

3 **Forecast Verification for Land-Sea Breeze and**
4 **Boundary Layer Mixing Processes**

5 **Ewan Short¹ | Ben Price² | Derryn Griffiths³ |**
Alexei Hider³

¹ARC Centre of Excellence for Climate

Extremes, School of Earth Sciences,
University of Melbourne, Parkville, VIC,
3010, Australia

²Bureau of Meteorology, Casuarina, NT,
0810, Australia

³Bureau of Meteorology, Melbourne, VIC,
3208, Australia

Correspondence

Ewan Short, ARC Centre of Excellence for
Climate Extremes, School of Earth Sciences,
University of Melbourne, Parkville, VIC,
3010, Australia
Email: ewan.short@unimelb.edu.au

Funding information

ARC Centre of Excellence for Climate
System Science

This study presents a methodology for comparing the performance of Australian Bureau of Meteorology forecasts of the land-sea breeze with unedited model guidance products, such as those of the European Center for Medium-Range Weather Forecasting (ECMWF) and the Australian Community Climate and Earth System Simulation (ACCESS). The methodology is applied to the 8 Australian capital city airports. The results indicate that at some airports, human intervention to model guidance products adds value to land-sea breeze forecasts, whereas at other airports it does not.

KEYWORDS

land-sea breeze, forecast verification, Australia, ACCESS, ECMWF

1 | INTRODUCTION

Modern weather forecasts are produced by models in conjunction with human forecasters. For instance, a forecaster working for the Australian Bureau constructs a seven day forecast by first loading model data into the Graphical Forecast Editor (GFE) software package, then manually editing this model data as they see fit. Forecasters can choose which model to base their forecast on, and refer to this as a choice of *model guidance*. Edits are typically made to account for processes that are underresolved at synoptic scale model resolutions, or to address known biases of the models being used.

It is therefore important to assess not only the overall accuracy of weather forecasts, but also the contribution human forecaster edits make to this accuracy. If effective, but routine, editing procedures can be identified they can be automated, freeing forecasters up to focus on other tasks. One common edit involves changing the surface wind fields near coastlines to try to represent sea-breezes more realistically. Forecasters invest time in making sea-breeze edits because accurate predictions of near-surface winds are highly valued by a number of users, such as the aviation and energy (Smith et al., 2009) industries. Accurate sea-breeze forecasts are also valuable to environmental monitoring authorities, as these winds provide ventilation to coastal urban areas.

Assessing the accuracy of a weather forecast is a task far more nuanced than it might first appear. For instance, attempting to assess the accuracy of a precipitation forecast by comparing the rainfall amounts measured at an individual weather station to the closest grid point of a model prediction will often give poor results. Although the synoptic drivers of convection are usually well predicted, exactly where convective cells form, and where the most rain falls, is highly unpredictable. As such, it is often appropriate to use "fuzzy" verification metrics which measure the agreement between prediction and observation in a more indirect way. For instance, one approach known as "upscaling" is to first average forecast and observational data over a given spatial domain before calculating verification scores. Ebert (2008) provided a review of current "fuzzy verification" methodologies, and a framework for how they can be used to determine the spatial scales at which a given forecast has predictive skill.

Relatively few forecast verification studies have focused on near-surface winds, and the ones that have generally only considered wind speeds. Pinson and Hagedorn (2012) performed a verification study of the ECMWF 10 m wind

speeds across western Europe over December, January, February 2008/09. First, they interpolated ECMWF model data onto the locations of weather stations across Europe, then they compared the interpolated model data at these stations with the station observations themselves. They found that the worst performing regions were coastal and mountainous areas, and attributed this poor performance to the small scale processes, e.g. sea and mountain breezes, that are underresolved at ECMWF's coarse 50km spatial resolution. They noted that future work could better identify the effect of diurnal cycles on verification statistics by considering forecasts at different times of day.

Lynch et al. (2014) also performed a verification study of ECMWF 10 m wind speed data, with the goal of assessing skill at lead times of between 14 to 20 days. They compared ECMWF 32-day forecast model wind speeds with gridded ERA-Interim wind speeds between 2008-12, with both datasets analysed at a six hour temporal resolution. Before conducting the comparison, the wind speed data were transformed into wind-speed "anomaly" data by first calculating the mean wind speed at 0000, 0600, 1200 and 1800 UTC for each calendar day from the entire ERA-Interim record, and from a 20 year ECMWF 32-day model hindcast, then subtracting these means from the ERA-Interim and ECMWF 32-day model data respectively. Wind speed anomaly data was used so that stable seasonal and diurnal cycles did not contribute to verification scores. At the 14-20 day timescale around western Europe, the greatest skill was found in the boreal winter (austral summer) months of December, January and February.

Pinson and Hagedorn (2012) and Lynch et al. (2014) restricted their verification studies to wind speeds, but wind directions are also crucial to diagnosing whether land sea breezes - and the diurnal wind cycle more generally - are being forecast correctly. Furthermore, no previous published work has proposed a verification methodology to assess the accuracy of the diurnal wind cycle in forecasts, or of the contributions made to this accuracy by human forecaster edits of model output. Finally, no previously published work has considered the performance of ACCESS near surface winds, which together with ECMWF, are the model guidance products most widely used by Australian forecasters. Thus, the present study has two goals. First, to describe a methodology for comparing human edited forecasts of the land-sea breeze to unedited model guidance forecasts, in order to assess where and when human edits are producing an increase in accuracy. Second, to apply this methodology across Australia. The remainder of this paper is organised as follows. Section 2 describes the methodology in detail, section 3 provides results, and sections ?? and 5 provide a discussion and a conclusion, respectively.

2 | DATA AND METHODS

This study compares both edited and non-edited Australian Bureau of Meteorology forecast data with automatic weather station (AWS) data across Australia. The comparison is performed by first isolating the diurnal signals of each dataset, then comparing these signals on an hour-by-hour basis. If the diurnal cycle cannot be resolved correctly using wind perturbations, it cannot be resolved correctly in the overall wind fields, which are subject to additional synoptic scale errors between the models and observations.

2.1 | Data

Four datasets are considered in this study; they are the Australian Bureau of Meteorology's Official wind forecast data, model data from the European Center for Medium Range Weather Forecasting (ECMWF), model data from the Australian Community Climate and Earth System Simulator (ACCESS), and observational data from automatic weather stations. The Official, ECMWF and ACCESS data are at a $1^\circ, 1^\circ$ degree spatial resolution respectively. What are the resolutions of these datasets as they're used in Jive? Does the ACCESS model data in Jive Official, ACCESS and AWS data exists at each UTC hour. ECMWF data exists at a three hour resolution. To be consistent with the other data sets, ECMWF is therefore linearly interpolated to an hourly resolution: this is also what happens in practice when forecasters load ECMWF wind data into the GFE. Two time periods are considered, the austral summer months (December, January, February) of 2017/18, and the austral winter months (June, July, August) of 2018.

