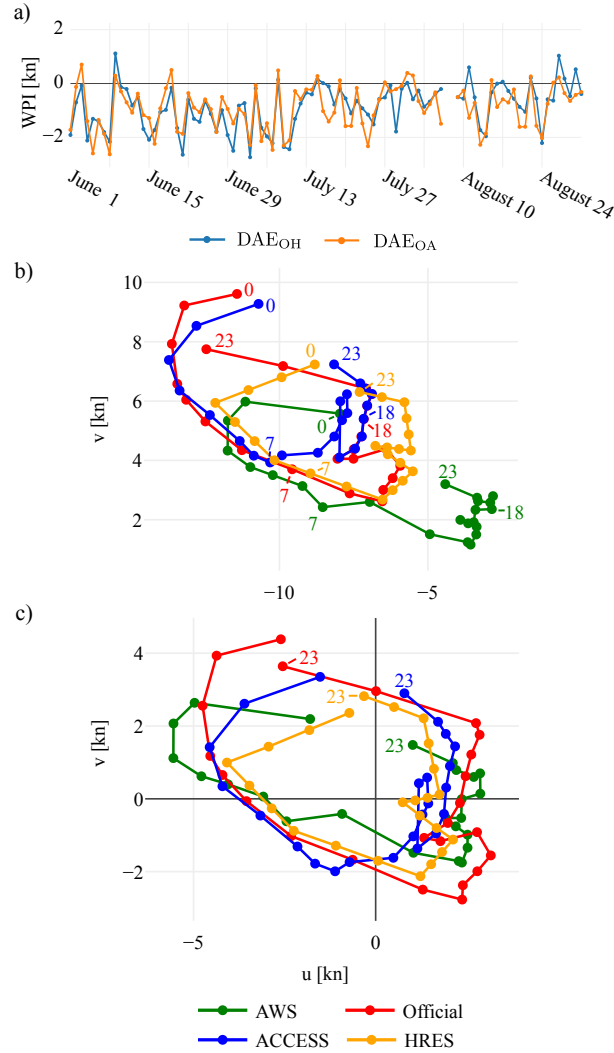


FIG. 3. Time series, a), of the difference in absolute error DAE defined in equation (1) for the official forecast versus ACCESS, DAE_{OA} , and official forecast versus HRES, DAE_{OH} , for the NT coastal station group shown in Fig. 1 at 23:00 UTC. Also, temporal hodographs in hours UTC showing hourly changes in winds, b), and wind perturbations from a 24 hour running mean, c), at the NT coastal station group on the 3rd of July 2018.

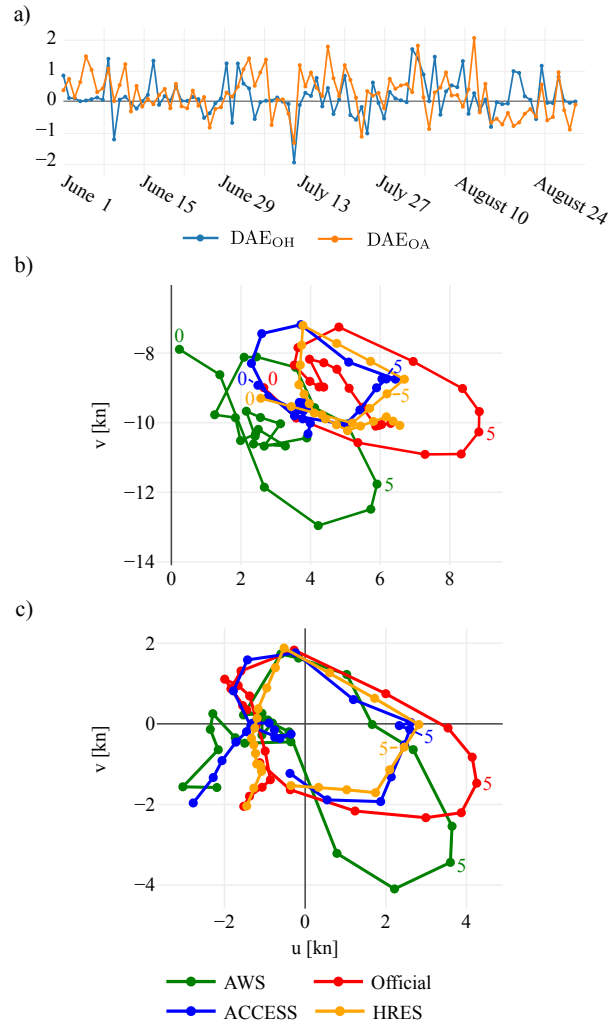


FIG. 4. As in Fig. 3, but for, a), the South WA coastal station group at 05:00 UTC, and b) and c), the winds
and wind perturbations, respectively, over the South WA coastal station group on the 9th June 2018.

