

OPERATIONAL PREDICTION SYSTEM NOTES

Evaluation of Global Wave Probabilities Consistent with Official Forecasts

CHARLES R. SAMPSON,^a EFREN A. SERRA,^b JOHN A. KNAFF,^c AND JOSHUA H. COSSUTH^d

^aNaval Research Laboratory, Monterey, California

^bDeVine Consulting, Naval Research Laboratory, Monterey, California

^cNOAA/STAR, CIRA, Fort Collins, Colorado

^dOffice of Naval Research, Arlington, Virginia

(Manuscript received 12 March 2021, in final form 22 June 2021)

ABSTRACT: The U.S. Navy is keenly interested in analyses and predictions of waves at sea due to their effects on important tasks such as shipping, base preparedness, and disaster relief. U.S. Tropical Cyclone (TC) Forecast Centers routinely disseminate wind probabilities consistent with official TC forecasts worldwide, but do not do the same for wave forecasts. These probabilities are especially important at longer leads where TC forecast accuracy diminishes. This work describes global wave probabilities consistent with both the official TC forecasts and their wind probabilities. Real-time runs for 84 TCs between May 2018 and March 2019, with probabilities generated for 12- and 18-ft significant wave heights are used to calculate verification statistics. This results in 347, 319, 261, 214, 155, and 112 verification cases at lead times of 1, 2, 3, 4, and 5 days where each verification case consists of a $20^\circ \times 20^\circ$ latitude–longitude grid around the verifying TC position. When compared with wave probabilities generated solely by a global numerical weather prediction model, the wind probability–based algorithm demonstrates improved consistency with official forecasts and provides additional benefits. Those benefits include an improved capability to discriminate between 12- and 18-ft significant wave events and nonevents. The verification statistics also shows that the wind probability–based algorithm has a consistent high bias. How these biases can be reduced in future efforts is also discussed.


SIGNIFICANCE STATEMENT: The extreme wave heights associated with tropical cyclones are difficult to accurately forecast deterministically or probabilistically. To exacerbate matters, existing global ensemble systems cannot resolve the strongest winds in hurricanes and typhoons, and so they provide input to wave models that is inconsistent with official forecasts. This paper describes an algorithm that provides ensemble wind wave products that are both more realistic and consistent with official forecasts from tropical cyclone forecast centers. We show that this method provides improved identification of extreme wave events, which should provide improved input for ship navigation and hazard avoidance that saves both lives and property.

KEYWORDS: Forecast verification/skill; Forecasting techniques; Operational forecasting; Probability forecasts/models/distribution

1. Introduction

U.S. Navy operations are adversely impacted by high seas, especially those from tropical cyclones (TCs). In particular, the U.S. Navy is concerned about significant wave heights and their effects on safely routing ships, routine and emergency ship sorties, and Human Assistance Disaster Relief activities. Traditionally, wave model ensembles are run with Numerical Weather Prediction (NWP) model surface winds to produce significant wave heights and wave height probabilities around TCs. However, the NWP models are generally inconsistent with official forecasts from the U.S. TC forecast centers and lack the resolution to adequately capture large gradients in TC structure specified in the official forecasts (e.g., Tolman et al. 2005). This is problematic for forecasters and downstream

applications as the inconsistencies add confusion to an already stressful situation. To address this issue, the U.S. Navy's Fleet Numerical Meteorology and Oceanography Center (FNMO) implemented a deterministic global wave model forecast that uses postprocessed winds from U.S. TC forecast centers as input to WAVEWATCH III (WW3; Tolman 1991; Tolman et al. 2002; NCEP 2020). This algorithm is named for the WAVEWATCH III model (WW3) and its input TC winds from the U.S. TC forecast centers (OFCL), thus named WW3TCOFCL (Sampson et al. 2013). Faced with deficiencies in both the forcing winds and resolution for forecasting TC generated waves in the Northwest Australian region, the Australian Bureau of Meteorology (Zieger et al. 2018; Aijaz et al. 2019) designed a postprocessing method that correct wind distribution biases associated with TCs in the NWP model ensembles used to force their high resolution (8 km) wave model. For each ensemble member, the method constructs a synthetic vortex to replace the existing one, keeping the asymmetric flow in the numerical model. An evaluation of operational real-time runs found improvements in both TC wind and TC-generated wave probabilities, and importantly they had consistency between the winds from the NWP

 Denotes content that is immediately available upon publication as open access.

Corresponding author: Charles R. Sampson, Buck.Sampson@nrlmry.navy.mil

DOI: 10.1175/WAF-D-21-0037.1

For information regarding reuse of this content and general copyright information, consult the [AMS Copyright Policy \(www.ametsoc.org/PUBSReuseLicenses\)](http://www.ametsoc.org/PUBSReuseLicenses).

ensemble and the waves. These consistency and resolution issues are important to operations, and as yet there is no operational global wave model ensemble consistent with U.S. TC forecast center forecasts, wind probabilities associated with TC forecasts (DeMaria et al. 2013), and deterministic wave forecasts derived from U.S. forecast center forecasts (Sampson et al. 2013).

To address both consistency and resolution issues, a post-processing algorithm has been developed that constructs and inserts realistic wind structure in the vicinity of TCs out to 120 h. These winds are consistent with the forecasts from the U.S. TC forecast centers, which are frequently quite different in track, intensity and/or structure from the NAVGEM or other numerical model forecasts. These differences between official U.S. TC forecast and NWP forecasts can cause confusion for forecasters, warning managers and the general public in a time when coordinated and clear communication is of the utmost importance. The postprocessed winds can then be used in the Navy Global Environmental Model (NAVGEM, Hogan et al. 2014) global wave model ensemble to produce wave probability fields that are consistent with deterministic TC forecasts and wind probabilities generated at the U.S. TC forecast centers. The current incarnation of this algorithm is designed to run as a 20-member ensemble on a 0.25° global WW3 grid, the same as currently used at FNMOC. This is an intentional design to be consistent with the current NAVGEM global wave model ensemble so that implementation is simplified, extra computational resources are minimal, and the wind postprocessing algorithm can be run independently of the NAVGEM global wave model ensemble. Sampson et al. (2016) demonstrated that more ensemble members would be beneficial, but computational restrictions may not allow for expanding the ensemble. NRL has implemented the post-processing algorithm with the WW3 ensemble, executed in real time for over a year, and gathered runs for this evaluation. The algorithm, hereafter referred to as WW3TCOFCL ensemble, is described in section 2. Section 3 provides a description of how the data are used to conduct our evaluation. The result of the evaluations is provided in section 4, where individual cases and probabilistic verification is presented followed by conclusions and discussion of future work.

2. Algorithm description

The WW3TCOFCL ensemble follows the algorithm published in Sampson et al. (2016), except that the number of ensemble members has been reduced to 20 (the same number as in the FNMOC operational WW3 ensemble run using NAVGEM ensemble surface winds, hereafter referred to as the WW3NAVGEM ensemble) from 128. The WW3TCOFCL ensemble grid has also been expanded to a global $0.25^\circ \times 0.25^\circ$ grid to match the operational WW3NAVGEM ensemble. These changes are made so that the algorithm adheres to computing and other resource constraints at FNMOC, and so that the algorithm could also be implemented within the current WW3NAVGEM ensemble job instead of as a completely separate algorithm. Expanding the application to a global grid and reducing the number of ensemble members to 20 introduced

major changes to the algorithm with potentially adverse effects. Also, there have been important changes (new sensors and new methods) in wind structure analysis that occurred at the Joint Typhoon Warning Center since the original evaluation that could potentially change the performance of the WW3TCOFCL ensemble. And finally, the global grid allows waves to propagate around the world as they do in the real world while the limited domains in Sampson et al. (2016) did not. All these changes require vetting since their overall effects on performance are uncertain.

To summarize the current WW3TCOFCL ensemble algorithm: First, 20 forecast ensemble members from the original 1000 generated using the wind speed probability (WSP) algorithm (DeMaria et al. 2013) are randomly selected. Each WSP ensemble member is made available to the WW3TCOFCL deterministic model (Sampson et al. 2013) independently to create each ensemble member. The ensemble member is essentially the same as an official forecast defined at 0, 12, 24, 36, 48, 72, 96, and 120 h with the extent of the circulation extending to 20 kt ($1 \text{ kt} = 0.514 \text{ m s}^{-1}$) at the radius of outermost closed isobar specified in the TC analysis. Hourly TC forecast wind fields are created and interpolated to high-resolution hourly storm-scale gridded fields using O'Reilly and Guza (1993) tessellation. Then, NAVGEM ensemble surface wind fields are postprocessed by removing the NWP model TC vortex from each member's set of forecast fields. Location is determined by using predicted centers from the National Centers for Environmental Prediction (NCEP) vortex tracker (Marchok 2002). The entire area out to the analyzed radius of outermost closed isobar is removed at all forecast times. This is done to remove geographical displaced and structurally different NAVGEM ensemble forecasts so that only the background field remains. The removed TC vortex is replaced with bilinear interpolated data from the borders of the removed area. The final step of the gridded surface wind processing is inserting the hourly storm-scale gridded fields (one for each active TC) into the NAVGEM 10 m winds (originally at 1° resolution) to a $0.25^\circ \times 0.25^\circ$ global grid for WW3 v5.16—the operational version at FNMOC during 2018 and 2019. Even this resolution is insufficient to resolve the highest winds and waves, especially with TCs that have small eyes. The resultant set of gridded surface wind field forecasts at 1-h forecast intervals provide the wind forcing for WW3 to generate ocean wave forecasts for each ensemble member. Ensemble member wave forecasts are then combined to produce probability fields of significant wave height exceeding a threshold (e.g., 12 or 18 ft) on a 1° resolution grid. This resolution is consistent with the current operational WW3NAVGEM ensemble probabilities available from FNMOC. An example of the 12-ft significant wave height probabilities on the right side of Fig. 1. Since we are only running 20 members of the WW3 ensemble, the probability fields are generated on a $1^\circ \times 1^\circ$ global grid to reduce graininess noted in Sampson et al. (2016). Still, this graininess is visible at longer lead forecast times such as the 96-h WW3TCOFCL ensemble forecast probabilities shown in Fig. 1.