Only station data from the seven Australian capital city airport automatic weather stations are considered; Official, ECMWF and ACCESS data is (linearly?) interpolated to the coordinates of the airport weather stations. Capital city airports have been chosen as the focus of this study for a number of reasons. Automatic weather stations located at airports tend to provide the most accurate wind data, and wind forecasts at airports are important to the aviation industry. Moreover, the capital city airports are all reasonably close to coastlines, resulting in a clear diurnal signal. Finally, these airports are also all close to their respective capital cities, which are high priority regions for accurate forecasting. The datasets are hosted on the Bureau's Jive database, but are not currently generally available, although

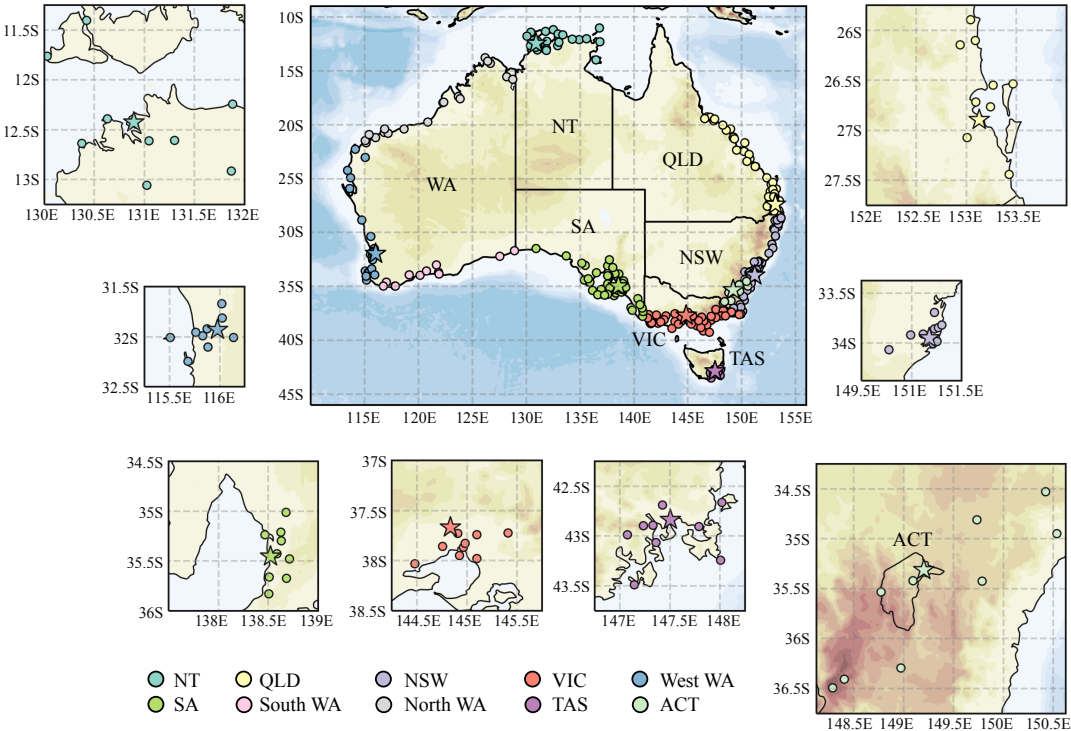


FIGURE 1 Locations of the automatic weather stations used in this study. Stars indicate capital city airport stations. Height and depth shading intervals every 200 and 1000 m, respectively.

the long term plan is for this to change. Can I extract and host the data I need myself? Can I obtain copies of the relevant Jive Functions so that I can post complete code online?

As described above, the Australian Bureau of Meteorology's official wind forecast is constructed out of model data, which is then edited by human forecasters using the Graphical Forecast Editor (GFE) software package. Australian forecasters typically construct wind forecasts out of model data either from the European Center for Medium Range Weather Forecasting (ECMWF), or the Australian Community Climate and Earth System Simulator (ACCESS). Testing whether the official forecast data conforms more closely to the AWS observations than ECMWF or ACCESS therefore provides a way to assess the extra accuracy gained by forecaster edits.

2.2 | Assessing Diurnal Cycles

Although close to coastlines the land-sea breeze is generally the dominant diurnal wind process, the overall diurnal signal may also include mountain-valley breezes, boundary layer mixing processes, atmospheric tides, and urban heat island circulations. Forecasters typically edit model output to account for *both* unresolved sea-breezes *and* unresolved boundary layer mixing; attempting to focus solely on sea-breezes without examining the entire diurnal cycle therefore risks erroneous conclusions, with the effects of one category of edit mistaken for another. In general it is hard to separate boundary layer mixing edits from sea-breeze edits in the diurnal cycle composites, so this point maybe needs to be reworked. Or could simply comment on this in the discussion.

Sea-breezes are therefore analysed by examining the overall diurnal signal in each dataset, with the assumption that close to coastlines the land-sea breeze is the dominant diurnal process. The diurnal signal is identified by subtracting a twenty hour centred running mean *background wind* from each zonal and meridional hourly wind data point. This provides a collection of zonal and meridional wind *perturbation* datasets. Note that thinking of land-sea breezes in terms of perturbations from a background wind may require a conceptual shift from the usual operational definitions. A forecaster would likely define a sea-breeze to be a reversal in wind direction from a primarily offshore flow during the night and morning, to an onshore flow in the afternoon and evening. However, even if the wind is offshore the entire day, sea-breeze *perturbations* are generally still detectable as a weakening of the offshore flow throughout the

afternoon and evening.

Note that subtracting background winds may raise concerns, because perturbations obviously depend on background winds. However, the forecaster does not have knowledge of the observations when they make the diurnal process edits. They are implicitly assuming that the true mean winds will be close enough to the predicted mean state - however this prediction is produced - to justify making diurnal edits on the basis of the predicted mean state.

Once the wind perturbation datasets have been constructed, the accuracy of the Official, ACCESS and ECMWF diurnal cycles are quantified by first calculating the Euclidean distances of the perturbations at each hour from the corresponding AWS perturbations. For instance, to quantify how closely the Official forecast perturbations match the AWS observations, we calculate the Euclidean distances $|u_{AWS} - u_O|$ at each time step. The accuracy with which the Official and ACCESS datasets resolve the diurnal cycle can then be compared by defining the *Wind Perturbation Index* (WPI)

$$WPI_{OA} \equiv |u_{AWS} - u_A| - |u_{AWS} - u_O|. \quad (1)$$

At a given time, the Official forecast wind perturbation is closer to the AWS perturbation than that of ACCESS if and only if $WPI > 0$. Similarly, the WPI can be used to provide a comparison of the Official and ECMWF datasets, or a comparison of the two model guidance datasets ACCESS and ECMWF.

To assess which dataset provides, in general, the most accurate representation of the diurnal cycle, we then take means of the WPI on an hourly basis; i.e. all the 00:00 UTC WPI values are averaged, all the 01:00 UTC values are averaged, and so forth. The sampling distributions of these means can then be modelled as Student's t -distributions, and from this we can calculate the probability that $\overline{WPI} > 0$ at each hour, where the bar denotes a temporal average. Temporal autocorrelations of WPI, i.e. correlations between WPI values at a particular hour from one day to the next, are accounted for using the standard method of reducing the "effective" sample size to $n(1 - \rho_1)/(1 + \rho_1)$, where n is the actual sample size and ρ_1 is the lag-1 autocorrelation (Zwiers and von Storch, 1995; Wilks, 2011), although in practice temporal autocorrelations of WPI are either non-existent or very small. To assess how well the diurnal

perturbations of an overall region are predicted, for instance those of the Victorian coastal station group (see Fig. 1), the perturbations across each station group are averaged before WPI values calculated. The temporal means and sampling distributions of the WPI are then calculated as before, with each value of WPI calculated from the spatially averaged perturbations treated as a single observation. This provides a conservative method for dealing with spatial correlation in the perturbations.

The advantage of the WPI method is it's clarity and simplicity: we are essentially just comparing the magnitudes of vector differences, then applying a two sided t -test to determine whether one dataset's perturbations are consistently closer to observations than another's. One factor that complicates interpretation of statistics of WPI, is that the near surface winds observed in AWS data are consistently noisier than those of the Official, ECMWF and ACCESS forecasts. This is likely due to unresolved subgrid scale turbulence in the Official, ECMWF and ACCESS model datasets. It would be unreasonable to expect forecasters to be able to predict this essentially random additional observed variability, and so a direct comparison of observed and modelled diurnal cycles is overly stringent.