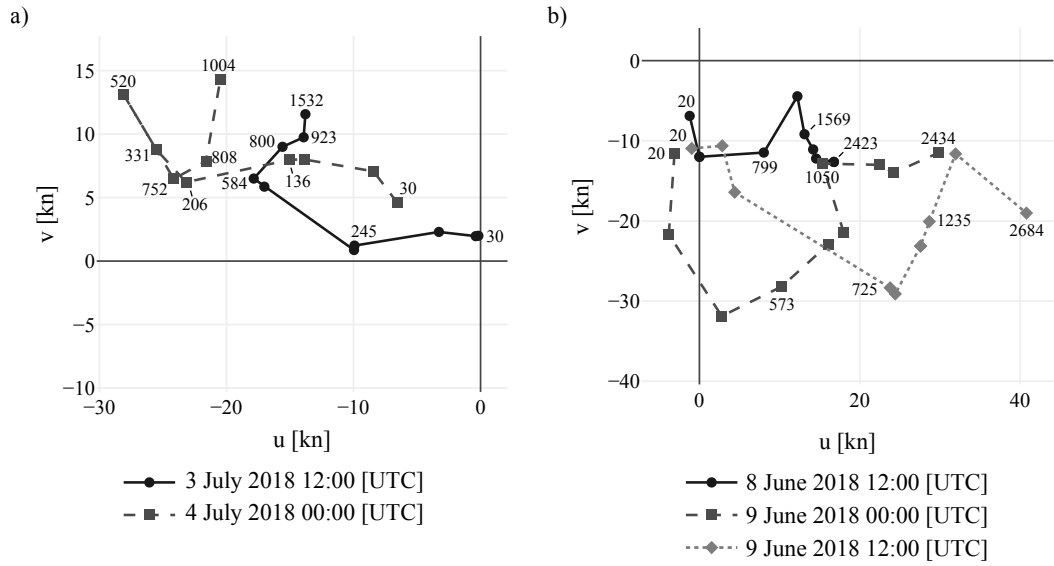


FIG. 5. Vertical wind soundings at, a), Darwin Airport, and b), Perth Airport, with heights given in metres.

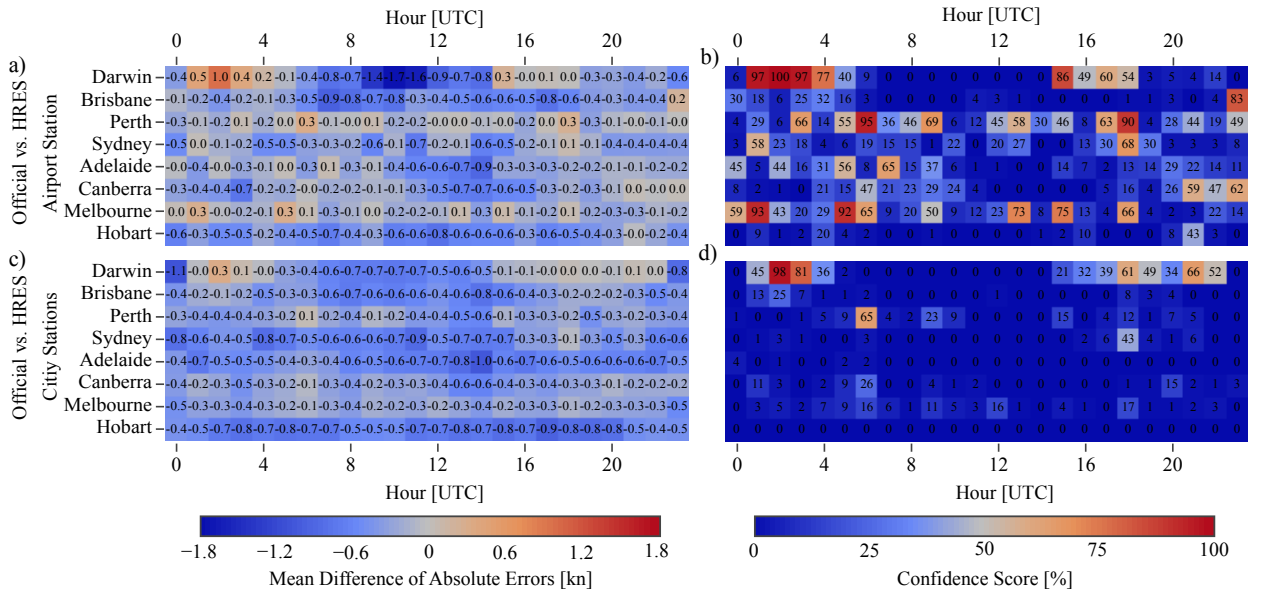


FIG. 6. As in Fig. 2, but for the official versus HRES mean difference of absolute error $\overline{\text{DAE}}_{\text{OH}}$ values, a) and

c), and confidence scores, b) and d), for the airport stations, a) and b), and city station groups, c) and d).