The entire 10-m wind field preparation process takes just a few minutes on a Cray XC-30, and an estimated 1 h of wall-time to run both the wind field preparation and the 20 WW3

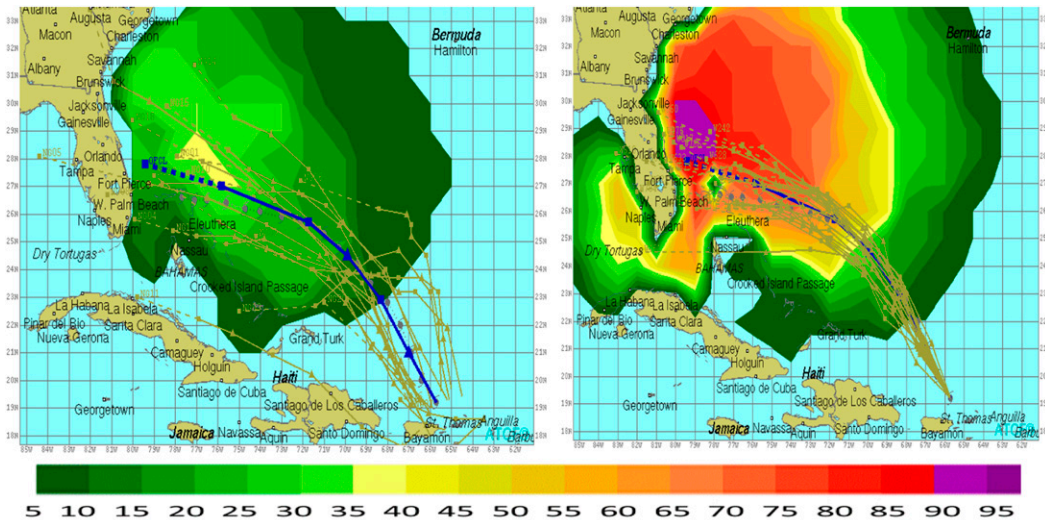


FIG. 1. (left) WW3NAVGEN ensemble and (right) WW3TCOFCL ensemble 96-h forecast 12-ft significant wave height probabilities for Dorian (AL052019) at 0000 UTC 29 Aug 2019. National Hurricane Center forecast track (blue) is shown for reference. Also, NAVGEN ensemble TC tracks and wind probability realizations generated by the U.S. TC forecast center wind probabilities (brown) are included. Probability (%) color bar is shown at the bottom.

ensemble members using 16 processors per ensemble member. Although attempts are made to warm start the WW3TCOFCL ensemble every 12 h using the previous 12-h forecast, this was not feasible when NRL computer resources became unavailable for extended periods of time. In these instances, the WW3TCOFCL ensemble was cold started with potentially adverse effects on seas and swell in the early forecast times. These effects become less important beyond 24 h, but they are worth noting as they are plainly visible in visual inspection.

3. Evaluation data

The WW3TCOFCL ensemble was run in real time on 84 TCs that existed between May 2018 and March 2019. NRL was able to produce forecast data in the vicinity of TCs in all regions of the globe. As with most nonoperational real-time NWP systems, NRL had issues with data acquisition and unscheduled computer downtime. As a result of this computer downtime, the evaluation set has periodic gaps resulting in some artifacts from the many WW3 cold starts, some of which are visible in our evaluation. Since the WW3TCOFCL ensemble was run on the same grid and has the same number of members as the WW3NAVGEN ensemble, verification of head-to-head cases will provide insight into both ensembles.

For ground truth we use the WW3TCOFCL deterministic model analysis of significant wave height in feet (ft; 1 ft = 0.3048 m), as that is the parameter most commonly used in Navy operations. Noting again that the WW3TCOFCL deterministic model uses postprocessed winds forecasted by U.S. TC forecast centers. Since the U.S. Navy is most concerned about significant wave heights in ship routing, we chose to evaluate significant wave height probabilities. We present results using WW3TCOFCL deterministic model significant

wave height analyses, but we also evaluated results against WW3NAVGEN deterministic analyses. The WW3NAVGEN deterministic model analyses assimilate altimeter data (Cummins and Wittmann 2009), but little difference was found between results using the WW3TCOFCL and WW3NAVGEN deterministic model analyses as ground truth. The 12- and 18-ft thresholds chosen for evaluation are not necessarily the thresholds used for operational forecasting, but span a reasonable range of significant wave heights associated with TCs and are routinely available for the WW3NAVGEN ensemble.

To gather data with 12- and 18-ft significant wave heights, which are not common in the tropics, our verification was limited to a $20^\circ \times 20^\circ$ box surrounding the verifying TC position. This area is likely larger than the TC wind field (Frank 1977) and also generally encompasses the extreme waves associated with TCs. In most cases a $20^\circ \times 20^\circ$ box will include many cases of zero probabilities in both the forecast and verification data (null cases), which affects results and their interpretation. The verification impacts of null cases are discussed in section 3. We also attempted this evaluation using a $10^\circ \times 10^\circ$ box around the verifying TC location, and found that this smaller area did not always encompass the TC-driven waves and highest significant wave height probabilities at longer forecast leads. At these longer leads, the area of high significant wave height probabilities can be both larger and dislocated from the $10^\circ \times 10^\circ$ box around the verifying position. Our evaluation was also limited to TCs with verifying intensities of 35 kt or greater intensity, which results in limiting the false alarm rates for both algorithms.

Although we verify WW3TCOFCL ensemble probabilities against WW3NAVGEN deterministic model significant wave height analyses (which assimilate altimeter wave heights), we do not attempt verification ensemble runs against buoys

and/or altimetry data explicitly, other than anecdotally. These observations have coverage issues that hinder verification of steep gradients and rare events, and can yield misleading results (see Sampson et al. 2013).

Table 1 provides a summary of the cases used in the verification. Each $20^\circ \times 20^\circ$ verification area represents 400 potential paired forecast and verification points, so the values in Table 1 are effectively 1/400th of the paired forecast points evaluated (minus an estimated 10% that verified over land and were removed from verification). Grid differences also accounted for minor differences in the matched pairs over water, 1 or 2 paired forecasts in approximately 10% of the cases. This represents differences of less than 0.1% and is ignored.

Summary statistics at the end of the results section are provided with significance using a two-tailed Student's t test. To remove correlation issues within the data, each $20^\circ \times 20^\circ$ (each with potentially 400 paired forecasts) is treated as a single case. Then the t tests are provided for the summary statistics—discrimination distance, relative/receiver operating characteristic (ROC) area under curve (AUC), and Brier score. No effort is made to account for the effects of serial correlation in the summary data, but the degrees of freedom are conservatively estimated using the number of cases rather than the number of matched pairs (i.e., counting every point in the $20^\circ \times 20^\circ$ box as a case).

4. Results

To demonstrate significant wave height forecasts we present results in three ways. We first present two cases that exemplify our real-time assessment of the differences between WW3TCOFCL ensemble and WW3NAVGEN ensemble significant wave height probabilities. We then verify WW3TCOFCL ensemble and WW3NAVGEN ensemble against WW3TCOFCL deterministic model significant wave height analyses, and for completeness, against WW3NAVGEN deterministic significant wave height analyses. For objective probabilistic verification statistics generation, we use the Model Evaluation Tools (MET; Development Testbed Center 2020) grid verification tools. We employ MET parameters reliability, likelihood, calibration, ROC, ROC AUC, and Brier score to obtain a reasonably complete summary of performance characteristics of each ensemble. Each of these metrics is described in section 4c.

a. Typhoon Maria (WP102018)—Intensifying to 140 kt

To highlight differences in the two algorithms (WW3 run with/without postprocessing) in an intensifying TC, we choose the Maria (WP102018). Maria, the eighth named storm of the 2018 typhoon season, was a powerful tropical cyclone that affected Guam, the Ryukyu Islands, Taiwan, and East China in early July 2018. Here we examine 96-h forecasts valid at 0000 UTC 9 July 2018, initiated at 0000 UTC 5 July when the storm was located southeast of Guam and forecast to intensify as it moved toward Okinawa. Figure 2 shows details of the WW3TCOFCL ensemble (left column) and WW3NAVGEN ensemble (right column) forecasts of 12-ft seas. Consistent among the TCs inspected (approximately 30 cases) are that the

