1  Title: Evaluation of Global Wave Probabilities Consistent with Official Forecasts

2  Authors: Charles R. Sampson[1], Efren A. Serra[2], John A. Knaff[3], Joshua H. Cossuth[4]

3  Affiliations:    [1]Naval Research Laboratory, Monterey, CA, USA

4                   [2]DeVine Consulting, Naval Research Laboratory, Monterey, CA, USA

5                   [3]NOAA STAR, CIRA, Fort Collins, CO, USA

6                   [4]Office of Naval Research, Arlington, VA, USA

7

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

9

10

11                                        Abstract

12         The U.S. Navy is keenly interested in analyses and predictions of waves at sea due to

13    their effects on important tasks such as shipping, base preparedness and disaster relief.  U.S.

14    Tropical Cyclone (TC) Forecast Centers routinely disseminate wind probabilities consistent with

15    official TC forecasts worldwide, but do not do the same for wave forecasts.  These probabilities

16    are especially important at longer leads where TC forecast accuracy diminishes.  This work

17    describes global wave probabilities consistent with both the official TC forecasts and their wind

18    probabilities.  Real-time runs for 84 TCs between May 2018 and March 2019, with probabilities

19    generated for 12-ft and 18-ft significant wave heights are used to calculate verification

20    statistics.  This results in 347, 319, 261, 214, 155, and 112 verification cases at lead times of 1, 2,

21    3, 4, and 5 days where each verification case consists of a 20x20 degree latitude longitude grid

22    around the verifying TC position.  When compared with wave probabilities generated solely by

23    a global numerical weather prediction model, the wind probability-based algorithm

24    demonstrates improved consistency with official forecasts and provides additional benefits.

25    Those benefits include an improved capability to discriminate between 12-ft and 18-ft

26    significant wave events and non-events.  The verification statistics also shows that the wind

27    probability-based algorithm has a consistent high bias.  How these biases can be reduced in

28    future efforts is also discussed.

29

30

31                                Significance Statement

32      The extreme wave heights associated with tropical cyclones are difficult to accurately forecast

33      deterministically or probabilistically. To exacerbate matters, existing global ensemble systems

34      cannot resolve the strongest winds in hurricanes and typhoons and provide input to wave

35      models that is inconsistent with official forecasts. This paper describes an algorithm that

36      provides ensemble winds wave products that are both more realistic and consistent with

37      official forecasts from tropical cyclone forecast centers.  We show that this method provides

38      improved identification of extreme wave events, which should provide improved input for ship

39      navigation and hazard avoidance that saves both lives and property.

40

41      1.      Introduction

42      U.S. Navy operations are adversely impacted by high seas, especially those from tropical

43      cyclones (TCs).   In particular, the U.S. Navy is concerned about significant wave heights and

44      their effects on safely routing ships, routine and emergency ship sorties, and Human Assistance

45      Disaster Relief activities. Traditionally, wave model ensembles are run with Numerical Weather

46      Prediction (NWP) model surface winds to produce significant wave heights and wave height

47      probabilities around TCs.  However, the NWP models are generally inconsistent with official

48      forecasts from the U.S. TC forecast centers and lack the resolution to adequately capture large

49      gradients in TC structure specified in the official forecasts (e.g., Tolman et al. 2005).  This is

50      problematic for forecasters and downstream applications as the inconsistencies add confusion

                                                                                              3

51  to an already stressful situation.  To address this issue, the U.S. Navy's Fleet Numerical

52  Meteorology and Oceanography Center (FNMOC) implemented a deterministic global wave

53  model forecast that uses post-processed winds from U.S. TC forecast centers as input to

54  WAVEWATCH III ® (WW3; Tolman 1991, Tolman et al. 2002, NCEP 2020).  This algorithm is

55  named for the WAVEWATCH III model (WW3) and its input TC winds from the U.S. TC forecast

56  centers (OFCL), thus named WW3TCOFCL (Sampson et al. 2013).  Faced with deficiencies in

57  both the forcing winds and resolution for forecasting TC generated waves in the Northwest

58  Australian region, the Australian Bureau of Meteorology (Zieger et al. 2018, Aijaz et al. 2019)

59  designed a post-processing method that correct wind distribution biases associated with TCs in

60  the NWP model ensembles used to force their high resolution (8 km) wave model.  For each

61  ensemble member, the method constructs a synthetic vortex to replace the existing one,

62  keeping the asymmetric flow in in the numerical model.  An evaluation of operational real-time

63  runs found improvements in both TC wind and TC-generated wave probabilities, and

64  importantly they had consistency between the winds from the NWP ensemble and the waves.

65  These consistency and resolution issues are important to operations, and as yet there is no

66  operational global wave model ensemble consistent with U.S. TC forecast center forecasts,

67  wind probabilities associated with TC forecasts (DeMaria et al. 2013), and deterministic wave

68  forecasts derived from U.S. forecast center forecasts (Sampson et al. 2013).

69      To address both consistency and resolution issues, a post-processing algorithm has been

70  developed that constructs and inserts realistic wind structure in the vicinity of TCs out to 120 h.

71  These winds are consistent with the forecasts from the U.S. TC forecast centers, which are

72  frequently quite different in track, intensity and/or structure from the NAVGEM or other

4

73    numerical model forecasts.  These differences between official U.S. TC forecast and NWP

74    forecasts can cause confusion for forecasters, warning managers and the general public in a

75    time when coordinated and clear communication is of the utmost importance.  The post-

76    processed winds can then be used in the Navy Global Environmental Model (NAVGEM, Hogan

77    et al. 2015) global wave model ensemble to produce wave probability fields that are consistent

78    with deterministic TC forecasts and wind probabilities generated at the U.S. TC forecast centers.

79    The current incarnation of this algorithm is designed to run as a 20-member ensemble on a 0.25

80    degree global WW3 grid, the same as currently used at FNMOC.  This is an intentional design to

81    be consistent with the current NAVGEM global wave model ensemble so that implementation is

82    simplified, extra computational resources are minimal, and the wind post-processing algorithm

83    can be run independently of the NAVGEM global wave model ensemble.  Sampson et al. (2016)

84    demonstrated that more ensemble members would be beneficial, but computational

85    restrictions may not allow for expanding the ensemble.  NRL has implemented the post-

86    processing algorithm with the WW3 ensemble, executed in real-time for over a year, and

87    gathered runs for this evaluation. The algorithm, hereafter referred to as WW3TCOFCL

88    Ensemble, is described in section 2.  Section 3 provides a description of how the data is used to

89    conduct our evaluation. The result of the evaluations is provided in section 4, where individual

90    cases and probabilistic verification is presented followed by conclusions and discussion of

91    future work.

92

93    2.      Algorithm Description

94    The WW3TCOFCL Ensemble follows the algorithm published in Sampson et al. (2016),

95    except that the number of ensemble members has been reduced to 20 (the same number as in

96    the FNMOC operational WW3 ensemble run using NAVGEM Ensemble surface winds, hereafter

97    referred to as the WW3NAVGEM Ensemble) from 128.  The WW3TCOFCL Ensemble grid has

98    also been expanded to a global 0.25x0.25 degree grid to match the operational WW3NAVGEM

99    Ensemble.  These changes are made so that the algorithm adheres to computing and other

100   resource constraints at FNMOC, and so that the algorithm could also be implemented within

101   the current WW3NAVGEM Ensemble job instead of as a completely separate algorithm.

102   Expanding the application to a global grid and reducing the number of ensemble members to

103   20 introduced major changes to the algorithm with potentially adverse effects.  Also, there

104   have been important changes (new sensors and new methods) in wind structure analysis that

105   occurred at the Joint Typhoon Warning Center since the original evaluation that could

106   potentially change the performance of the WW3TCOFCL Ensemble.  And finally, the global grid

107   allows wave to propagate around the world as they do in the real world while the limited

108   domains in Sampson et al. (2016) did not.   All these changes require vetting since their overall

109   effects on performance are uncertain.