To reduce the significance of unpredictable noise, we also compare temporal averages of the perturbations for each dataset. These comparisons have less operational significance: people generally care how well the actual weather forecast performed, not whether the average of a predicted quantity matched the average of an observed quantity. However, comparisons of averages arguably better represent what we can realistically expect from human forecaster edits, and from weather forecasts overall, particularly in regards to small scale processes like sea-breezes. Furthermore, when temporal averages of perturbations are considered, the diurnal signal becomes dramatically clearer, and structural differences become much easier to diagnose.

To quantify how closely the temporally averaged Official forecast perturbations match those of the AWS observations, we calculate $|\bar{u}_{AWS} - \bar{u}_O|$ for each hour. To assess the performance of the Official temporally averaged perturbations against those ACCESS, we define the *Climatological Wind Perturbation Index* (CWPI)

$$CWPI_{OA} \equiv |\bar{u}_{AWS} - \bar{u}_O| - |\bar{u}_{AWS} - \bar{u}_A|. \quad (2)$$

As with the WPI, the CWPI can also be used to provide a comparison of the Official and ECMWF datasets, or a comparison of the two model guidance datasets ACCESS and ECMWF. Uncertainty in the CWPI is estimated through bootstrapping (Efron, 1979). This is done by performing resampling with replacement on the underlying perturbation datasets, and calculating the CWPI multiple times using these resampled datasets. This provides a distribution of CWPI values, from which the probability that $CWPI > 0$ can be calculated. Similarly to with the WPI, performance over a particular region can be assessed by first averaging perturbation values over multiple stations before the CWPI is calculated.

Although the WPI and CWPI provide quantitative information on the accuracy of the diurnal cycle at different times of day, they do not provide much information about the structure of the diurnal wind cycles of each dataset, or provide insight into the reason one dataset is outperforming another. Gille et al. (2005) obtained summary statistics on the observed structure of temporally averaged diurnal wind cycles across the globe by using linear regression to calculate the coefficients u_i, v_i $i = 0, 1, 2$, for the elliptical fit

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

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

where ω is the angular frequency of the earth and t is the local solar time in seconds. Descriptive quantities - like the angle the semimajor axis of the ellipse makes with the horizontal - were then calculated directly from the coefficients u_1, u_2, v_1 and v_2 .

Gille et al. (2005) applied this fit to satellite scatterometer wind observations, which after temporal averaging provided only four temporal datapoints at each $0.25^\circ \times 0.25^\circ$ spatial grid cell. As such, their fit was very good, explaining over 90% of the wind variability in each spatial gridcell. However, the choice of ellipse parametrisation in equations 5 and 6 assumes that datapoints lie on the ellipse at equal intervals of time t . When observational or model data with an hourly or smaller timestep is considered, this assumption becomes too stringent, as heating asymmetries imply

that wind perturbations evolve much more rapidly during the day than at night (see Fig. XX). Note I'm also basing this point on knowledge of the land vs sea breeze, and knowledge of heating vs cooling asymmetries (Brown et al., 2017, e.g.).

Thus, we model the climatological diurnal cycles with the equations

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

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

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

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

where t is time in units of hours UTC, and ψ gives to the time when the wind perturbations vary least with time. Need to confirm whether least or most! 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.

Descriptive quantities can then be calculated from these parameters. The value of α at which the winds align with the semimajor axis, α_M , satisfies

$$\alpha_M = \frac{1}{2} \arctan \left(\frac{2(u_1 u_2 + v_1 v_2)}{u_1^2 + v_1^2 - u_2^2 - v_2^2} \right) \bmod \pi, \quad (8)$$

The time at which the perturbations align with the major axis t_M can then be calculated by inverting equation (7), fixing ψ to the value obtained from the nonlinear regression. The lengths of the semimajor and semiminor axes, and the

angle the semimajor axis makes with lines of latitude ϕ , can then be calculated from α_M using the same expressions as Gille et al. (2005).

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 (austral winter) 2018. First, error is assessed on a daily basis using the Wind Perturbation Index (WPI) at three different spatial scales. Second, overall seasonal biases during this time period are assessed using the Climatological Wind Perturbation Index CWPI, and by comparing quantities derived from ellipses fitted to the climatological wind perturbations. Unless otherwise stated, values throughout this section are provided to two significant figures.

3.1 | Daily Comparison

Figure 2 provides the mean wind perturbation index values \overline{wpi} and confidence scores $P(\overline{WPI} > 0)$ for the coastal station groups for \overline{wpi}_{OA} , \overline{wpi}_{OE} and \overline{wpi}_{EA} , which represent the the Official versus ACCESS, Official versus ECMWF, and ECMWF versus ACCESS comparisons, respectively. Values of \overline{wpi}_{OA} and \overline{wpi}_{OE} are negative for the majority of station groups and hours, and often both $P(\overline{WPI}_{OA} > 0) < 5\%$ and $P(\overline{WPI}_{OE} > 0) < 5\%$. This implies that at this level of spatial aggregation, there is often high confidence that both the unedited ACCESS and ECMWF models outperform the Official forecast. The lowest \overline{wpi} values of -0.9 kn occur for the NT station group at 23:00 and 00:00 UTC for both \overline{wpi}_{OA} and \overline{wpi}_{OE} , with $\overline{wpi}_{EA} = 0$ kn. Comparatively low values also occur at 08:00 UTC with $\overline{wpi}_{OA} = \overline{wpi}_{OE} = -0.6$ kn, but $\overline{wpi}_{EA} = 0$ kn. This suggests the Official forecast may be performing particularly poorly over the NT station group.

Although Official outperforms at least one of ACCESS or ECMWF with high confidence at a few dozen times and station groups, there is only one group and time where it outperforms both. At 05:00 UTC over the South WA station