TABLE 1. Numbers of WW3TCOFCL ensemble and WW3NAVGEN ensemble cases (each being a $20^\circ \times 20^\circ$ grid) gathered from real-time execution from 0000 UTC 26 May 2018 to 0000 UTC 18 Mar 2019 with 84 TCs occurring around the world during that period. Each 12- and 18-ft case was required to have both ensemble forecasts and verifying WW3TCOFCL deterministic model analysis.

| Tau | 0 | 24 | 48 | 72 | 96 | 120 |
|-------|-----|-----|-----|-----|-----|-----|
| 12 ft | 347 | 319 | 261 | 214 | 155 | 112 |
| 18 ft | 347 | 319 | 261 | 214 | 155 | 112 |

WW3NAVGEN ensemble input forecast tracks (Fig. 2 top row) and intensities both have reasonably large spread, but that ensemble member intensities tend to be too low, with intensities, unrealistically peaking near 70 kt for all members (Fig. 2 second row). In comparison, the WSP tracks and intensities appear to be well-calibrated with individual forecasts encompassing the forecast, and thus provide more realistic wind forcing input to WW3. In the case of Maria, this results in large areas relatively weak wind forcing input to the WW3NAVGEN ensemble, and much lower 12-ft significant height probabilities when compared to those from the WW3TCOFCL ensemble (Fig. 2, third row)—the issue is even more pronounced for higher significant wave height thresholds (not shown). These differences are not isolated, but seen throughout the dataset, especially for developing TCs.

b. Hurricane Ileana (EP112018)—Maintaining intensity at 40–45 kt

The majority of TCs are not forecast to intensify beyond 70 kt. To highlight differences between a weaker TC that is not forecast to intensify, we choose Hurricane Ileana's 48-h forecast valid at 0000 UTC 8 August 2018, initiated at 0000 UTC 6 August. Ileana was a remarkably small TC and the ninth tropical storm in the east Pacific in 2018 and during its life cycle tracked parallel to the Mexican coast. At this time, NHC forecasted Ileana to remain weak as it approached the Baja California Peninsula. In this case, the initial intensities used in the WW3NAVGEN ensemble encapsulate the initial estimate from NHC (Fig. 3, second row). The forecast track (Fig. 3, top row) and intensity spreads (Fig. 3, second row) are larger than those produced from the WSP algorithm. The 12-ft seas probabilities forecasts (Fig. 3, third row) from WW3TCOFCL ensemble are still noticeably higher probabilities in the vicinity of the highest observed wave heights (Fig. 3, bottom row). Much of the difference in 12-ft significant wave height probabilities generated from the WW3TCOFCL ensemble and WW3NAVGEN ensemble can be explained by larger forecast track spread in the WW3NAVGEN ensemble input.

c. Objective scores

Once the analyses are limited to $20^\circ \times 20^\circ$ boxes centered on the TC best track position, the probability forecasts can be intercompared using standard probability metrics such as reliability (Fig. 4), discrimination (Fig. 5), relative/receiver operating characteristic (ROC; Fig. 6), and summary or

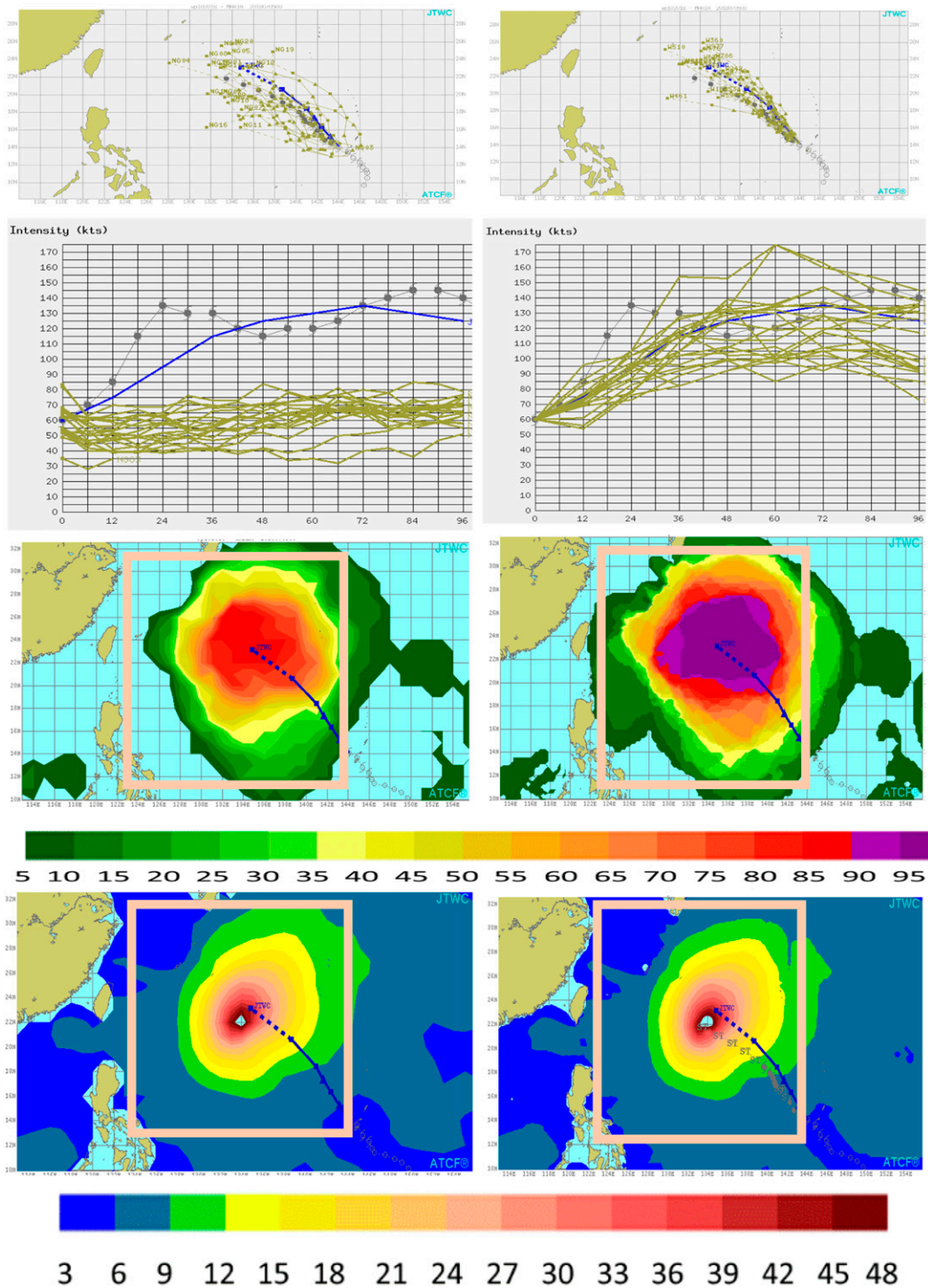


FIG. 2. (left) WW3NAVGENM ensemble and (right) WW3TCOFCL ensemble. (first row) Input 96-h forecast tracks, (second row) input forecast and verifying intensities (brown lines and black typhoon symbols), (third row) 96-h forecast 12-ft significant wave height probabilities, and (fourth row) verifying significant wave height (ft) analyses for WW3NAVGENM deterministic model in the left panels and WW3TCOFCL deterministic model in the right panels. Forecasts and analyses valid at 0000 UTC 9 Jul 2018 for Maria (WP102018). Significant wave heights for this case are above the end of the color bar (48 ft). Joint Typhoon Warning Center forecast track and intensity (blue) is shown for reference. Verifying track labeled “ST” for Super Typhoon is shown (brown) in bottom-right panel.

