

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.*

Benjamin ?. Price and Nicholas ?. Loveday

Bureau of Meteorology, Casuarina, Northern Territory, Australia

Alexei Hider

Bureau of Meteorology, 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

13 Forecasts issued by the Australian Bureau of Meteorology (BoM) are based
14 on model data that is edited by human forecasters. Two types of edits are
15 commonly made to the wind fields. These edits aim to improve how bound-
16 ary layer mixing and land-sea breeze processes are resolved in the forecast. In
17 this study we compare the diurnally varying component of the BoM's Official
18 edited wind forecast, with that of station observations and unedited model
19 datasets, to assess changes to error and bias resulting from these edits. We
20 consider coastal locations across Australia over June, July and August 2018,
21 aggregating data over three spatial scales. The edited forecast generally only
22 produces a lower mean absolute error than model guidance at the coarsest
23 spatial scale (1000 - 2000 km), but can achieve lower seasonal biases over all
24 spatial scales. However, the edited forecast only reduces errors or biases at
25 particular times and locations, and rarely produces lower errors or biases than
26 all model guidance products simultaneously. This suggests that forecaster
27 skill lies mostly in making the choice of model guidance, rather than in mak-
28 ing edits. To better understand the biases in the diurnal wind cycles, we fit
29 modified ellipses to the temporal hodographs of seasonally averaged diurnal
30 wind cycles. Biases in the Official forecast diurnal cycle vary with location
31 for multiple reasons, including biases in the directions sea-breezes approach
32 coastlines, amplitude and shape biases in the hodographs, and disagreement
33 as to whether sea-breeze or boundary layer mixing processes contribute most
34 to the diurnal cycle.

35 1. Introduction

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

47 Australian forecasters generally make two types of edits to the surface wind fields on a routine
48 daily basis. The first involves modifying the surface winds after sunrise at locations where the
49 forecaster believes the model guidance is providing a poor representation of boundary layer mixing
50 processes. Boundary layer mixing occurs as the land surface heats up, producing an unstable
51 boundary layer which transports momentum downward to the surface layer. Before this mixing
52 occurs, winds are typically both weaker and ageostrophically oriented due to surface friction (Lee
53 2018), and so mixing can affect both the speed and direction of the surface winds. Australian
54 forecasters perform these edits using a GFE tool which allows them to specify a region over which
55 to apply the edit, a height z and a percentage p , with the tool then calculating a weighted average
56 of the surface winds and winds at z , weighted by p .

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

63 Forecasters, and the weather services that employ them, have good reasons for ensuring the
64 diurnally varying component of their wind forecasts are as accurate as possible. In addition to the
65 significant contribution diurnal wind cycles make to overall wind fields (e.g. Dai and Deser 1999),
66 diurnal wind cycles are important for the ventilation of pollution, with sea-breezes transporting
67 clean maritime air inland, where it helps flush polluted air out of the boundary layer (Miller et al.
68 2003; Physick and Abbs 1992). Furthermore, diurnal wind cycles affect the function of wind
69 turbines (Englberger and Dörnbrack 2018) and the design of wind farms (Abkar et al. 2016), as
70 daily patterns of boundary layer stability affect turbine wake turbulence, and the losses in wind
71 power that result.

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

80 Other studies have assessed near-surface wind forecasts, verifying the total wind speeds, not
81 just the diurnal component. Pinson and Hagedorn (2012) studied the 10 m wind speeds from the
82 European Centre for Medium Range Weather Forecasting (ECMWF) operational model ensemble
83 across western Europe over December, January, February 2008/09. They found that the worst
84 performing regions were coastal and mountainous areas, and attributed this to the small scale
85 processes, e.g. sea and mountain breezes, that are under-resolved by the ensembles coarse 50 km
86 spatial resolution.

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

94 **2. Data and Methods**

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

99 *a. Data*

100 Four datasets are considered in this study; the human edited Official BoM wind forecast data that
101 is issued to the public, observational data from automatic weather stations (AWS) across Australia,

unedited data from the ECMWF’s high resolution 10-day forecast model (HRES), and unedited model data from the Australian Community Climate and Earth System Simulator (ACCESS), noting that HRES and ACCESS are two of the model guidance products most commonly used by Australian forecasters for winds. We consider just the lead-day one forecasts of Official, HRES and ACCESS, for reasons discussed below.

This study primarily considers the austral winter months of June, July and August 2018. This short time period was chosen to reduce the effect of changing seasonal and climatic conditions, changing forecasting practice and staff, and of developments to the ACCESS and HRES models. Results for December, January and February 2017/18 are occasionally mentioned to strengthen conclusions or provide a seasonal contrast.

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

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

The Official forecast dataset is produced on a state by state basis at forecasting centres located in most state capitals. To construct the Official forecast dataset, forecasters make a choice of model guidance in the GFE, which then interpolates or upscales the model data onto a standard 3 km

126 spatial grid for Victoria and Tasmania, or a 6 km grid for the rest of the country. GFE displays
127 model data at hourly intervals by taking the model guidance output at each hour UTC, with the
128 exception of the HRES model data which is only provided to the BoM at 3 hourly intervals, and is
129 therefore linearly interpolated to hourly intervals by the GFE. Forecasters then make edits to these
130 3 or 6 km hourly grids to produce the Official forecast datasets.

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

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

142 *b. Assessing Diurnal Cycles*

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

147 To compare errors in the Official, ACCESS and HRES diurnal cycles we calculate the Euclidean
148 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 perturbations, \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 and ECMWF perturbations, and of the ACCESS and ECMWF 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 climatology, or model predictions of background conditions, 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

171 contribute to differences in the perturbations. To minimise the importance of background state
172 differences, this study focuses exclusively on lead-day one forecasts.

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

181 Although temporal autocorrelations of DAE, i.e. correlations between DAE values at a particular
182 hour from one day to the next, are in practice small or non-existent, they are still accounted for
183 by reducing the “effective” sample size to $n(1 - \rho_1) / (1 + \rho_1)$, where n is the actual sample size
184 and ρ_1 is the lag-1 autocorrelation (Zwiers and von Storch 1995; Wilks 2011). In the language of
185 statistical hypothesis testing, the null hypothesis that $E = 0$ would be rejected at significance level
186 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
187 value of $\Pr(E > 0)$, referring to this as a *confidence score*, and noting $\Pr(E < 0) = 1 - \Pr(E > 0)$.
188 We say Official outperforms model guidance with “high confidence” if $\Pr(E > 0) \geq 95\%$, or
189 that model guidance outperforms Official with “high confidence” if $\Pr(E > 0) \leq 5\%$, with high
190 confidence implicit whenever it is not explicitly mentioned.

191 Following the “fuzzy verification” agenda outlined by Ebert (2008), forecast and observational
192 perturbation datasets are compared not only at individual stations, but are also averaged over two
193 coarser spatial scales before being compared. The individual stations we consider are the 8 capital
194 city *airport stations*, marked by stars in Fig. 1, as their high operational significance means that

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

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

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

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

with DB_{OH} and DB_{HA} defined analogously, where the over-bars denote temporal averages of the perturbations at a particular hour, over June, July and August 2018. These temporally averaged perturbations can be viewed as the climatological diurnal wind cycles over the three month study period for each dataset. Biases over the city and coastal station groups are calculated by taking the spatial average before the temporal average. Uncertainty in the DB is estimated through bootstrap-

ping (Efron 1979). This is done by performing resampling with replacement on the underlying
 perturbation datasets, and calculating the DB multiple times using these resampled datasets. This
 provides a distribution of DB values, which analogously to with DAE, we treat as a sample from
 a random variable B , and use this to estimate $\Pr(B > 0)$.

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

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

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

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

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

comparatively stagnant at night (e.g. Fig. 9). Thus, we instead fit the equations

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

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

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 nonlinear regression, performed using the `least_squares` function from the `scipy.optimize` python module (SciPy 2019).

