

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 Forecasters working for Australia’s Bureau of Meteorology (BoM) produce
10 a seven day forecast in two key steps: first they choose a model guidance
11 dataset to base the forecast on, then they use graphical software to manually
12 edit this data. Two types of edits are commonly made to the wind fields that
13 aim to improve how the influences of boundary layer mixing and land-sea
14 breeze processes are represented in the forecast. In this study I compare the
15 diurnally varying component of the BoM’s official wind forecast, with that of
16 station observations and unedited model guidance datasets. I consider coastal
17 locations across Australia over June, July and August 2018, aggregating data
18 over three spatial scales. The edited forecast generally only produces a lower
19 mean absolute error than model guidance at the coarsest spatial scale (over
20 fifty thousand square kilometres), but can achieve lower seasonal biases over
21 all spatial scales. However, the edited forecast only reduces errors or biases
22 at particular times and locations, and rarely produces lower errors or biases
23 than all model guidance products simultaneously. To better understand phys-
24 ical reasons for biases in the mean diurnal wind cycles, I fit modified ellipses
25 to the seasonally averaged diurnal wind temporal hodographs. Biases in the
26 official forecast diurnal cycle vary with location for multiple reasons, includ-
27 ing biases in the directions sea-breezes approach coastlines, amplitude biases,
28 and disagreement in the relative contribution of sea-breeze and boundary layer
29 mixing processes to the mean diurnal cycle.

30 1. Introduction

31 Modern weather forecasts are typically produced by models in conjunction with human forecast-
32 ers. Operational forecasters working for the Australian Bureau of Meteorology (BoM) undertake
33 two key steps to construct a seven day forecast.

34 First, they choose a *model guidance* dataset on which to base the official forecast. Datasets from
35 both the BoM and international modelling centres are available to Australia forecasters, with the
36 BoM’s Operational Consensus Forecast (OCF) an increasingly common choice. In the second step,
37 the forecaster uses GFE to *manually edit* the model guidance data. Such edits aim to incorporate
38 processes that are under-resolved at the resolutions of the model guidance products, or to correct
39 for perceived biases of the model guidance being used. Forecasters working for the United States
40 National Weather Service also use GFE, and utilise a similar approach.

41 Australian forecasters regularly make two types of edits to the surface wind fields. The first
42 involves modifying the surface winds after sunrise at locations where the forecaster believes the
43 model guidance is providing a poor representation of boundary layer mixing processes. Boundary
44 layer mixing occurs as the land surface heats up, producing an unstable boundary layer which
45 transports momentum downward to the surface layer. Before this mixing occurs, winds are typi-
46 cally both weaker and ageostrophically oriented due to surface friction (Lee 2018), and so mixing
47 can affect both the speed and direction of the surface winds. Australian forecasters perform bound-
48 ary layer mixing edits using a GFE tool which allows them to specify a region over which to apply
49 the edit, a height z and a percentage p , with the tool then calculating a weighted average of the
50 surface winds and winds at z , weighted by p .

51 The second type of edit involves changing the afternoon and evening surface winds around those
52 coastlines where the forecaster believes the model guidance is resolving the sea-breeze poorly.

53 Similarly to with boundary layer mixing, these edits are performed using a GFE tool that allows
54 forecasters to trace out the relevant coastline graphically, choose a wind speed and a time, with
55 the tool then smoothly blending in winds of the given speed perpendicular to the traced coastline
56 at the given time. In Australia, the official gridded forecast datasets resulting from a forecaster's
57 choice of model guidance and subsequent edits are then provided to the public through the BoM's
58 online MetEye data browser (Bureau of Meteorology 2019b), and are also translated into text and
59 icon forecasts algorithmically.

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

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

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

84 The present study has two goals. First, to describe a method for comparing the diurnal wind
85 signals of human edited forecasts to those of unedited model guidance forecasts, in order to assess
86 where and when human choice of model guidance and edits produce a reduction in error or bias.
87 Second, to apply this methodology across Australian coastal locations. The remainder of this paper
88 is organised as follows. Section 2 describes the methodology, and datasets to which it is applied,
89 section 3 provides results, and sections 4 and 5 provide a synthesis and conclusion, respectively.

90 **2. Data and Methods**

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

95 *a. Data*

96 Five datasets are considered in this study (Bureau of Meteorology 2019a); the human edited
97 official BoM wind forecast data that is issued to the public, observational data from automatic
98 weather stations (AWS) across Australia, unedited data from the ECMWF’s high resolution 10-day

99 forecast model (HRES), unedited data from the operational Australian Community Climate and
100 Earth System Simulator (ACCESS) regional model, and gridded Operational Consensus Forecast
101 (OCF) data, which blends output from multiple operational models. HRES, ACCESS and OCF
102 are three of the model guidance products commonly used by Australian forecasters for winds. I
103 consider just the lead-day one forecasts of the official forecast, HRES, ACCESS and OCF, for
104 reasons discussed below.

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

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

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

121 The BoM's gridded Operational Consensus Forecast (OCF) is based on the work of Woodcock
122 and Engel (2005) and Engel and Ebert (2007). OCF first corrects biases in model data, then

123 forms a weighted average of an ensemble of models in a way that minimises error with recent
124 observations. The methodology was expanded by the BoM in order to produce gridded datasets
125 that could be used by forecasters within the GFE, with 10 m horizontal winds added in June 2012
126 (Bureau of Meteorology 2005, 2008, 2012). For the time period of this study, the OCF ensemble
127 was comprised of the ACCESS and HRES datasets described above, and 5 other model datasets
128 (Bureau of Meteorology 2018).

129 To form a consensus wind forecast, OCF works with wind speed and direction, as taking av-
130 erages of u and v wind components can suppress wind speeds (Glahn and Lowry 1972), and this
131 is viewed as undesirable in an operational context. Speeds are calculated from each ensemble
132 member, bias corrected, then a weighted average calculated, with weights chosen based on the
133 performance of each member over the previous 20 days. Consensus wind direction is chosen as
134 the median wind direction from the members (Bureau of Meteorology 2012). Because data from
135 some members are only provided to the BoM at 3 hourly time intervals, interpolation and post-
136 processing is applied to produce an hourly OCF dataset that forecasters can use in GFE (Bureau of
137 Meteorology 2008). Gridded OCF is an objective alternative to the forecaster's subjective choice
138 of model guidance. When the overall wind field is assessed at six hourly intervals, gridded OCF
139 produces lower errors in both wind speed and direction than all the model guidance products that
140 comprise it (Bureau of Meteorology 2012).

141 The Bureau's official forecast dataset is produced on a state by state basis at forecasting centres
142 located in most state capitals. To construct the official forecast dataset, forecasters make a choice
143 of model guidance in the GFE, which then downscales the model data, or in the case of high spatial
144 resolution mesoscale model guidance, upscales the model data, onto a standard 3 km spatial grid
145 for Victoria and Tasmania, or a 6 km grid for the rest of the country. GFE displays model data
146 at hourly intervals by taking the model guidance output at each hour UTC. An exception is the

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

150 I therefore compare the official forecast and model guidance datasets as they appear in the GFE,
151 i.e. I compare the upscaled or downscaled datasets on the standardised 3 or 6 km, hourly grids.
152 This both ensures a consistent comparison between model guidance products of different spatial
153 resolutions, and an assessment of how the official forecast compares to the model guidance prod-
154 ucts as they actually appear to forecasters in the GFE. This is the standard approach the BoM takes
155 when comparing the performance of the official forecast to unedited model guidance (e.g. Griffiths
156 et al. 2017).

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

162 *b. Assessing Diurnal Variability*

163 Forecasters edit model guidance wind data to account for under-resolved sea-breeze and bound-
164 ary layer mixing processes. Instead of attempting to assess each type of edit individually, I study
165 the overall diurnal signal by subtracting a twenty four hour centred running mean *background*
166 *wind* from each zonal and meridional hourly wind data point, to create wind *perturbation* datasets.
167 Because records are not kept as to which model guidance product was used for the official fore-
168 cast on a given day, nor of what kinds of edits where performed, I compare the official forecast

169 on a pairwise basis with three unedited model guidance datasets commonly used by Australian
170 forecasters for winds, ACCESS, HRES and OCF.

171 The first metric I consider is the *difference of absolute errors* (DAE) in the perturbations, with
172 Fig. 1 illustrating how DAE is calculated. To compare errors in the diurnal signals of the official
173 forecast and model guidance, I calculate the Euclidean distances between the official or model
174 guidance perturbation vectors at each hour UTC, and the corresponding AWS perturbation vectors
175 at each hour UTC, and take their difference, viewing the Euclidean distance as a measure of
176 absolute error.

177 For example, to assess whether the official forecast perturbations, \mathbf{u}_O , or model guidance pertur-
178 bations, \mathbf{u}_M , produce lower absolute errors when compared with the observed AWS perturbations,
179 \mathbf{u}_{AWS} , I calculate

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

180 I then calculate statistics from the DAE values on an hourly basis, in particular, I calculate the
181 arithmetic mean of all the 00:00 UTC DAE values, denoting such an average by $\overline{\text{DAE}}$, and repeat
182 this for each hour of the day. If $\overline{\text{DAE}} > 0$ at a particular hour, then the official forecast perturbations
183 at that hour are, on average, closer to the observed perturbations than model guidance, and vice
184 versa if $\overline{\text{DAE}} < 0$.

185 Diurnal processes like the sea-breeze and boundary layer mixing depend on the background
186 atmospheric conditions in which they occur. By comparing wind perturbations rather than the
187 overall wind fields I am not claiming these background conditions are irrelevant to these processes.
188 However, when a forecaster makes an edit of a wind forecast to better resolve these processes, they
189 are implicitly assuming that future background conditions will be close enough to the preceding
190 24 hour mean state, or to model predictions of the mean state, to justify making the edit. Thus, it
191 makes sense to compare forecast perturbations to observed perturbations, as long as differences are

192 interpreted as a consequence not only of how the forecaster or model resolves diurnal processes,
193 but of how differences in the background state contribute to differences in the perturbations. To
194 minimise the importance of background state differences, this study focuses exclusively on lead-
195 day one forecasts.

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

204 Although temporal autocorrelations of DAE, i.e. correlations between DAE values at a particular
205 hour from one day to the next, are in practice small or non-existent, they are still accounted for
206 by reducing the “effective” sample size to $n(1 - \rho_1) / (1 + \rho_1)$, where n is the actual sample size
207 and ρ_1 is the lag-1 autocorrelation (Zwiers and von Storch 1995; Wilks 2011). In the language of
208 statistical hypothesis testing, the null hypothesis that $E = 0$ would be rejected at significance level
209 α if $\Pr(E > 0) > 1 - \frac{\alpha}{2}$ or $\Pr(E < 0) > 1 - \frac{\alpha}{2}$. However, in this study I simply state the value of
210 $\Pr(E > 0)$, referring to this as a *confidence score*, and noting $\Pr(E < 0) = 1 - \Pr(E > 0)$. I say the
211 official forecast outperforms model guidance with “high confidence” if $\Pr(E > 0) \geq 95\%$, or that
212 model guidance outperforms the official forecast with “high confidence” if $\Pr(E > 0) \leq 5\%$, with
213 high confidence implicit whenever it is not explicitly mentioned.