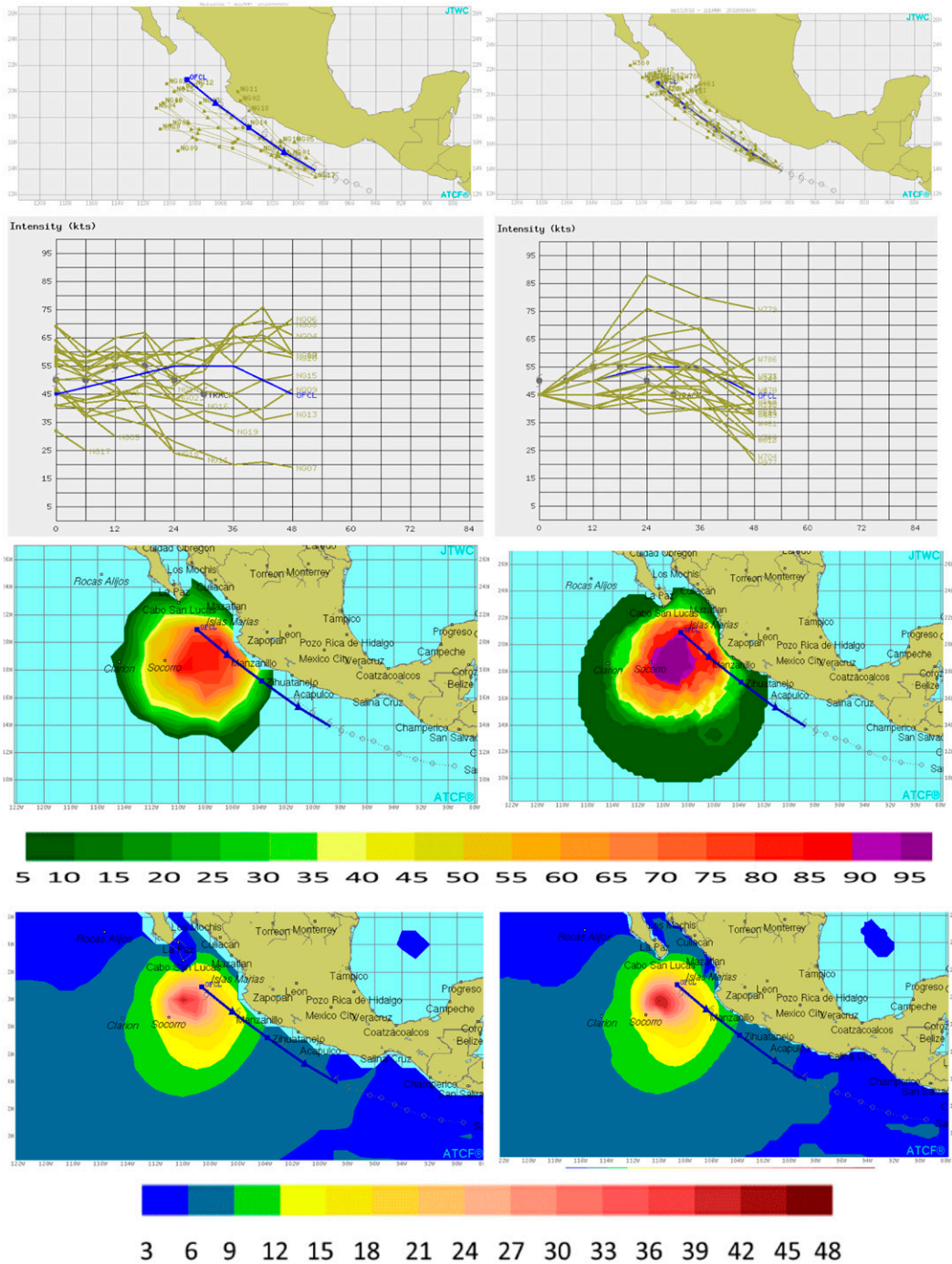


FIG. 3. (left) WW3NAVGEM ensemble and (right) WW3TCOFCL ensemble. (first row) Input 96-h forecast tracks, (second row) input forecast and verifying intensities (brown lines and black typhoon symbols), (third row) 96-h forecast 12-ft significant wave height probabilities, and (fourth row) verifying significant wave height (ft) analyses for WW3NAVGEM deterministic model in the left panels and WW3TCOFCL deterministic model in the right panels. Forecasts and analyses for Ileana (EP112018) 48-h forecast valid at 0000 UTC 8 Aug 2018. Significant wave heights for this case are above the end of the color bar (48 ft). National Hurricane Center forecast track and intensity (blue) is shown for reference.

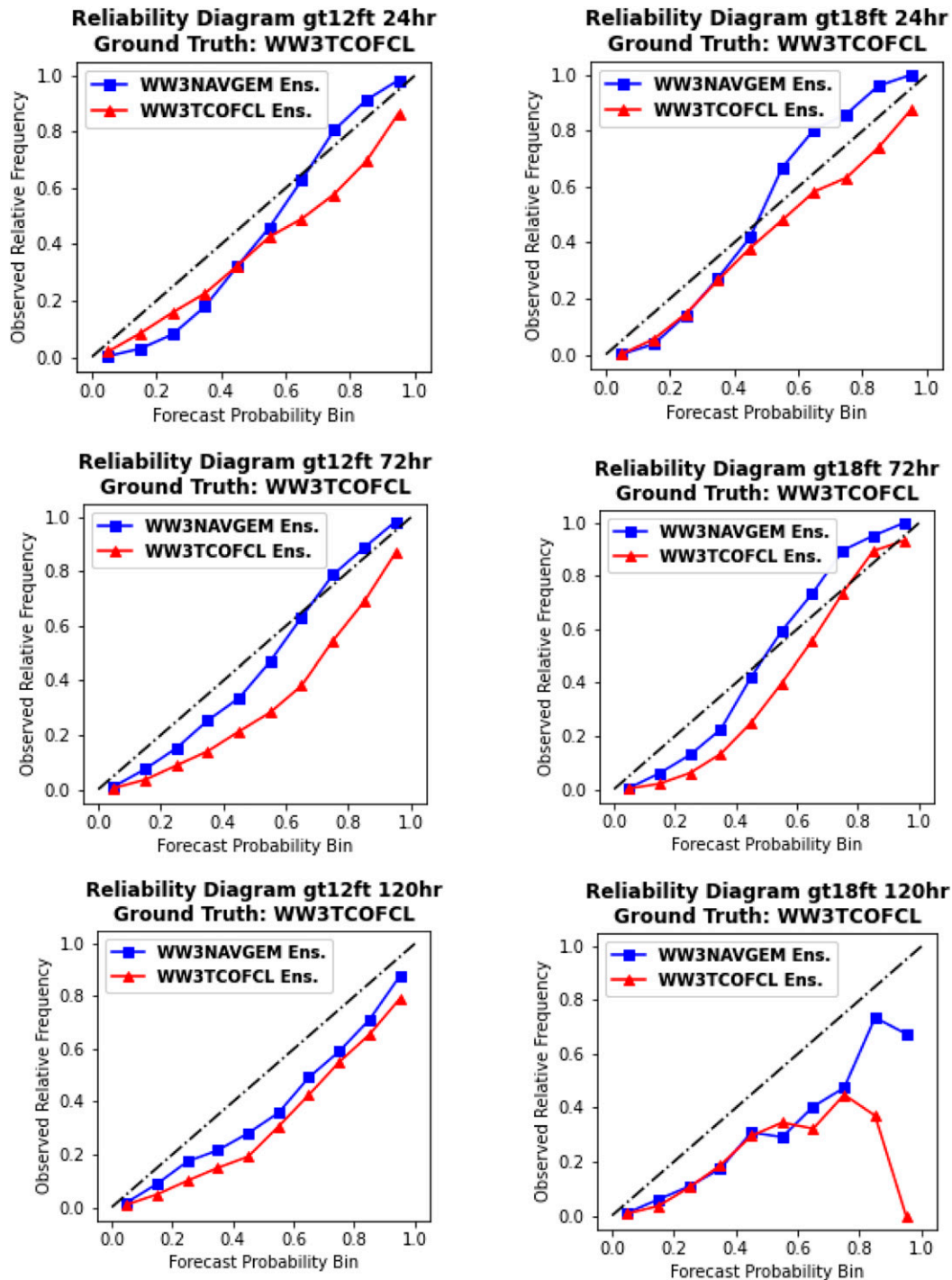


FIG. 4. Reliability diagrams for WW3TCOFCL ensemble and WW3NAVGEM ensemble (left) 12- and (right) 18-ft significant wave height with WW3TCOFCL deterministic model analysis employed as ground truth. Sequence progresses from (top) 24-, (middle) 72-, to (bottom) 120-h forecast. See Table 3 for numbers of head-to-head cases. Dashed lines represent perfect reliability.