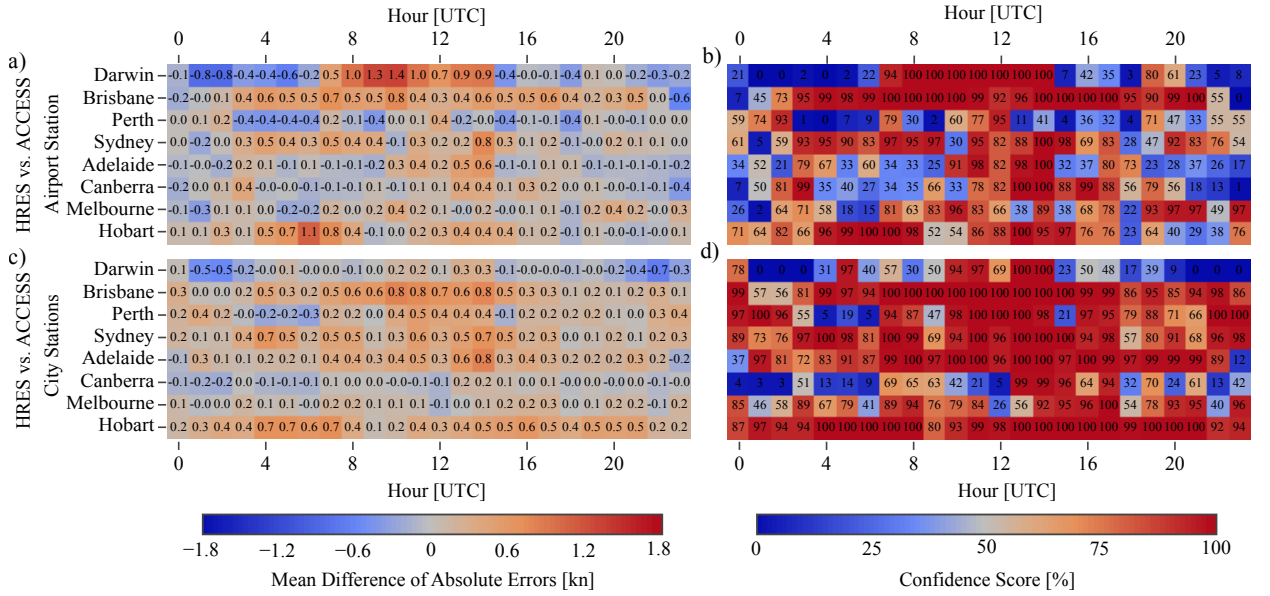


FIG. 7. As in Fig. 6, but for the HRES versus ACCESS mean difference in absolute error $\overline{\text{DAE}}_{\text{HA}}$ values and

confidence scores.