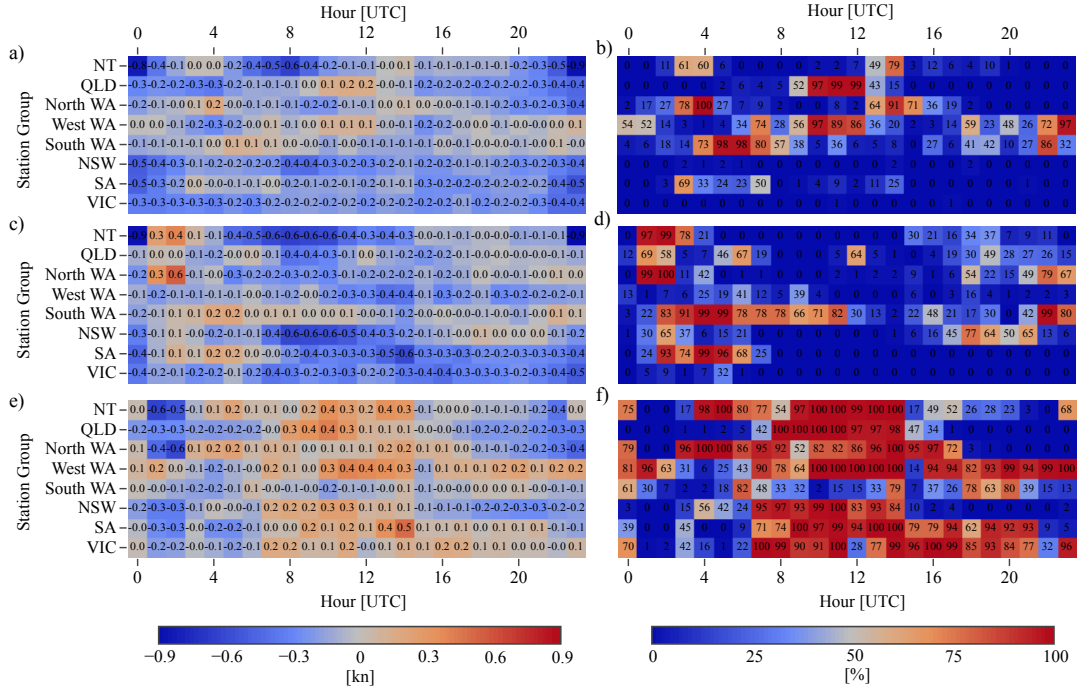


FIGURE 2 Heatmaps of \overline{WPI} values and confidence scores for each coastal station group and hour of the day: a) and b), Official versus ACCESS, c) and d) Official versus ECMWF, e) and f) ECMWF versus ACCESS. Positive \overline{WPI} values mean that the former dataset in each pair is on average closer to observations than the latter dataset. Confidence scores provide the probability the population \overline{WPI} is greater than zero. Values within the heatmaps are accurate to two significant figures.

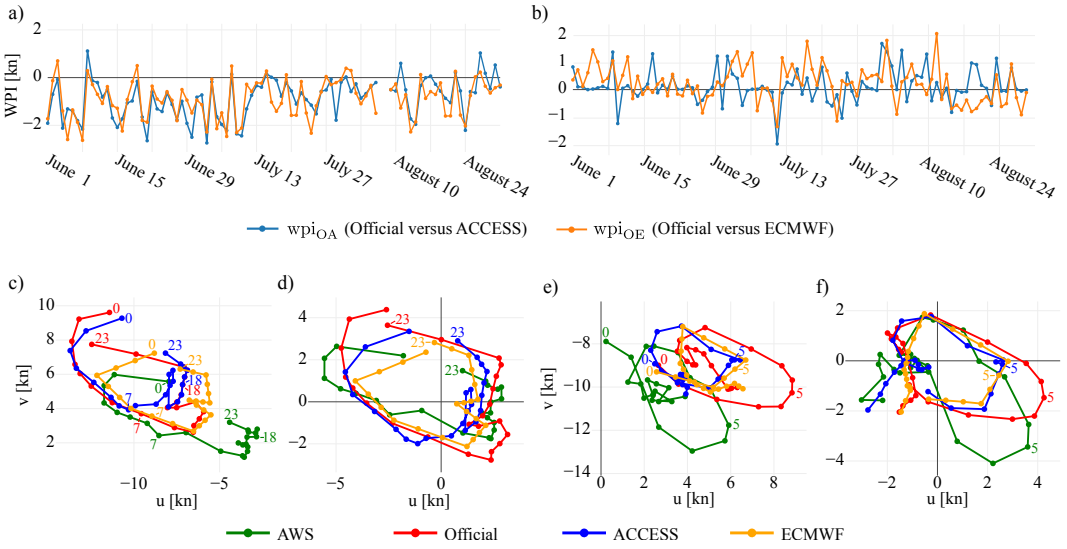


FIGURE 3 Time series, a) and b), of \overline{wpi}_{OA} and \overline{wpi}_{OE} for a), the NT station group at 23:00 UTC, and b), the south WA station group at 05:00 UTC. Hodographs, c) to f), showing change in winds, c) and e), and wind perturbations, d) and f), for the NT station group, c) and d), and south WA station group, e) and f).

group, $\overline{wpi}_{OA} = 0.2$ kn and $\overline{wpi}_{OE} = 0.1$ kn, both with confidence scores $\geq 95\%$, although the actual \overline{wpi} values are comparatively small. Note that ECMWF generally outperforms ACCESS from 10:00 - 14:00 UTC, with the South WA station group being the main exception.

Using the NT and South WA station groups as case studies, Figures 3 a) and b) provide time series of \overline{wpi}_{OA} and \overline{wpi}_{OE} for a), the NT station group at 23:00 UTC, and b), the South WA station group at 05:00 UTC. The \overline{wpi}_{OA} and \overline{wpi}_{OE} values for the NT station group show significant temporal variability over the three month period, exceeding -2 kn on at least 10 days each, and occasionally becoming positive. The \overline{wpi} values for the South WA station at 05:00 UTC also show significant temporal variability, with \overline{wpi}_{OA} and \overline{wpi}_{OE} each exceeding 1 kn on at least 9 separate days, despite \overline{wpi}_{OA} and \overline{wpi}_{OE} being small.

Fig. 3 a) shows that there are four days where \overline{wpi}_{OA} and \overline{wpi}_{OE} are both less than -2 kn: the 8th of June and the 3rd, 9th and 10th of July. Figures 3 c) and d) show hodographs of the winds and wind perturbations, respectively, at each hour UTC for the AWS observations, Official forecast, and ACCESS and ECMWF model datasets on the 3rd of July, which provides an interesting example. Figure 3 e) shows that the Official wind forecast on this day was likely based on edited ACCESS from 00:00 to 06:00 UTC, then edited ECMWF from 07:00 to 13:00 UTC, then unedited

ACCESS from 15:00 to 21:00 UTC. The final two hours of the forecast show the Official winds acquiring a stronger east-northeasterly component than either the AWS observations, ACCESS, or ECMWF; this rapid, exaggerated change is even clearer in the perturbation hodograph shown in Fig. 3 f). Note that at this time of year the prevailing winds throughout the NT are east-southeasterly, and 22:00 UTC corresponds to \approx 08:30 LST in this region, so the rapid departure of the Official forecast from ACCESS at this time likely represents an edit made by a forecaster to capture boundary layer mixing processes. Figure 4 a) shows the first ten values from wind soundings at Darwin Airport - the nearest station to issue vertical wind soundings - at 12:00 UTC on July 3rd and 00:00 UTC on July 4th. In both instances the winds are indeed east-southeasterly, and so the rapidly changing wind perturbations at 22:00 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. The 8th of June and 9th and 10th of July examples are all similar in this respect.

Considering now the South WA station group, Fig. 3 b) shows that wpi_{OA} and wpi_{OE} both exceed 1 kn on the 9th of June and the 3rd of August. Figures 3 c) and d) show hodographs of the winds and wind perturbations, respectively, at each hour UTC for the AWS observations, Official forecast, and ACCESS and ECMWF model datasets on the 9th of June, which is the more interesting example. The perturbation hodograph shows both ECMWF and ACCESS underpredicting the amplitude of the diurnal wind cycle on this day. In each dataset the 05:00 UTC perturbations are westerly to northwesterly, and given the orientation of the South WA coastline (see Fig. ??) and the fact that 05:00 UTC corresponds to \approx 13:00 local solar time (LST) in this region, the perturbations likely indicate boundary layer mixing processes, rather than the land-sea breeze. Furthermore, the AWS perturbations rapidly become northwesterly between 01:00 and 02:00 UTC, \approx 09:00 - 10:00 LST, which would be about three hours after the sun has risen, consistent with a boundary layer mixing mechanism.