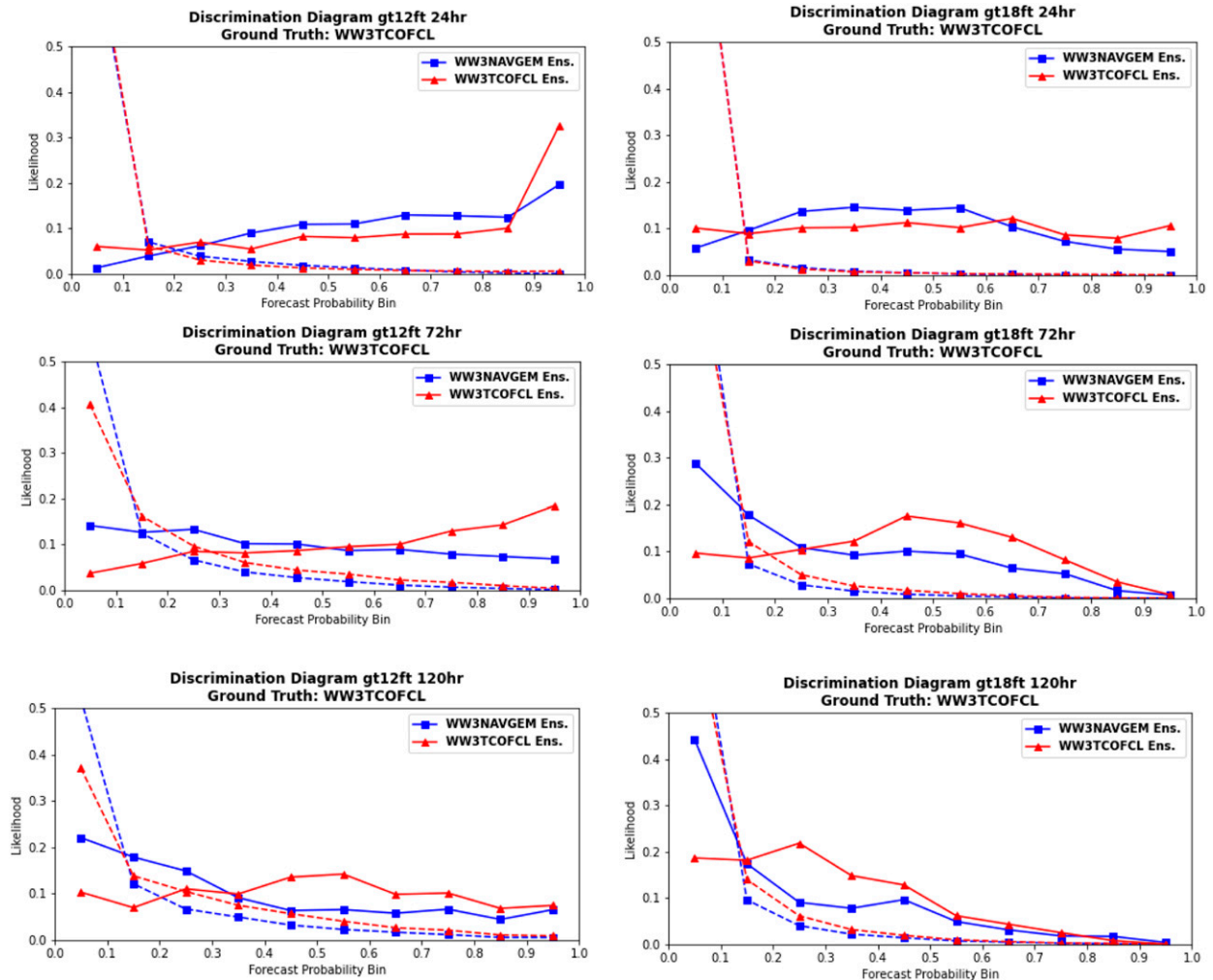


FIG. 5. Discrimination diagrams for WW3TCOFCL ensemble and WW3NAVGEN ensemble (left) 12- and (right) 18-ft significant wave height with WW3TCOFCL deterministic model analysis employed as ground truth. Solid lines indicate observed yes, dashed lines indicate observed no distributions. Sequence progresses from (top) 24-, (middle) 72-, to (bottom) 120-h forecast. See Table 3 for numbers of head-to-head cases.

derivative metrics such as discrimination distance, area under ROC curve, and Brier score (Fig. 7). Each of these metrics answers a specific question that we discuss below. Again, our evaluation uses MET, which in turn cites Wilks (2011) for most of its statistical algorithms. Results shown here are for a homogeneous dataset, meaning that the scores from the two different algorithms can be compared since they are for the same TCs on the same dates. For ground truth we again use analyzed significant wave heights from the WW3TCOFCL deterministic model (Sampson et al. 2013) as these have been shown to have realistic TC structure. We also performed the same tests using WW3NAVGEN deterministic model analyzed significant wave heights for verification, but somewhat surprisingly found consistent results in both statistical analyses for the metrics chosen. Finally, the evaluation was conducted for 0, 1, 2, 3, 4, and 5 day forecasts, but we limit presentation of the reliability, discrimination and ROC

charts to 1, 3, and 5 days and the results to those using the WW3TCOFCL model deterministic analysis as ground truth for brevity.

1) RELIABILITY

Reliability determines how well the probabilities compare to observed frequencies. On a reliability diagram, perfect reliability is a diagonal (1:1) line from lower left to upper right, biases are indicated by model reliability being below (high bias) and above (low bias) the 1:1 line, and forecast confidence is provided by the slope of model reliability relative to the 1:1 line, that is underconfident when the slope is less than and overconfidence when the slope is greater than one (Wilks 2011). Reliability for both 12- and 18-ft significant wave height probabilities is shown in Fig. 4. The reliability for WW3TCOFCL ensemble 12-ft significant wave height appears high biased (overforecasting in Wilks 2011) throughout. The

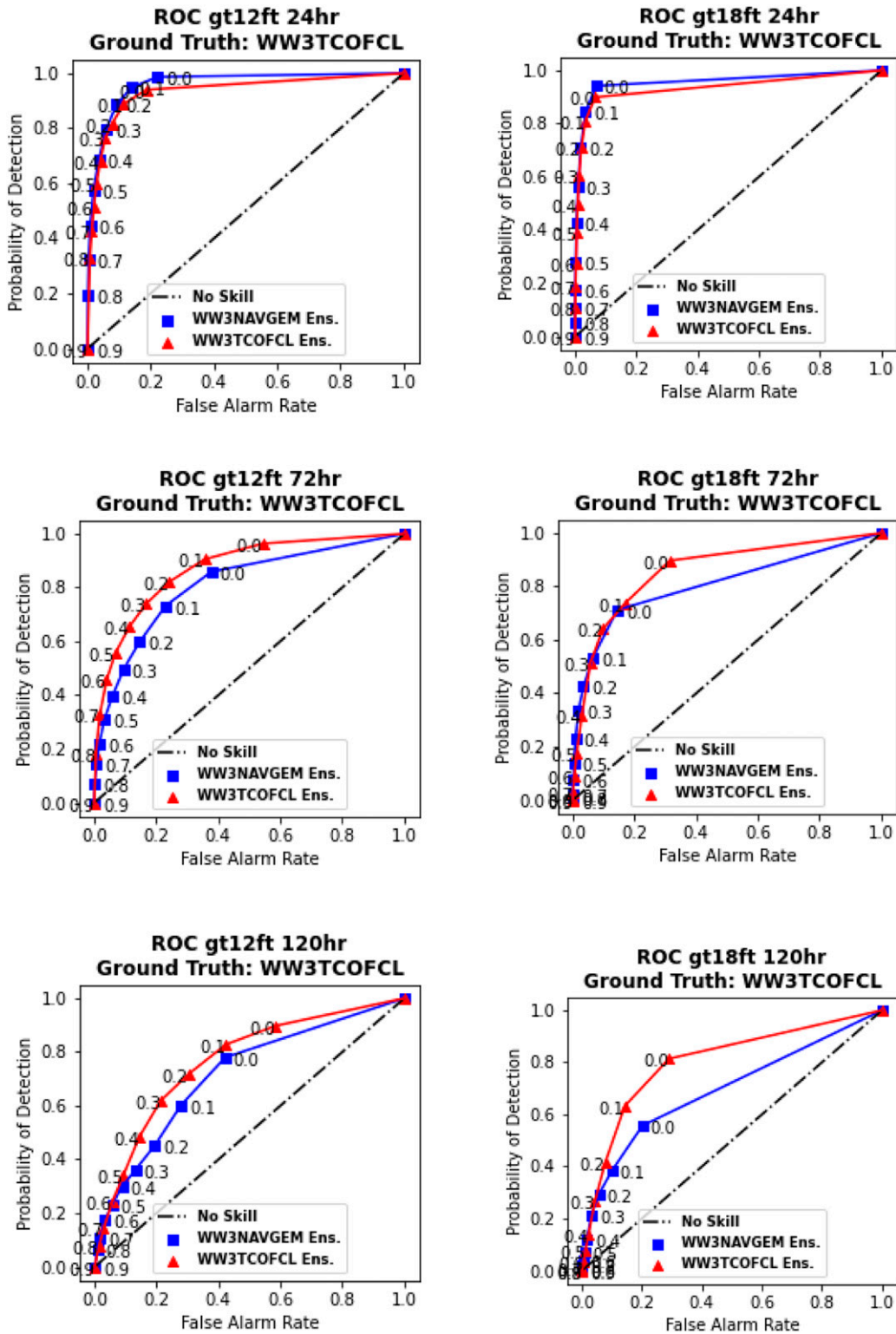


FIG. 6. ROC diagrams for (left) WW3NAVGEN ensemble and (right) WW3TCOFCL ensemble with WW3TCOFCL deterministic model analysis employed as ground truth. Sequence progresses from (top) 24-, (middle) 72-, to (bottom) 120-h forecast. Dashed line indicates no skill. See Table 3 for numbers of head-to-head cases.