110   To summarize the current WW3TCOFCL Ensemble algorithm: First, 20 forecast ensemble

111   members from the original 1000 generated using the Wind Speed Probability (WSP) algorithm

112   (DeMaria et al. 2013) are randomly selected. Each WSP ensemble member is made available to

113   the WW3TCOFCL deterministic model (Sampson et al. 2013) independently to create each

114   ensemble member. The ensemble member is essentially the same as an official forecast defined

115   at 0, 12, 24, 36, 48, 72, 96, and 120 h with the extent of the circulation extending to 20 kt at the

6

116    radius of outermost closed isobar specified in the TC analysis.  Hourly TC forecast wind fields

117    are created and interpolated to high-resolution hourly storm-scale gridded fields using O'Reilly

118    and Guza (1993) tessellation.   Then, NAVGEM Ensemble surface wind fields are post-processed

119    by removing the NWP model TC vortex from each member's set of forecast fields.  Location is

120    determined by using predicted centers from the National Centers for Environmental Prediction

121    (NCEP) vortex tracker (Marchok 2002).  The entire area out to the analyzed radius of outermost

122    closed isobar is removed at all forecast times.  This is done to remove geographical displaced

123    and structurally different NAVGEM Ensemble forecasts so that only the background field

124    remains. The removed TC vortex is replaced with bilinear interpolated data from the borders of

125    the removed area.  The final step of the gridded surface wind processing is inserting the hourly

126    storm-scale gridded fields (one for each active TC) into the NAVGEM 10 m winds (originally at 1

127    degree resolution) to a 0.25x0.25 degree global grid for WW3 v5.16 — the operational version

128    at FNMOC during 2018 and 2019. Even this resolution is insufficient to resolve the highest

129    winds and waves, especially with TCs that have small eyewalls.  The resultant set of gridded

130    surface wind field forecasts at 1-h forecast intervals provide the wind forcing for WW3 to

131    generate ocean wave forecasts for each ensemble member.  Those ensemble wave forecasts

132    are then combined to yield significant wave height probability fields exceeding a threshold (e.g.,

133    12 or 18 ft) on a 1 degree resolution grid, which has a resolution consistent with the current

134    WW3NAVGEM Ensemble probabilities available from FNMOC for evaluation purposes.  An

135    example of the 12-ft significant wave height probabilities on the right side of Fig. 1.  Since we

136    are only running 20 members of the WW3 ensemble, the probability fields are generated on a

137    1x1 degree global grid to reduce graininess noted in Sampson et al. (2016).  Still, this graininess

7

138     is visible at longer lead forecast times such as the 96-h WW3TCOFCL Ensemble forecast

139     probabilities shown in Fig. 1.

140         The entire 10-m wind field preparation process takes just a few minutes on a Cray XC-

141     30, and an estimated 1 hour of wall-time to run both the wind field preparation and the 20

142     WW3 ensemble members using 16 processors per ensemble member.  Although attempts are

143     made to warm start the WW3TCOFCL Ensemble every 12 hours using the previous 12-h

144     forecast, this was not feasible when NRL computer resources became unavailable for extended

145     periods of time.  In these instances, the WW3TCOFCL Ensemble was cold started with

146     potentially adverse effects on seas and swell in the early forecast times.  These effects become

147     less important beyond 24 h, but they are worth noting as they are plainly visible in visual

148     inspection.

149

150     3.      Evaluation Data

151         The WW3TCOFCL Ensemble was run in real-time on 84 TCs that existed between May

152     2018 and March 2019.  NRL was able produce forecast data in the vicinity of TCs in all regions of

153     the globe.  As with most non-operational real-time NWP systems, NRL had issues with data

154     acquisition and unscheduled computer downtime.  As a result of this computer downtime, the

155     evaluation set has periodic gaps resulting in some artifacts from the many WW3 cold starts,

156     some of which are visible in our evaluation.  Since the WW3TCOFCL Ensemble was run on the

157     same grid and has the same number of members as the WW3NAVGEM Ensemble, verification

158     of head-to-head cases will provide insight into both ensembles.

8

159    For ground truth we use the WW3TCOFCL deterministic model analysis of significant

160    wave height in feet (ft; 1 ft = 0.3048m), as that is the parameter most commonly used in Navy

161    operations. Noting again that the WW3TCOFCL deterministic model uses post-processed winds

162    forecasted by U.S. TC forecast centers.   Since the U.S. Navy is most concerned about significant

163    wave heights in ship routing, we chose to evaluate significant wave height probabilities.  We

164    present results using WW3TCOFCL deterministic model significant wave height analyses, but we

165    also evaluated results against WW3NAVGEM deterministic analyses.  The WW3NAVGEM

166    deterministic model analyses assimilate altimeter data (Cummings and Wittmann 2009), but

167    little difference was found between results using the WW3TCOFCL and WW3NAVGEM

168    deterministic model analyses as ground truth.  The 12-and 18-ft thresholds chosen for

169    evaluation are not necessarily the thresholds used for operational forecasting, but span a

170    reasonable range of significant wave heights associated with TCs and are routinely available for

171    the WW3NAVGEM Ensemble.

172

173    To gather data with 12-and 18-ft significant wave heights, which are not common in the

174    tropics, our verification was limited to a 20x20 degree box surrounding the verifying TC

175    position.  This area is likely larger than the TC wind field (Frank 1977) and also generally

176    encompasses the extreme waves associated with TCs.  In most cases a 20x20 degree box will

177    include many cases of zero probabilities in both the forecast and verification data (null cases),

178    which affects results and their interpretation. The verification impacts of null cases are

179    discussed section 3.   We also attempted this evaluation using a 10x10 degree box around the

180    verifying TC location, and found that this smaller area did not always encompass the TC-driven

9

181   waves and highest significant wave height probabilities at longer forecast leads.  At these longer

182   leads, the area of high significant wave height probabilities can be both larger and  dislocated

183   from the 10x10 degree box around the verifying position.  Our evaluation was also limited to

184   TCs with verifying intensities of 35 knot (kt; 1 kt =0.514 m s$^{-1}$) or greater intensity, which results

185   in limiting the false alarm rates for both algorithms.

186

187        Although we verify WW3TCOFCL Ensemble probabilities against WW3NAVGEM

188   deterministic model significant wave height analyses (which assimilate altimeter wave heights),

189   we do not to attempt verification ensemble runs against buoys and/or altimetry data explicitly,

190   other than anecdotally. These observations have coverage issues that hinder verification of

191   steep gradients and rare events, and can yield misleading results (see Sampson et al. 2013).

192

193        Table 1 provides a summary of the cases used in the verification.  Each 20x20 degree

194   verification area represents 400 potential paired forecast and verification points, so the values

195   in Table 1 are effectively 1/400$^{th}$ of the paired forecast points evaluated (minus an estimated

196   10% that verified over land and were removed from verification).  Grid differences also

197   accounted for minor differences in the matched pairs over water, 1 or 2 paired forecasts in

198   approximately 10% of the cases.  This represents differences of less than 0.1% and is ignored.

199        Summary statistics at the end of the Results Section are provided with significance using

200   a 2-tailed Student's t-test.  To remove correlation issues within the data, each 20x20 degree

201   (each with potentially 400 paired forecasts) is treated as a single case.  Then the t-tests are

202   provided for the summary statistics –Discrimination Distance, ROC AUC, and Brier Score.  No

10

203 effort is made to account for the effects of serial correlation in the summary data, but the

204 degrees of freedom are conservatively estimated using the number of cases rather than the

205 number of matched pairs (i.e., counting every point in the 20x20 degree box as a case).

206

207   4.      Results

208       To demonstrate significant wave height forecasts we present results in three ways.  We

209 first present two cases that exemplify our real-time assessment of the differences between

210 WW3TCOFCL Ensemble and WW3NAVGEM Ensemble significant wave height probabilities.  We

211 then verify WW3TCOFCL Ensemble and WW3NAVGEM Ensemble against WW3TCOFCL

212 deterministic model significant wave height analyses, and for completeness, against

213 WW3NAVGEM deterministic significant wave height analyses.  For objective probabilistic

214 verification statistics generation, we use the Model Evaluation Tools (MET; Development Test

215 Center 2020) grid verification tools.  We employ MET parameters Reliability, Likelihood,

216 Calibration, ROC, ROC AUC, and Brier Score to obtain a reasonably complete summary of

217 performance characteristics of each ensemble. Each of these metrics is described in section 4c.