214 Following the “fuzzy verification” approach outlined by Ebert (2008), forecast and observational
215 perturbation datasets are compared not only at individual stations, but are also averaged over two

216 coarser spatial scales before being compared. The individual stations I consider are the 7 capital
 217 city *airport stations*, marked by stars in Fig. 2, as their high operational significance means that
 218 they are typically the most well maintained. An intermediate spatial scale is formed by averag-
 219 ing perturbation data over the 10 stations closest to each capital city airport station, with some
 220 flexibility allowed to ensure stations are roughly parallel to the nearest coastline. These station
 221 groups are referred to as the *city station groups*. The coarsest spatial scale is formed by averaging
 222 over all stations within 100 km of the nearest coastline, and grouping these by state. The West-
 223 ern Australian coastline (see Fig. 2) is subdivided into three pieces, and stations along the Gulf
 224 of Carpentaria, north Queensland Peninsula, and Tasmanian coastlines are neglected, in order to
 225 ensure each station group corresponds to an approximately linear segment of coastline to better re-
 226 solve the land-sea breeze after spatial averaging (e.g. Vincent and Lane 2016). These eight station
 227 groups are referred to as the *coastal station groups*.

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

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

$$DB = |\bar{\mathbf{u}}_{AWS} - \bar{\mathbf{u}}_O| - |\bar{\mathbf{u}}_{AWS} - \bar{\mathbf{u}}_M|, \quad (2)$$

235 where the over-bars denote temporal averages of the perturbations at a particular hour, over June,
 236 July and August 2018. These temporally averaged perturbations can be viewed as the mean diurnal
 237 wind cycle over the three month study period for each dataset. Biases over the city and coastal sta-

tion groups are calculated by taking the spatial average before the temporal average. Uncertainty in the DB is estimated through bootstrapping (Efron 1979). This is done by performing resampling with replacement on the underlying perturbation datasets, and calculating the DB 1000 times using these resampled datasets. This provides a distribution of DB values, which analogously to with DAE, I treat as a sample from a random variable B , and use this to estimate $\Pr(B > 0)$.

Note that on a given day, at a given location, wind perturbations do not necessarily reflect genuinely diurnal processes. There is a large degree of random turbulence in AWS wind observations, and convective cold pools or synoptic fronts can produce rapid changes in background winds that induce large perturbations. However, averaging multiple perturbations at a given hour over many days cancels out much of the diurnal variability not genuinely associated with diurnal processes. When this is repeated for each hour of the day, the signal that remains reflects the mean diurnal cycle (e.g. Figs. 10 and 11). Similar ideas apply to the DAE metric. Note that spatially averaging perturbations accomplishes a similar thing to temporal averaging, helping to cancel out random variability. These ideas can be explored with synthetic data, and some preliminary work to this end is available online (Short 2020).

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

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

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

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. 11). Thus, I 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 mean diurnal wind cycle, I solve for the seven parameters $u_0, u_1, u_2, v_0, v_1, v_2$ and ψ using non-linear regression.

Importantly, the metrics defined in this section compare just *some aspects* of the official forecast with model guidance: they do not, for instance, assess whether diurnal variance of the official forecast is more realistic than that of model guidance. Thus, any statements about performance made throughout this paper refer solely to the metrics defined here, 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. Furthermore, comparing results at different locations is *not* intended as a “ranking” of forecasting centres in different states because, for instance, station density varies significantly with location so it is hard to define station groups at a given spatial scale in a completely consistent way across locations.

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, mean 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. Unless otherwise noted, times are given in UTC.

a. Absolute Errors

Figure 3 provides the mean difference of absolute error values and confidence scores defined in section 2 for the coastal station groups shown in Fig. 2. Results are given for the official forecast versus ACCESS, official forecast versus HRES, and official forecast versus OCF comparisons. The results indicate that for the majority of station groups and hours, the unedited ACCESS, HRES and OCF datasets outperform the official forecast. The lowest \overline{DAE} values occur at the Northern Territory (NT) station group at 23:00 and 00:00 for both the official forecast versus ACCESS, and official forecast versus HRES comparisons, and at 22:00 and 23:00 for the official forecast versus OCF comparison. Although the official forecast outperforms at least one of ACCESS, HRES and OCF at multiple times and station groups, the only group and time where it outperforms all three is 05:00 UTC over the South Western Australia (WA) station group.

300 Figures 4 and 5 provide case studies of the Northern Territory (NT) and South Western Australia
301 (WA) station groups, respectively. Figure 4 a) provides a time series of DAE for the NT station
302 group at 23:00. The time series shows significant temporal variability, with DAE frequently dropping
303 below -2 kt. Figures 4 b) and c) show hodographs of the winds and wind perturbations,
304 respectively, at each hour UTC on the 3rd of July, which provides an interesting example. Note
305 that care must be taken when interpreting perturbations and DAE scores on individual days physically,
306 as discussed in section 2.

307 Figure 4 b) shows that the official wind forecast on this day was likely based on edited ACCESS
308 from 00:00 to 06:00, then edited HRES from 07:00 to 13:00 UTC, then unedited ACCESS from
309 15:00 to 21:00. At 22:00 and 23:00, the official forecast winds acquire stronger east-southeasterly
310 components than the other datasets. For comparison, Fig. 6 a) shows the first ten values from
311 wind soundings at Darwin Airport at 12:00 on July 3rd and 00:00 on July 4th. In both instances
312 the winds are east-southeasterly, and so the rapidly changing wind perturbations at 22:00 in the
313 official forecast may reflect a boundary layer mixing edit that has been applied either too early,
314 or has strengthened the southeasterly component of the winds too much. Similar issues appear to
315 create the low DAE values on the 8th of June and 9th and 10th of July.

316 Figure 5 a) provides a time series of DAE for the South WA station group at 05:00. As with
317 the NT station group there is significant temporal variability, with DAE frequently exceeding 1 kt.
318 Figures 5 b) and c) provide hodographs of the winds and wind perturbations, respectively, on the
319 9th of June, another interesting example. Both the raw winds and the perturbations appear to show
320 both HRES and ACCESS under-predicting the amplitude of the diurnal wind cycle on this day,
321 with OCF performing better in this regard. Figure 6 b) shows wind soundings at Perth Airport,
322 the nearest station to provide wind soundings, between 12:00 on the 8th June and 12:00 on the
323 9th June. The 8th June 12:00 sounding shows surface northerlies of around 6 kt, becoming west

324 to northwesterlies of over 20 kt 2.4 km above the surface. However, the subsequent sounding at
325 00:00 on the 9th of June shows that the winds acquire a strong northerly component of 30 kt in the
326 first 500 m of the atmosphere, with the final sounding indicating a strong northwesterly wind at
327 725 m persisting until 12:00.

328 In Fig. 5 c), the OCF and official forecast perturbations from 04:00 to 07:00 show stronger
329 westerly perturbations than either ACCESS or HRES, improving the magnitude of both dataset's
330 perturbations. However, the AWS perturbations are more northerly than those of the official fore-
331 cast or OCF. Possible explanations for this discrepancy are that the official forecast has been edited
332 based on the June 8th 12:00 sounding, with the winds above the surface changing direction in the
333 subsequent 12 hours, or that the official forecast has been based on OCF, which underestimates
334 the northerly component of the perturbations.

335 Fig. 7 presents the $\overline{\text{DAE}}$ values and confidence scores for the city station groups, for the official
336 forecast versus HRES and official forecast versus OCF comparisons; the official forecast versus
337 ACCESS comparisons (not shown) are similar to those for HRES and have been omitted for space.
338 Both HRES and OCF outperform the official forecast almost uniformly, with the Darwin city sta-
339 tion group the main exception. At Darwin, the official forecast outperforms both HRES and OCF
340 at 02:00 UTC, and there is ambiguity at some other times of day. The OCF comparison shows
341 less ambiguity at Darwin, but more at Melbourne and Brisbane. The city station group results for
342 December, January, February 2017/18 (not shown) are similar but slightly more ambiguous, par-
343 ticularly for ACCESS. These results were replicated using alternative city station groups, defined
344 by taking all stations within $100 \text{ km} \times 100 \text{ km}$ boxes centred on each capital city airport: the
345 results (not shown) were very similar, with both HRES and OCF almost uniformly outperforming
346 the official forecast.

Fig. 8 presents the comparisons for the airport stations. Here the results are noisier than at both the city and coastal spatial scales, but similarities also exist. For instance, the official forecast outperforms both OCF and HRES at 02:00 at Darwin Airport, the Darwin city station group, and the NT coastal station group with at least 90% confidence. There are four other instances where the official forecast outperforms HRES with at least 90% confidence, although this could simply be occurring by chance due repeated testing (Wilks 2011, p. 178). By contrast, the official forecast outperforms OCF over four hour intervals at both Perth and Brisbane airports.

b. Seasonal Biases

Figure 9 provides the difference of biases (DB) and confidence scores defined in section 2, for the coastal station groups, for the official forecast versus ACCESS, official forecast versus HRES, and official forecast versus OCF comparisons. At the NT station group at 03:00, the official forecast outperforms both ACCESS and HRES with confidence $\geq 93\%$. However, ACCESS, HRES and OCF each outperform the official forecast at 23:00 and 00:00, and from 06:00 to 10:00, consistent with the $\overline{\text{DAE}}$ results of Fig. 3. Figure 10 c) shows that these DB results reflect amplitude biases in the official forecast's mean diurnal cycle.

At the South WA station group from 01:00 to 05:00, the official forecast outperforms HRES with confidence scores of at least 88%. Figure 11 a) shows that HRES underestimates the westerly perturbations at these times, with these perturbations potentially associated with boundary layer mixing processes, as discussed in section 3 a. The official forecast, ACCESS and HRES all underestimate the amplitude of the diurnal cycle between 02:00 and 10:00, including both the westerly perturbations and the southerly sea-breeze perturbations. OCF better approximates the amplitude of the diurnal cycle between 02:00 and 05:00, but shows the greatest underestimation of the southerly perturbations between 06:00 and 10:00.

370 At the South Australia (SA) station group, the official forecast slightly outperforms ACCESS
371 and HRES from 02:00 to 05:00 and 09:00 to 12:00, although confidence scores do not exceed
372 64% and 90% respectively. The official forecast also slightly outperforms OCF between 00:00
373 and 02:00, and between 08:00 and 09:00, although confidence scores do not exceed 74%. Figure
374 11 b) shows that although the official forecast captures the amplitude of the perturbations from
375 01:00 to 05:00 almost perfectly, its mean diurnal cycle is out of phase with that of AWS during
376 this period, explaining the only slightly positive DB values.

377 For comparison, Figs. 12 and 13 present the DB values and confidence scores for the official
378 forecast versus HRES and official forecast versus OCF comparisons, for the city station groups
379 and airport stations, respectively. Some regions exhibit consistent results across all three spatial
380 scales. For example, the official forecast outperforms HRES between 14:00 and 18:00, with at
381 least 83% confidence, at Sydney Airport, the Sydney city station group, and the NSW coastal
382 station group.

383 Other results are markedly different between spatial scales. For instance, the official forecast
384 outperforms OCF for most of the day at Darwin Airport, but the opposite is true at the Darwin
385 city and NT coastal station groups. Figure 10 a) shows that the mean AWS diurnal cycle is highly
386 asymmetric, with a sharp peak occurring at 06:00. This peak is captured well by HRES and the
387 official forecast, but not by OCF or ACCESS. Figures 10 b) and c) show that over the Darwin city
388 and NT coastal station groups, the mean diurnal cycles are much smoother, with the amplitudes of
389 the official forecast diurnal cycles exaggerated relative to AWS and OCF.