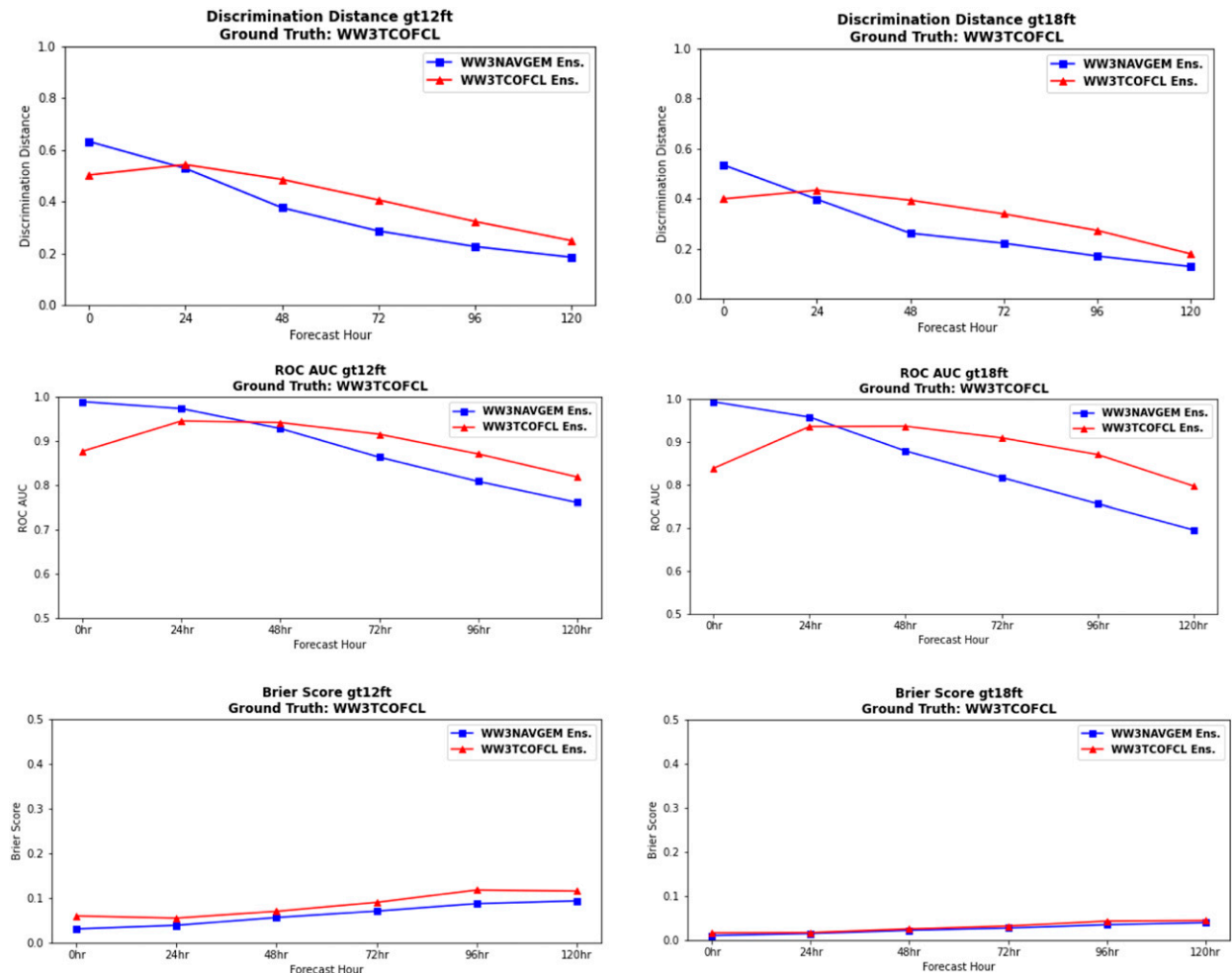


FIG. 7. (top) Discrimination distances, (middle) ROC AUC, and (bottom) Brier scores for WW3TCOFCL ensemble and WW3NAVGEN ensemble. (left) 12- and (right) 18-ft significant wave height shown with WW3TCOFCL deterministic model analysis employed as ground truth. See Table 3 for numbers of head-to-head cases.

WW3NAVGEN ensemble appears to overestimate low probabilities and underestimate higher probabilities in shorter forecast leads (underconfident), and overestimate probabilities like the WW3TCOFCL ensemble does at longer forecast leads. The number of cases drops precipitously for the 120-h 18-ft significant wave height probabilities above 80%, dropping to 400 head-to-head cases or one grid (SH112019 verifying at 1200 UTC 7 March 2019). So the reliability diagrams at 120 h for 18-ft significant wave height at the highest probability thresholds have few verification cases, reflected in the erratic changes in the reliability.

In the case of the WW3NAVGEN ensemble (underconfident in short-term forecast leads, overforecasting at longer-term forecast leads), the authors suspect that the ensemble is challenged by resolution in that circulations tend to be too large at longer forecast leads. In the case of WW3TCOFCL ensemble, the authors suspect several potential issues. The first is that WW3 is likely more appropriately run with 10-min mean wind speeds since it is developed to use

NWP fields. This is in contrast to U.S. official forecast center specified TC winds and wind probability realizations, which are both considered 1-min wind speed estimates. Operational forecasters use conversion rates such as 0.93 (Harper et al. 2010) to convert the 1-min wind speeds to 10-min wind speeds, and this conversion would likely reduce the high bias. Another potential source of bias is the statistical wind radius model (DRCL; Knaff et al. 2007, 2018) used in the wind probabilities. DRCL wind radii become more symmetric as the forecast progresses in time, and these symmetric forecasts could provide unrealistic durations for TC winds. DRCL will never emulate the large symmetry fluctuations seen in nature. A more appropriate treatment of the asymmetries, especially at longer forecast periods, could provide more realistic changes in fetch and duration of winds around TCs.

2) DISCRIMINATION

Discrimination is the relative frequency with which a forecast can discriminate between events and nonevents, where

perfect discrimination would entail no overlap between distributions of forecast probabilities for events and nonevents. Discrimination diagrams show these frequencies, where superior discrimination is indicated by separation between the events and nonevents. Figure 5 shows discrimination for probabilities from our two algorithms at 1, 3, and 5 days. One obvious trend is that the separation between events and nonevents becomes smaller as forecast length increases, as seen by the lines of the same color converging toward each other. The ability to discriminate between events and nonevents drops with forecast lead time for both algorithms.

3) DISCRIMINATION DISTANCE

An easier way to visualize and summarize the discrimination is to graph the discrimination distance (the difference between the average of the event and nonevents) for all forecast leads on one graph (Fig. 7). The discrimination distances for the WW3TCOFCL ensemble are lower than for WW3NAVGEN ensemble probabilities out to approximately 24 h, then remain approximately 10% higher for the longer leads. Significant differences using a two-tailed t test at the 5% level are present at all but the 24-h time period for 12-ft probabilities, and at all but 24- and 120-h time periods for the 18-ft probabilities. Discrimination distances for 12-ft are about 10% higher than for 18-ft significant wave heights at all forecast leads, indicating more skill in discrimination of 12-ft significant wave heights. The discrimination distances also decay at longer leads, indicating less skill in discrimination between events and nonevents at these forecast lead times.

4) ROC

ROC is another measure of the ability of the forecast to discriminate between two alternative outcomes, thus measuring resolution. It is not sensitive to bias in the forecast, so says nothing about reliability. A biased forecast may still have good resolution and produce a good ROC curve, which means that it may be possible to improve the forecast through calibration (e.g., correcting the bias). ROC can thus be considered as a measure of potential usefulness (Development Testbed Center 2020). A perfect ROC curve follows the y axis from 0 to 1, then across the top of the diagram to 1, 1. The ROC degrades for both algorithms as forecast time increases (Fig. 6). This is true for both the 12- and 18-ft thresholds.

5) ROC AUC

ROC AUC is a convenient way to summarize how a forecast discriminates between event/nonevent (Wilks 2011). Values can theoretically go from 0 to 1. A perfect score is 1, describing the area under a curve that passes from $x = 0, y = 0$, through $x = 0, y = 1$, to $x = 1, y = 1$. The ROC AUC for the no-skill diagonal is 0.5 (the area under a diagonal from $x = 0, y = 0$ to $x = 1, y = 1$ on a ROC diagram). As expected, the ROC AUC (Fig. 7) for the WW3TCOFCL ensemble probabilities is relatively low at analysis time due to the many cold starts in our dataset. The WW3TCOFCL ensemble ROC AUC improves until about the 48-h forecast time, then gradually drops off through 120 h. The WW3NAVGEN ensemble ROC AUC drops gradually

TABLE 2. Contingency table for WW3NAVGEN ensemble and WW3TCOFCL ensemble greater than 12-ft significant wave height probabilities for the 96-h forecast case shown in Fig. 2. Observed yes and observed no for the $20^\circ \times 20^\circ$ grid encompassing the verifying TC position in Fig. 2.

| Probability | WW3NAVGEN ensemble matched pairs | | WW3TCOFCL ensemble matched pairs | |
|-------------|-------------------------------------|----------------|--|----------------|
| | Observed yes | Observed no | Observed yes | Observed no |
| 0.05 | 0 | 82 | 0 | 56 |
| 0.15 | 0 | 58 | 0 | 37 |
| 0.25 | 0 | 51 | 0 | 38 |
| 0.35 | 0 | 42 | 0 | 37 |
| 0.45 | 10 | 32 | 1 | 33 |
| 0.55 | 16 | 18 | 1 | 35 |
| 0.65 | 14 | 1 | 5 | 29 |
| 0.75 | 28 | 2 | 14 | 12 |
| 0.85 | 32 | 0 | 15 | 5 |
| 0.95 | 14 | 0 | 78 | 2 |

through the forecast and is approximately 15% lower than the WW3TCOFCL ensemble between 72 and 120 h. Differences in the ROC AUC pass significance tests at all forecast periods except at 48 h for 12-ft and at 0 h for 18-ft significant wave height. The numbers of cases (each case representing an entire $20^\circ \times 20^\circ$ grid) for this ROC AUC at 48, 72, 96, and 120 h are all well below 200, so conclusions on significance tests 18-ft significant wave height should await more cases. Recall that the high bias in the WW3TCOFCL ensemble is not penalized in either the ROC or the ROC AUC, and that the ROC AUC is only used to discriminate between the event and nonevent. It is encouraging that the WW3TCOFCL ensemble probabilities maintain high ROC AUC out to 120 h since high bias, not depicted in either the ROC or ROC AUC, can be corrected through adjustments in the algorithm.