Gille et al. (2005) fit equations (3) and (4) to the temporally averaged wind fields, so that (u_0, v_0) could be interpreted as the mean wind over the study's time period, and the remaining terms providing the average diurnal perturbations. In this study we fit equations (5) and (6) to the average perturbations themselves, with (u_0, v_0) now necessary to offset the asymmetry introduced by α , i.e. to ensure the time integral of the fitted values is approximately zero.

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.

254 *a. Absolute Errors*

255 Figure 2 provides the mean difference of absolute error $\overline{\text{DAE}}$ values and confidence scores de-
256 fined in section 2 for the coastal station groups shown in Fig. 1, for $\overline{\text{DAE}}_{\text{OA}}$, $\overline{\text{DAE}}_{\text{OH}}$ and $\overline{\text{DAE}}_{\text{HA}}$,
257 which represent the the Official versus ACCESS, Official versus HRES, and HRES versus AC-
258 CESS comparisons, respectively. The results indicate that for the majority of station groups and
259 hours, both the unedited ACCESS and HRES models outperform the Official forecast. The lowest
260 $\overline{\text{DAE}}$ values occur at the NT station group at 23:00 and 00:00 UTC for both $\overline{\text{DAE}}_{\text{OA}}$ and $\overline{\text{DAE}}_{\text{OH}}$.
261 Although Official outperforms at least one of ACCESS or HRES at multiple times and station
262 groups, the only group and time where it outperforms both is 05:00 UTC over the South WA sta-
263 tion group. HRES generally outperforms ACCESS from 10:00 - 14:00 UTC, with the South WA
264 station group being the main exception.

265 Figures 3 and 4 provide case studies of the NT and South WA station groups, respectively. Figure
266 3 a) provides a time series of DAE for the NT station group at 23:00 UTC. The time series shows
267 significant temporal variability, with DAE frequently dropping below -2 kn. Figures 3 b) and c)
268 show hodographs of the winds and wind perturbations, respectively, at each hour UTC on the 3rd
269 of July, which provides an interesting example.

270 Figure 3 b) shows that the Official wind forecast on this day was likely based on edited ACCESS
271 from 00:00 to 06:00 UTC, then edited HRES from 07:00 to 13:00 UTC, then unedited ACCESS
272 from 15:00 to 21:00 UTC. At 22:00 and 23:00 UTC, the Official winds acquire stronger east-
273 northeasterly components than the other datasets. Figure 5 a) shows the first ten values from
274 wind soundings at Darwin Airport at 12:00 UTC on July 3rd and 00:00 UTC on July 4th. In both
275 instances the winds are east-southeasterly, and so the rapidly changing wind perturbations at 22:00
276 UTC in the Official forecast likely reflect a boundary layer mixing edit that has been applied either

too early, or has strengthened the southeasterly component of the winds too much. Similar issues create low DAE values on the 8th of June and 9th and 10th of July.

Figure 4 a) provides a time series of DAE for the South WA station group at 05:00 UTC. As with the NT station group there is significant temporal variability, with DAE frequently exceeding 1 kn. Figures 4 b) and c) provide hodographs of the winds and wind perturbations, respectively, on the 9th of June, another interesting example. The perturbation hodograph shows both HRES and ACCESS under-predicting the amplitude of the diurnal wind cycle on this day. Figure 5 shows wind soundings at Perth Airport, the nearest station to provide wind soundings, between 12:00 UTC on the 8th June and 12:00 UTC on the 9th June. The 8th June 12:00 UTC sounding shows surface northerlies of around 6 kn, becoming west to northwesterlies of over 20 kn 2.4 km above the surface. However, the subsequent sounding at 00:00 UTC on the 9th of June shows that the winds acquire a strong northerly component of 30 kn in the first 500 m of the atmosphere, with the final sounding indicating a strong northwesterly wind at 725 m persisting until 12:00 UTC.

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

Fig. 6 presents the \overline{DAE} values and confidence scores for the airport stations, and city station groups, for the Official versus HRES comparison, i.e. \overline{DAE}_{OH} . The results for the airport stations are noisier than the results for the coastal station groups in Figs. 2 c) and d), although they share

301 some similarities. For instance, Official outperforms HRES at 01:00 and 02:00 UTC at both the
302 Darwin airport station and the NT coastal station group. There are four other instances where
303 Official outperforms HRES with at least 90% confidence, although this could simply be occurring
304 by chance due repeated testing (Wilks 2011, p. 178).

305 For the city station groups, HRES outperforms Official almost uniformly. The main exception
306 is the Darwin city station group, where Official outperforms HRES at 02:00 UTC, and there is
307 ambiguity as to whether Official or HRES performs better at 01:00, 03:00 and 04:00 UTC, and
308 from 15:00 to 22:00 UTC. The analogous \overline{DAE}_{OA} Official versus ACCESS comparisons (not
309 shown) are similar, with the airport station results noisy, but ACCESS outperforming Official over
310 the city station groups for the vast majority of times and locations. Over the December, January,
311 February 2017/18 season, HRES also outperforms Official almost uniformly over the city station
312 groups, although the Official versus ACCESS comparisons are more ambiguous.

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

320 *b. Seasonal Biases*

321 Figure 8 provides the difference of biases (DB) and confidence scores defined in section 2, for
322 the coastal station groups for DB_{OA} , DB_{OH} and DB_{HA} , which represent the the Official versus
323 ACCESS, Official versus HRES, and HRES versus ACCESS comparisons, respectively. At the

NT station at 03:00 UTC, Official outperforms both ACCESS and HRES with confidence $\geq 93\%$. However, both ACCESS and HRES outperform Official at 23:00 and 00:00 UTC, and from 05:00 to 11:00 UTC, consistent with the $\overline{\text{DAE}}$ results of Fig. 2. Figure 9 a) shows that these biases are mostly a consequence of amplitude biases in Official's diurnal cycle.

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

At the SA station group from 02:00 to 05:00 UTC and 09:00 to 12:00 UTC, Official 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 Official only slightly outperforms ACCESS in the results of Figures 8 a) and b).

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

352 *c. Ellipse Fits*

353 The hodographs in Fig. 9 are roughly elliptical in shape, suggesting that descriptive quantities
354 can be estimated by fitting equations (5) and (6) to the zonal and meridional climatological per-
355 turbations, as described in section 2. Figure 11 gives the R^2 values for the fits of the zonal and
356 meridional perturbations to equations (5) and (6), respectively. The fit performs best at the coastal
357 station group spatial scale, with R^2 generally above 95%.

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

369 To check whether this just represents a direction bias of the Brisbane Airport weather station,
 370 Fig. 13 c) shows the climatological perturbations at the nearby Spitfire Channel station (see Fig. 1).
 371 While the amplitude bias is smaller at Spitfire Channel than Brisbane Airport, the directional bias
 372 is at least as high. A similar directional bias is evident at the nearby Inner Beacon station (not
 373 shown), although the bias is smaller than at Spitfire Channel and Brisbane Airport. Similar biases
 374 are also evident at these stations during December, January and February 2017/18, with the semi-
 375 major axis of Official's ellipse fit oriented 29° clockwise from AWS's at Brisbane airport. Figure
 376 1 shows there are two small islands to the east of Brisbane airport; the more north-northeasterly
 377 orientation of the Brisbane Airport sea-breeze suggests these islands may be redirecting winds
 378 between the east coast of Brisbane and the west coasts of these islands, and that this local effect is
 379 not being captured in Official, ACCESS or HRES.

380 Another example is the Hobart Airport station. Figure 12 c) shows that the semi-major axis of
 381 the AWS ellipse fit is oriented 31, 35 and 62 degrees anti-clockwise from the semi-major axes of
 382 the HRES, Official and ACCESS ellipse fits, respectively. Figures 11 a) and b) show that the ellipse
 383 fit for the AWS perturbations at Hobart airport only achieve R^2 values of 59% and 68% for the u
 384 and v components, respectively, but figures 13 d) and e) show that the fit still captures orientations
 385 accurately, although it underestimates the maximum AWS perturbation. Figure 13 f) provides the
 386 climatological perturbations at the Hobart (city) station, which also show a large difference in ori-
 387 entation between ACCESS and AWS. Given the timing of the westerly perturbations in ACCESS,
 388 and the fact that the prevailing winds around Tasmania are westerly, these results suggest that AC-
 389 CESS is exaggerating the boundary layer mixing processes involved in the diurnal cycle around
 390 Hobart. These biases are not present during December, January and February 2017/18, as strong
 391 south to southeasterly sea-breeze perturbations are now dominant all four datasets, although the
 392 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 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.