390 *c. Ellipse Fits*

391 The hodographs in Figs. 10 and 11 are roughly elliptical in shape, suggesting that descriptive
392 quantities can be estimated by fitting equations (5) and (6) to the zonal and meridional mean

393 perturbations, as described in section 2. Figure 14 gives the R^2 values for the fits of the zonal and
394 meridional perturbations to equations (5) and (6), respectively. The fit performs best at the coastal
395 station group spatial scale, with R^2 generally above 95%.

396 Figure 15 provides four descriptive quantities based on the fits of equations (5) and (6) to the
397 mean perturbations: these are maximum perturbation speed, eccentricity of the fitted ellipse, angle
398 the semi-major axis makes with lines of latitude, and the time at which the maximum perturbation
399 speed is achieved.

400 Figure 15 a) shows OCF has mean diurnal cycle amplitude biases at the airport station scale, with
401 the exception of Hobart. These biases persist, but are smaller, at the city station group scale, but
402 are absent at the coastal station group scale, with the exception of Queensland (QLD). Given that
403 OCF represents a blended average of multiple model guidance datasets (Engel and Ebert 2007),
404 and that OCF's gridding process involves additional interpolation steps (Bureau of Meteorology
405 2008, 2012), this result is perhaps not surprising: at the individual station scale OCF has undergone
406 more smoothing than ACCESS or HRES, but at the coarser spatial scales this lessens in importance
407 as all datasets undergo comparable smoothing. Note that this does *not* mean OCF's overall wind
408 speeds or directions are biased at the individual station scale, only the amplitude of OCF's mean
409 diurnal cycle, subject to how mean diurnal cycles are treated in this study.

410 Considering specific locations, Brisbane provides an interesting example, as Fig. 15 a) shows
411 that at Brisbane Airport the maximum AWS perturbation is at least 1 kt greater than the official
412 forecast, ACCESS and HRES, and 3.5 kt greater than that of OCF. Furthermore Fig. 15 c) shows
413 that the orientation of the AWS fitted ellipse is at least 20 degrees anti-clockwise from that of the
414 other datasets.

415 Figures 16 a) and b) show hodographs of the Brisbane Airport mean perturbations and ellipse
416 fits, respectively. Although the ellipse fits suppress some of the asymmetric details, they capture

417 the amplitudes and orientations of the real mean diurnal cycles well. In this case the results show
418 that the mean AWS sea-breeze approaches from the northeast, whereas the official forecast, HRES,
419 ACCESS and OCF sea-breezes approach more from the east-northeast. The amplitude of OCF's
420 mean diurnal cycle is significantly weaker than those of the other datasets.

421 To check whether these results just represent a direction bias of the Brisbane Airport weather
422 station, Fig. 16 c) shows the mean diurnal cycle at the nearby Spitfire Channel station (see Fig. 2).
423 While the amplitude biases are slightly smaller at Spitfire Channel than Brisbane Airport, the
424 directional bias is at least as high. A similar directional bias is evident at the nearby Inner Beacon
425 station (not shown), although the bias is smaller than at Spitfire Channel and Brisbane Airport.
426 Similar biases are also evident at these stations in analogous figures for December, January and
427 February 2017/18 (not shown), with the semi-major axis of the official forecast's ellipse fit oriented
428 29° clockwise from AWS's at Brisbane Airport. Figure 2 shows there are two small islands to the
429 east of Brisbane Airport; the more north-northeasterly orientation of the Brisbane Airport sea-
430 breeze suggests these islands may be redirecting winds between the east coast of Brisbane and the
431 west coasts of these islands, and that this local effect is not being captured in the official forecast,
432 ACCESS, HRES or OCF.

433 The South WA station group provides another interesting example, as Fig. 15 shows the semi-
434 major axes of the ACCESS, OCF and official forecast ellipse fits are oriented at least 48 degrees
435 anti-clockwise from those of the AWS and HRES ellipse fits, and the HRES perturbations peak
436 between 1.2 and 4 hours after the other datasets. Figure 11 a) shows that these differences occur
437 because the westerly perturbations, potentially associated with boundary layer mixing, are weaker
438 for HRES than for the other datasets, resulting in HRES's semimajor axis being oriented more
439 meridionally. Analogously, the southerly perturbations, potentially associated with the sea-breeze,
440 are stronger for AWS than the other datasets, with a similar effect on orientation and timing as

with HRES. Similar points can be made for the Victorian (VIC) and NT coastal station groups, and at Darwin Airport.

4. Synthesis

For land-sea breeze and boundary layer mixing edits to reduce absolute errors in the subsequent day's wind forecast, these edits should reduce the absolute errors in the diurnal component of the wind fields. However, Figs. 3, 7 and 8 indicate that this is generally only possible when absolute error is considered at coarse spatial scales, as at individual airport stations results are generally noisy and ambiguous, and over the intermediate city station group scale model guidance outperforms the official forecast almost uniformly.

Taking the effective resolutions of the models considered in this study to be approximately $7\Delta x$ (e.g. Skamarock 2004; Abdalla et al. 2013), where Δx is the horizontal grid spacing, the effective resolutions of ACCESS and HRES are ≈ 84 km and ≈ 63 km, respectively. From resolution considerations alone, one might expect that forecaster edits would be able to reduce errors at the individual airport station scale, and the intermediate city station group scale (see Fig. 2), as motion at these scales is unresolved or only partially resolved by ACCESS and HRES.

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

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

where var, cov and the over-bar denote the sample variance, covariance and mean respectively.

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

square of the bias between u_{AWS} and u_O . Equation (8) can also be applied to the mean square errors (MSEs) of ACCESS, HRES and OCF. Note that the MSE is closely related to \overline{DAE} and the squared bias components of the MSEs are closely related to DB.

Figure 17 shows the terms of equation (8) for both the official forecast and OCF, for Brisbane Airport, the Brisbane city station group, and the QLD coastal station group. At all three scales the official forecast varies more than OCF. The official forecast also generally varies more than ACCESS and HRES (not shown), and this is also true for the other stations and station groups considered in this study.

At Brisbane Airport the variance of AWS is significantly larger than either the official forecast or OCF. This additional variability is mostly uncorrelated to either dataset. Although the covariance between the official forecast and AWS increases between 20:00 and 08:00, the increase is not sufficient to offset the official forecast's additional variance, and the error variances are thus of comparable magnitude for both the official forecast and OCF.

The larger AWS variances are unsurprising from representation considerations alone (e.g. Zaron and Egbert 2006), as the official forecast and OCF data represent averages over 6 km spatial grid-cells, whereas the AWS data represent point values. As a result, error variance terms are generally much larger than the squared bias terms at this scale. The exception is OCF at 04:00, where the squared bias is ≈ 6 kt, while error variance is ≈ 15 kt. This results in a higher MSE for OCF than the official forecast around 04:00, consistent with the airport station \overline{DAE} results of Figs. 8 c) and d).

At the intermediate Brisbane city station group scale, the AWS variances are again larger than those of OCF, but of comparable magnitude to those of the official forecast, with the official forecast's additional variability again mostly uncorrelated to AWS. This results in larger error variance terms for the official forecast, consistent with OCFs almost complete outperformance

486 of the official forecast in Figs. 7 c) and d). However, OCF's squared bias terms remain larger
487 than the official forecast's, resulting in OCF's MSE slightly exceeding the official forecast's at
488 around 04:00. These results are consistent with Figs. 7 c) and d), where the official forecast slightly
489 outperforms OCF at 04:00 with a confidence score of 79%.

490 Over the coarse QLD coastal station group scale, variances in all three datasets are small enough
491 that the error variance terms are less dominant over the bias terms. Although the error variance of
492 the official forecast is still larger than that of OCF, OCF's zonal biases around 04:00 UTC are again
493 sufficient to result in larger MSEs around this time. When considered with the analogous plots for
494 the meridional perturbations (not shown), for which OCFs squared bias terms peak slightly later,
495 the results are consistent with Figs. 3 c) and d).

496 Analogous points can be made for the other locations and datasets considered in this study. At
497 the airport station scale, AWS variance is generally significantly higher than that of the official
498 forecast and model guidance, producing high error variance and likely explaining why the airport
499 station DAE results of Fig. 8 are comparatively noisier than those of the city or coastal station
500 group scales. Interesting exceptions include OCF at Brisbane and Perth airports, where amplitude
501 biases in OCF's diurnal cycle are sufficient to affect airport station DAE scores.

502 At the city station group scale, the official forecast is generally outperformed by HRES and
503 OCF in the $\overline{\text{DAE}}$ results of Fig. 7, and in the analogous comparisons with ACCESS (not shown).
504 This occurs because the official forecast is generally more variable than model guidance, and
505 this additional variability is mostly random, in the sense of being uncorrelated with AWS. At the
506 coastal station group scale, random variability in each dataset is reduced, and biases are sufficiently
507 large relative to error variance to affect the $\overline{\text{DAE}}$ results of Fig. 3.

508 These results suggest that switching model guidance products or performing edits can add more
509 random noise to the diurnal component of the official forecast than what can be offset by reductions

510 in bias, or improved correlations with AWS. Because the official forecast is built from multiple
511 model guidance datasets, switching between datasets with different means will tend to produce
512 greater variance than any of the component datasets. If the choice of model guidance is made
513 primarily on which model best captures more slowly evolving, larger amplitude synoptic scale
514 features, then switching model guidance may add random variability to the diurnal component of
515 the official forecast. Furthermore, unless all forecasters follow identical thought processes when
516 making edits, the edits will also add random variability.

517 These results could have implications for forecasting practice. Model guidance products are
518 indeed biased in how they resolve diurnal wind cycles (e.g. Fig. 16), and there is therefore scope
519 for forecaster edits to reduce these biases. However, editing model guidance generally fails to
520 reduce error in the forecast diurnal signal, even at scales finer than the effective resolutions of the
521 models, as at these scales diurnal cycles are significantly masked by random variability. Averaging
522 over large areas reduces this random variability, better revealing the diurnal cycle, and so biases
523 have a greater impact on forecast error. However, even at large scales Fig. 3 shows model guidance
524 still outperforms the official forecast more often than not.

525 Reducing the random variability of the official forecast, or the model guidance datasets that
526 comprise it, could therefore improve the capacity of these types of edits to reduce error in the
527 diurnal cycle. One way to accomplish this would be to use an ensemble average model guidance
528 product like OCF, another would be to further post process model guidance products, such as by
529 averaging multiple time steps around the hour, before including them in the GFE.

530 **5. Conclusion**

531 In this study I have presented methods for assessing the diurnal component of wind forecasts,
532 with the intended application being the assessment of the edits Australian forecasters make to

533 model guidance datasets to better resolve land-sea breeze and boundary layer mixing processes.
534 I considered both errors and seasonal biases at each hour UTC, over three spatial scales, but the
535 methods are immediately generalisable to other spatiotemporal scales. Crucially, the results of
536 this study depend on the metrics and methods chosen, and no claim is being made that these are
537 sufficient to completely describe the overall accuracy, or value to the user, of competing forecasts.

538 When the methods are applied to Australian forecast data, the results indicate that the official
539 edited forecast only produces lower absolute errors in the diurnal wind signal when wind perturba-
540 tion data is averaged over the coarse “coastal station group” spatial scale (see Fig. 2) of 500×100
541 km^2 to $2000 \times 100 \text{ km}^2$. Even at these scales, reductions in error are isolated to particular loca-
542 tions and times of day, and the official forecast rarely has lower mean absolute error than the three
543 model guidance products considered in this study simultaneously.