Figure 4 provides hodographs of wind with height throughout the first two km of the atmosphere between 12:00 UTC on the 8th June and 12:00 UTC on the 9th June; the soundings were taken at Perth Airport, which is the nearest station to the South WA station group to provide wind soundings. The 8th June 12:00 UTC hodograph shows surface northerlies of \approx 6 kn, becoming west to northwesterlies of over 20 kn 2.4 km above the surface. A forecaster basing a model edit of the following days winds on this sounding would therefore gradually strengthen the westerly compo-

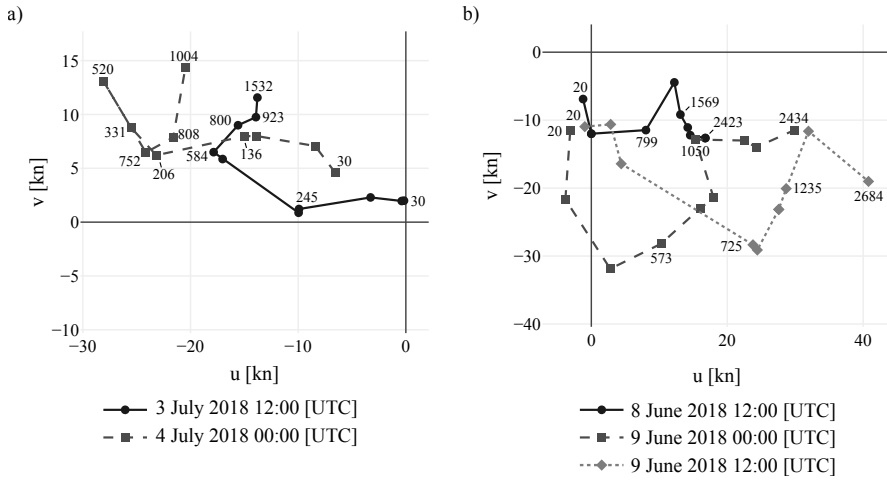


FIGURE 4 Hodographs showing change in winds with height at a), Darwin Airport, and b), Perth Airport.

233 nent of the surface winds in the hours after sunrise. However, the subsequent sounding at 00:00 UTC on the 9th of
 234 June shows that the winds acquire a strong northerly component of 30 kn in the first 500 m of the atmosphere, with
 235 the final sounding indicating a strong northwesterly wind at 725 m persisting until 12:00 UTC. In Fig. 3 d), the Official
 236 perturbations from 04:00 to 07:00 UTC show stronger westerly perturbations than either ACCESS or ECMWF, im-
 237 proving the amplitude of Official's diurnal wind cycle. However, the AWS perturbations are more northerly than those
 238 of Official, and so the Official forecast winds have been strengthened in a slightly incorrect direction. An explanation
 239 for this discrepancy is that the Official forecast for the southwest region of WA has been edited based on the June 8th
 240 12:00 UTC Perth Airport sounding, with the winds above the surface changing direction in the subsequent 12 hours.
 241 Note that the 3rd of August example is similar, although in this case the Official forecast slightly improves both the
 242 magnitude and direction of the 05:00 UTC wind perturbations.

243 Figure ?? presents the \overline{wpi} values and confidence scores for the Official versus ECMWF comparisons, i.e. \overline{wpi}_{OE}
 244 and $P(\overline{wpi}_{OE} > 0)$, for the airport stations, and airport station groups. The results for the airport stations are noisier
 245 than the analogous results for the coastal station groups in Figures 2 c) and d), although they do share some similarities.
 246 Official outperforms ECMWF at 01:00 and 02:00 UTC at both the Darwin airport station and the NT station group,
 247 although ECMWF outperforms Official between 08:00 and 14:00 UTC at Darwin and Brisbane airports, and the
 248 corresponding NT and QLD station groups, with the exception of the QLD station group at 12:00 UTC where $\overline{wpi}_{OE} =$

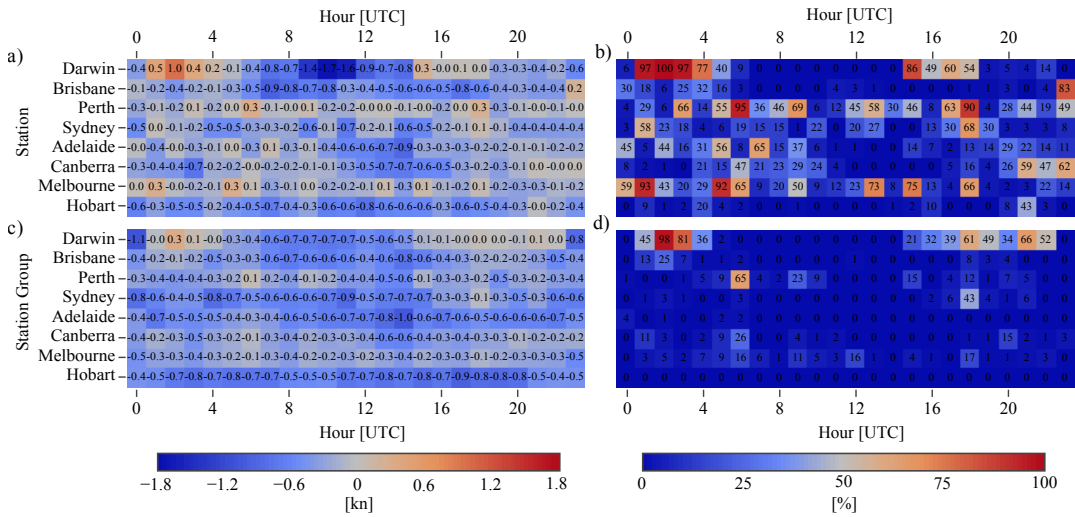
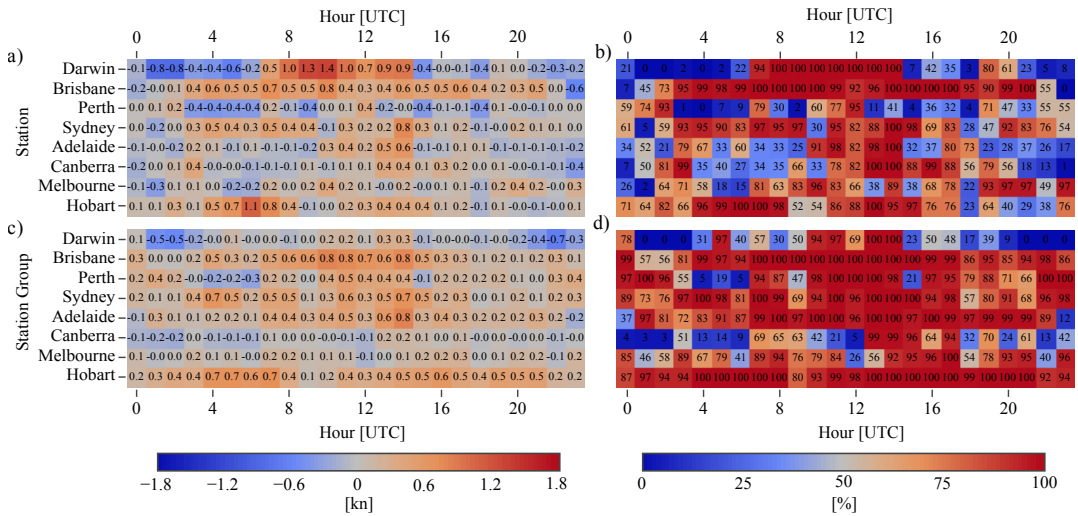


FIGURE 5 The \overline{wpi}_{OE} (Official versus ECMWF comparison) values, a) and c), and confidence scores, b) and d), for the airport stations, a) and b), and airport station groups, c) and d), respectively.

0. ECMWF also outperforms Official at Hobart airport at almost all hours of the day, and at Adelaide and Canberra airports from 11:00 to 14:00 UTC.