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

4. Discussion

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

416 To investigate the effect of spatial scale on error, consider first just the zonal components of the
 417 AWS and Official wind perturbations, denoted by u_{AWS} and u_{O} respectively. Considering just the
 418 values at a particular hour UTC, over the entire June, July, August time period, the mean square
 419 error $\text{mse}(u_{\text{AWS}}, u_{\text{O}}) = \overline{(u_{\text{AWS}} - u_{\text{O}})^2}$ can be decomposed $\text{mse}(u_{\text{AWS}}, u_{\text{O}}) =$

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

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

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

436 At the intermediate Adelaide city station group scale, the AWS variances are of similar magni-
 437 tudes to those of ECMWF, but smaller than those of Official, with Official's additional variability

mostly uncorrelated to AWS. This results in larger error variance terms for Official, consistent with ECMWFs almost complete outperformance of Official 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 Official is still larger than that of ECMWF, ECMWF'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 Official'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 Official being more variable than ECMWF, with this additional variability mostly random, in the sense of being uncorrelated with AWS. Similarly, Official is generally more variable than ACCESS, explaining why Official also struggles to outperform ACCESS at these scales, and ACCESS is generally more variable than ECMWF, explaining why ECMWF 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 Official 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 ECMWF: 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, 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 Official more often than not.

Reducing the random variability of Official, 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 very coarse spatial scales of 500 to 2000 km. Even at these scales, reductions in error are isolated

485 to particular locations and times of day, and Official rarely produces lower mean absolute error
486 than both commonly used model guidance products simultaneously. This suggests that forecaster
487 skill lies more in making the choice of model guidance than in making edits.

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

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

503 Future research could extend this study in multiple directions. One important question is whether
504 the random variability in the Official forecast, or the model guidance products that comprise it,
505 could be reduced through ensemble forecasting or post-processing, as reducing random variability
506 would both decrease errors, and increase the value of land-sea breeze and boundary layer mixing
507 edits. Another goal could be to identify precisely the spatiotemporal scales at which diurnal wind

508 cycles can be identified against background noise, so as to better understand the scales at which
509 land-sea breeze and boundary layer mixing edits can add value to a forecast.

510 *Acknowledgments.* Funding for this study was provided for Ewan Short by the Australian Re-
511 search Council’s Centre of Excellence for Climate Extremes (CE170100023). Datasets and soft-
512 ware were generously provided by the Australian Bureau of Meteorology’s Evidence Tasked Au-
513 tomation team. Thanks are due to Michael Foley and Deryn Griffiths for providing support at the
514 Bureau of Meteorology’s Melbourne office, and to Prof. Craig Bishop for some helpful conversa-
515 tions. The code written for this study is freely available online (Short 2019).

516 **References**

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

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

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

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

529 Dai, A., and C. Deser, 1999: Diurnal and semidiurnal variations in global surface wind
 530 and divergence fields. *Journal of Geophysical Research*, **104**, 31 109–31 125, doi:10.1029/
 531 1999JD900927.

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

534 Efron, B., 1979: Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7** (1),
 535 1–26, doi:10.1214/aos/1176344552.

536 Englberger, A., and A. Dörnbrack, 2018: Impact of the diurnal cycle of the atmospheric bound-
 537 ary layer on wind-turbine wakes: a numerical modelling study. *Boundary-Layer Meteorology*,
 538 **166** (3), 423–448, doi:10.1007/s10546-017-0309-3.

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

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

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

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

549 Lock, A. P., A. R. Brown, M. R. Bush, G. M. Martin, and R. N. B. Smith, 2000: A new bound-
 550 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 cycle 41r2. ECMWF, [Available online at <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.

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

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

SciPy, 2019: Optimization and root finding (scipy.optimize). SciPy, [Available online at <https://docs.scipy.org/doc/scipy/reference/optimize.html>].

572 Short, E., 2019: eshort0401/forecast_verification_paper. GitHub, [Available online at [https://](https://github.com/eshort0401/forecast_verification_paper)
573 github.com/eshort0401/forecast_verification_paper].

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

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

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

582 Zaron, E. D., and G. D. Egbert, 2006: Estimating open-ocean barotropic tidal dissipation: The
583 hawaiian ridge. *Journal of Physical Oceanography*, **36** (6), 1019–1035, doi:10.1175/JPO2878.
584 1.

585 Zwiers, F. W., and H. von Storch, 1995: Taking serial correlation into account in tests of the mean.
586 *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 *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. 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 Official versus ACCESS, a) and b), Official 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 Official versus ACCESS, DAE_{OA} , and Official 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. 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 9th 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 $\overline{(u_{AWS} - u_H)^2}$, a), e), and i), decomposed into the error variance $\text{var}(u_{AWS} - u_H)$ and squared bias $(\bar{u}_{AWS} - \bar{u}_H)^2$ terms of equation (8). Also, the decomposed mean square error between the AWS and Official zonal perturbations, b), f) and j). Additionally, the HRES and AWS error variance term $\text{var}(u_{AWS} - u_H)$ decomposed into the $\text{var}(u_{AWS})$, $\text{var}(u_H)$ and $-2 \cdot \text{cov}(u_{AWS}, u_H)$ terms, c), g) and k), and analogously for the Official and AWS error variance term $\text{var}(u_{AWS} - u_O)$, 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