544 By contrast, the official forecast can produce lower seasonal biases than model guidance at
545 all three spatial scales, but again, it rarely produces lower biases than the three model guidance
546 products considered here simultaneously. Reduced seasonal biases do not translate into reduced
547 errors at the two smaller spatial scales because the diurnal cycle is mostly masked by the random
548 variability in each dataset. Furthermore, because the official forecast generally exhibits much
549 greater random variability than model guidance, model guidance almost uniformly outperforms
550 the official forecast over the intermediate $50 \times 50 \text{ km}^2$ to $200 \times 200 \text{ km}^2$ city station group spatial
551 scale.

552 I also compare structural features of the mean diurnal wind cycles of each dataset by fitting
553 modified ellipses to their temporal hodographs, then deriving metrics from these ellipses. This ap-
554 proach reveals structural biases in the official forecast, including directional biases in the approach
555 of the sea-breeze at Brisbane Airport, and amplitude biases along the southwest coast of Western
556 Australia.

557 Future research could extend this study in multiple directions. One approach would be to study
558 how the difference of absolute errors (DAE) metric defined in this study responds to synthetic,
559 or idealised model data, so that the influence of random and synoptic variability can be better
560 understood: some preliminary work to this end is available online (Short 2020). Another important
561 question is whether the random variability in the official forecast, or the model guidance products
562 that comprise it, can be reduced through ensemble forecasting or post-processing, as reducing
563 random variability would both decrease errors, and could increase the value of land-sea breeze and
564 boundary layer mixing edits. The BoM's Operational Consensus Forecast (OCF) is an effective
565 way to accomplish this, and future work could assess whether it is possible, or desirable, to adjust
566 OCF's wind algorithm to reduce the amplitude biases identified in OCF's mean diurnal cycle,
567 noting that these biases are subject to how the mean diurnal cycle has been defined in this study.
568 Another goal could be to identify precisely the spatiotemporal scales at which diurnal wind cycles
569 can be separated from random variability, so as to better understand the scales at which land-sea
570 breeze and boundary layer mixing edits can reduce error in a forecast.

571 *Acknowledgments.* Funding for this study was provided for Ewan Short by the Australian Re-
572 search Council's Centre of Excellence for Climate Extremes (CE170100023). Datasets and soft-
573 ware were generously provided by the Australian Bureau of Meteorology's Evidence Targeted
574 Automation team, with additional code available online (Short 2019). Thanks are due to Michael
575 Foley, Deryn Griffiths, Nicholas Loveday, Ben Price and Alexei Hider for providing support at
576 the Bureau of Meteorology's Melbourne and Darwin offices, and to Craig Bishop, Todd Lane
577 and Claire Vincent from the University of Melbourne, and Carly Kovacik from the United States'
578 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 11 December 2019].
- Abkar, M., A. Sharifi, and F. Porté-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, 2005: Analysis and prediction operations bulletin no. 60. Tech. rep., Bureau of Meteorology, Melbourne, Victoria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/APOB74.pdf> - Accessed 4 February 2020].
- Bureau of Meteorology, 2008: Analysis and prediction operations bulletin no. 74. Tech. Rep. 74, Bureau of Meteorology, Melbourne, Victoria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/APOB74.pdf> - Accessed 4 February 2020].
- 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/nwp/doc/bulletins/apob83.pdf> - Accessed 11 December 2019].
- Bureau of Meteorology, 2012: NMOC operations bulletin number 91. Tech. Rep. 91, Bureau of Meteorology, Melbourne, Victoria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/apob91.pdf> - Accessed 4 February 2020].
- Bureau of Meteorology, 2016: APS2 upgrade to the ACCESS-R numerical weather prediction system. Tech. Rep. BNOC Operations Bulletin No. 104, Bureau of Meteorology, Melbourne, Victoria.

600 ria. [Available online at <http://www.bom.gov.au/australia/charts/bulletins/apob107-external.pdf>
601 - Accessed 11 December 2019].

602 Bureau of Meteorology, 2018: BNOC operations bulletin number 113. Tech. rep., Bureau of
603 Meteorology, Melbourne, Victoria. [Available online at [http://www.bom.gov.au/australia/charts/](http://www.bom.gov.au/australia/charts/bulletins/BNOC_Operations_Bulletin_113.pdf)
604 [bulletins/BNOC_Operations_Bulletin_113.pdf](http://www.bom.gov.au/australia/charts/bulletins/BNOC_Operations_Bulletin_113.pdf) - Accessed 4 February 2020].

605 Bureau of Meteorology, 2019a: Datasets used in “Verifying operational forecasts of land-sea
606 breeze and boundary layer mixing processes”. Zenodo, [Available online at [http://doi.org/10.](http://doi.org/10.5281/zenodo.3570002)
607 [5281/zenodo.3570002](http://doi.org/10.5281/zenodo.3570002) - Accessed 11 December 2019], doi:10.5281/zenodo.3570002.

608 Bureau of Meteorology, 2019b: Meteye. Bureau of Meteorology, [Available online at [http://www.](http://www.bom.gov.au/australia/meteye/)
609 [bom.gov.au/australia/meteye/](http://www.bom.gov.au/australia/meteye/) - Accessed 11 December 2019].

610 Dai, A., and C. Deser, 1999: Diurnal and semidiurnal variations in global surface wind
611 and divergence fields. *Journal of Geophysical Research*, **104**, 31 109–31 125, doi:10.1029/
612 1999JD900927.

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

615 Efron, B., 1979: Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7** (1),
616 1–26, doi:10.1214/aos/1176344552.

617 Engel, C., and E. Ebert, 2007: Performance of hourly operational consensus forecasts
618 (OCFs) in the Australian region. *Weather and Forecasting*, **22** (6), 1345–1359, doi:10.1175/
619 2007WAF2006104.1.

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

European Center for Medium Range Weather Forecasting, 2018: *Part IV: Physical processes*, 223. No. 4, IFS Documentation, European Center for Medium Range Weather Forecasting, [Available online at <https://www.ecmwf.int/node/18714> - Accessed 11 December 2019].

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

Glahn, H. R., and D. A. Lowry, 1972: The use of model output statistics (MOS) in objective weather forecasting. *Journal of Applied Meteorology*, **11** (8), 1203–1211, doi:10.1175/1520-0450(1972)011<1203:TUOMOS>2.0.CO;2.

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> - Accessed 11 December 2019].

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 boundary 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.

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

643 Mass, C. F., D. Ovens, K. Westrick, and B. A. Colle, 2002: Does increasing horizontal resolution
644 produce more skillful forecasts? *Bulletin of the American Meteorological Society*, **83** (3), 407–
645 430, doi:10.1175/1520-0477(2002)083<0407:DIHRPM>2.3.CO;2.

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

648 Modigliani, U., and C. Maass, 2017: Detailed information of implementation of IFS cy-
649 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)
650 [information+of+implementation+of+IFS+cycle+41r2](https://confluence.ecmwf.int/display/FCST/Detailed+information+of+implementation+of+IFS+cycle+41r2) - Accessed 11 December 2019].

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

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

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

658 Short, E., 2019: eshort0401/forecast_verification_paper. GitHub, [Available online at [https://](https://github.com/eshort0401/forecast_verification_paper)
659 github.com/eshort0401/forecast_verification_paper - Accessed 11 December 2019].

660 Short, E., 2020: DAE synthetic data tests. [Available online at [https://github.com/eshort0401/](https://github.com/eshort0401/forecast_verification_paper/blob/master/code/DAE%20Synthetic%20Data%20Tests.ipynb)
661 [forecast_verification_paper/blob/master/code/DAE%20Synthetic%20Data%20Tests.ipynb](https://github.com/eshort0401/forecast_verification_paper/blob/master/code/DAE%20Synthetic%20Data%20Tests.ipynb) -
662 Accessed 7 February 2020].

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

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

669 Vincent, C. L., and T. P. Lane, 2016: Evolution of the diurnal precipitation cycle with the passage
 670 of a Madden-Julian Oscillation event through the Maritime Continent. *Monthly Weather Review*,
 671 **144** (5), 1983–2005, doi:10.1175/MWR-D-15-0326.1.

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

674 Woodcock, F., and C. Engel, 2005: Operational consensus forecasts. *Weather and Forecasting*,
 675 **20** (1), 101–111, doi:10.1175/WAF-831.1.

676 Zaron, E. D., and G. D. Egbert, 2006: Estimating open-ocean barotropic tidal dissipation: The
 677 hawaiian ridge. *Journal of Physical Oceanography*, **36** (6), 1019–1035, doi:10.1175/JPO2878.
 678 1.

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

LIST OF FIGURES

681		
682	Fig. 1.	Illustration of the method for calculating the <i>difference of absolute errors</i> (DAE) in the
683		diurnal signals of an unedited model guidance dataset, and the human edited official forecast
684		dataset, when compared with automatic weather station (AWS) observations, at an example
685		time of 12:00 UTC. 34
686	Fig. 2.	Locations of the automatic weather stations, and the groupings of these stations, considered
687		in this study. The <i>coastal station groups</i> are indicated in a), with the <i>airport stations</i> shown
688		by stars. The Perth, Adelaide, Melbourne, Hobart, Darwin, Brisbane and Sydney <i>city station</i>
689		<i>groups</i> are shown shown by b) to h), respectively. 35
690	Fig. 3.	Heatmaps of mean difference of absolute error $\overline{\text{DAE}}$ values, a), c), e), and confidence scores,
691		b), d), f), for the <i>coastal station groups</i> (see Fig. 2). Results given for each hour of the day,
692		for the official forecast versus ACCESS, a) and b), official forecast versus HRES, c) and
693		d), and official forecast versus OCF, e) and f), comparisons. Positive $\overline{\text{DAE}}$ values indicate
694		that the former dataset in each pair is on average $\overline{\text{DAE}}$ kt closer to observations than the
695		latter dataset (see equation 1), where $1 \text{ kt} \approx 0.514 \text{ m s}^{-1}$. Confidence scores provide the
696		probability the population or “true” value of $\overline{\text{DAE}}$ is greater than zero (see section 2). . . . 36
697	Fig. 4.	Time series, a), of the difference in absolute error DAE defined in equation (1) for the
698		official forecast versus ACCESS, official forecast versus HRES, and official forecast versus
699		OCF comparisons, for the Northern Territory (NT) coastal station group shown in Fig. 2, at
700		23:00 UTC. Also, temporal hodographs in hours UTC showing hourly changes in winds, b),
701		and wind perturbations from a 24 hour running mean, c), at the NT coastal station group on
702		the 3 rd of July 2018. 37
703	Fig. 5.	As in Fig. 4, but for, a), the South Western Australia (WA) coastal station group at 05:00
704		UTC, and b) and c), the winds and wind perturbations, respectively, over the South WA
705		coastal station group on the 9 th June 2018. 38
706	Fig. 6.	Vertical wind soundings at, a), Darwin Airport, and b), Perth Airport, with heights given in
707		metres. 39
708	Fig. 7.	As in Fig. 3, but for the official forecast versus HRES comparison a) and b), and the official
709		forecast versus OCF comparison, c) and d), for the <i>city station groups</i> (see Fig. 2.) . . . 40
710	Fig. 8.	As in Fig. 3, but for the official forecast versus HRES comparison a) and b), and the official
711		forecast versus OCF comparison, c) and d), for the <i>airport stations</i> (see Fig. 2.) . . . 41
712	Fig. 9.	As in Fig. 3, but for the difference of biases (DB) values and confidence scores. . . . 42
713	Fig. 10.	Temporal hodographs in hours UTC of wind perturbations at, a), Darwin Airport, b), spa-
714		tially averaged over the Darwin city station group, and c), the NT coastal group (see Fig. 2),
715		then temporally averaged over June, July and August 2018. 43
716	Fig. 11.	Temporal hodographs in hours UTC of diurnal wind perturbations spatially averaged over
717		the, a), South Western Australia (WA), and b), South Australia (SA) coastal station groups
718		(see Fig. 2), and temporally averaged over June, July and August 2018. 44
719	Fig. 12.	As in Fig. 7, but for the difference of biases (DB) values and confidence scores. . . . 45
720	Fig. 13.	As in Fig. 8, but for the difference of biases (DB) values and confidence scores. . . . 46