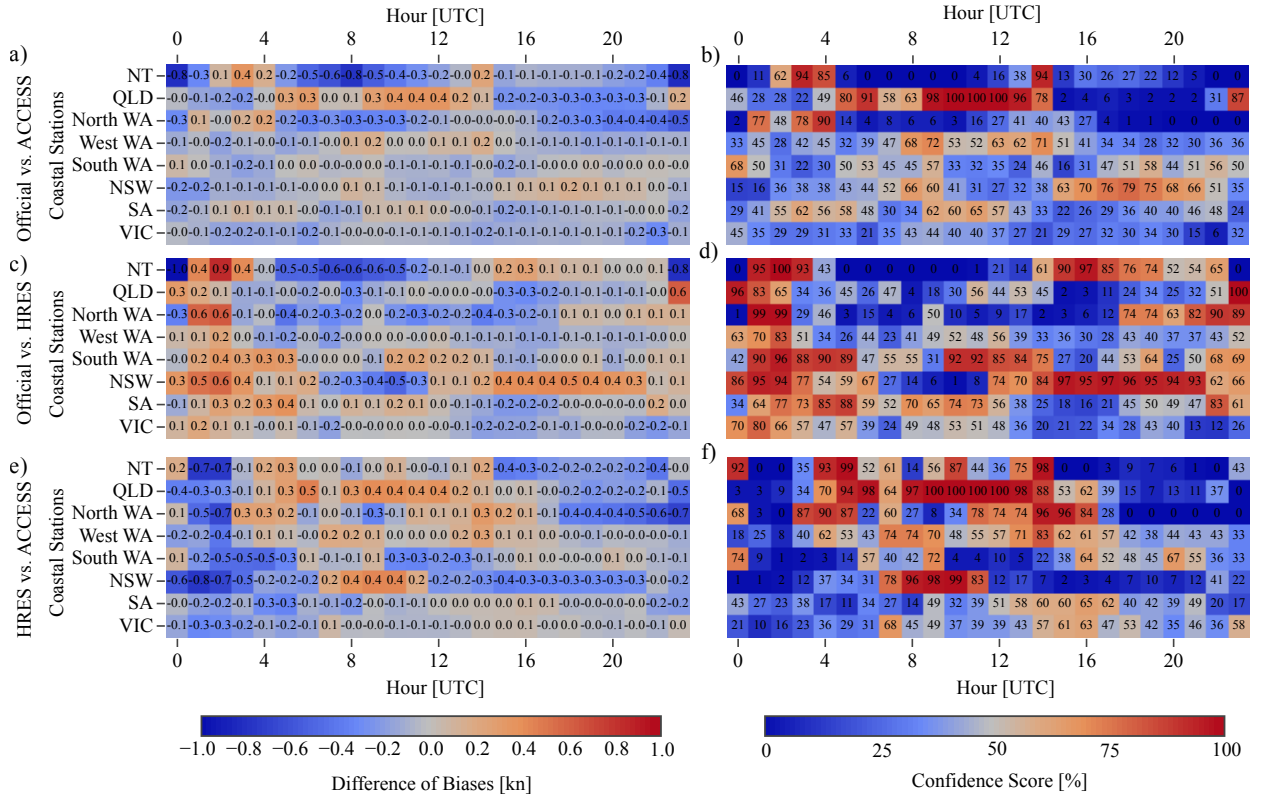


FIG. 8. As in Fig. 2, but for the difference of biases (DB) values and confidence scores.