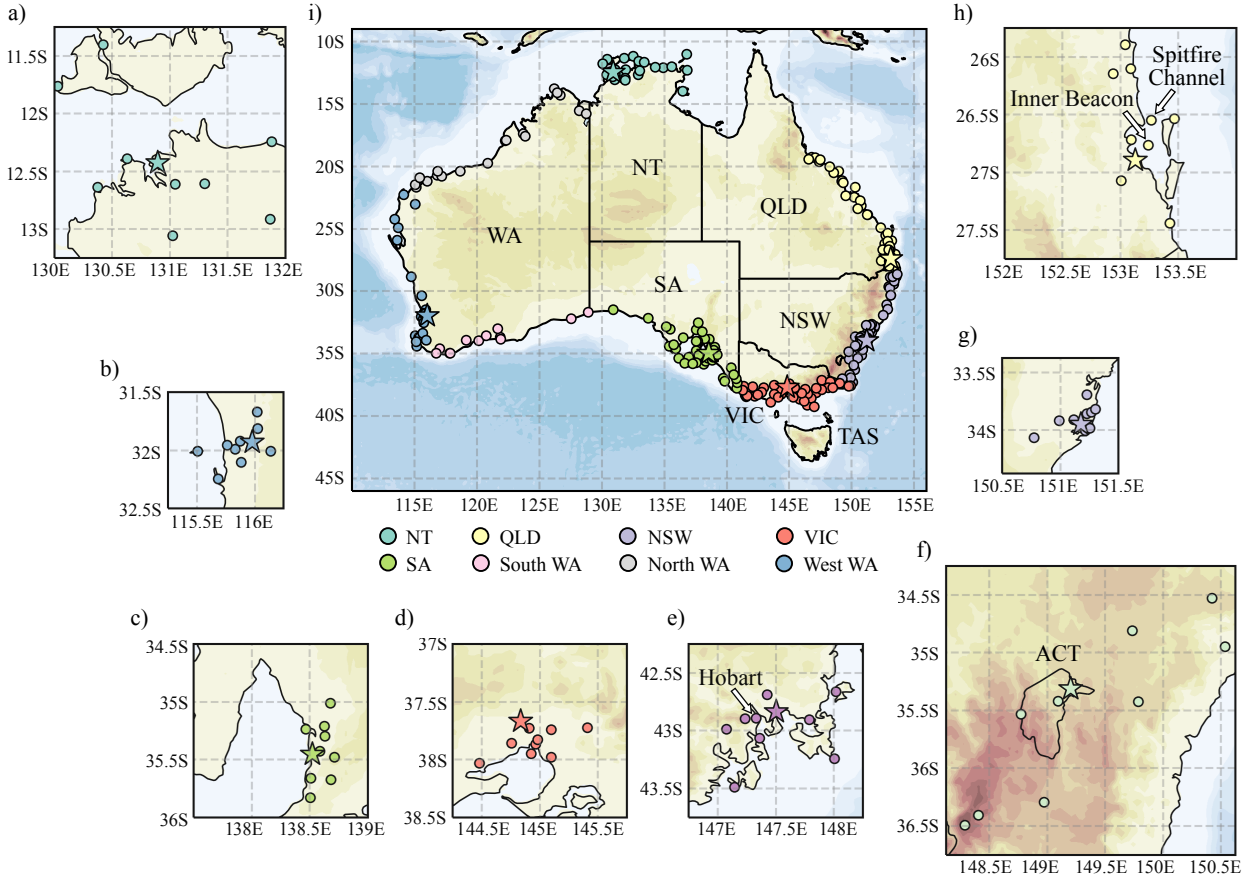


FIG. 1. Locations of the automatic weather stations considered in this study, where stars give the locations of the capital city *airport stations*. Stations are divided into the Darwin, Perth, Adelaide, Melbourne, Hobart, Canberra, Sydney and Brisbane *city station groups*, a) to h), respectively, and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

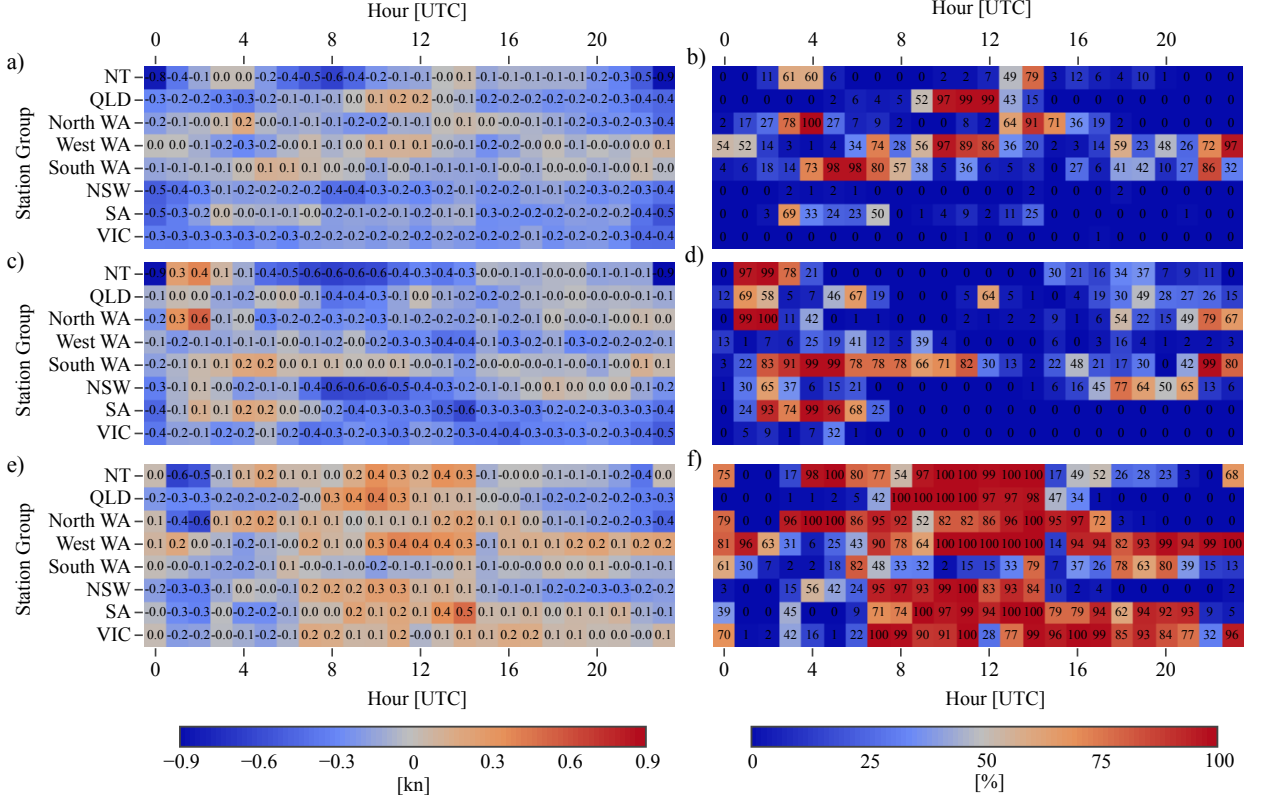


FIG. 2. Heatmaps of mean difference of absolute error $\overline{\text{DAE}}$ values, a), c), e), and confidence scores, b), d), f), for each coastal station group (see Fig. 1) and hour of the day, for Official versus ACCESS, a) and b), Official versus HRES, c) and d), HRES versus ACCESS, e) and f). 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 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).

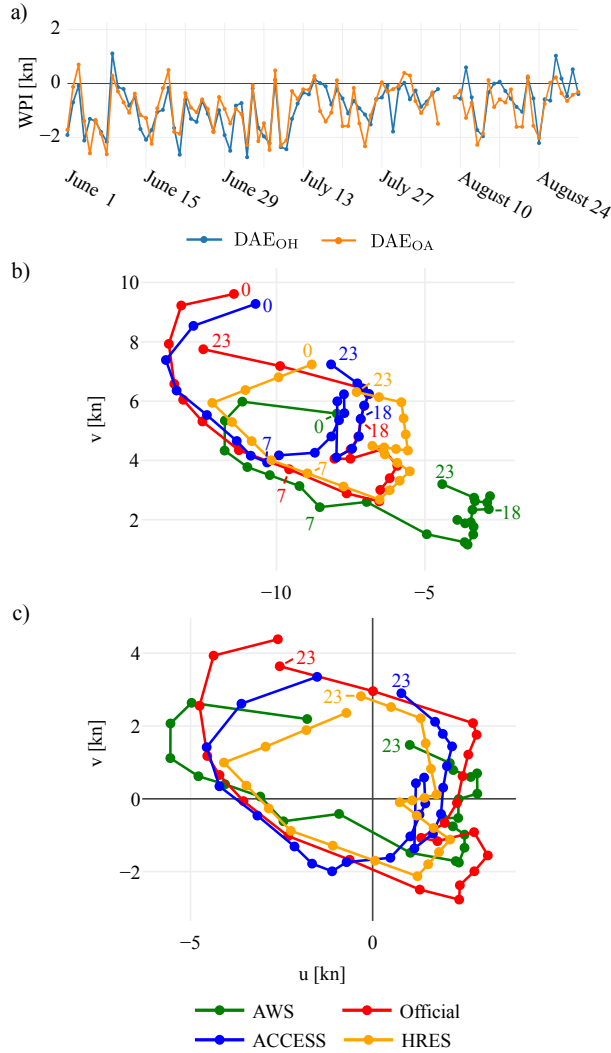


FIG. 3. Time series, a), of the difference in absolute error DAE defined in equation (1) for Official versus ACCESS, DAE_{OA}, and Official versus HRES, DAE_{OH}, for the NT coastal station group shown in Fig. 1 at 23:00 UTC. Also, temporal hodographs in hours UTC showing hourly changes in winds, b),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