721	Fig. 14.	R^2 values as percentages for the fit of equation (5) to the zonal perturbations, a), c) and e),	
722		and equation (6) to the meridional perturbations, b), d) and f), for the airport stations, a) and	
723		b), city station groups, c) and d), and coastal station groups, e) and f), shown in Fig. 2. . . .	47
724	Fig. 15.	Metrics derived from fitting ellipse equations (5) and (6) to wind perturbations at the Aus-	
725		tralian capital city airport stations, a) to d), and to wind perturbations spatially averaged	
726		over the city station groups and coastal station groups shown in Fig. 2, e) to h) and i) to l)	
727		respectively, with perturbations also temporally averaged over June, July and August 2018	
728		in each case. Metrics given are the maximum perturbation speed, a), e) and i), eccentricity	
729		of fitted ellipse, b), f) and j), orientation semi-major axis makes with lines of latitude, c), g)	
730		and k), and time of maximum perturbation, d), h) and l).	48
731	Fig. 16.	Temporal hodograph, a), and ellipse fit, b), of wind perturbations at each hour UTC averaged	
732		over June, July and August 2018 at Brisbane Airport. For comparison, c) provides the	
733		hodograph of the mean perturbations at the nearby Spitfire Channel station (see Fig. 2). . . .	49
734	Fig. 17.	Mean square error between the AWS and official forecast zonal perturbations $\overline{(u_{AWS} - u_O)^2}$,	
735		a), e), and i), decomposed into the error variance $\text{var}(u_{AWS} - u_O)$ and squared bias	
736		$(\bar{u}_{AWS} - \bar{u}_O)^2$ terms of equation (8). Also, the decomposed mean square error between	
737		the AWS and OCF zonal perturbations, b), f) and j). Additionally, the AWS and official	
738		forecast error variance term $\text{var}(u_{AWS} - u_O)$ decomposed into the $\text{var}(u_{AWS})$, $\text{var}(u_O)$ and	
739		$-2 \cdot \text{cov}(u_{AWS}, u_O)$ terms, c), g) and k), and analogously for the official forecast and OCF	
740		error variance term $\text{var}(u_{AWS} - u_O)$, d), h) and l). Decompositions given for Brisbane Air-	
741		port, a) to d), the Brisbane city station group, e) to h), and the Queensland coastal station	
742		group, i) to l). See Fig. 2 for station locations.	50

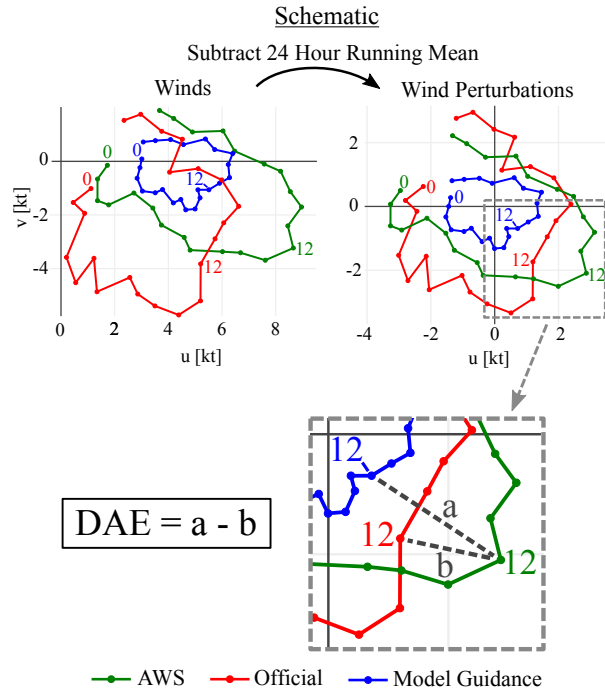


FIG. 1. Illustration of the method for calculating the *difference of absolute errors* (DAE) in the diurnal signals of an unedited model guidance dataset, and the human edited official forecast dataset, when compared with automatic weather station (AWS) observations, at an example time of 12:00 UTC.

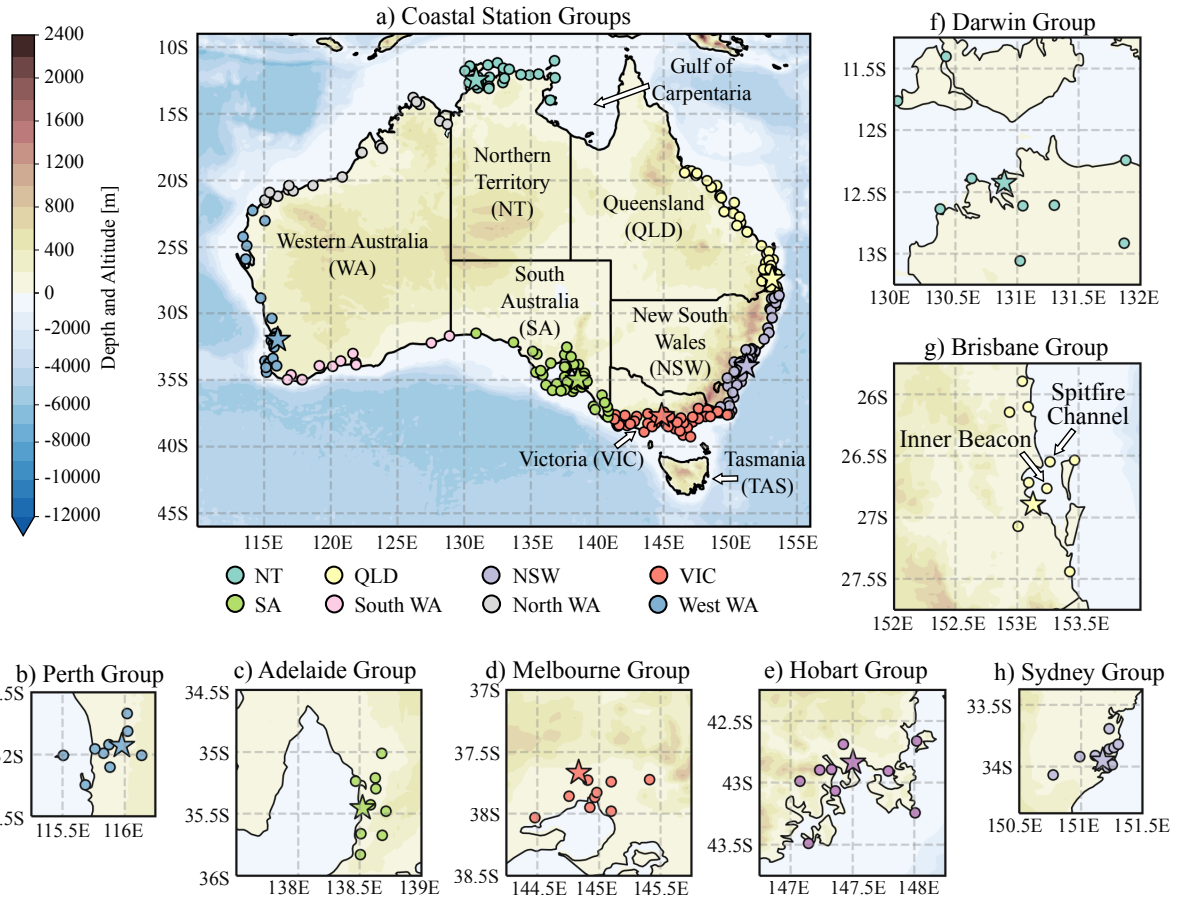


FIG. 2. Locations of the automatic weather stations, and the groupings of these stations, considered in this study. The *coastal station groups* are indicated in a), with the *airport stations* shown by stars. The Perth, Adelaide, Melbourne, Hobart, Darwin, Brisbane and Sydney *city station groups* are shown shown by b) to h), respectively.

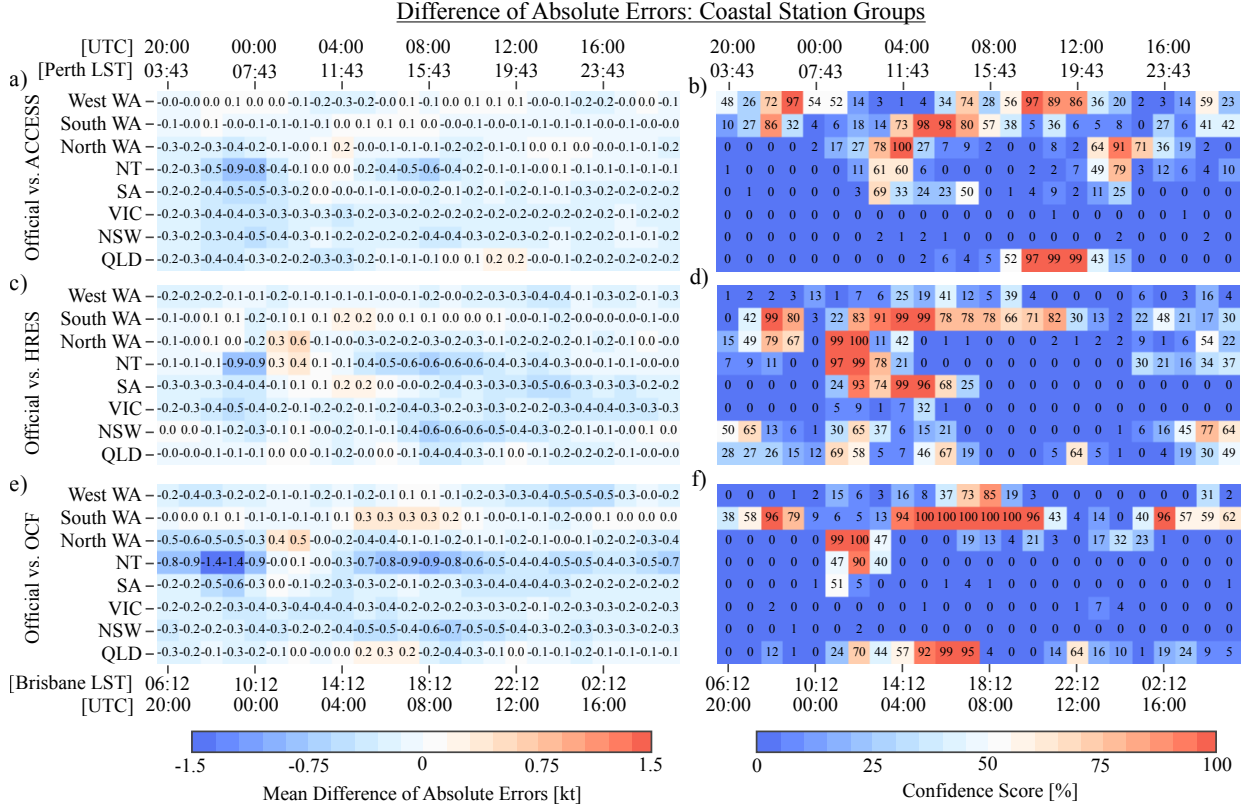


FIG. 3. Heatmaps of mean difference of absolute error $\overline{\text{DAE}}$ values, a), c), e), and confidence scores, b), d), f), for the *coastal station groups* (see Fig. 2). Results given for each hour of the day, for the official forecast versus ACCESS, a) and b), official forecast versus HRES, c) and d), and official forecast versus OCF, e) and f), comparisons. Positive $\overline{\text{DAE}}$ values indicate that the former dataset in each pair is on average $\overline{\text{DAE}}$ kt closer to observations than the latter dataset (see equation 1), where $1 \text{ kt} \approx 0.514 \text{ m s}^{-1}$. Confidence scores provide the probability the population or “true” value of $\overline{\text{DAE}}$ is greater than zero (see section 2).