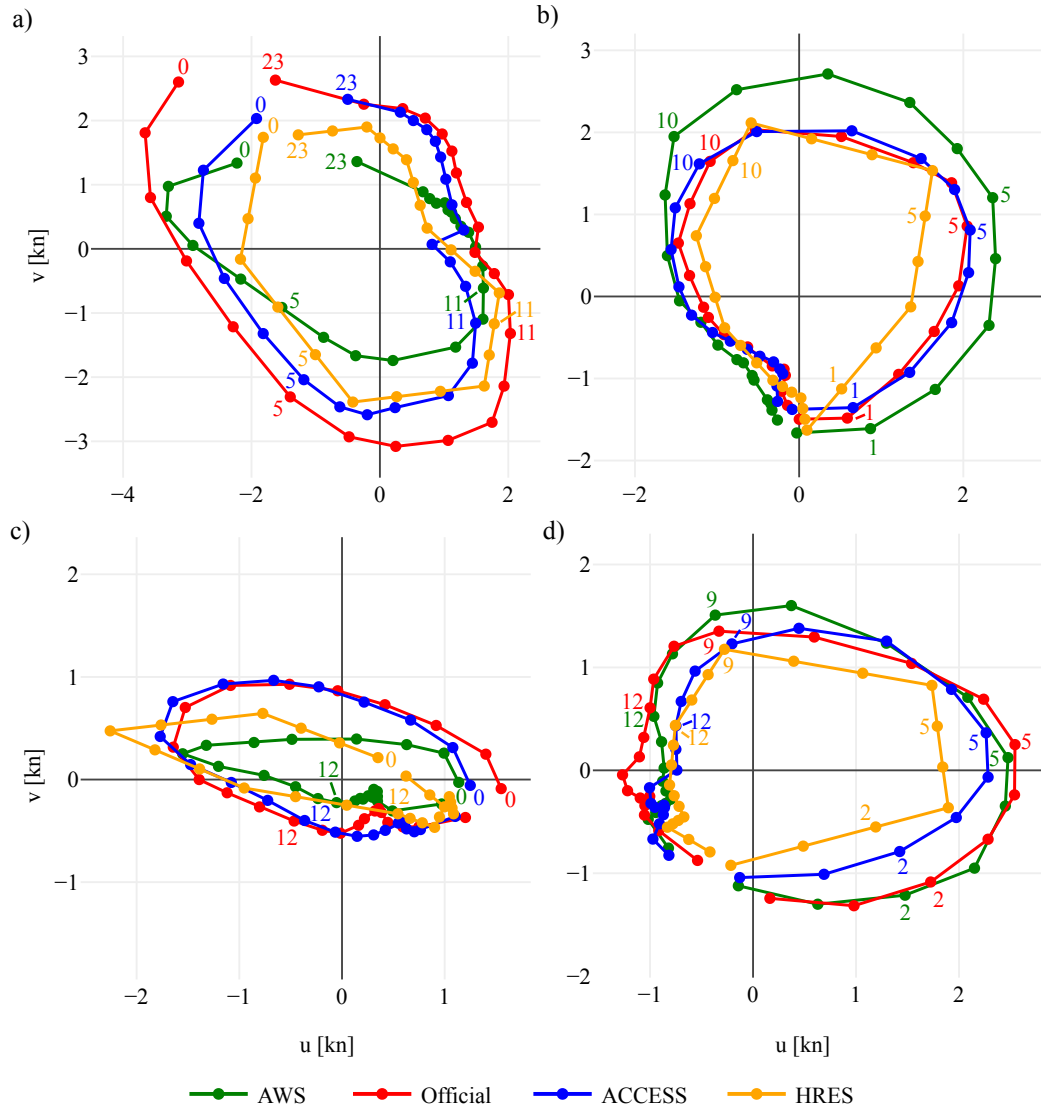


FIG. 9. Temporal hodographs in hours UTC of wind perturbations spatially averaged over the, a), NT, b) South WA, c) NSW and d), SA coastal station groups (see Fig. 1), and temporally averaged over June, July and August 2018.

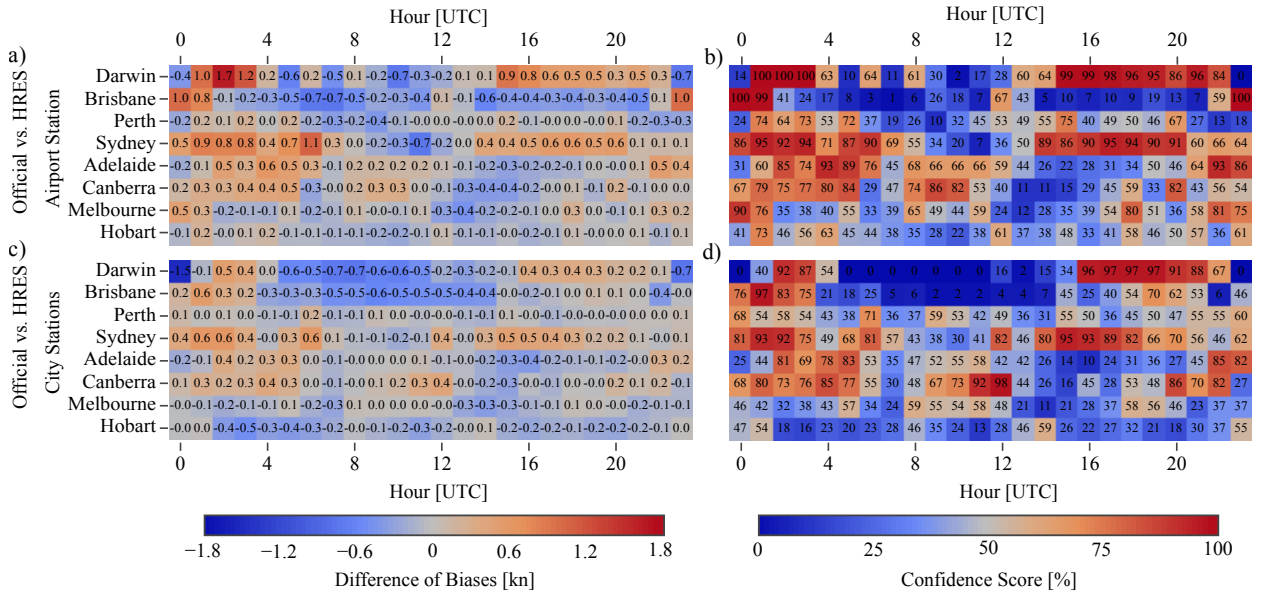


FIG. 10. As in Fig. 6, but for the difference of biases (DB) values and confidence scores.