6) BRIER SCORES

Brier scores are another standard skill score for probabilistic forecasts, and measure both reliability and resolution (the ability to distinguish an event from a nonevent). The Brier score measures the mean square error of probabilities. Here again we use the WW3TCOFCL deterministic model analyses as ground truth. Brier scores range from 0 to 1, 0 being a perfect score. Brier scores for both ensembles evaluated are shown in Fig. 7 and they are within 3% of each other for both 12- and 18-ft thresholds. These generally rise as forecast time increase, indicating skill drops with forecast lead. The uptick in the WW3TCOFCL ensemble at analysis time is expected as this ensemble was frequently cold started throughout the testing period and the WW3TCOFCL ensemble (and its input) has little spread at analysis time. The WW3NAVGEN ensemble probabilities have slightly lower Brier scores than the WW3TCOFCL ensemble probabilities at all forecast times for the 12-ft significant wave height threshold, and scores from the two algorithms are within 3% of each other. Differences for 12-ft probabilities are significant at all forecast periods. Brier

TABLE 3. No. of cases for reliability, discrimination, and ROC shown in Figs. 4–6.

| Probability | WW3NAVGEM ensemble | | WW3TCOFCL ensemble | | WW3NAVGEM ensemble | | WW3TCOFCL ensemble | |
|-------------|--------------------|-------------|--------------------|-------------|--------------------|-------------|--------------------|-------------|
| | Observed yes | Observed no | Observed yes | Observed no | Observed yes | Observed no | Observed yes | Observed no |
| | 24 h 12 ft | | 24 h 12 ft | | 24 h 18 ft | | 24 h 18 ft | |
| 0.05 | 102 | 76 240 | 1326 | 79 715 | 59 | 99 054 | 169 | 99 583 |
| 0.15 | 169 | 7871 | 565 | 6578 | 129 | 3689 | 193 | 3299 |
| 0.25 | 319 | 4397 | 609 | 3330 | 246 | 1804 | 267 | 1432 |
| 0.35 | 517 | 3087 | 596 | 2083 | 344 | 876 | 306 | 747 |
| 0.45 | 813 | 2093 | 744 | 1429 | 412 | 560 | 335 | 513 |
| 0.55 | 1049 | 1499 | 870 | 1058 | 432 | 238 | 323 | 330 |
| 0.65 | 1484 | 923 | 960 | 799 | 428 | 119 | 384 | 235 |
| 0.75 | 1821 | 461 | 1082 | 671 | 392 | 53 | 329 | 165 |
| 0.85 | 2189 | 204 | 1357 | 560 | 353 | 21 | 355 | 89 |
| 0.95 | 5363 | 64 | 5717 | 645 | 430 | 0 | 564 | 46 |
| | 120 h 12 ft | | 120 h 12 ft | | 120 h 18 ft | | 120 h 18 ft | |
| 0.05 | 438 | 20 507 | 253 | 14 547 | 316 | 30 250 | 212 | 27 093 |
| 0.15 | 603 | 4695 | 327 | 5347 | 298 | 3549 | 229 | 5250 |
| 0.25 | 542 | 2512 | 369 | 4099 | 208 | 1494 | 280 | 2306 |
| 0.35 | 457 | 1750 | 390 | 2803 | 189 | 879 | 261 | 1184 |
| 0.45 | 427 | 1187 | 560 | 2200 | 213 | 579 | 244 | 743 |
| 0.55 | 522 | 839 | 744 | 1481 | 147 | 313 | 153 | 390 |
| 0.65 | 557 | 648 | 650 | 985 | 95 | 195 | 118 | 349 |
| 0.75 | 688 | 438 | 819 | 790 | 61 | 98 | 78 | 118 |
| 0.85 | 621 | 216 | 784 | 408 | 61 | 42 | 26 | 71 |
| 0.95 | 1149 | 217 | 1108 | 354 | 13 | 13 | 0 | 23 |
| | 72 h 12 ft | | 72 h 12 ft | | 72 h 18 ft | | 72 h 18 ft | |
| 0.05 | 440 | 39 933 | 134 | 28 743 | 326 | 59 769 | 123 | 52 648 |
| 0.15 | 622 | 9051 | 278 | 11 742 | 332 | 5526 | 140 | 8876 |
| 0.25 | 818 | 4929 | 548 | 7149 | 324 | 2156 | 216 | 3807 |
| 0.35 | 907 | 3019 | 661 | 4478 | 314 | 1132 | 320 | 1995 |
| 0.45 | 983 | 2031 | 798 | 3265 | 365 | 620 | 469 | 1259 |
| 0.55 | 1048 | 1324 | 878 | 2604 | 387 | 363 | 520 | 720 |
| 0.65 | 1261 | 751 | 1213 | 1631 | 347 | 170 | 477 | 336 |
| 0.75 | 1362 | 450 | 1498 | 1224 | 307 | 41 | 371 | 120 |
| 0.85 | 1463 | 252 | 1782 | 680 | 118 | 6 | 188 | 33 |
| 0.95 | 1934 | 71 | 3048 | 309 | 46 | 0 | 42 | 3 |

scores for 18-ft significant wave height thresholds are within 1% of each other with the WW3NAVGEM ensemble scoring lower (better). Differences are significant at 24 and 96 h, but just barely pass significance tests. In the case shown in Fig. 2, the Brier score for WW3NAVGEM ensemble (0.082098) is lower than for WW3TCOFCL ensemble (0.13089). This may seem counterintuitive as the WW3TCOFCL ensemble probabilities “look” to capture the 12-ft significant wave heights in the WW3TCOFCL deterministic model analysis from 96 h later. But upon further inspection (Table 2), the distribution of probability forecasts for WW3NAVGEM ensemble is skewed to lower probabilities so that it scores much higher in the large number of nonevents than the WW3TCOFCL ensemble probabilities for this case. The Brier score becomes inadequate for very rare (or very frequent) events because it does not sufficiently discriminate between small changes in forecast that are significant for rare events (Benedetti 2010). Thus, Brier score unfairly penalizes extremely rare (or common) event forecasts and can actually leads to conclusions that disagree

with our intuition (Jewson 2008), such as indicating that the WW3NAVGEM ensemble outperforms the WW3TCOFCL ensemble for the case in Fig. 2. The Brier scores are still useful in our evaluation as they confirm high bias in the WW3TCOFCL ensemble that, if corrected, could decrease the Brier scores. However, tuning specifically to Brier scores is not advised as that could result in undesired reduction in extreme event prediction (described as underconfident in Wilks 2011). An analog to this would be tuning a TC wind intensity consensus (e.g., see Sampson et al. 2008) to minimize mean forecast error when the most impactful errors are associated with rare and difficult to forecast rapid intensification events.

5. Conclusions and future work

A postprocessing algorithm for insertion of real-time operational TC surface wind forecasts into a $0.25^\circ \times 0.25^\circ$ global 20-member ensemble surface wind field is described. This algorithm was run twice a day (at 0000 and 1200 UTC) for

approximately one year and included active TCs from all basins. Each set of postprocessed wind fields was then used as wind input to WW3 in order to generate a 20-member ensemble of forecasted significant wave height fields out to 5 days. The resultant significant wave height fields from each ensemble member were then compiled to create significant wave height probabilities on a $1^\circ \times 1^\circ$ global grid.

Evaluation was performed using $20^\circ \times 20^\circ$ boxes around verifying positions of the TCs at each forecast day using the MET statistics package. Both WW3NAVGEN and WW3TCOFCL deterministic model analyses were used as ground truth for evaluation of the probabilities and little difference was found between evaluations with the two ground truth datasets. Case studies indicated large discrepancies frequently existed between input winds from the two algorithms. NAVGEM ensemble tracks and intensities generally had large spreads, and certainly larger than those generated by the WSP algorithm that are used in the WW3TCOFCL ensemble for weaker TCs. WW3NAVGEN ensemble input intensities were generally low-biased for intense TCs as the NAVGEM ensemble resolution was challenged to represent steep wind gradients in relatively small TCs. Large discrepancies also existed between significant wave height probabilities generated by each of the ensemble forecasts. The WW3NAVGEN ensemble significant wave height probabilities tended to be more widespread and lower in magnitude than those from the WW3TCOFCL ensemble.