For the remaining stations and times, only the Perth airport station at 06:00 UTC and the Melbourne airport station at 01:00 UTC exhibit $\overline{wpi}_{OE} > 0$ with $P(\overline{WPI}_{OE} > 0) \geq 95\%$. However, in both cases $\overline{wpi}_{OE} = 0.3$, which is small compared to the maximum value of 1.0 which occurs at the Darwin airport station at 02:00 UTC. Furthermore, in both cases there is no clear pattern to the \overline{wpi}_{OE} values over the rest of the day. Given the random appearance of the \overline{wpi}_{OE} values, the *multiplicity problem* (Wilks, 2011, p. 178) requires care be taken before giving meaning to these two examples: i.e., given that we are calculating twenty four confidence scores for eight stations, then assuming WPI were uncorrelated across each station and hour we would expect to find $0.05 \times 24 \times 8 \approx 10$ instances where $P(\overline{WPI}_{OE} > 0) \geq 95\%$, even if \overline{WPI}_{OE} was in fact equal to zero. **Comment on performance versus ACCESS.**

For the airport station groups, ECMWF outperforms Official for the majority of station groups and times. The main exception is the Darwin airport station group, where Official outperforms ECMWF at 02:00 UTC, and there is ambiguity as to whether Official or ECMWF performs better at 01:00, 03:00 and 04:00 UTC, and from 15:00 to 22:00 UTC. In the analogous comparisons of Official and ACCESS (not shown), the airport station results are similarly noisy,



although the airport station group results are slightly more favourable to Official, with Official outperforming ACCESS from 10:00 to 12:00 UTC at the Brisbane station group, and fewer occasions overall where ACCESS outperforms Official than ECMWF does.

Figure ?? shows the \overline{wpi} values and confidence scores for the ECMWF versus ACCESS comparisons, i.e. \overline{wpi}_{EA} and $P(\overline{wpi}_{EA} > 0)$, for the airport stations, and airport station groups. As with the Official versus ECMWF comparison in Fig. ??, the results for the airport stations are noisy, but more often than not show that ECMWF outperforms ACCESS. The results for the airport station group show ECMWF usually outperforms ACCESS, the main exceptions being the Darwin and Canberra airport station groups.

At face value, the fact that ECMWF generally outperforms ACCESS at these scales is surprising, as ACCESS runs at a higher spatiotemporal resolution than ECMWF, and is calibrated for Australian conditions, and so one would expect ACCESS would better resolve small scale processes like the land-sea breeze and boundary layer mixing processes. However, these results are unsurprising if one considers the scales at which predictable atmospheric motion occurs, and the scales being resolved by AWS, ACCESS and ECMWF. The AWS data resolves motion with time scales as low as 10 minutes, and arbitrarily small spatial scales: it therefore includes highly unpredictable eddy turbulence. This explains why the results for the airport stations are noisier than for the airport station groups or coastal station groups. Furthermore, because ACCESS runs at a higher resolution than ECMWF, it includes additional scales of motion,

and therefore adds additional variability to the wind fields. Unless this additional variability in ACCESS is perfectly correlated with observations, the average of $|u_{AWS} - u_A|$ will therefore increase, unless this additional variability is compensated for by a reduction in bias, i.e. $|\bar{u}_{AWS} - \bar{u}_A|$ decreases. These ideas are discussed in greater detail in section 4. Note finally that the results for the Official versus ECMWF comparison in Fig. ?? largely mirror those of the ECMWF versus ACCESS comparison in Fig. ??, e.g. for the Darwin airport station and station group, Official outperforms ECMWF at the same times that ACCESS does, suggesting that either the Official forecast at these spatial scales is largely based on ACCESS, or that ECMWF is highly biased at these scales and times.

3.2 | Seasonal Comparison

Figure ?? provides the climatological wind perturbation index values, $cwpi$, and confidence scores, $P(CWPI > 0)$, for the coastal station groups for $cwpi_{OA}$, $cwpi_{OE}$ and $cwpi_{EA}$, which represent the the Official versus ACCESS, Official versus ECMWF, and ECMWF versus ACCESS comparisons, respectively. At the NT station group Official outperforms both ACCESS and ECMWF at 03:00 UTC with $cwpi_{OA} = cwpi_{OE} = 0.4$, $P(cwpi_{OA} > 0) = 94\%$ and $P(cwpi_{OE} > 0) = 93\%$. However, both ACCESS and ECMWF outperform Official at 23:00 and 00:00 UTC, consistent with the \overline{wpi} results in Fig. 2. The NT station group results are discussed in more detail in section 4.

At the North WA station group at 01:00, 03:00 and 04:00, Official outperforms ACCESS with confidence scores of 77, 78 and 90%, respectively; Official also outperforms ECMWF at 01:00 and 02:00 UTC with confidence scores above 99%. Figure 6 a) shows that ECMWF's poor performance at 01:00 and 02:00 UTC is simply due to its linear interpolation at these times, whereas Official's outperformance of ACCESS at 01:00, 03:00 and 04:00 is due to ACCESS's climatological diurnal cycle being slightly out of phase with that of the AWS observations, and the Official forecast appearing to correct for this somewhat. Both Official and ECMWF slightly exaggerate the magnitude of the climatological sea-breeze with ACCESS doing a good job in this regard.

At the South WA station group from 01:00 to 05:00 UTC, $cwpi_{OE}$ is positive with confidence scores of at least 88%, although $cwpi_{OA}$ is negative or zero at these times. Figure 6 b) shows that ECMWF underestimates the westerly perturbations at these times, with these perturbations likely associated with boundary layer mixing processes, as

discussed in section 3.1. Each of Official, ACCESS and ECMWF 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, $cwpi_{OA}$ and $cwpi_{OE}$ are at least 0.4 and 0.1 kn, respectively, with confidence scores of at least 95% and 75%, respectively. Figure 6 c) shows that these times correspond to a strange "dimple" in perturbation hodograph that is present in all four datasets. The Official hodograph closely resembles that of ACCESS, except for this dimple, which has been exaggerated relative to ACCESS. **Don't know what is going on here.** Figure 6 c) also shows that although ECMWF 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 01:00 to 05:00 UTC and 09:00 to 11:00 UTC both $cwpi_{OA}$ and $cwpi_{OE}$ are positive, with maximum values of 0.4 and 0.1 kn, although confidence scores do not exceed 88% and 65% respectively. Figure 6 shows that the Official forecast captures the amplitude of the perturbations from 01:00 to 05:00 UTC almost perfectly, matching the amplitude of the AWS perturbations better than both ACCESS and ECMWF. However, the Official diurnal cycle is slightly out of phase with the AWS cycle during this period, explaining why Official only slightly outperforms ACCESS in the results of Figures ?? a) and b).

4 | DISCUSSION

The methods developed in this study can be readily extended to analyse *just* the sea-breezes satisfying the operational definition above. For instance, to study the sea-breezes at a station near a coastline with inward pointing normal vector \hat{n} , the wind perturbation datasets could be restricted to just those days where the corresponding raw wind vector \mathbf{u} satisfies $\hat{n} \cdot \mathbf{u} > 0$ for at least one of the hours of that day.

How much time should forecasters spend on sea-breeze edits (if any)? What is the value of an improved diurnal cycle climatology? Improving the accuracy of forecast climatologies will have little value to the typical forecast user. Are there applications where a higher performing climatological forecast yields better outcomes, even if errors increase or even get worse?