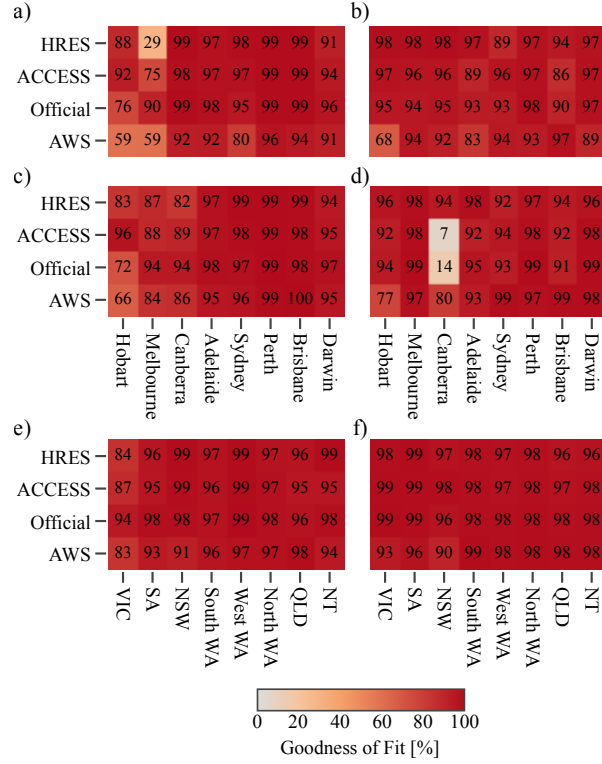


FIG. 11. R^2 values as percentages for the fit of equation (5) to the zonal perturbations, a), c) and e), and equation (6) to the meridional perturbations, b), d) and f), for the airport stations, a) and b), city station groups, c) and d), and coastal station groups, e) and f), shown in Fig. 1.