218       a)  Typhoon Maria (WP102018)—intensifying to 140 kt

219       To highlight differences in the two algorithms (WW3 run with/without post-processing)

220 in an intensifying TC, we choose the Maria (WP102018). Maria, the eighth named storm of the

221 2018 typhoon season, was a powerful tropical cyclone that affected Guam, the Ryukyu Islands,

222 Taiwan, and East China in early July 2018.  Here we examine 96-h forecasts valid July 9, 2018 at

223 00 UTC, initiated on 5 July 00 UTC when the storm was located southeast of Guam and forecast

11

224　to intensify as it moved toward Okinawa.  Figure 2 shows details of the WW3TCOFCL Ensemble

225　(left column) and WW3NAVGEM Ensemble (right column) forecasts of 12-ft seas. Consistent

226　among the TCs inspected (approximately 30 cases) are that the WW3NAVGEM Ensemble input

227　forecast tracks (Fig. 2 top row) and intensities both have reasonably large spread, but that

228　ensemble member intensities tend to be too low, with intensities, unrealistically peaking near

229　70 kt for all members (Fig.2 second row). In comparison, the WSP tracks and intensities appear

230　to be well-calibrated with individual forecasts encompassing the forecast, and thus provide

231　more realistic wind forcing input to WW3. In the case of Maria, this results in large areas

232　relatively weak wind forcing input to the WW3NAVGEM Ensemble, and much lower 12-ft

233　significant height probabilities when compared to those from the WW3TCOFCL Ensemble (Fig.2,

234　third row)—the issue is even more pronounced for higher significant wave height thresholds

235　(not shown). These differences are not isolated, but seen throughout the data set, especially for

236　developing TCs.

237

238　　　b)  Hurricane Ileana (EP112018)—maintaining intensity at 40-45 kt

239　　　The majority of TCs are not forecast to intensify beyond 70 kt.  To highlight differences

240　between a weaker TC that is not forecast to intensify, we choose Hurricane Ileana's 48-h

241　forecast valid August 8, 2018 at 00 UTC, initiated on 6 August 00UTC. Ileana was a remarkably

242　small TC and the ninth tropical storm in the East Pacific in 2018 and during its lifecycle tracked

243　parallel to the Mexican coast.  At this time, NHC forecasted Ileana to remain weak as it

244　approached the Baja California Peninsula.  In this case, the initial intensities used in the

12

245     WW3NAVGEM Ensemble encapsulate the initial estimate from NHC (Fig. 3, second row). The

246     forecast track (Fig 3, top row) and intensity spreads (Fig. 3, second row) are larger than those

247     produced from the WSP algorithm. The 12-ft seas probabilities forecasts (Fig.3, third row) from

248     WW3TCOFCL Ensemble are still noticeably higher probabilities in the vicinity of the highest

249     observed wave heights (Fig. 3, bottom row). Much of the difference in 12-ft significant wave

250     height probabilities generated from the WW3TCOFCL Ensemble and WW3NAVGEM Ensemble

251     can be explained by larger forecast track spread in the WW3NAVGEM Ensemble input.

252

253         c) Objective Scores

254         Once the analyses are limited to 20x20 degree boxes centered on the TC best track

255     position, the probability forecasts can be inter-compared using standard probability metrics

256     such as Reliability (Fig. 4), Discrimination (Fig. 5), Relative/Receiver Operating Characteristic

257     (ROC; Fig. 6), and summary or derivative metrics such as Discrimination Distance, Area Under

258     ROC Curve, and Brier Score (Fig. 7). Each of these metrics answers a specific question that we

259     discuss below. Again, our evaluation uses MET, which in turn cites Wilks (2011) for most of its

260     statistical algorithms. Results shown here are for a homogeneous data set, meaning that the

261     scores from the two different algorithms can be compared since they are for the same TCs on

262     the same dates. For ground truth we again use analyzed significant wave heights from the

263     WW3TCOFCL deterministic model (Sampson et al. 2013) as these have been shown to have

264     realistic TC structure. We also performed the same tests using WW3NAVGEM deterministic

265     model analyzed significant wave heights for verification, but somewhat surprisingly found

13

266    consistent results in both statistical analyses for the metrics chosen.  Finally, the evaluation was

267    conducted for 0, 1, 2, 3, 4, and 5 day forecasts, but we limit presentation of the Reliability,

268    Discrimination and ROC charts to 1, 3, and 5 days and the results to those using the

269    WW3TCOFCL model deterministic analysis as ground truth for brevity.

270        i)      Reliability

271        Reliability determines how well the probabilities compare to observed frequencies.  On

272    a Reliability Diagram, perfect reliability is a diagonal (1:1) line from lower left to upper right,

273    biases are indicated by model reliability being below (high bias) and above (low bias) the 1:1

274    line, and forecast confidence is provided by the slope of model reliability relative to the 1:1 line,

275    that is under-confident when the slope is less than and overconfidence when the slope is

276    greater than one (Wilks 2011).    Reliability for both 12-ft and 18-ft significant wave height

277    probabilities is shown in Fig. 4.  The reliability for WW3TCOFCL Ensemble 12-ft significant wave

278    height appears high biased (over-forecasting in Wilks 2011) throughout.  The WW3NAVGEM

279    Ensemble appears to overestimate low probabilities and underestimate higher probabilities in

280    shorter forecast leads (under-confident), and overestimate probabilities like the WW3TCOFCL

281    Ensemble does at longer forecast leads.  The number of cases drops precipitously for the 120-h

282    18-ft significant wave height probabilities above 80%, dropping to 400 head-to-head cases or

283    one grid (SH112019 verifying Mar 7 2019 at 12:00 UTC).  So the Reliability Diagrams at 120 h for

284    18-ft significant wave height at the highest probability thresholds have few verification cases,

285    reflected in the erratic changes in the reliability.

14

286    In the case of the WW3NAVGEM Ensemble (under-confident in short-term forecast

287    leads, over-forecasting at longer-term forecast leads), the authors suspect that the ensemble is

288    challenged by resolution in that circulations tend to be too large at longer forecast leads.  In the

289    case of WW3TCOFCL Ensemble, the authors suspect several potential issues.  The first is that

290    WW3 is likely more appropriately run with 10-minute mean wind speeds since it is developed

291    to use NWP fields.  This is in contrast to U.S. official forecast center specified TC winds and wind

292    probability realizations, which are both considered 1-minute wind speed estimates.

293    Operational forecasters use conversion rates such as .93 (Harper et al. 2010) to convert the 1-

294    minute wind speeds to 10-minute wind speeds, and this conversion would likely reduce the

295    high bias.  Another potential source of bias is the statistical wind radius model (DRCL; Knaff et

296    al. 2007 and Knaff et al. 2018) used in the wind probabilities.  DRCL wind radii become more

297    symmetric as the forecast progresses in time, and these symmetric forecasts could provide

298    unrealistic durations for TC winds.  DRCL will never emulate the large symmetry fluctuations

299    seen in nature.  A more appropriate treatment of the asymmetries, especially at longer forecast

300    periods, could provide more realistic changes in fetch and duration of winds around TCs.

301    ii)    Discrimination

302    Discrimination is the relative frequency with which a forecast can discriminate between

303    events and non-events, where perfect discrimination would entail no overlap between

304    distributions of forecast probabilities for events and non-events.  Discrimination Diagrams show

305    these frequencies, where superior discrimination is indicated by separation between the events

306    and non-events.  Figure 5 shows discrimination for probabilities from our two algorithms at 1,

307    3, and 5 days.  One obvious trend is that the separation between events and non-events

15

308     becomes smaller as forecast length increases, as seen by the lines of the same color converging

309     towards each other.  The ability to discriminate between events and non-events drops with

310     forecast lead time for both algorithms.