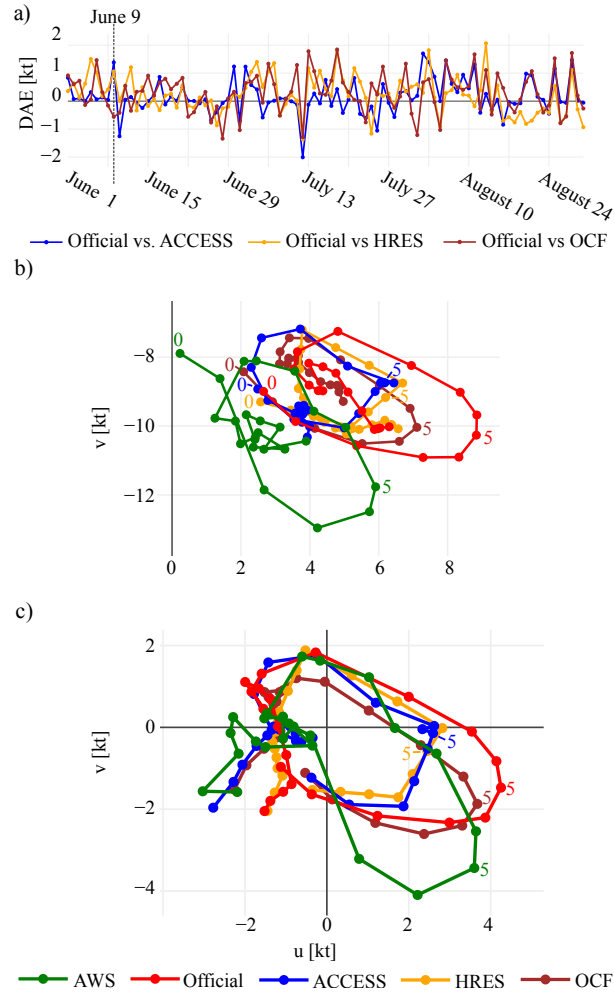


FIG. 5. As in Fig. 4, but for, a), the South Western Australia (WA) coastal station group at 05:00 UTC, and b) and c), the winds and wind perturbations, respectively, over the South WA coastal station group on the 9th June 2018.

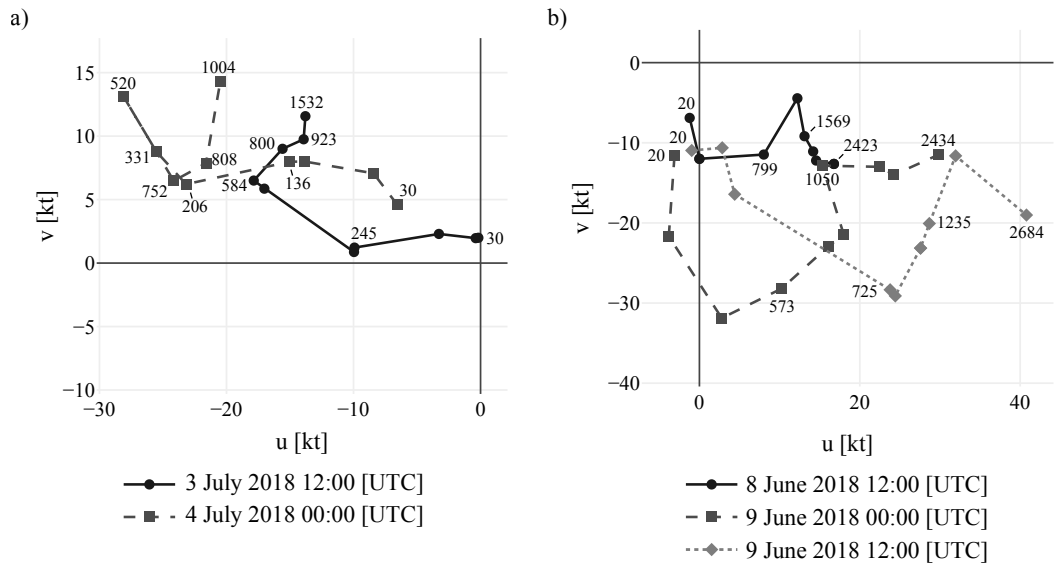


FIG. 6. Vertical wind soundings at, a), Darwin Airport, and b), Perth Airport, with heights given in metres.

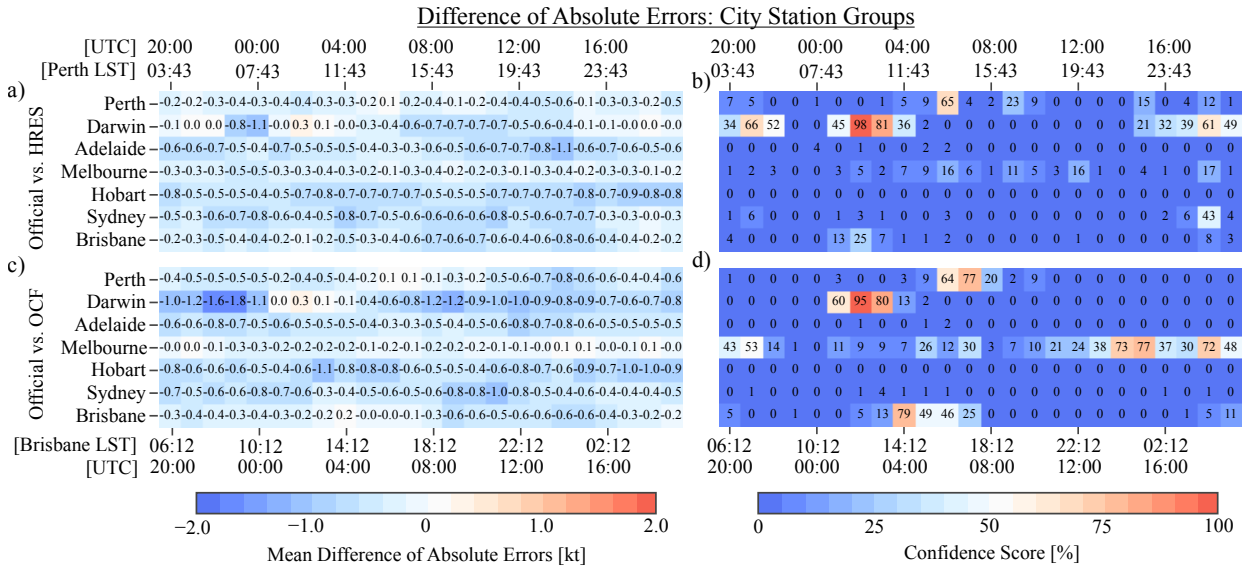


FIG. 7. As in Fig. 3, but for the official forecast versus HRES comparison a) and b), and the official forecast versus OCF comparison, c) and d), for the *city station groups* (see Fig. 2.)

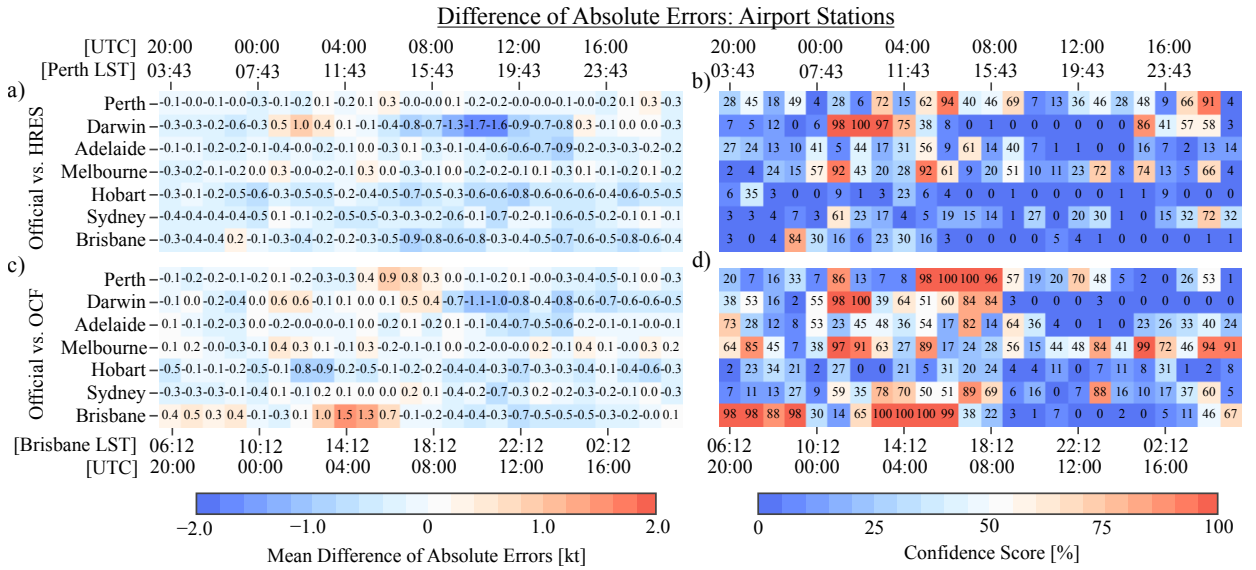


FIG. 8. As in Fig. 3, but for the official forecast versus HRES comparison a) and b), and the official forecast versus OCF comparison, c) and d), for the *airport stations* (see Fig. 2.)