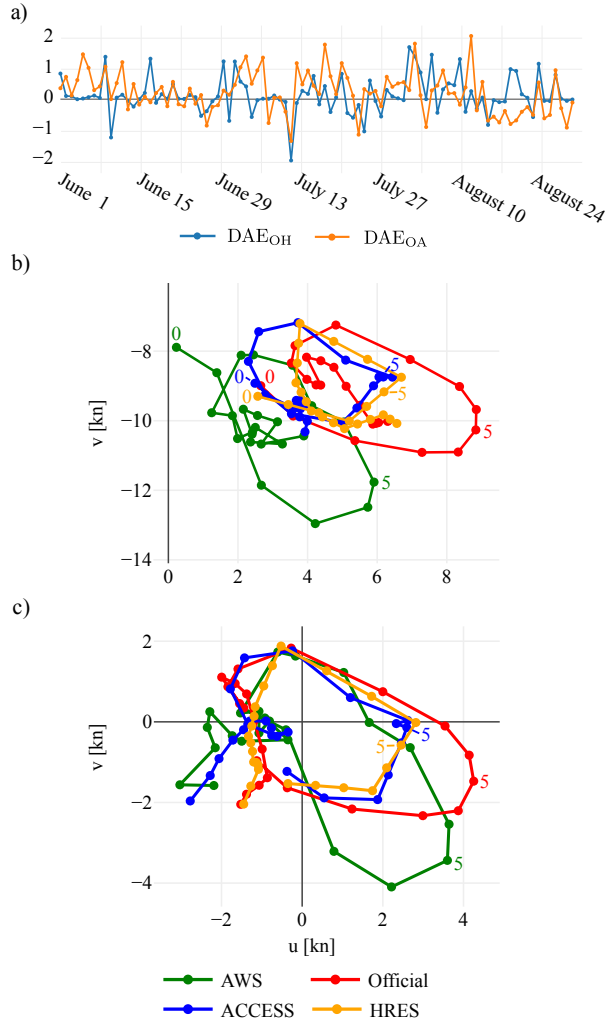


FIG. 4. As in Fig. 3, but for, a), the South WA coastal station group at 05:00 UTC, and b) and c), the winds and wind perturbations, respectively, over the South WA coastal station group on the 9th June 2018.

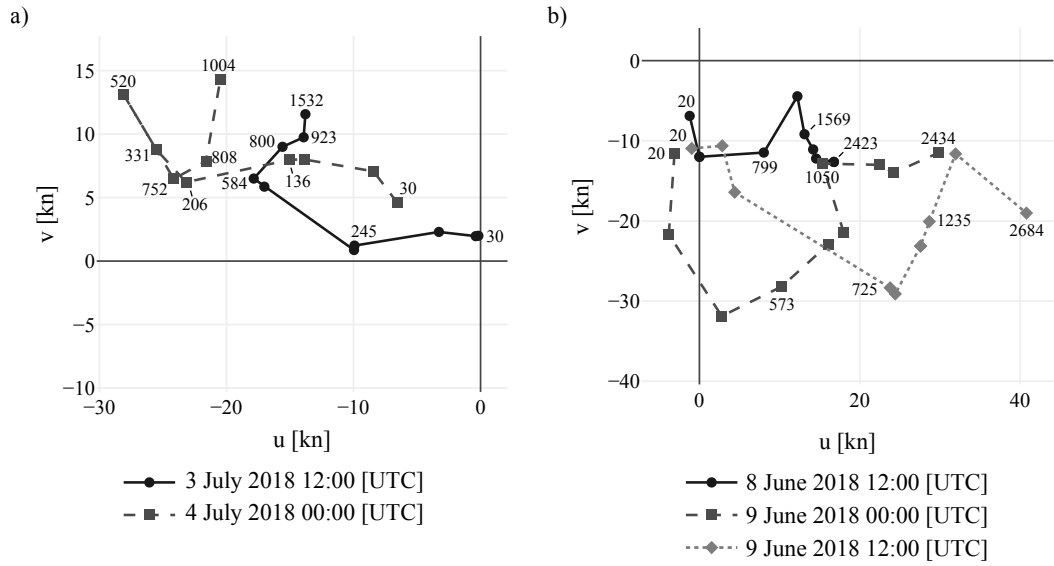


FIG. 5. Vertical wind soundings at, a), Darwin Airport, and b), Perth Airport, with heights given in metres.

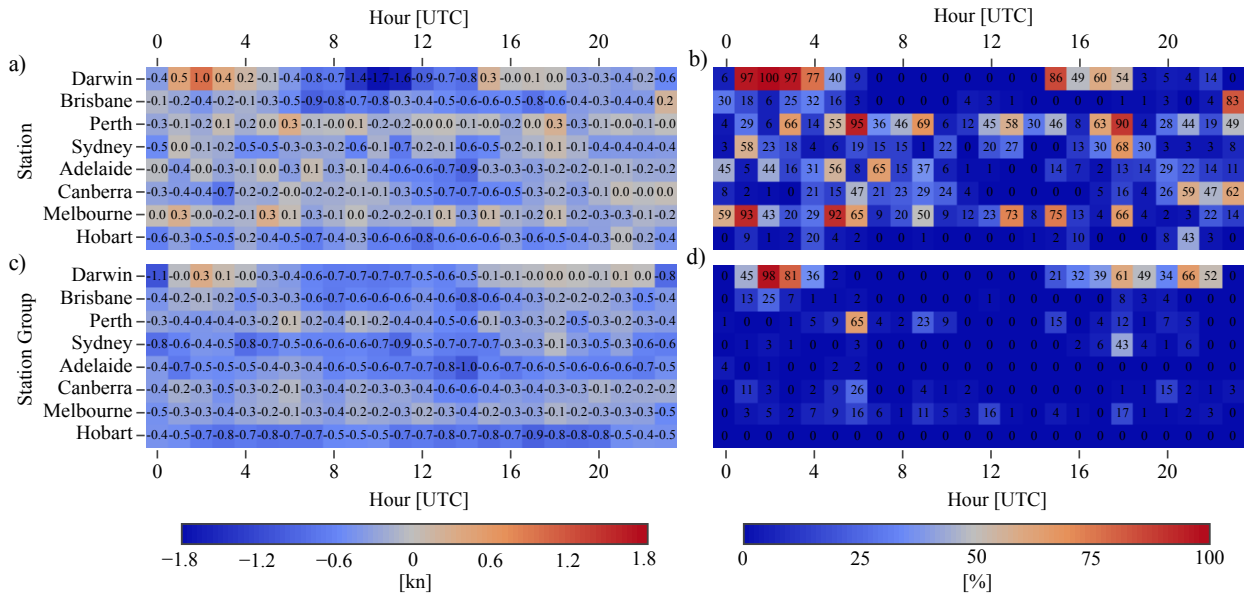


FIG. 6. As in Fig. 2, but for the Official versus HRES mean difference of absolute error \overline{DAE}_{OH} values, a) and c), and confidence scores, b) and d), for the airport stations, a) and b), and city station groups, c) and d).

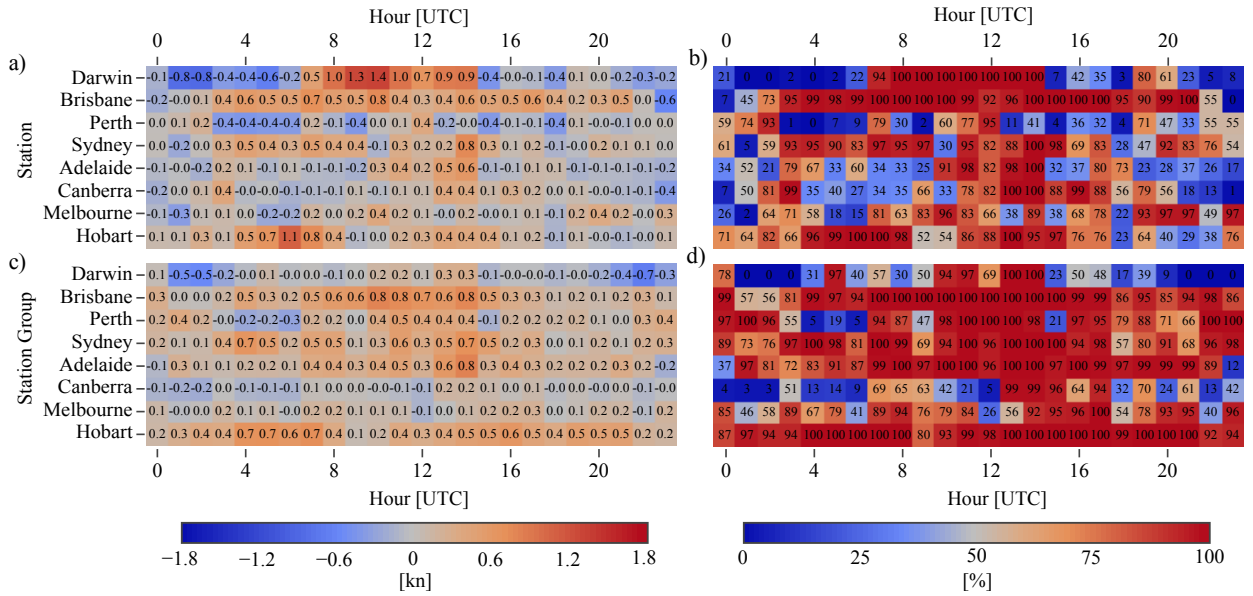


FIG. 7. As in Fig. 6, but for the HRES versus ACCESS mean difference in absolute error \overline{DAE}_{HA} values and confidence scores.