In objective evaluation, reliability diagrams show that WW3NAVGEN ensemble overestimated low probabilities and underestimated higher probabilities in short-range forecasts, then generally overestimated probabilities by 5 days. WW3TCOFCL ensemble generally overestimated all probabilities throughout the entire forecast. Brier scores for WW3NAVGEN ensemble were a few percent better than WW3TCOFCL ensemble at 12-ft significant wave height forecasting at all forecast lengths, but inspection of individual cases indicated that those scores were heavily influenced by forecasts of very low probability for nonevents (no 12- or 18-ft significant wave height in ground truth). Brier scores for 18-ft significant wave height were within about 1% at all forecast lengths. ROC curves and ROC AUC indicated that discrimination between events and nonevents degrades with forecast period for both sets of probabilities, but that WW3TCOFCL ensemble forecast generally appeared better at discriminating events from nonevents beyond 24 h. These results are confirmed by the discrimination diagrams, discrimination distances, and significance tests for discrimination distances.

The WW3TCOFCL ensemble high bias noted in the reliability diagrams is likely correctable. Whether by converting the WW3TCOFCL ensemble input 1–10-min mean winds that are more representative of NWP model winds, by replacing the Wind Radii CLIPER Model (DRCL) with more realistic wind distribution realizations, or by applying some combination of the above, the high bias can be addressed. Also, the validation package developed in this work could be modified to validate whether changes in algorithms upstream of the WW3 ensembles (e.g., the WSP algorithm and the NAVGEM ensemble) adversely affect the significant wave height probabilities.

Operational forecasts are certain to improve in the future through use of new sensors, improved NWP representation of the vortex, and more advanced postprocessing in the wind probability algorithm—all of which can affect these ensembles. Construction of TC-specific significant wave height probability verification was time-consuming, but the process to achieve this is in place and could be used as is or improved upon to validate TC-specific wave probabilities in the future. And addition of Object-Based Diagnostic Evaluation (MODE) verification available in MET may compliment the evaluation done within this work as it follows features (e.g., TCs) and reports statistics different than those here when comparing the features. That evaluation would be similar to and hopefully more rigorous than the 12-ft sea radii evaluation against operational NHC estimates as done in Sampson et al. (2016).

Acknowledgments. This effort is dedicated to our late friend and colleague Paul Wittmann, WAVEWATCH expert and all around great person. Publication of this work was graciously funded by the Office of Naval Research, Program Elements 0602435N and 0603207N and NOAA/NESDIS base funding. We thank Chuck Skupniewicz at Fleet Numerical Meteorology and Oceanography Center, who provided both encouragement and funding to keep the WW3TCOFCL ensemble project going. Thanks also to Chris Landsea at the National Hurricane Center and his advocacy for this type of effort—it is greatly appreciated. We also thank John Gotway and the entire MET Team, without their help and software this evaluation would be significantly more difficult. Finally, kudos to Jim Hansen for suggesting expanding the deterministic WW3TCOFCL model to generate probabilities for sortie decisions. The views, opinions, and findings contained in this report are those of the authors and should not be construed as an official National Oceanic and Atmospheric Administration or U.S. government position, policy, or decision.

Data availability statement. Both sets of ensemble significant wave height probabilities have been archived and are available on request; however, a nondisclosure agreement public release approval may be required to provide data.

REFERENCES

- Aijaz, S., J. D. Kepert, H. Ye, Z. Huang, and A. Hawksford, 2019: Bias correction of tropical cyclone parameters in the ECMWF ensemble prediction system in Australia. *Mon. Wea. Rev.*, **147**, 4261–4285, <https://doi.org/10.1175/MWR-D-18-0377.1>.
- Benedetti, R., 2010: Scoring rules for forecast verification. *Mon. Wea. Rev.*, **138**, 203–211, <https://doi.org/10.1175/2009MWR2945.1>.
- Cummings, J. A., and P. Wittmann, 2009: Navy implements data assimilation capability for its wave forecasting model. *JCSDA Quarterly*, No. 28, Joint Center for Satellite Data Assimilation, Camp Springs, MD, 2–3, accessed 6 October 2021, <https://static1.squarespace.com/static/5bad1a12c2ff616821035c9ff/t/5d1bc190f87d390001d7f83e/1562100113337/200909JCSDAQarterly.pdf>.
- DeMaria, M., and Coauthors, 2013: Operational tropical cyclone wind speed probabilities. Part I: Recent model improvements

- and verification. *Wea. Forecasting*, **28**, 586–602, <https://doi.org/10.1175/WAF-D-12-00116.1>.
- Development Testbed Center, 2020: Model evaluation tools (MET). Development Testbed Center, accessed 10 June 2021, <https://dtcenter.org/community-code/model-evaluation-tools-met>.
- Frank, W. M., 1977: The structure and energetics of the tropical cyclone. I: Storm structure. *Mon. Wea. Rev.*, **105**, 1119–1135, [https://doi.org/10.1175/1520-0493\(1977\)105<1119: TSAEOT>2.0.CO;2](https://doi.org/10.1175/1520-0493(1977)105<1119: TSAEOT>2.0.CO;2).
- Harper, B. A., J. D. Kepert, and J. D. Ginger, 2010: Guidelines for converting between various wind averaging periods in tropical cyclone conditions. WMO Tech. Doc. WMO/TD-1555, World Meteorological Society, 64 pp., https://library.wmo.int/index.php?lvl=notice_display&id=135#.X6gp8mR7mUk.
- Hogan, T. F., and Coauthors, 2014: The Navy Global Environmental Model. *Oceanography*, **27**, 116–125, <https://doi.org/10.5670/oceanog.2014.73>.
- Jewson, S., 2008: The problem with the Brier score. <https://arxiv.org/abs/physics/0401046v1>.
- Knaff, J. A., C. R. Sampson, M. DeMaria, T. P. Marchok, J. M. Gross, and C. J. McAdie, 2007: Statistical tropical cyclone wind radii prediction using climatology and persistence. *Wea. Forecasting*, **22**, 781–791, <https://doi.org/10.1175/WAF1026.1>.
- , —, and K. D. Musgrave, 2018: Statistical tropical cyclone wind radii prediction using climatology and persistence: Updates for the western North Pacific. *Wea. Forecasting*, **33**, 1093–1098, <https://doi.org/10.1175/WAF-D-18-0027.1>.
- Marchok, T. P., 2002: How the NCEP tropical cyclone tracker works. Preprints, *25th Conf. on Hurricanes and Tropical Meteorology*, San Diego, CA, Amer. Meteor. Soc., P1.13, https://ams.confex.com/ams/25HURR/techprogram/paper_37628.htm.
- NCEP, 2020: WAVEWATCH III model. NCEP, accessed 12 June 2020, <https://polar.ncep.noaa.gov/waves/wavewatch/>.
- O'Reilly, W. C., and R. T. Guza, 1993: A comparison of two spectral wave models in the Southern California Bight. *Coast. Eng.*, **19**, 263–282, [https://doi.org/10.1016/0378-3839\(93\)90032-4](https://doi.org/10.1016/0378-3839(93)90032-4).
- Sampson, C. R., J. L. Franklin, J. A. Knaff, and M. DeMaria, 2008: Experiments with a simple tropical cyclone intensity consensus. *Wea. Forecasting*, **23**, 304–312, <https://doi.org/10.1175/2007WAF2007028.1>.
- , P. A. Wittmann, E. A. Serra, H. L. Tolman, J. Schauer, and T. Marchok, 2013: Evaluation of wave forecasts consistent with tropical cyclone wind forecasts. *Wea. Forecasting*, **28**, 287–294, <https://doi.org/10.1175/WAF-D-12-00060.1>.
- , J. Hansen, P. A. Wittmann, J. A. Knaff, and A. Schumacher, 2016: Wave probabilities consistent with official tropical cyclone forecasts. *Wea. Forecasting*, **31**, 2035–2045, <https://doi.org/10.1175/WAF-D-15-0093.1>.
- Tolman, H. L., 1991: A third-generation model for wind waves on slowly varying, unsteady, and inhomogeneous depths and currents. *J. Phys. Oceanogr.*, **21**, 782–797, [https://doi.org/10.1175/1520-0485\(1991\)021<0782:ATGMFW>2.0.CO;2](https://doi.org/10.1175/1520-0485(1991)021<0782:ATGMFW>2.0.CO;2).
- , B. Balasubramanian, L. D. Burroughs, D. V. Chalikov, Y. Y. Chao, H. S. Chen, and V. M. Gerald, 2002: Development and implementation of wind generated ocean surface wave models at NCEP. *Wea. Forecasting*, **17**, 311–333, [https://doi.org/10.1175/1520-0434\(2002\)017<0311:DAIOWG>2.0.CO;2](https://doi.org/10.1175/1520-0434(2002)017<0311:DAIOWG>2.0.CO;2).
- , J. G. M. Alves, and Y. Y. Chao, 2005: Operational forecasting of wind-generated waves by Hurricane Isabel at NCEP. *Wea. Forecasting*, **20**, 544–557, <https://doi.org/10.1175/WAF852.1>.
- Wilks, D., 2011: *Statistical Methods in the Atmospheric Sciences*. 3rd ed. International Geophysics Series, Vol. 100, Academic Press, 704 pp.
- Zieger, S., D. Greenslade, and J. D. Kepert, 2018: Wave ensemble forecast system for tropical cyclones in the Australian region. *Ocean Dyn.*, **68**, 603–625, <https://doi.org/10.1007/s10236-018-1145-9>.