311          iii)     Discrimination Distance

312          An easier way to visualize and summarize the discrimination is to graph the

313     Discrimination Distance (the difference between the average of the event and non-events) for

314     all forecast leads on one graph (Fig. 7).  The Discrimination Distances for the WW3TCOFCL

315     Ensemble are lower than for WW3NAVGEM Ensemble probabilities out to approximately 24 h,

316     then remain approximately 10% higher for the longer leads.  Significant differences using a 2-

317     tailed t-test at the 5% level are present at all but the 24-h time period for 12-ft probabilities,

318     and at all but 24-h and 120-h time periods for the 18-ft probabilities.  Discrimination Distances

319     for 12-ft are about 10% higher than for 18-ft significant wave heights at all forecast leads,

320     indicating more skill in discrimination of 12-ft significant wave heights.  The Discrimination

321     Distances also decay at longer leads, indicating less skill in discrimination between events and

322     non-events at these forecast leads times.

323          iv)     ROC

324          ROC is another measure of the ability of the forecast to discriminate between two

325     alternative outcomes, thus measuring resolution. It is not sensitive to bias in the forecast, so

326     says nothing about reliability. A biased forecast may still have good resolution and produce a

327     good ROC curve, which means that it may be possible to improve the forecast through

328     calibration (e.g., correcting the bias).  ROC can thus be considered as a measure of potential

16

329　usefulness (Development Test Center, 2020).  A perfect ROC curve follows the y axis from 0 to

330　1, then across the top of the diagram to 1, 1.  The ROC degrades for both algorithms as forecast

331　time increases (Fig. 6).  This is true for both the 12-ft and 18-ft thresholds. .

332　　　　v)　　ROC Area Under Curve

333　　　　The area under a ROC Curve (ROC AUC) is a convenient way to summarize how a

334　forecast discriminates between event/non-event (Wilks 2011).  Values can theoretically go from

335　0 to 1. A perfect score is 1, describing the area under a curve that passes from x=0, y=0, through

336　x=0, y=1, to x=1, y=1).  The ROC AUC for the no-skill diagonal is .5 (the area under a diagonal

337　from x=0, y=0 to x=1, y=1 on a ROC Diagram).  As expected, the ROC AUC (Fig. 7) for the

338　WW3TCOFCL Ensemble probabilities is relatively low at analysis time due to the many cold

339　starts in our data set.  The WW3TCOFCL Ensemble ROC AUC improves until about the 48-h

340　forecast time, then gradually drops off through 120 h.  The WW3NAVGEM Ensemble ROC AUC

341　drops gradually through the forecast and is approximately 15% lower than the WW3TCOFCL

342　Ensemble between 72 and 120 h.  Differences in the ROC AUC pass significance tests at all

343　forecast periods except at 48 h for 12-ft, and at 0 h for 18-ft significant wave height.  The

344　numbers of cases (each case representing an entire 20x20 degree grid) for this ROC AUC at 48,

345　72, 96, and 120 h are all well below 200, so conclusions on significance tests 18-ft significant

346　wave height should await more cases.  Recall that the high bias in the WW3TCOFCL Ensemble is

347　not penalized in either the ROC or the ROC AUC, and that the ROC AUC is only used to

348　discriminate between the event and non-event.  It is encouraging that the WW3TCOFCL

349　Ensemble probabilities maintain high ROC AUC out to 120 h since high bias, not depicted in

350　either the ROC or ROC AUC, can be corrected through adjustments in the algorithm.

17

351    vi)    Brier Scores

352        Brier Scores are another standard skill score for probabilistic forecasts, and measure

353    both reliability and resolution (the ability to distinguish an event from a non-event). The Brier

354    Score measures the mean square error of probabilities.  Here again we use the WW3TCOFCL

355    deterministic model analyses as ground truth.   Brier Scores range from 0 to 1, 0 being a perfect

356    score.  Brier Scores for both ensembles evaluated are shown in Fig. 7 and they are within 3% of

357    each other for both 12- and 18-ft thresholds.  These generally rise as forecast time increase,

358    indicating skill drops with forecast lead.  The uptick in the WW3TCOFCL Ensemble at analysis

359    time is expected as this ensemble was frequently cold started throughout the testing period

360    and the WW3TCOFCL Ensemble (and its input) has little spread at analysis time.  The

361    WW3NAVGEM Ensemble probabilities have slightly lower Brier Scores than the WW3TCOFCL

362    Ensemble probabilities at all forecast times for the 12-ft significant wave height threshold, and

363    scores from the two algorithms are within 3% of each other.  Differences for 12-ft probabilities

364    are significant at all forecast periods.  Brier Scores for 18-ft significant wave height thresholds

365    are within 1% of each other with the WW3NAVGEM Ensemble scoring lower (better).

366    Differences are significant at 24 and 96 h, but just barely pass significance tests.  In the case

367    shown in Fig. 2, the Brier Score for WW3NAVGEM Ensemble (0.082098) is lower than for

368    WW3TCOFCL Ensemble (0.13089).  This may seem counterintuitive as the WW3TCOFCL

369    Ensemble probabilities "look" to capture the 12-ft significant wave heights in the WW3TCOFCL

370    deterministic model analysis from 96 hours later.  But upon further inspection (Table 2), the

371    distribution of probability forecasts for WW3NAVGEM Ensemble is skewed to lower

372    probabilities so that it scores much higher in the large number of non-events than the

18

373      WW3TCOFCL Ensemble probabilities for this case.   The Brier score becomes inadequate for

374      very rare (or very frequent) events because it does not sufficiently discriminate between small

375      changes in forecast that are significant for rare events (Benedetti 2010).   Thus, Brier Score

376      unfairly penalizes extremely rare (or common) event forecasts and can actually leads to

377      conclusions that disagree with our intuition (Jewson 2008), such as indicating that the

378      WW3NAVGEM Ensemble outperforms the WW3TCOFCL Ensemble for the case in Fig. 2. The

379      Brier Scores are still useful in our evaluation as they confirm high bias in the WW3TCOFCL

380      Ensemble that, if corrected, could decrease the Brier Scores.  However, tuning specifically to

381      Brier Scores is not advised as that could result in undesired reduction in extreme event

382      prediction (described as under-confident in Wilks 2011).  An analog to this would be tuning a TC

383      wind intensity consensus (e.g., see Sampson et al. 2008) to minimize mean forecast error when

384      the most impactful errors are associated with rare and difficult to forecast rapid intensification

385      events.

386

387      5.      Conclusions and future work

388          A post-processing algorithm for insertion of real-time operational TC surface wind

389      forecasts into a .25x.25 degree global 20-member ensemble surface wind field is described.

390      This algorithm was run twice a day (at 00 and 12 UTC) for approximately one year and included

391      active TCs from all basins.  Each set of post-processed wind fields was then used as wind input

392      to WW3 in order to generate a 20-member ensemble of forecasted significant wave height

19

393    fields out to 5 days.  The resultant significant wave height fields from each ensemble member

394    were then compiled to create significant wave height probabilities on a 1x1 degree global grid.

395        Evaluation was performed using 20x20 degree boxes around verifying positions of the

396    TCs at each forecast day using the MET statistics package.  Both WW3NAVGEM and

397    WW3TCOFCL deterministic model analyses were used as ground truth for evaluation of the

398    probabilities and little difference was found between evaluations with the two ground truth

399    datasets.  Case studies indicated large discrepancies frequently existed between input winds

400    from the two algorithms.  NAVGEM Ensemble tracks and intensities generally had large

401    spreads, and certainly larger than those generated by the WSP algorithm that are used in the

402    WW3TCOFCL Ensemble for weaker TCs.  WW3NAVGEM Ensemble input intensities were

403    generally low-biased for intense TCs as the NAVGEM Ensemble resolution was challenged to

404    represent steep wind gradients in relatively small TCs.  Large discrepancies also existed

405    between significant wave height probabilities generated by each of the ensemble forecasts.

406    The WW3NAVGEM Ensemble significant wave height probabilities tended to be more

407    widespread and lower in magnitude than those from the WW3TCOFCL Ensemble.