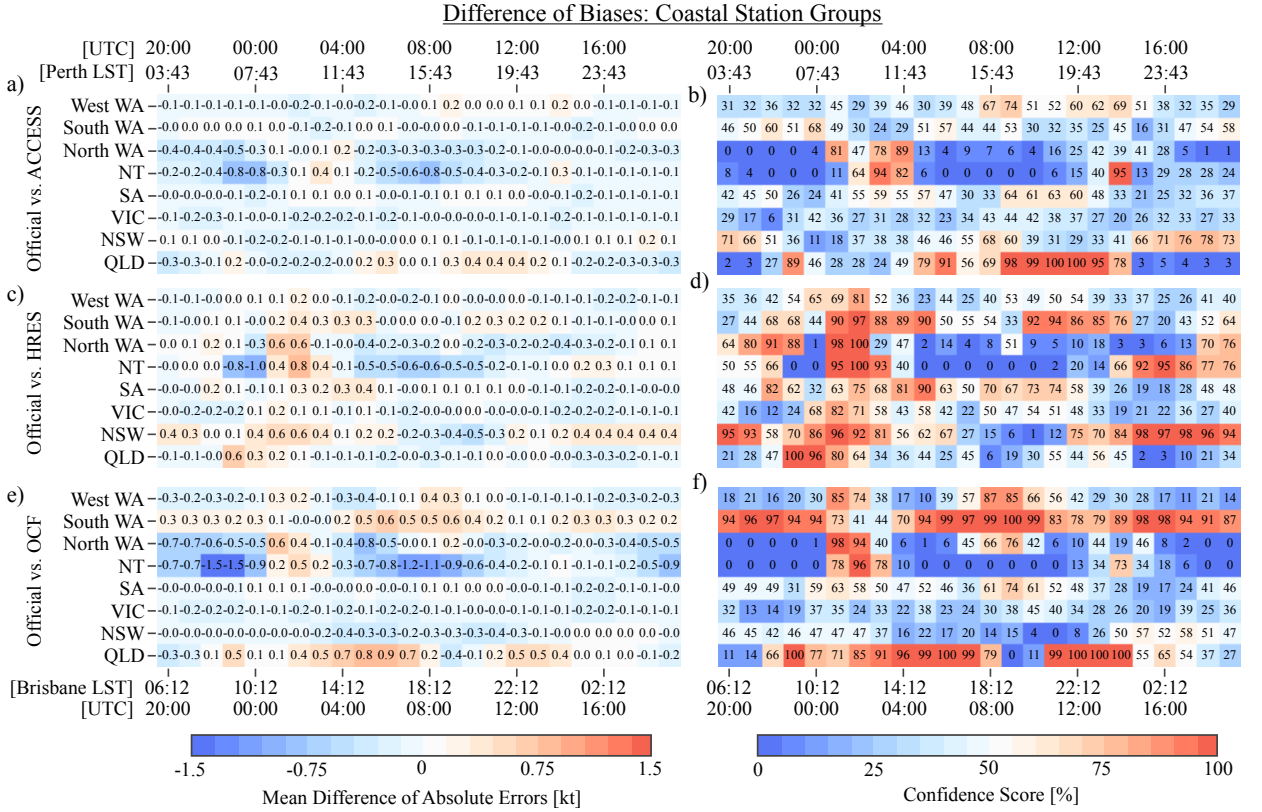


FIG. 9. As in Fig. 3, but for the difference of biases (DB) values and confidence scores.

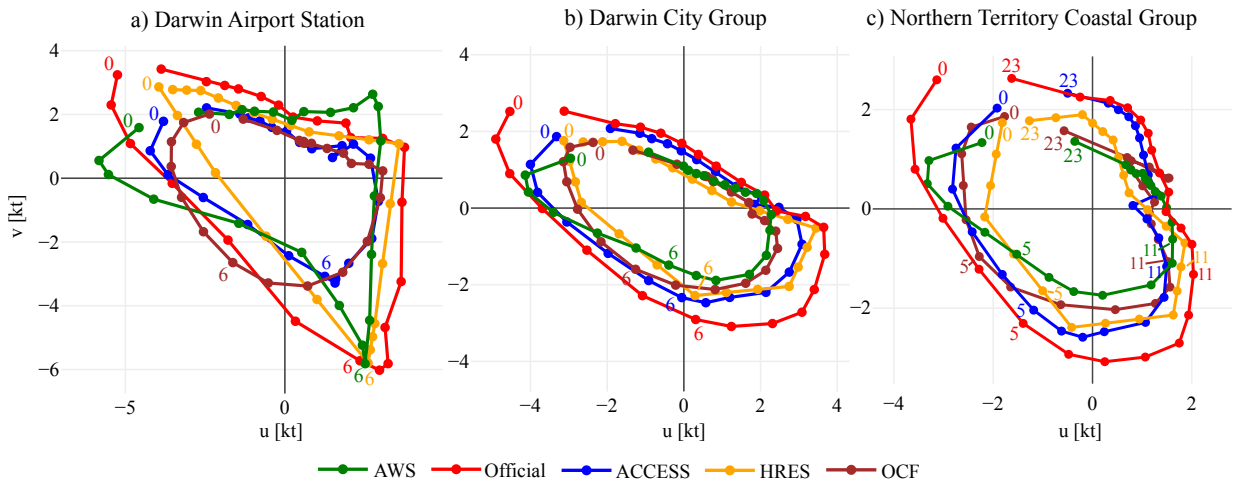


FIG. 10. Temporal hodographs in hours UTC of wind perturbations at, a), Darwin Airport, b), spatially averaged over the Darwin city station group, and c), the NT coastal group (see Fig. 2), then temporally averaged over June, July and August 2018.

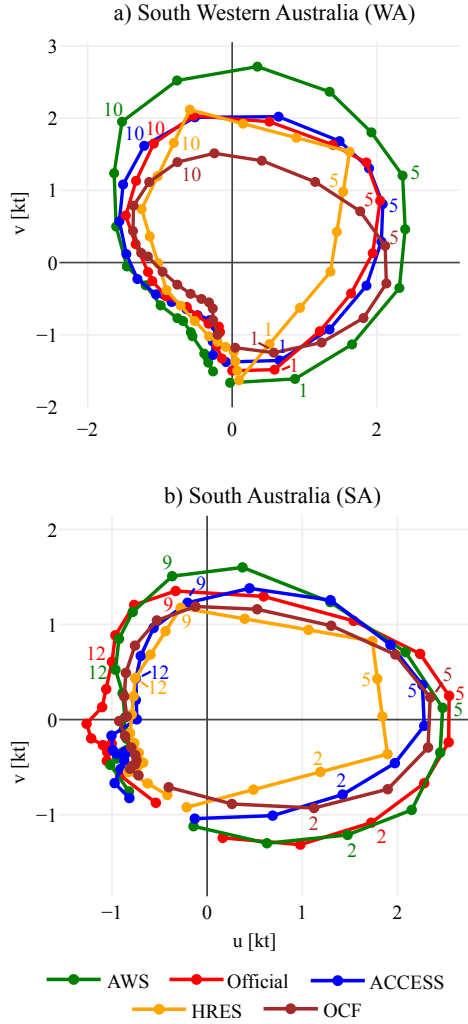


FIG. 11. Temporal hodographs in hours UTC of diurnal wind perturbations spatially averaged over the, a),
 South Western Australia (WA), and b), South Australia (SA) coastal station groups (see Fig. 2), and temporally
 averaged over June, July and August 2018.

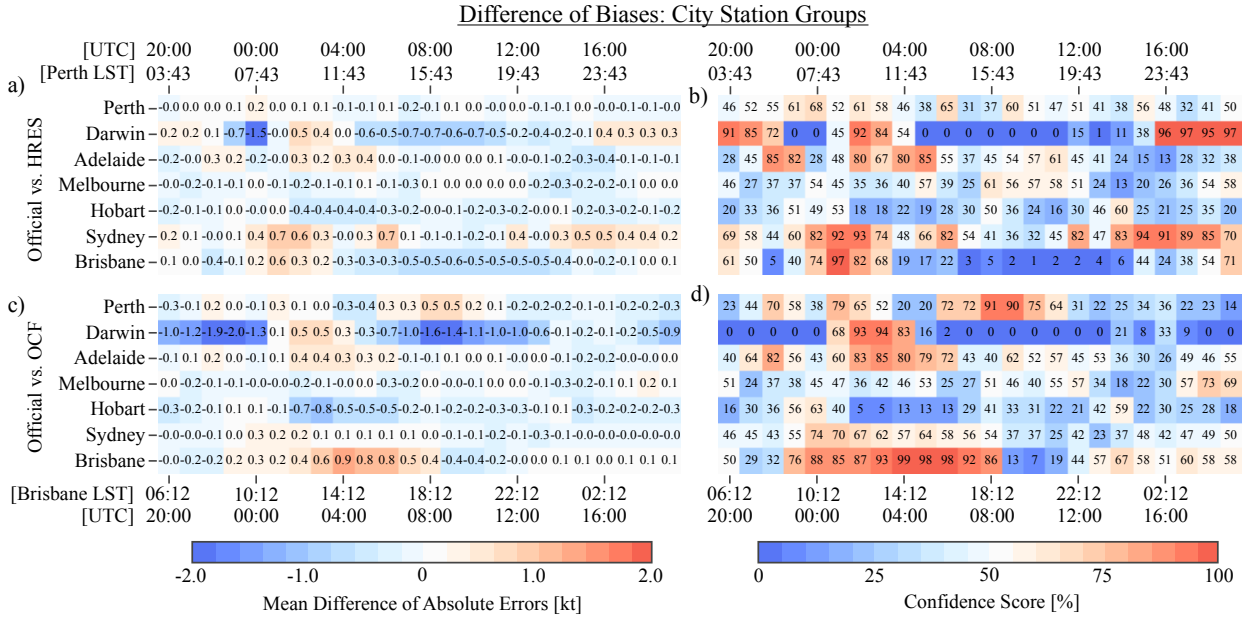


FIG. 12. As in Fig. 7, but for the difference of biases (DB) values and confidence scores.

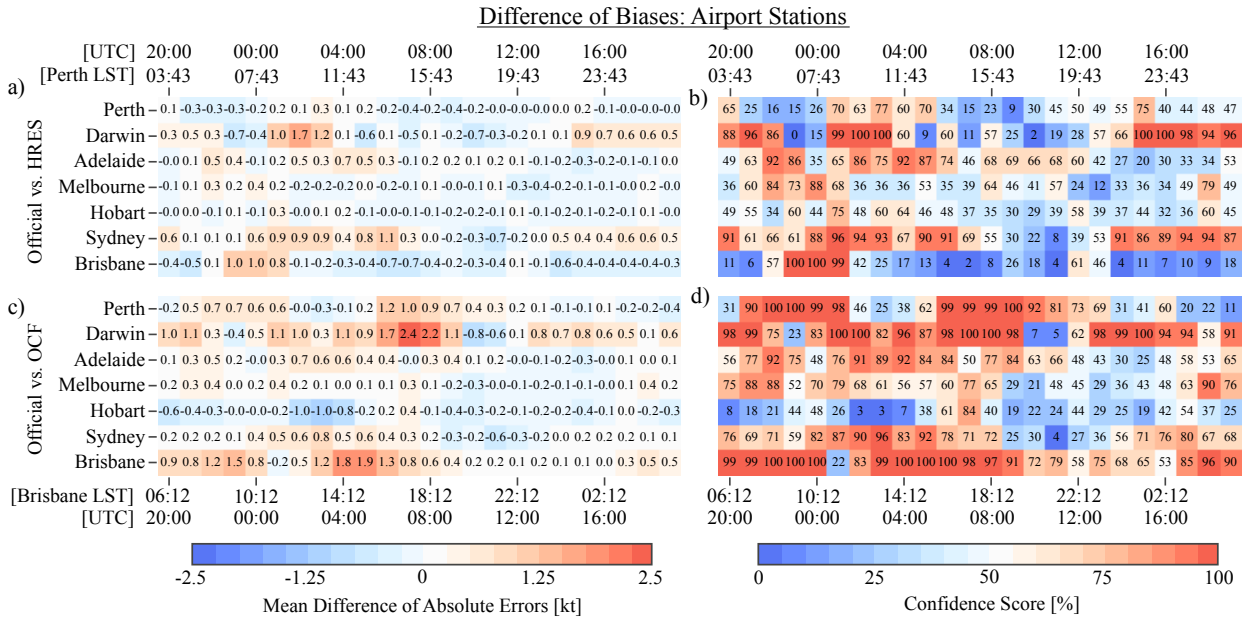


FIG. 13. As in Fig. 8, but for the difference of biases (DB) values and confidence scores.