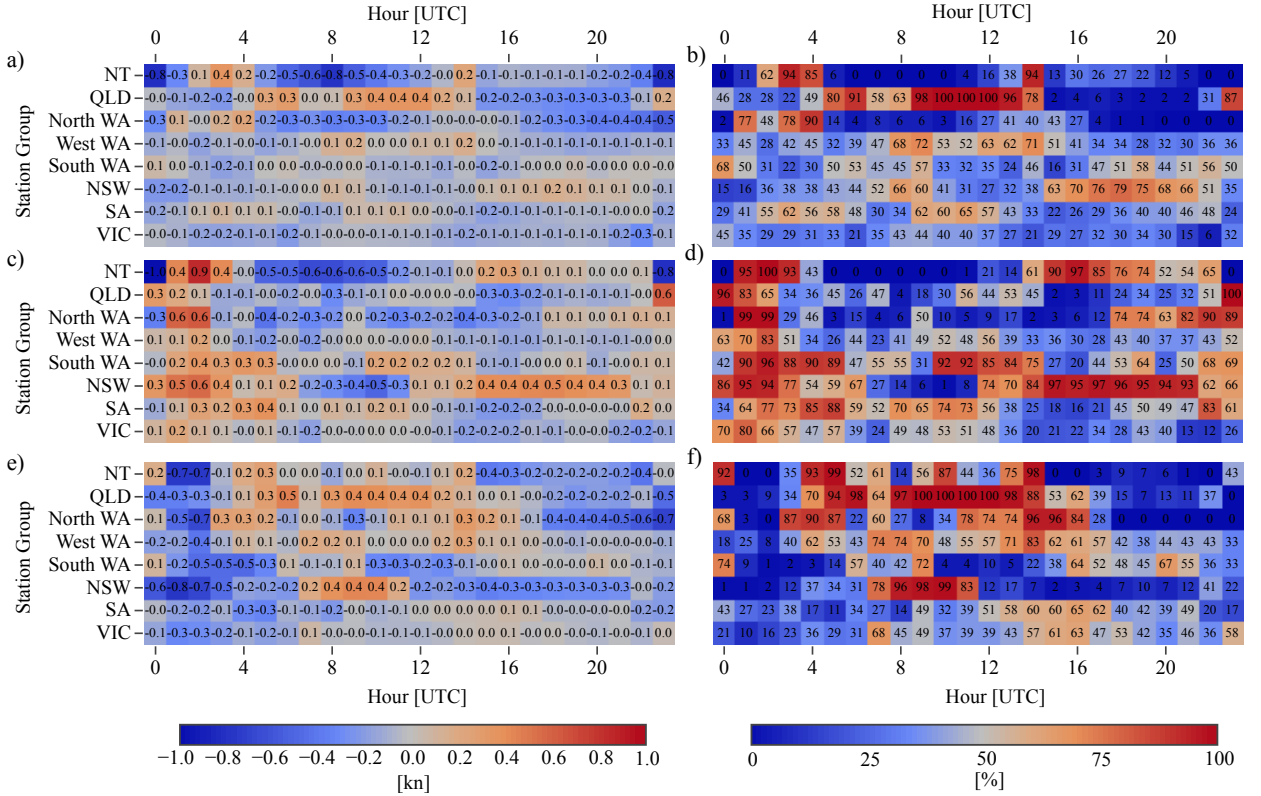


FIG. 8. As in Fig. 2, but for the difference of biases (DB) values and confidence scores.

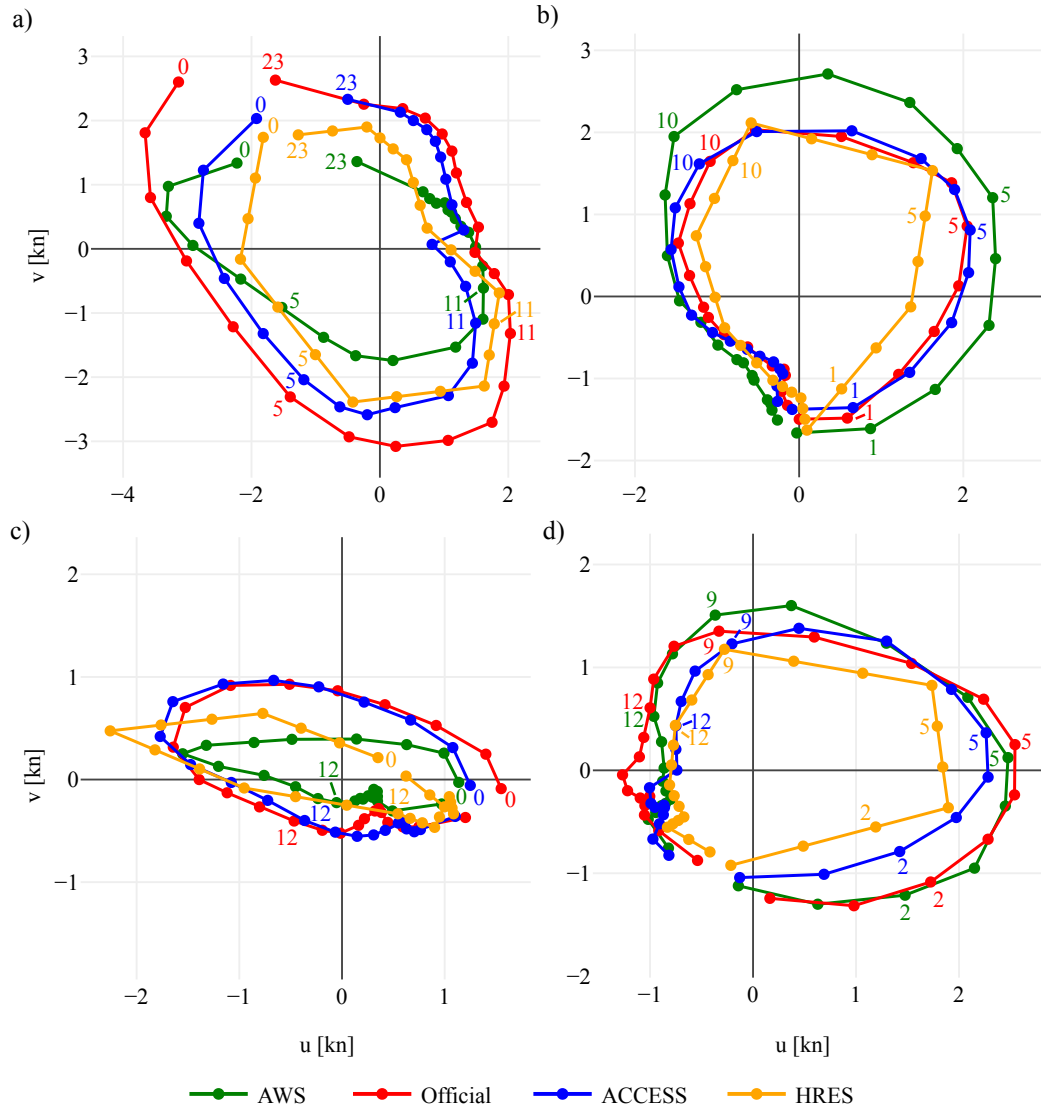


FIG. 9. Temporal hodographs in hours UTC of wind perturbations spatially averaged over the, a), NT, b) South WA, c) NSW and d), SA coastal station groups (see Fig. 1), and temporally averaged over June, July and August 2018.