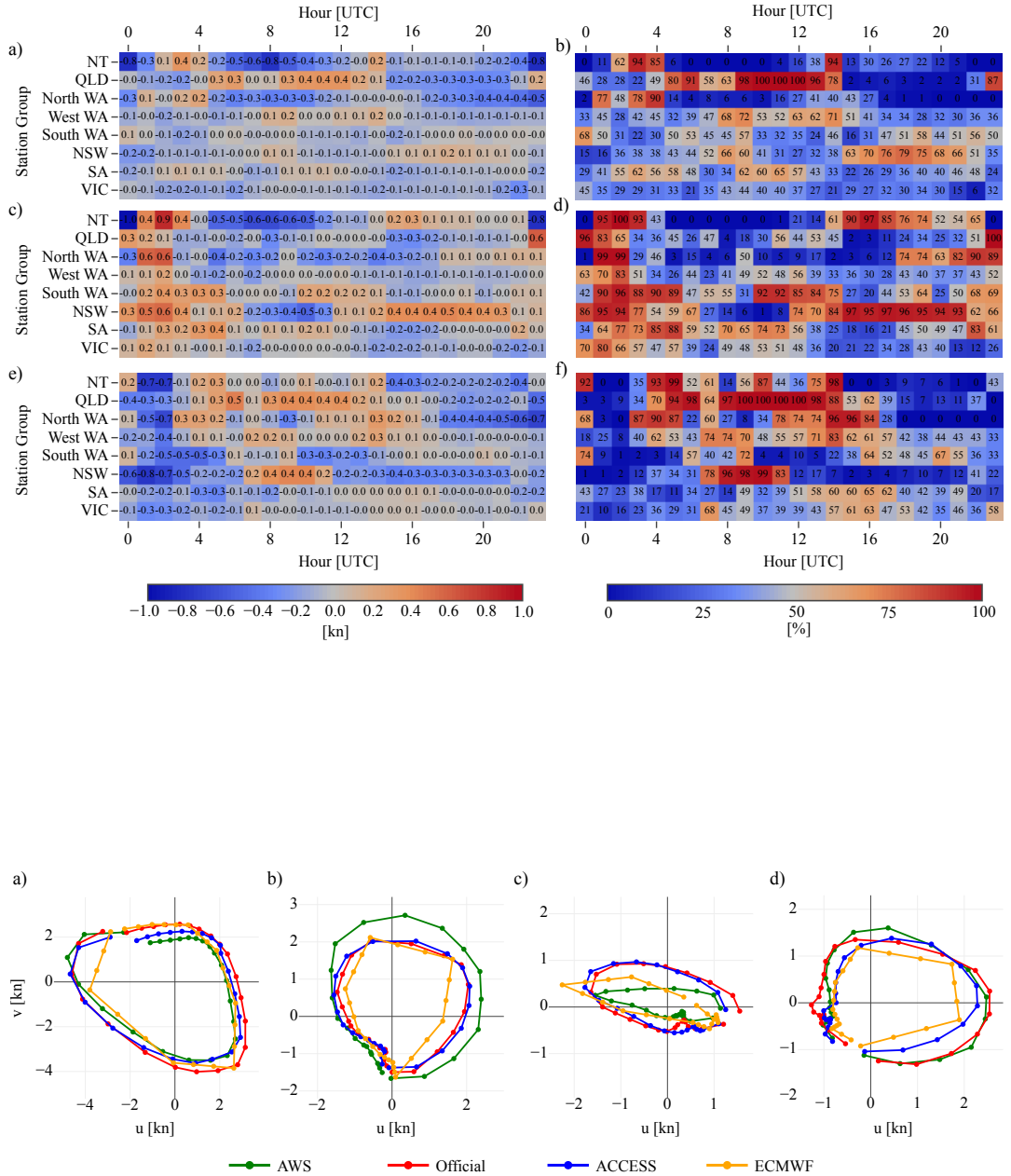


FIGURE 6 Climatological hodographs.

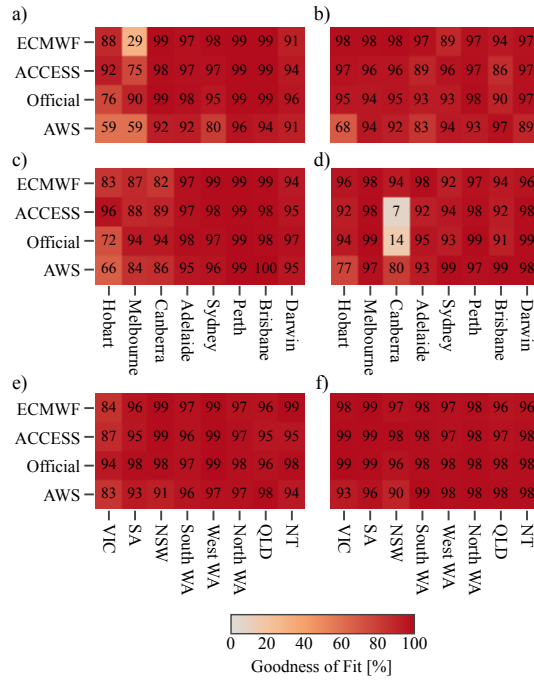


FIGURE 7 r squared.

Increasing the resolution of a forecast may reduce bias, but increase error.

Error, not bias, that generally matters for the forecast user. Standard methods for "improving" forecasts (adding parametrisations, increasing resolution) reduce bias, but actually increase errors!

Although they have similar definitions, \overline{WPI} and \overline{CWPI} measure different things. They do not converge as the length of the time period grows - they don't even necessarily approach the same sign. As a simple example, suppose that for each day, the observed and Official wind perturbations are given by $p_{AWS} = (5 \cos \omega t, 5 \sin \omega t)$ and $p_O = (6 \cos \omega t, 6 \sin \omega t)$, respectively. Furthermore, suppose that the ACCESS perturbations alternate between $p_A = (7 \cos \omega t, 7 \sin \omega t)$ and $p_A = (3 \cos \omega t, 3 \sin \omega t)$ from one day to the next. Then for any contiguous period of n days, $\overline{WPI} = 2 - 1 = 1$, but $\overline{CWPI} \approx -1$, with the approximation becoming exact for even n . Moreover $\overline{WPI} = 1$ with a confidence of 1, and using the bootstrapping procedure described above, the confidence that $\overline{CWPI} = -1$ approaches 1 as $n \rightarrow \infty$. This example shows that while the WPI and CWPI are sensitive both to random error and consistent biases between the different datasets, the CWPI becomes increasingly less sensitive to random error as