408        In objective evaluation, Reliability Diagrams show that WW3NAVGEM Ensemble

409    overestimated low probabilities and underestimated higher probabilities in short-range

410    forecasts, then generally overestimated probabilities by 5 days.  WW3TCOFCL Ensemble

411    generally overestimated all probabilities throughout the entire forecast.  Brier Scores for

412    WW3NAVGEM Ensemble were a few percent better than WW3TCOFCL Ensemble at 12-ft

413    significant wave height forecasting at all forecast lengths, but inspection of individual cases

414    indicated that those scores were heavily influenced by forecasts of very low probability for non-

20

415    events (no 12-ft or 18-ft significant wave height in ground truth).  Brier Scores for 18-ft

416    significant wave height were within about 1% at all forecast lengths.  ROC curves and ROC AUC

417    indicated that discrimination between events and non-events degrades with forecast period for

418    both sets of probabilities, but that WW3TCOFCL Ensemble forecast generally appeared better

419    at discriminating events from non-events beyond 24 h.  These results are confirmed by the

420    Discrimination Diagrams, Discrimination Distances, and significance tests for Discrimination

421    Distances.

422        The WW3TCOFCL Ensemble high bias noted in the Reliability Diagrams is likely

423    correctable.  Whether by converting the WW3TCOFCL Ensemble input 1-minute to 10-minute

424    mean winds that are more representative of NWP model winds, by replacing the Wind Radii

425    CLIPER Model (DRCL) with more realistic wind distribution realizations, or by applying some

426    combination of the above, the high bias can be addressed.  Also, the validation package

427    developed in this work could be modified to validate whether changes in algorithms upstream

428    of the WW3 ensembles (e.g., the WSP algorithm and the NAVGEM Ensemble) adversely affect

429    the significant wave height probabilities.  Operational forecasts are certain to improve in the

430    future through use of new sensors, improved NWP representation of the vortex, and more

431    advanced post-processing in the wind probability algorithm — all of which can affect these

432    ensembles.  Construction of TC-specific significant wave height probability verification was

433    time-consuming, but the process to achieve this is in place and could be used as is or improved

434    upon to validate TC-specific wave probabilities in the future.  And addition of Object-Based

435    Diagnostic Evaluation (MODE) verification available in MET may compliment the evaluation

436    done within this work as it follows features (e.g., TCs ) and reports statistics different than

437  those here when comparing the features.  That evaluation would be similar to and hopefully

438  more rigorous than the 12-ft sea radii evaluation against operational NHC estimates as done in

439  Sampson et al. (2016).

440

454

455  Data Availability: Both sets of ensemble significant wave height probabilities have been

456  archived and are available on request; however, a non-disclosure agreement public release

457  approval may be required to provide data.

458

459

460

461    8.      References

462 Aijaz, S., J. D. Kepert, H. Ye, Z. Huang, and A. Hawksford, 2019: Bias correction of tropical

463 cyclone parameters in the ECMWF Ensemble Prediction System in Australia. *Mon. Wea.*

464 *Rev.,* **147**, 4261–4285, doi: https://doi.org/10.1175/MWR-D-18-0377.1

465 Benedetti, R., 2010: Scoring rules for forecast verification. *Mon. Wea. Rev.*, **138**, 203–211, doi:

466 https://doi.org/10.1175/2009MWR2945.1

467 Cummings, J. A.,  and P. Wittmann, 2009: Navy implements data assimilation capability for its

468 wave forecasting model. *JCSDA Quarterly,* No. 28, Joint Center for Satellite Data Assimilation,

469 Camp Springs, MD, 2–3. [Available online

470 at https://static1.squarespace.com/static/5bad1a12c2ff616821035c9f/t/5d1bc190f87d390001

471 d7f83e/1562100113337/200909JCSDAQuarterly.pdf ]. Accessed 6/10/2021.

472 DeMaria, M., J. A. Knaff, M. Brennan, D. Brown, C. Lauer, R. T. DeMaria, A. Schumacher, R. D.

473 Knabb, D. P. Roberts, C. R. Sampson, P. Santos, D. Sharp, K. A. Winters, 2013: Operational

474 tropical cyclone wind speed probabilities part I: Recent model improvements and verification,

475 *Wea. Forecasting*, **28**, 586–602, doi: http://dx.doi.org/10.1175/WAF-D-12-00116.1

476 Development Testbed Center, 2020:  Model Evaluation Tools (MET) [available on line at

477 https://dtcenter.org/community-code/model-evaluation-tools-met] Accessed 6/10/2021.

478    Frank, W. M., 1977: The structure and energetics of the tropical cyclone. I: Storm structure.

479    *Mon. Wea. Rev.*, **105**, 1119–1135, doi:10.1175/1520-0493(1977)105<1119:TSAEOT>2.0.CO;2

480    Harper, B. A., J. D. Kepert, and  J. D. Ginger, 2010: Guidelines for converting between various

481    wind averaging periods in tropical cyclone conditions, World Meteorological Society, 64 pp.

482    [Available online at

483    https://library.wmo.int/index.php?lvl=notice_display&id=135#.X6gp8mR7mUk] Accessed

484    6/10/2021.

485    Jewson, S., cited. 2008: The problem with the Brier score. arXiv:physics/0401046v1 [physics.ao-

486    ph]. [Available online at http://arxiv.org/abs/physics/0401046v1] Accessed 6/10/2021.

487    Klotz, B. W., and D. S. Nolan, 2019: SFMR Surface wind undersampling over the tropical cyclone

488    life cycle. *Mon. Wea. Rev.*, **147**, 247–268, doi: https://doi.org/10.1175/MWR-D-18-0296.1

489    Knaff, J. A., C. R. Sampson, M. DeMaria, T. P. Marchok,J. M. Gross,and C. J. McAdie, 2007b:

490    Statistical tropical cyclone wind radii prediction using climatology and persistence. *Wea.*

491    *Forecasting*, **22**, 781–791, doi: https://doi.org/10.1175/WAF1026.1

492    Knaff, J. A., C. R. Sampson, and K. D. Musgrave, 2018: Statistical tropical cyclone wind radii

493    prediction using climatology and persistence: Updates for the western North Pacific. *Wea.*

494    *Forecasting*, **33**, 1093–1098, doi: https://doi.org/10.1175/WAF-D-18-0027.1

495    Marchok, T. P., 2002: How the NCEP Tropical Cyclone Tracker works. *Preprints, 25th Conf. on*

496    *Hurricanes and Tropical Meteorology*, San Diego, CA, Amer. Meteor. Soc., P1.13. [Available

497    online at http://ams.confex.com/ams/pdfpapers/37628.pdf] Accessed 6/10/2021.

498    Meissner, T., L Ricciardulli, and F.J. Wentz, 2017: Capability of the SMAP Mission to measure

499    ocean surface winds in storms, *Bull. Amer. Meteor. Soc.*, **98**, 1660-1677, doi: 10.1175/BAMS-D-

500    16-0052.1

501    Mouche, A., B. Chapron, J. Knaff, Y. L. Zhao, B. Zhang, and C. Combot, 2019:  Copolarized and

502    cross-colarized SAR measurements for high-resolution description of major hurricane wind

503    structures: Application to Irma Category 5 hurricane, *J Geophys Res-Oceans*, **124**, 3905–3922.

504    NCEP, 2020: WAVEWATCH III® Model. [Available on line at

505    https://polar.ncep.noaa.gov/waves/wavewatch/] Accessed 12/06/2020.

506    O'Reilly, W. C., and R. T. Guza, 1993: A comparison of two spectral wave models in the Southern

507    California Bight. *Coastal Eng.*, **19**, 263–282.

508    Reul, N., and co-authors (2017): A new generation of tropical cyclone size measurements from

509    space, *Bull. Amer. Meteorol. Soc.*, **98**, 2367–2386, doi: 10.1175/BAMS-D-15-00291.1

510    Sampson, C. R., J. L. Franklin, J. A. Knaff, and M. DeMaria, 2008: Experiments with a simple

511    tropical cyclone intensity consensus. *Wea. Forecasting*, **23**, 304–312.