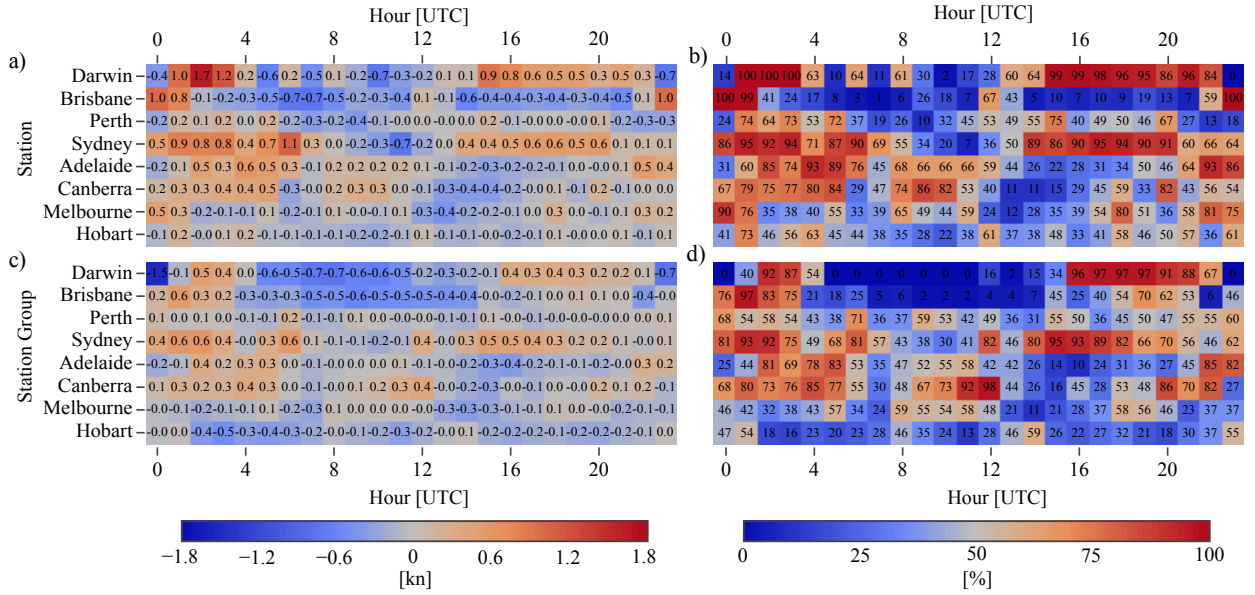


FIG. 10. As in Fig. 6, but for the difference of biases (DB) values and confidence scores.

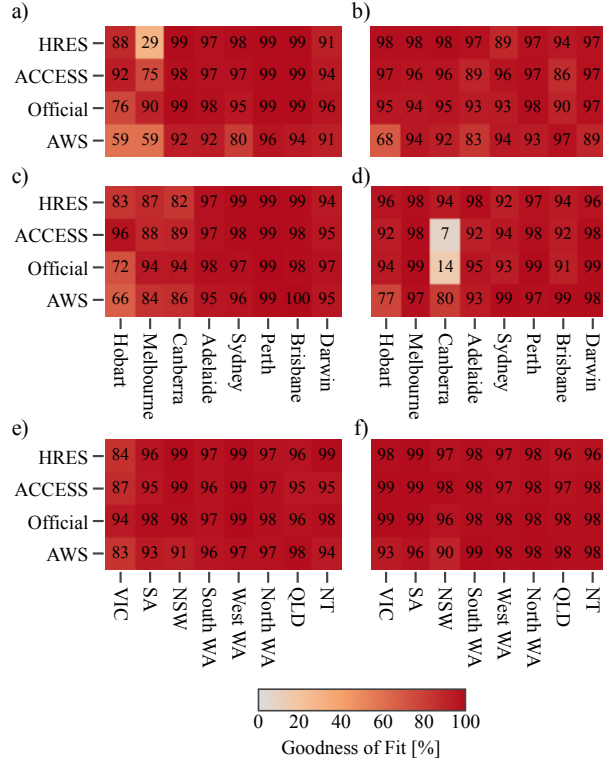


FIG. 11. R^2 values as percentages for the fit of equation (5) to the zonal perturbations, a), c) and e), and equation (6) to the meridional perturbations, b), d) and f), for the airport stations, a) and b), city station groups, c) and d), and coastal station groups, e) and f), shown in Fig. 1.

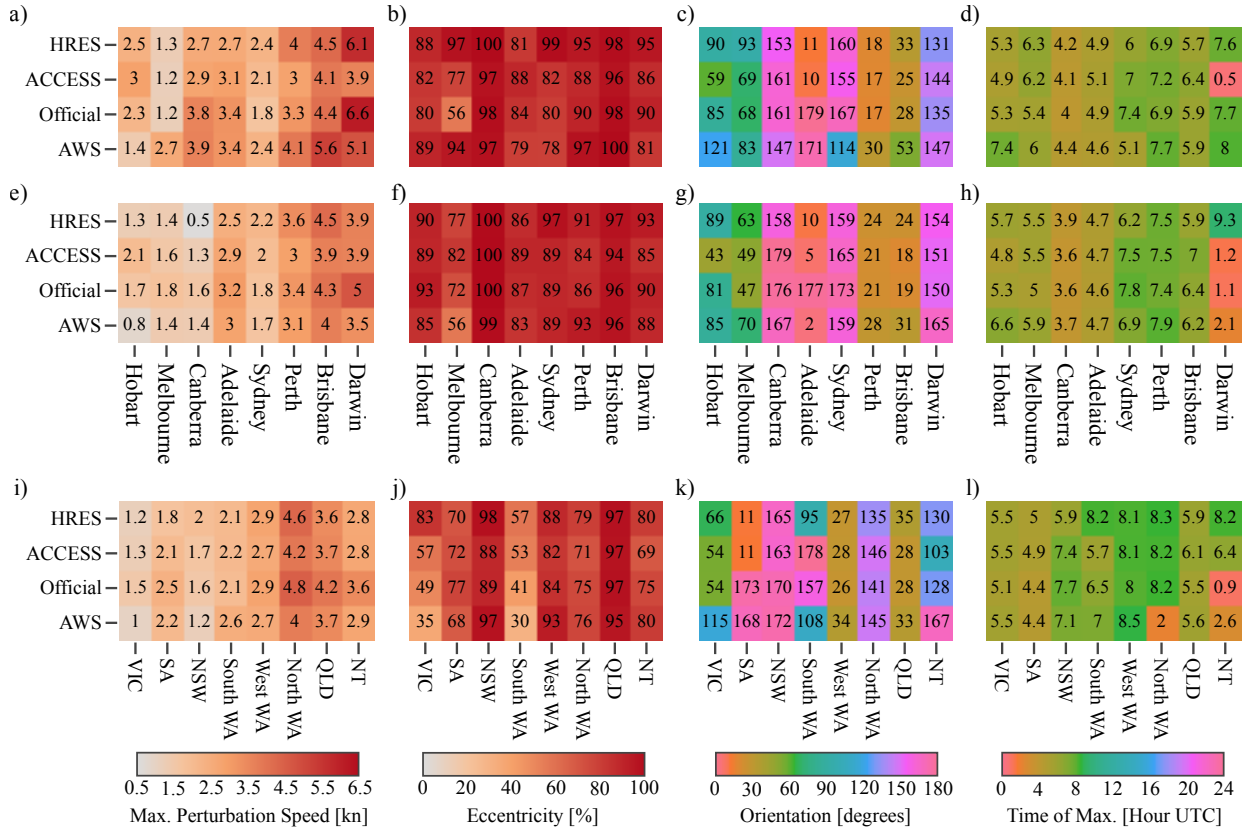


FIG. 12. Metrics derived from fitting ellipse equations (5) and (6) to wind perturbations at the Australian capital city airport stations, a) to d), and to wind perturbations spatially averaged over the city station groups and coastal station groups shown in Fig. 1, e) to h) and i) to l) respectively, 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).