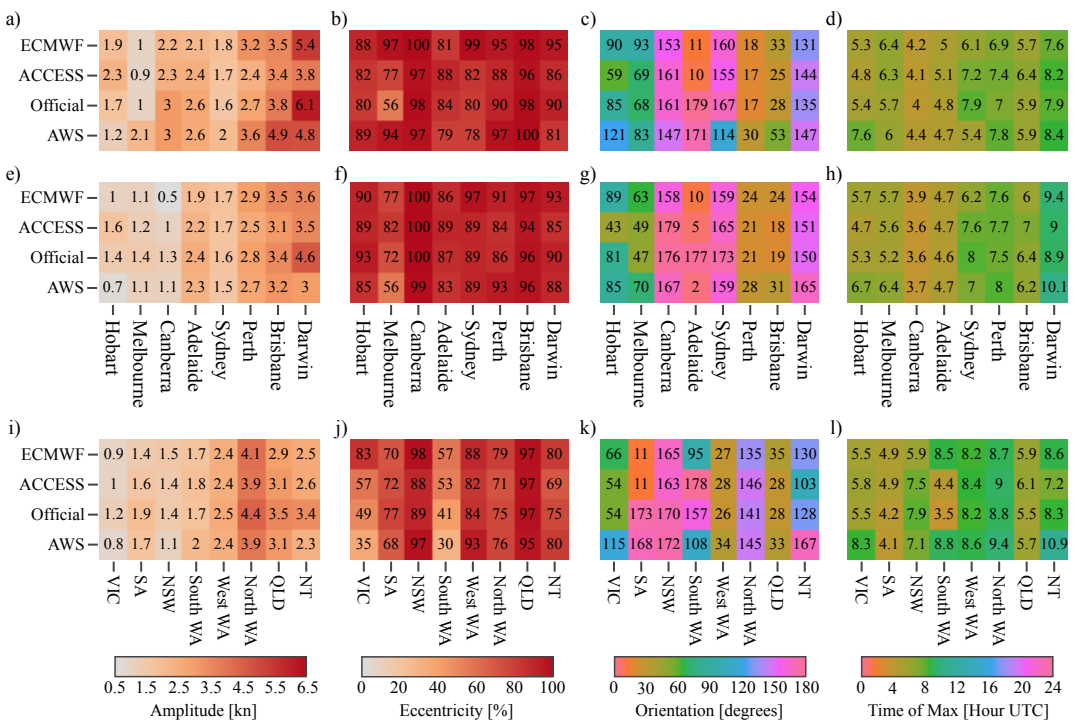


FIGURE 8 Ellipse fits.

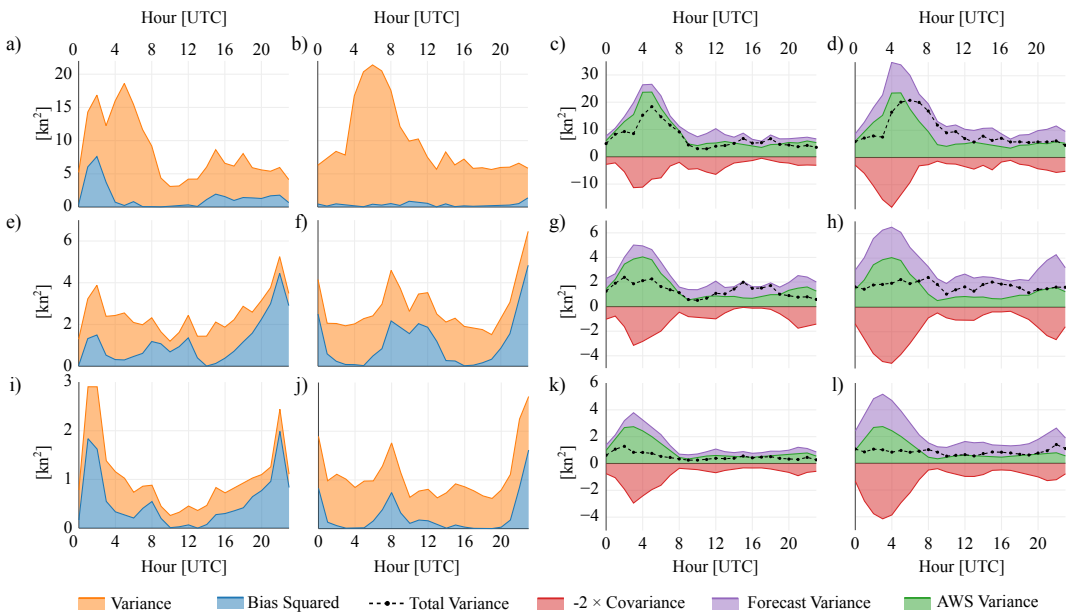


FIGURE 9 Actual perturbation standard deviation values. Note that official performs the worst at this scale!

the length of the time period being considered grows. Thus while the WPI arguably provides a more meaningful operational metric, as it measures the accuracy of actual forecast data, it may favour a more biased dataset over a less biased one, just because the internal variability of that dataset is lower. One consequence of this is that model data at a lower spatiotemporal resolution may outperform in \overline{WPI} model data of a higher resolution, purely because the internal variability is lower. In this way, the CWPI may actually provide more information about the performance of different forecasts.

Note that the Bureau has not yet moved to ensemble forecasting - and probabilistic verification methods are therefore not appropriate.

5 | CONCLUSION

In this report, a methodology for comparing the performance of Bureau forecasts of diurnal wind processes to unedited model guidance products has been developed and applied to a case study of the Darwin airport. The key results may

be summarised as follows.

1. During the dry season months of June, July and August 2017, the ECMWF sea-breeze is generally more accurate than that of the official forecast. However, during the wet season months of December, January and February 2017/18 this result is reversed, and the official forecast sea-breeze generally outperforms that of ECMWF.
2. In both seasons, boundary layer mixing processes are generally represented better in official forecasts than in ECMWF.
3. In the dry season, the climatological wind perturbations of the official forecast generally outperform those of ECMWF between 13:00 and 16:00 UTC. This is due to ECMWF not capturing the magnitude of the south-easterly mean perturbations.
4. During the wet season, the climatological wind perturbations of the official forecast generally outperform those of ECMWF at 11:00 UTC. This is due to ECMWF underestimating the magnitude of the mean land-breeze perturbation.

There a number of ways that this work could be extended. The most pressing would probably be to investigate whether the results presented here change when a more operational definition of the sea breeze is used in place of the entirely perturbation based definition used here: this could be done using the method described in section 2.

references

- Brown, A. L., Vincent, C. L., Lane, T. P., Short, E. and Nguyen, H. (2017) Scatterometer estimates of the tropical sea-breeze circulation near Darwin, with comparison to regional models. *Quart. J. Roy. Meteor. Soc.*
- Ebert, E. E. (2008) Fuzzy verification of high-resolution gridded forecasts: a review and proposed framework. *Meteor. Appl.*, **15**, 51–64. URL: <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/met.25>.
- Efron, B. (1979) Bootstrap methods: Another look at the jackknife. *The Annals of Statistics*, **7**, 1–26.

- 370 Gille, S. T., Llewellyn Smith, S. G. and Statom, N. M. (2005) Global observations of the land breeze. *Geophysical Research Letters*,
371 32. URL: <https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2004GL022139>.
- 372 Lynch, K. J., Brayshaw, D. J. and Charlton-Perez, A. (2014) Verification of european subseasonal wind speed forecasts. *Monthly*
373 *Weather Review*, **142**, 2978–2990. URL: <https://doi.org/10.1175/MWR-D-13-00341.1>.
- 374 Pinson, P. and Hagedorn, R. (2012) Verification of the ecmwf ensemble forecasts of wind speed against analyses and obser-
375 vations. *Meteor. Appl.*, **19**, 484–500. URL: <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/met.283>.
- 376 Smith, J. C., Thresher, R., Zavadil, R., DeMeo, E., Piwko, R., Ernst, B. and Ackermann, T. (2009) A mighty wind. *IEEE Power and*
377 *Energy Magazine*, **7**, 41–51.
- 378 Wilks, D. S. (2011) *Statistical methods in the atmospheric sciences*. [electronic resource]. International geophysics series: v. 100.
379 Elsevier.
- 380 Zwiers, F. W. and von Storch, H. (1995) Taking serial correlation into account in tests of the mean. *Journal of Climate*, **8**,
381 336–351. URL: [https://doi.org/10.1175/1520-0442\(1995\)008<0336:TSCIAI>2.0.CO;2](https://doi.org/10.1175/1520-0442(1995)008<0336:TSCIAI>2.0.CO;2).