512    Sampson, C. R., J. S. Goerss, J. A. Knaff, B. R. Strahl, E. M. Fukada, and E. A. Serra, 2018:  Tropical

513    Cyclone Gale Wind Radii Estimates, Forecasts, and Error Forecasts for the Western North

514    Pacific, *Wea. Forecasting*, **33**, 1081–1092.

515    Sampson, C. R., P. A. Wittmann, E. A. Serra, H. L. Tolman, J. Schauer, and T. Marchok, 2013:

516    Evaluation of wave forecasts consistent with tropical cyclone wind forecasts, *Wea. Forecasting*,

517    **28**, 287–294. doi: http://dx.doi.org/10.1175/WAF-D-12-00060.1

518     Sampson, C. R., J. Hansen, P. A. Wittmann, J. A. Knaff, and A. Schumacher, 2016: Wave

519     probabilities consistent with official tropical cyclone forecasts, *Wea. Forecasting*, **31**, 2035–

520     2045, doi: 10.1175/WAF-D-15-0093.1

521     Tolman, H. L., 1991: A third-generation model for wind waves on slowly varying, unsteady, and

522     inhomogeneous depths and currents. *J. Phys. Oceanogr.*, **21**, 782–797.

523     Tolman, H. L., B. Balasubramaniyan, L. D. Burroughs, D. V. Chalikov, Y. Y. Chao, H. S. Chen, and

524     V. M. Gerald, 2002: Development and implementation of wind generated ocean surface wave

525     models at NCEP. *Wea.Forecasting*, **17**, 311–333.

526     Tolman, H. L., J. G. M. Alves, and Y. Y. Chao, 2005: Operational forecasting of wind-generated

527     waves by Hurricane Isabel at NCEP. *Wea. Forecasting*, **20**, 544–557, doi:

528     https://doi.org/10.1175/WAF852.1

529     Wilks, D., 2011: Statistical methods in the atmospheric sciences, Elsevier, San Diego. 627 pp.

530     Zieger, S., D. Greenslade, J.D. Kepert, 2018: Wave ensemble forecast system for tropical

531     cyclones in the Australian region. *Ocean Dynamics* **68,** 603–625, doi: 10.1007/s10236-018-1145-

532     9

533

Table 1.  Numbers of WW3TCOFCL Ensemble and WW3NAVGEM Ensemble cases (each being a 20x20 grid) gathered from real-time execution from May 26, 2018 00:00 UTC to March 18, 2019 00:00 UTC with 84 TCs occurring around the world during that period.  Each 12-ft and 18-ft case 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 |

Table 2.  Contingency Table for WW3NAVGEM Ensemble (left) and WW3TCOFCL Ensemble (right) greater than 12-ft significant wave height probabilities for the 96-h forecast case shown in Figure 2.  Observed Yes and Observed No for the 20x20 degree grid encompassing the verifying TC position in Figure 2.

| | WW3NAVGEM Ensemble Matched Pairs | | WW3TCOFCL Ensemble Matched Pairs | |
|------|------|------|------|------|
| Prob | Ob Yes | Ob No | Ob Yes | Ob 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 |

27

Table 3. Numbers of cases for Reliability, Discrimination and ROC shown in Figures 4-6.

554

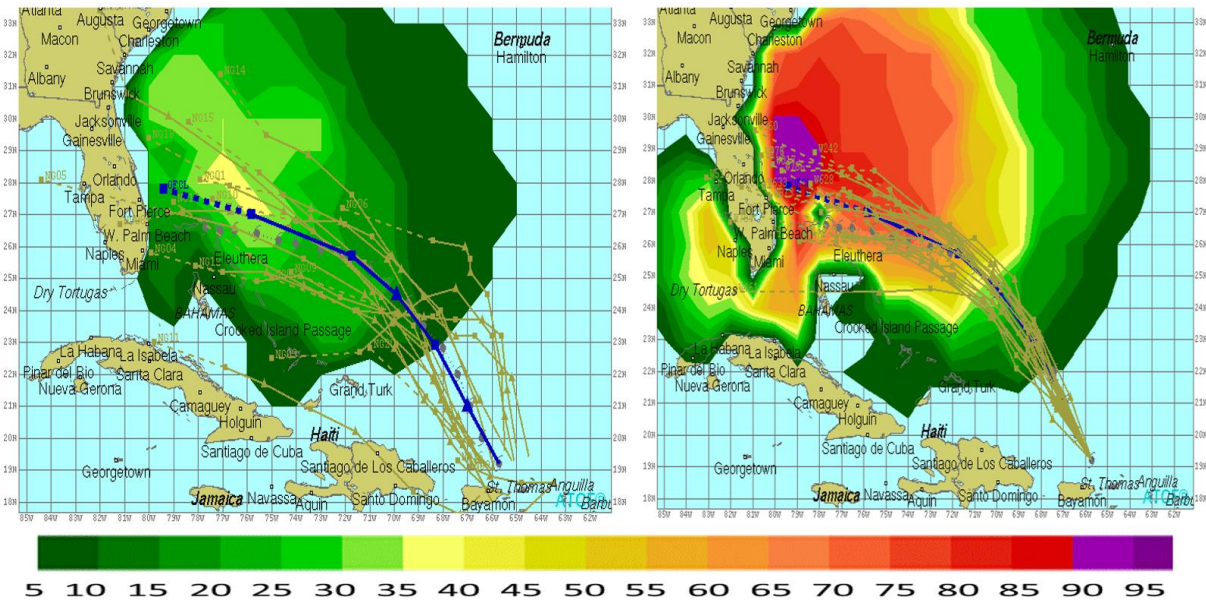| | WW3NAVGEM Ensemble | | WW3TCOFCL Ensemble | | WW3NAVGEM Ensemble | | WW3TCOFCL Ensemble | |
|---|---|---|---|---|---|---|---|---|
| | 24h 12ft | | 24h 12ft | | 24h 18ft | | 24h 18ft | |
| Prob | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No |
| 0.05 | 102 | 76240 | 1326 | 79715 | 59 | 99054 | 169 | 99583 |
| 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 |
| | 72h 12ft | | 72h 12ft | | 72h 18ft | | 72h 18ft | |
| Prob | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No |
| 0.05 | 440 | 39933 | 134 | 28743 | 326 | 59769 | 123 | 52648 |
| 0.15 | 622 | 9051 | 278 | 11742 | 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 |
| | 120h 12ft | | 120h 12ft | | 120h 18ft | | 120h 18ft | |
| Prob | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No | Ob Yes | Ob No |
| 0.05 | 438 | 20507 | 253 | 14547 | 316 | 30250 | 212 | 27093 |
| 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 |

555

556    Figure 1. (left) WW3NAVGEM Ensemble and (right) WW3TCOFCL Ensemble 96-h forecast 12-ft

557    sig wave ht probabilities for Dorian (AL052019) on Aug 29, 2019 at 00UTC. National Hurricane

558    Center forecast track (blue) is shown for reference. Also, (right) NAVGEM Ensemble TC tracks

559    and (right) wind probability realizations generated by the U.S. TC forecast center wind