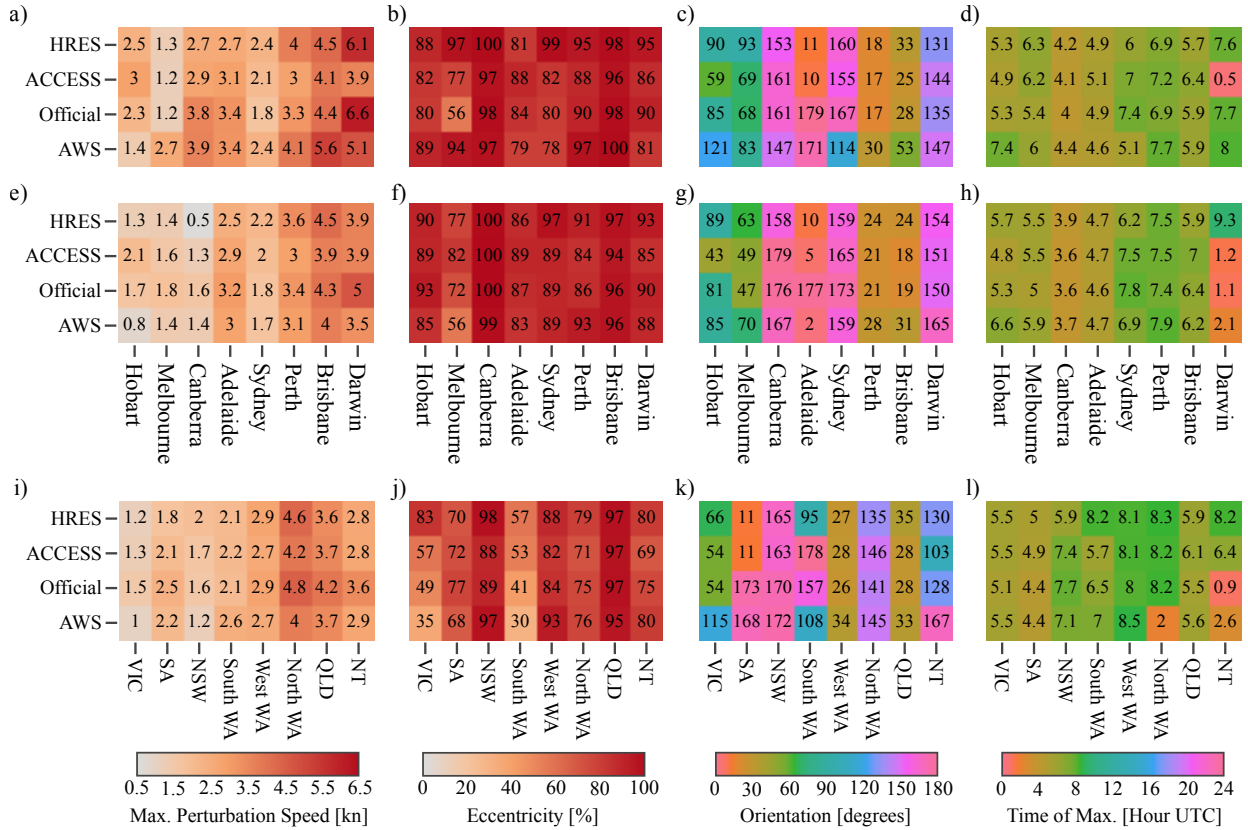


FIG. 12. Metrics derived from fitting ellipse equations (5) and (6) to wind perturbations at the Australian capital city airport stations, a) to d), and to wind perturbations spatially averaged over the city station groups and coastal station groups shown in Fig. 1, e) to h) and i) to l) respectively, with perturbations also temporally averaged over June, July and August 2018 in each case. Metrics given are the maximum perturbation speed, a), e) and i), eccentricity of fitted ellipse, b), f) and j), orientation semi-major axis makes with lines of latitude, c), g) and k), and time of maximum perturbation, d), h) and l).

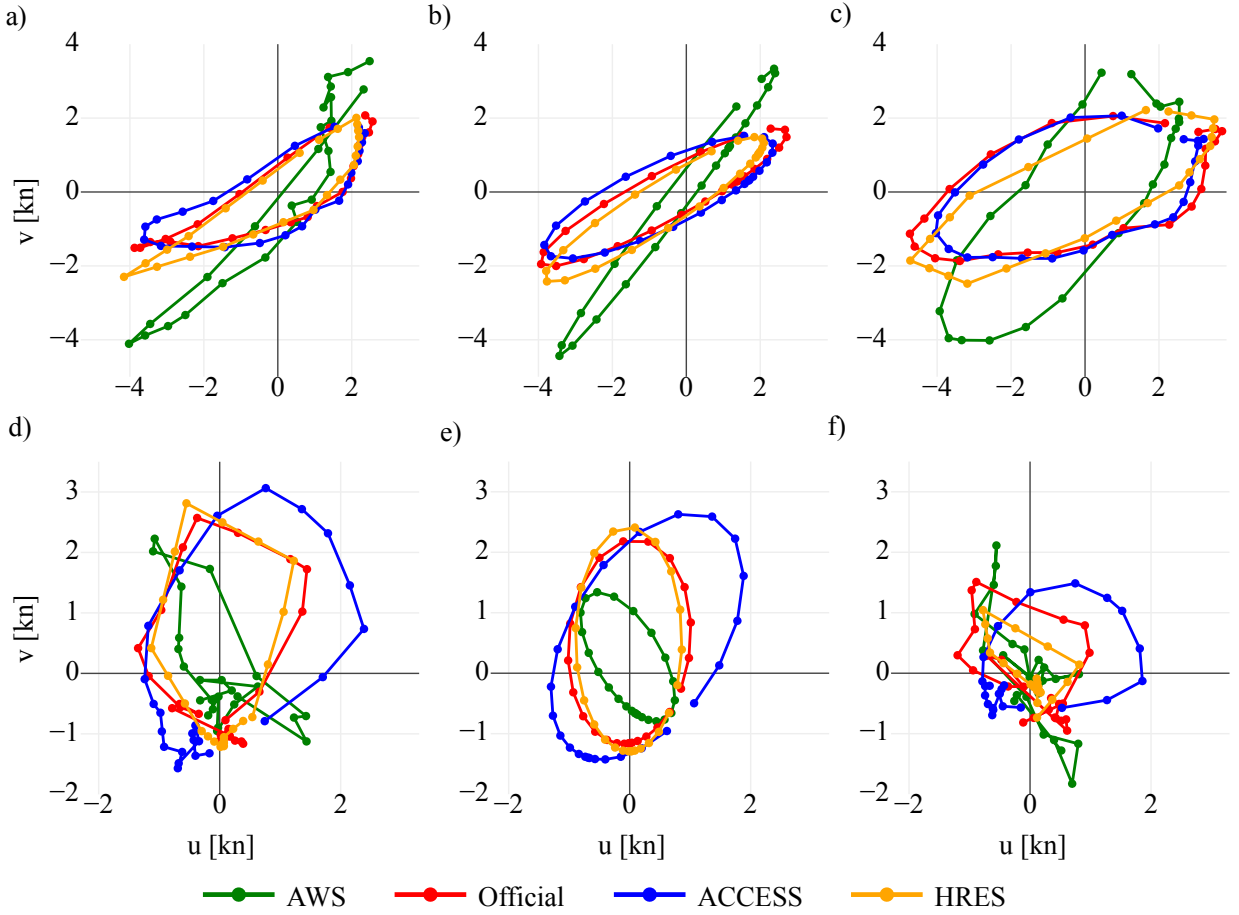


FIG. 13. Temporal hodographs of wind perturbations at each hour UTC averaged over June, July and August 2018, at Brisbane and Hobart airports, a) and d), and the associated ellipse fits, b) and e). For comparison, c) and f) provide the hodographs of the averaged perturbations at the Spitfire Channel and Hobart city stations, respectively (see Fig. 1).