Airport Stations	a)	Zonal Perturbations							b)	Meridional Perturbations						
	OCF	99	99	96	81	96	94	99	99	96	97	97	90	94	95	
	HRES	99	91	97	29	88	98	99	97	97	97	98	98	89	94	
	ACCESS	99	94	97	75	92	97	99	97	97	89	96	97	96	86	
	Official	99	96	98	90	76	95	99	98	97	93	94	95	93	90	
	AWS	96	91	92	59	59	80	94	93	89	83	94	68	94	97	
City Station Groups	c)								d)							
	OCF	99	96	95	89	95	98	99	94	95	92	98	97	99	99	
	HRES	99	99	83	87	97	94	99	94	92	96	98	98	96	97	
	ACCESS	98	98	96	88	97	95	99	92	94	92	98	92	98	98	
	Official	98	97	72	94	98	97	99	91	93	94	99	95	99	99	
	AWS	100	96	66	84	95	95	99	99	99	77	97	93	98	97	
Coastal Station Groups	e)								f)							
	OCF	99	98	97	94	94	89	99	97	98	98	98	97	97	96	
	HRES	99	97	97	99	96	84	99	96	97	98	96	99	98	97	
	ACCESS	99	96	97	95	95	87	99	95	97	98	98	99	99	98	
	Official	99	97	98	98	98	94	98	96	98	98	98	99	99	96	
	AWS	97	96	97	94	93	83	91	98	98	99	98	96	93	90	
		Perth	Darwin	Adelaide	Melbourne	Hobart	Sydney	Brisbane	Perth	Darwin	Adelaide	Melbourne	Hobart	Sydney	Brisbane	
		West WA	South WA	North WA	NT	SA	VIC	NSW	West WA	South WA	North WA	NT	SA	VIC	NSW	QLD

R^2 Goodness of Fit [%]

FIG. 14. 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. 2.

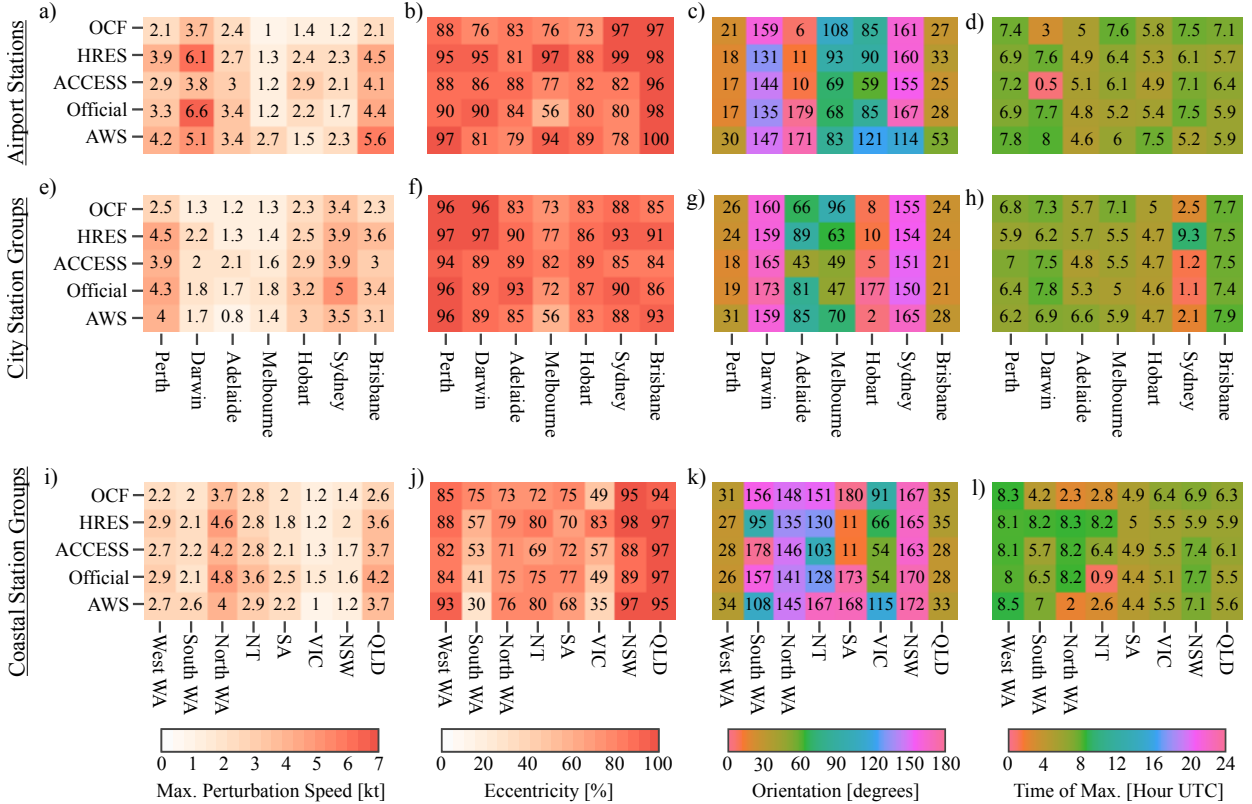


FIG. 15. Metrics derived from fitting ellipse equations (5) and (6) to wind perturbations at the Australian capital city airport stations, a) to d), and to wind perturbations spatially averaged over the city station groups and coastal station groups shown in Fig. 2, e) to h) and i) to l) respectively, 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).

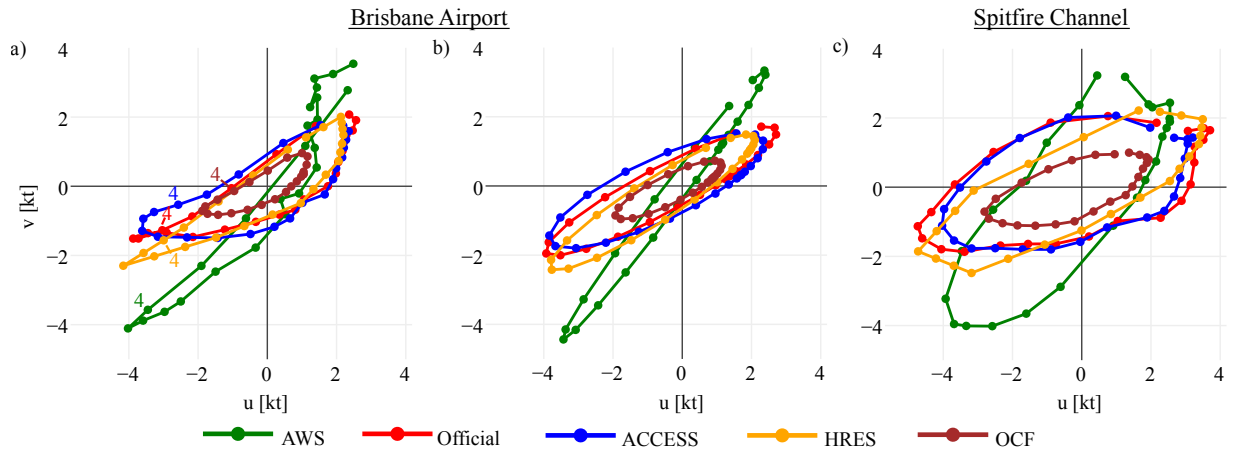


FIG. 16. Temporal hodograph, a), and ellipse fit, b), of wind perturbations at each hour UTC averaged over June, July and August 2018 at Brisbane Airport. For comparison, c) provides the hodograph of the mean perturbations at the nearby Spitfire Channel station (see Fig. 2).

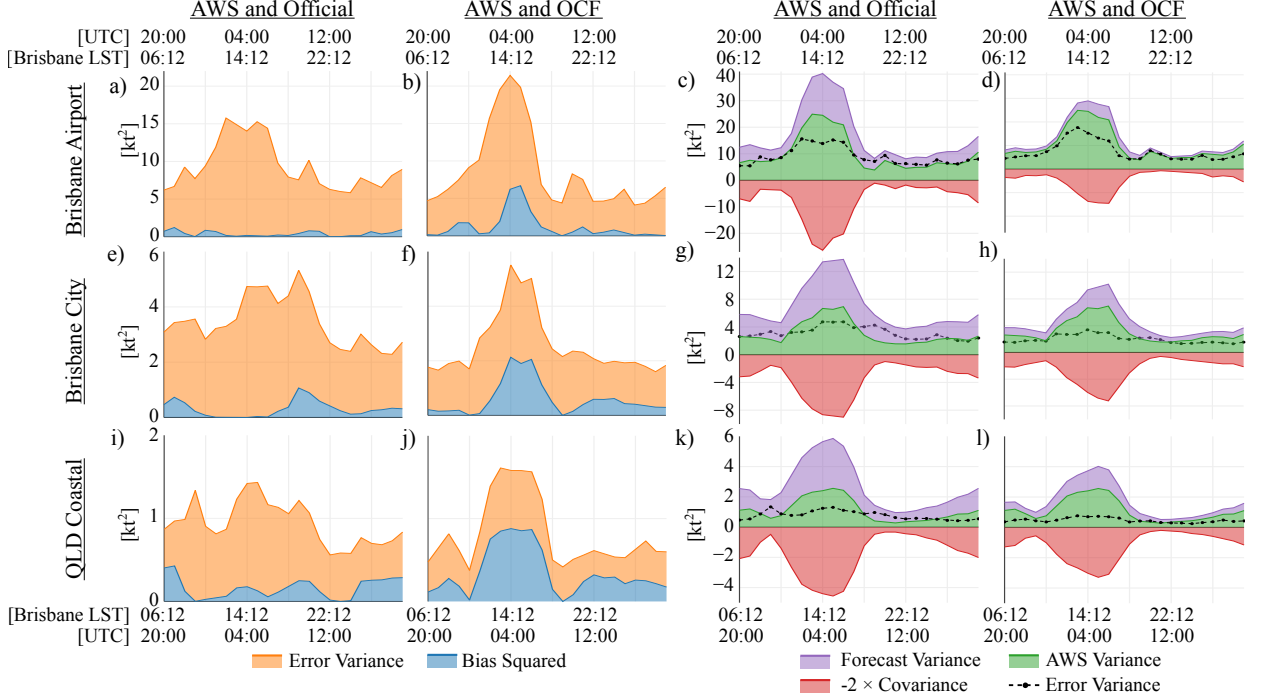


FIG. 17. Mean square error between the AWS and official forecast zonal perturbations $(u_{\text{AWS}} - u_{\text{O}})^2$, a), e), and i), decomposed into the error variance $\text{var}(u_{\text{AWS}} - u_{\text{O}})$ and squared bias $(\bar{u}_{\text{AWS}} - \bar{u}_{\text{O}})^2$ terms of equation (8). Also, the decomposed mean square error between the AWS and OCF zonal perturbations, b), f) and j). Additionally, the AWS and official forecast error variance term $\text{var}(u_{\text{AWS}} - u_{\text{O}})$ decomposed into the $\text{var}(u_{\text{AWS}})$, $\text{var}(u_{\text{O}})$ and $-2 \cdot \text{cov}(u_{\text{AWS}}, u_{\text{O}})$ terms, c), g) and k), and analogously for the official forecast and OCF error variance term $\text{var}(u_{\text{AWS}} - u_{\text{O}})$, d), h) and l). Decompositions given for Brisbane Airport, a) to d), the Brisbane city station group, e) to h), and the Queensland coastal station group, i) to l). See Fig. 2 for station locations.