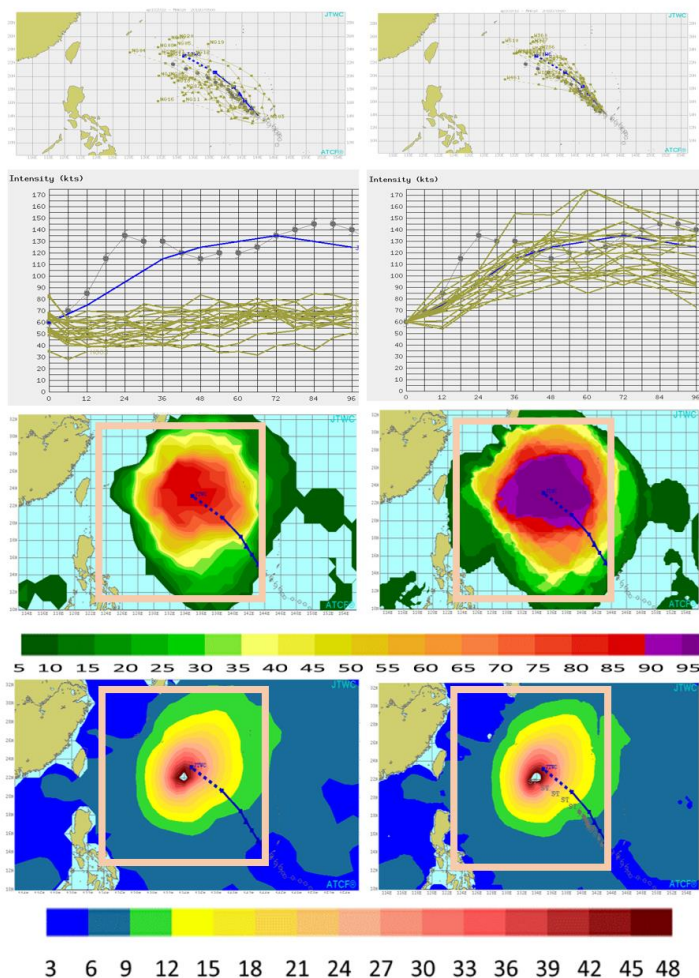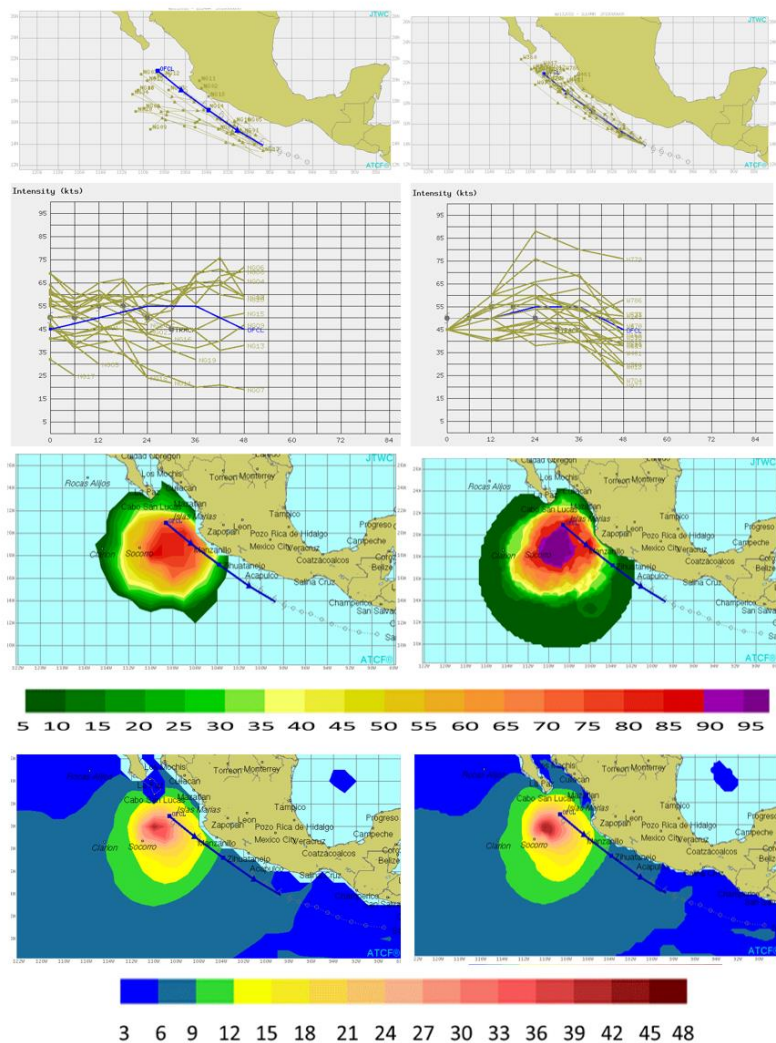560    probabilities (brown) included. Probability (%) colorbar is shown at the bottom.
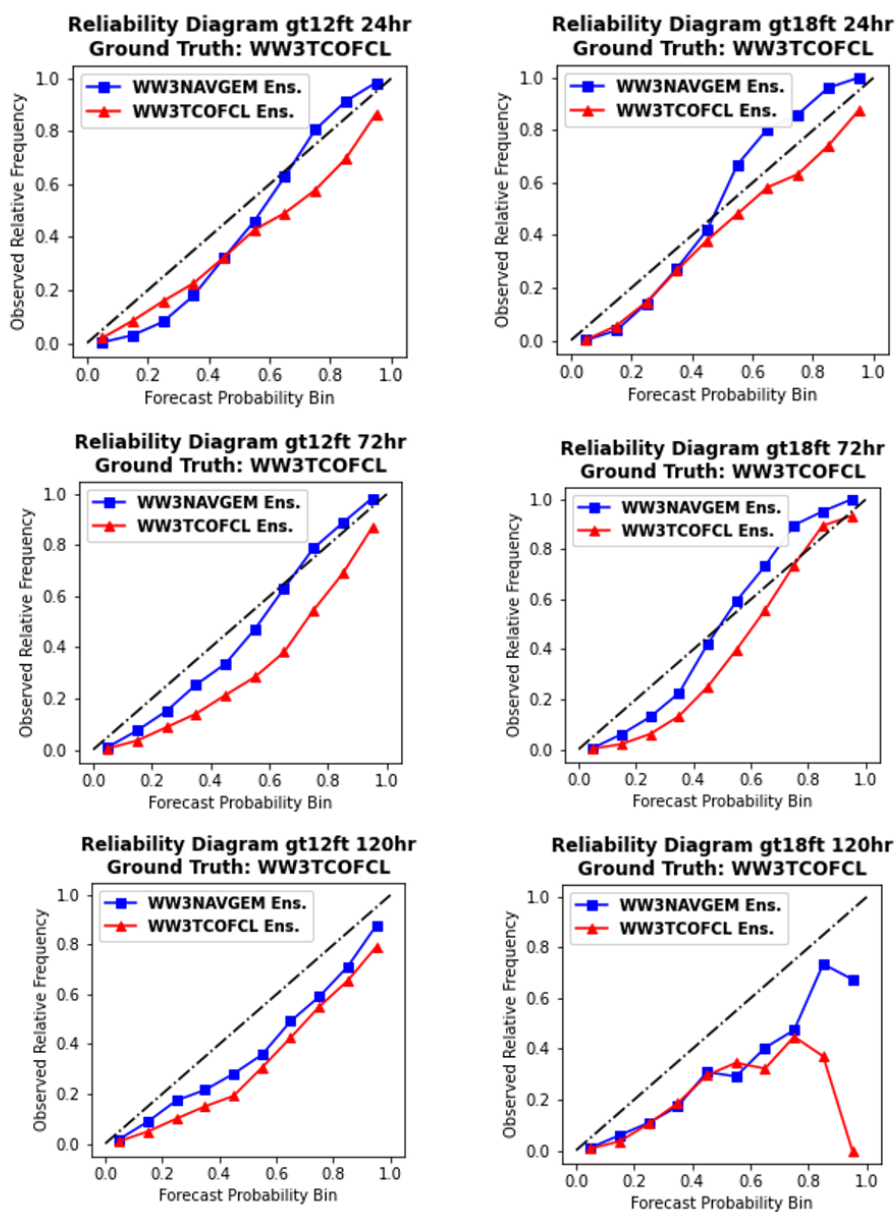
561

562

Figure 2.  (left) WW3NAVGEM Ensemble and (right) WW3TCOFCL Ensemble.  Input 96-h

forecast tracks (top row), input forecast and verifying intensities (brown lines and black

typhoon symbols, second row), 96-h forecast 12-ft sig wave ht probabilities (third row) and

verifying significant wave height (ft) analyses (fourth row) for WW3NAVGEM deterministic

model (left) and WW3TCOFCL deterministic model (right).  Forecasts and analyses valid July 9,

2018 at 00 UTC 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.

572

Figure 3. (left) WW3NAVGEM Ensemble and (right) WW3TCOFCL Ensemble.  Input 96-h

forecast tracks (top row), input forecast and verifying intensities (brown lines and black

typhoon symbols, second row), 96-h forecast 12-ft sig wave ht probabilities (third row)

and verifying significant wave height (ft) analyses (fourth row) for WW3NAVGEM

deterministic model (left) and WW3TCOFCL deterministic model (right).  Forecasts and

analyses for Ileana (EP112018) 48-h forecast valid August 8, 2018 at 00 UTC.  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.