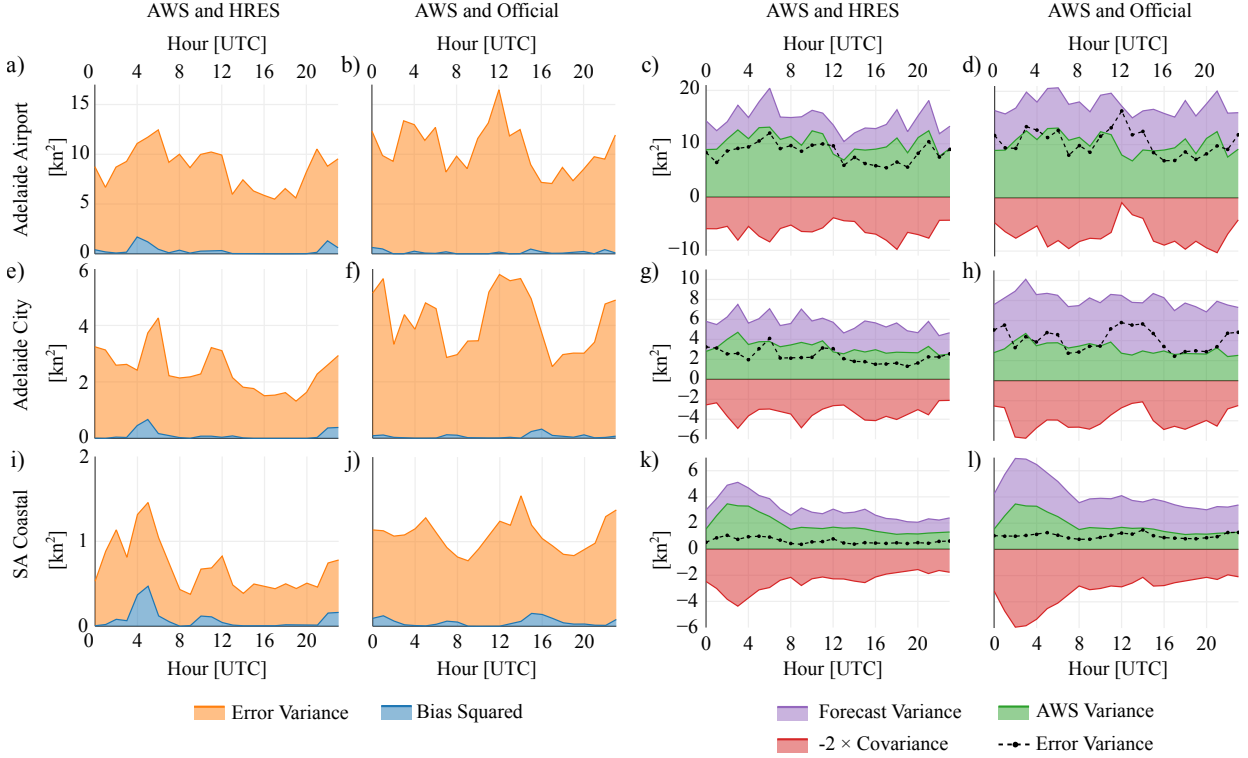


FIG. 14. Mean square error between the AWS and HRES zonal perturbations $\overline{(u_{\text{AWS}} - u_{\text{H}})^2}$, a), e), and i), decomposed into the error variance $\text{var}(u_{\text{AWS}} - u_{\text{H}})$ and squared bias $(\bar{u}_{\text{AWS}} - \bar{u}_{\text{H}})^2$ terms of equation (8). Also, the decomposed mean square error between the AWS and official forecast zonal perturbations, b), f) and j). Additionally, the HRES and AWS error variance term $\text{var}(u_{\text{AWS}} - u_{\text{H}})$ decomposed into the $\text{var}(u_{\text{AWS}})$, $\text{var}(u_{\text{H}})$ and $-2 \cdot \text{cov}(u_{\text{AWS}}, u_{\text{H}})$ terms, c), g) and k), and analogously for the official forecast and AWS error variance term $\text{var}(u_{\text{AWS}} - u_{\text{O}})$, d), h) and l). Decompositions given for Adelaide Airport, a) to d), the Adelaide city station group, e) to h), and the SA coastal station group, i) to l) (see Fig. 1.)