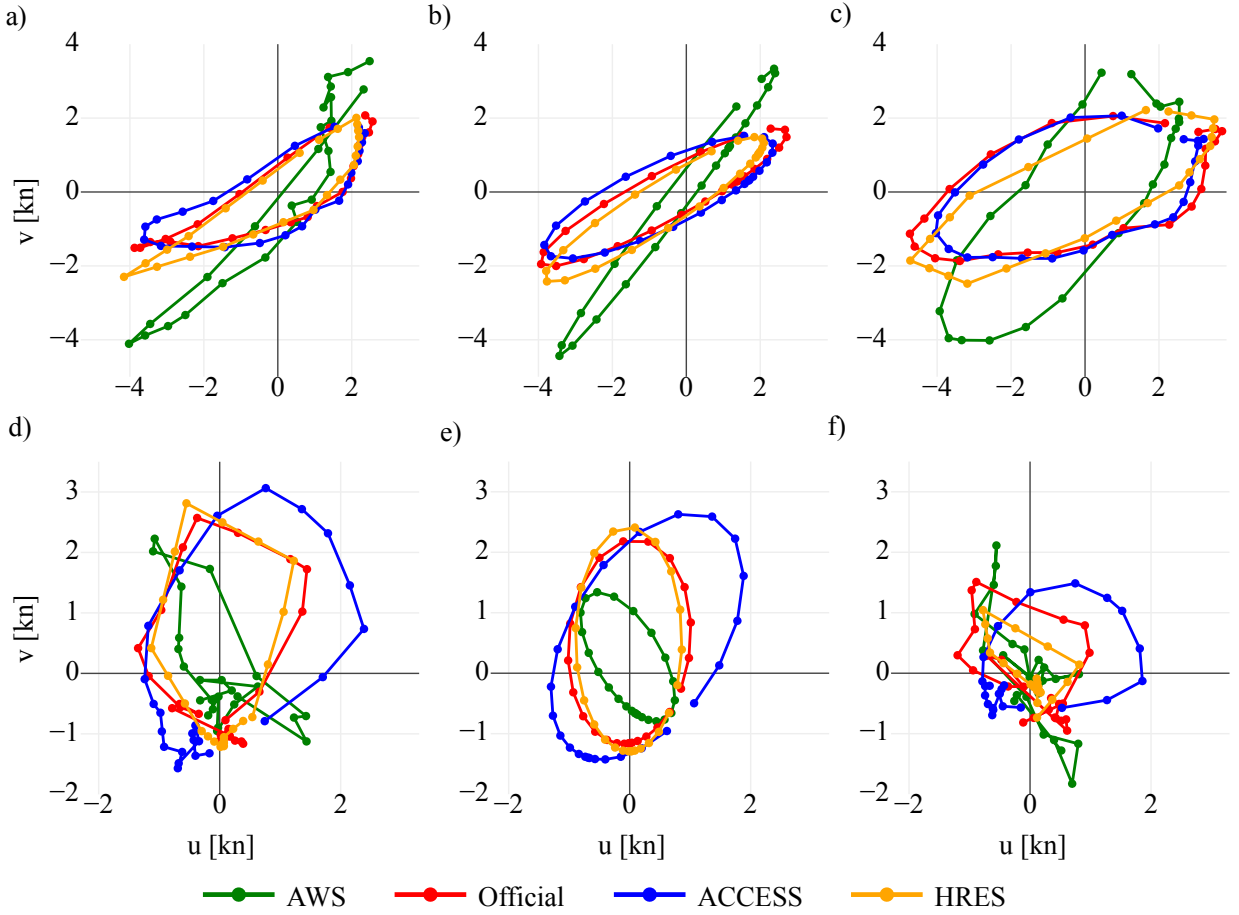


FIG. 13. Temporal hodographs of wind perturbations at each hour UTC averaged over June, July and August 2018, at Brisbane and Hobart airports, a) and d), and the associated ellipse fits, b) and e). For comparison, c) and f) provide the hodographs of the averaged perturbations at the Spitfire Channel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

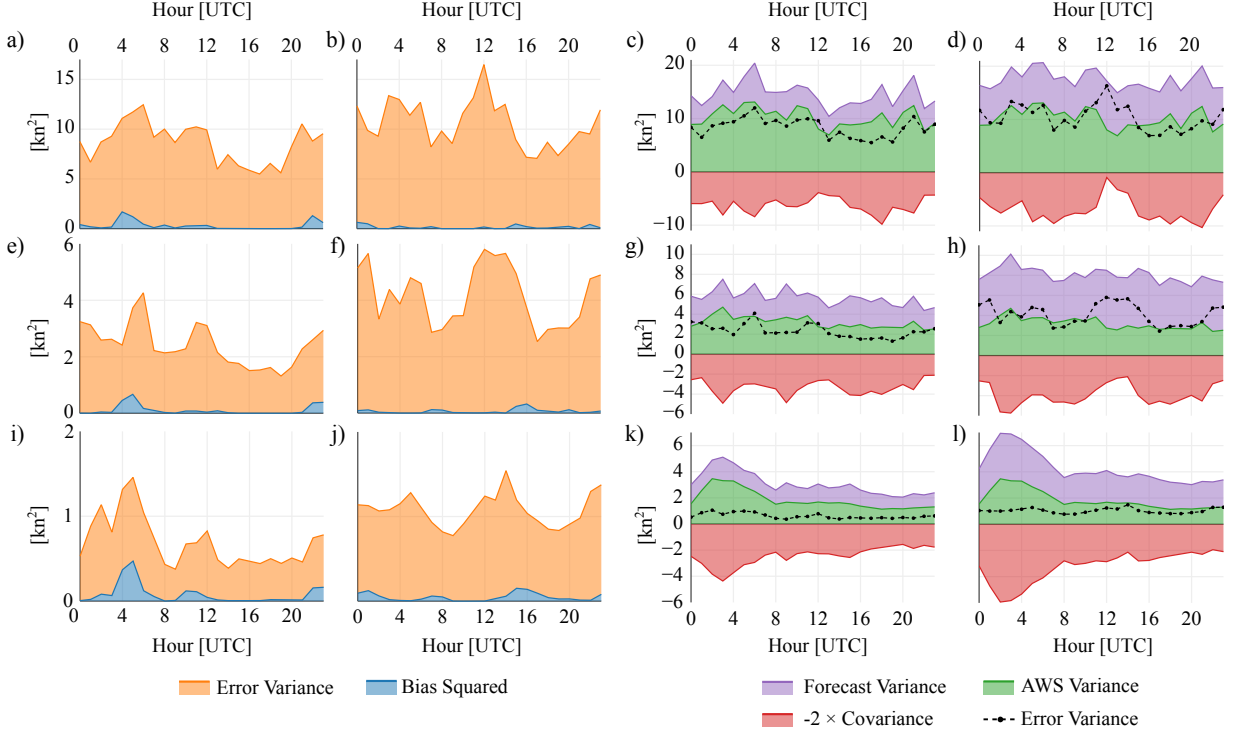


FIG. 14. Mean square error between the AWS and HRES zonal perturbations $(u_{AWS} - u_H)^2$, a), e), and i), decomposed into the error variance $\text{var}(u_{AWS} - u_H)$ and squared bias $(\bar{u}_{AWS} - \bar{u}_H)^2$ terms of equation (8). Also, the decomposed mean square error between the AWS and Official zonal perturbations, b), f) and j). Additionally, the HRES and AWS error variance term $\text{var}(u_{AWS} - u_H)$ decomposed into the $\text{var}(u_{AWS})$, $\text{var}(u_H)$ and $-2 \cdot \text{cov}(u_{AWS}, u_H)$ terms, c), g) and k), and analogously for the Official and AWS error variance term $\text{var}(u_{AWS} - u_O)$, 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.)