581



582

583    Figure 4.  Reliability Diagrams for WW3TCOFCL Ensemble and WW3NAVGEM Ensemble 12-ft

584    (left) and 18-ft significant wave height (right) with  WW3TCOFCL deterministic model analysis

585    employed as ground truth.   Sequence progresses from (top) 24-h to (middle) 72-h to (bottom)

586    120-h forecast.  See Table 3 for numbers of head-to-head cases.  Dashed lines represent perfect
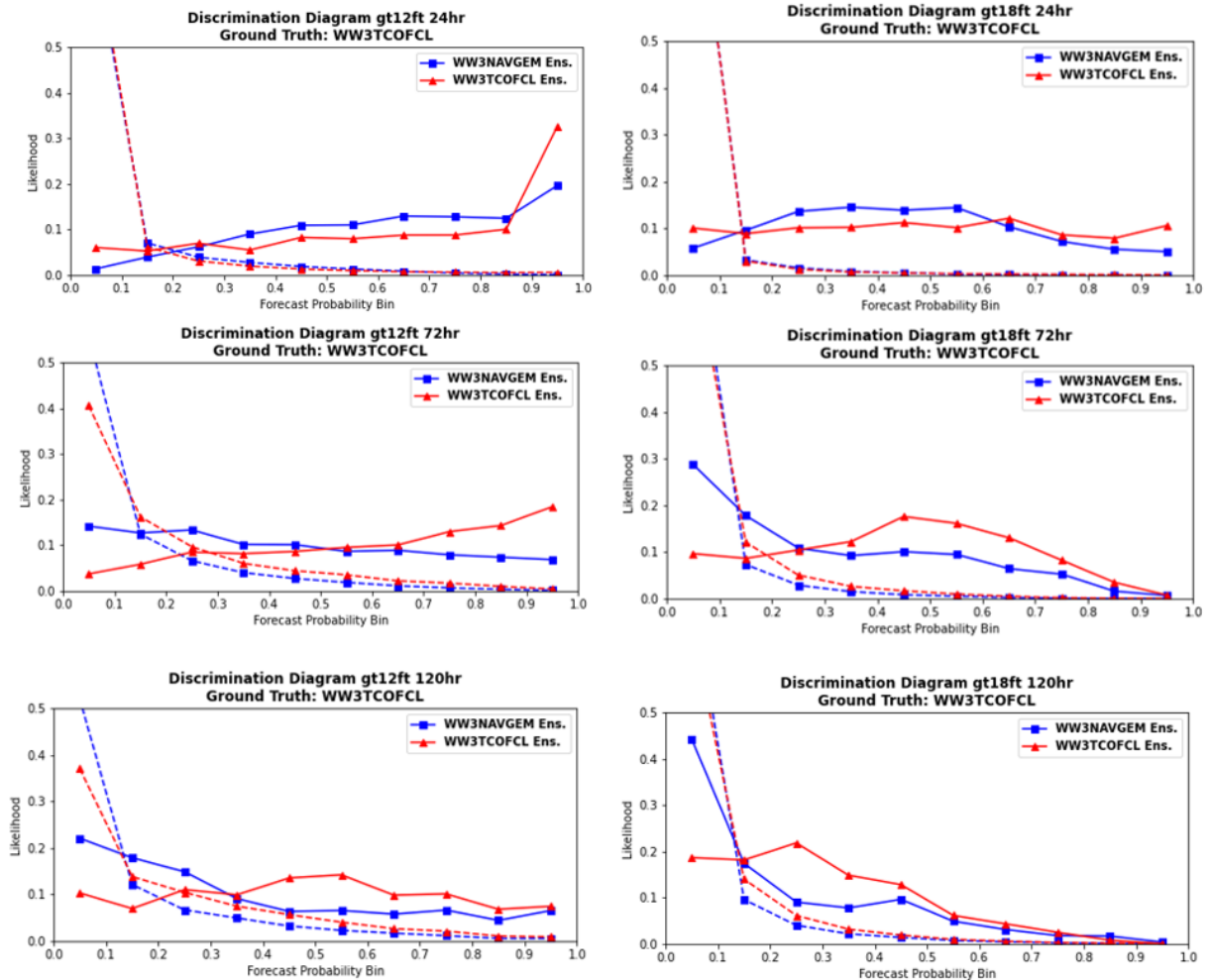
587    reliability.

588

589    Figure 5.  Discrimination Diagrams for WW3TCOFCL Ensemble and WW3NAVGEM Ensemble 12-

590    ft (left) and 18-ft significant wave height (right) with  WW3TCOFCL deterministic model analysis

591    employed as ground truth.   Solid lines indicate Observed Yes, dashed lines indicate Observed

592    No distributions.  Sequence progresses from (top) 24-h to (middle) 72-h to (bottom) 120-h

593    forecast.  See Table 3 for numbers of head-to-head cases.
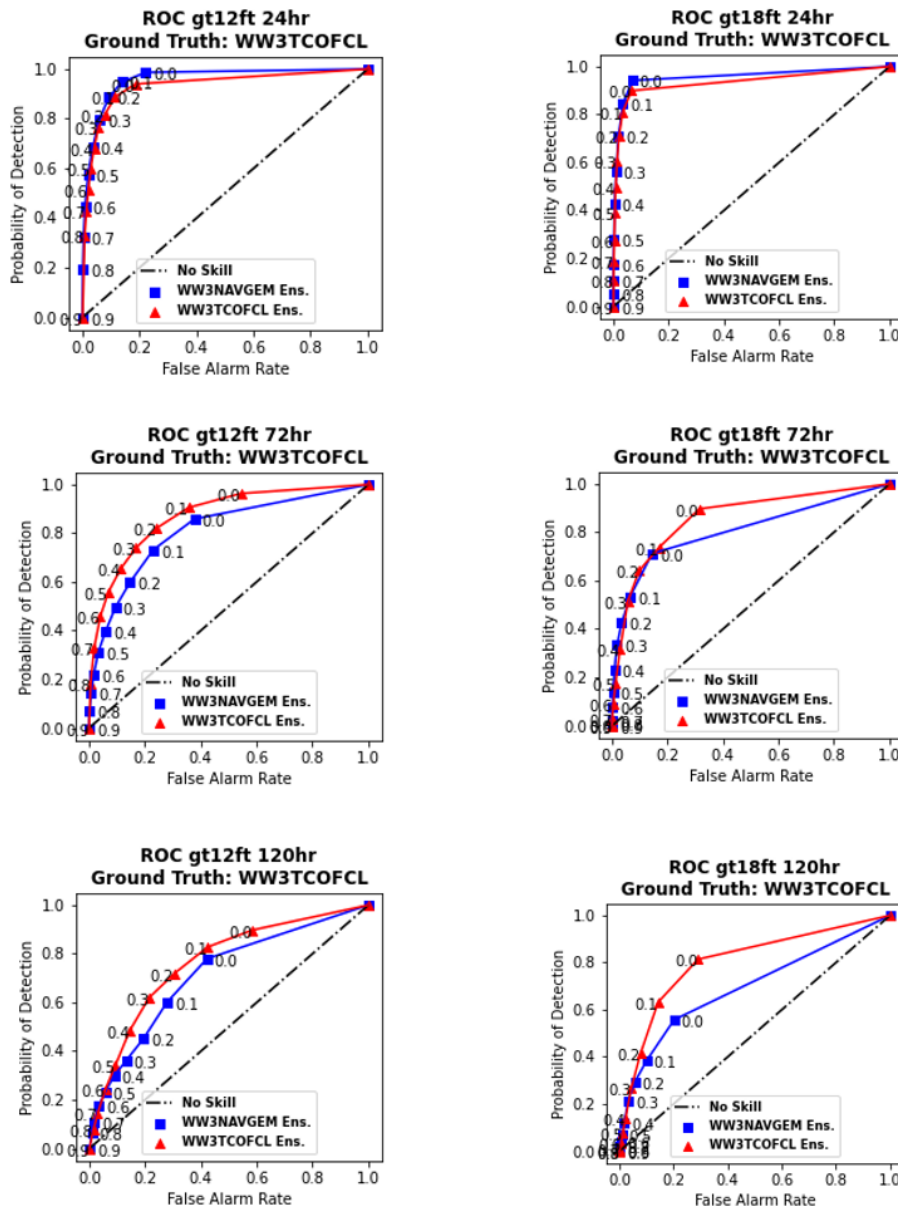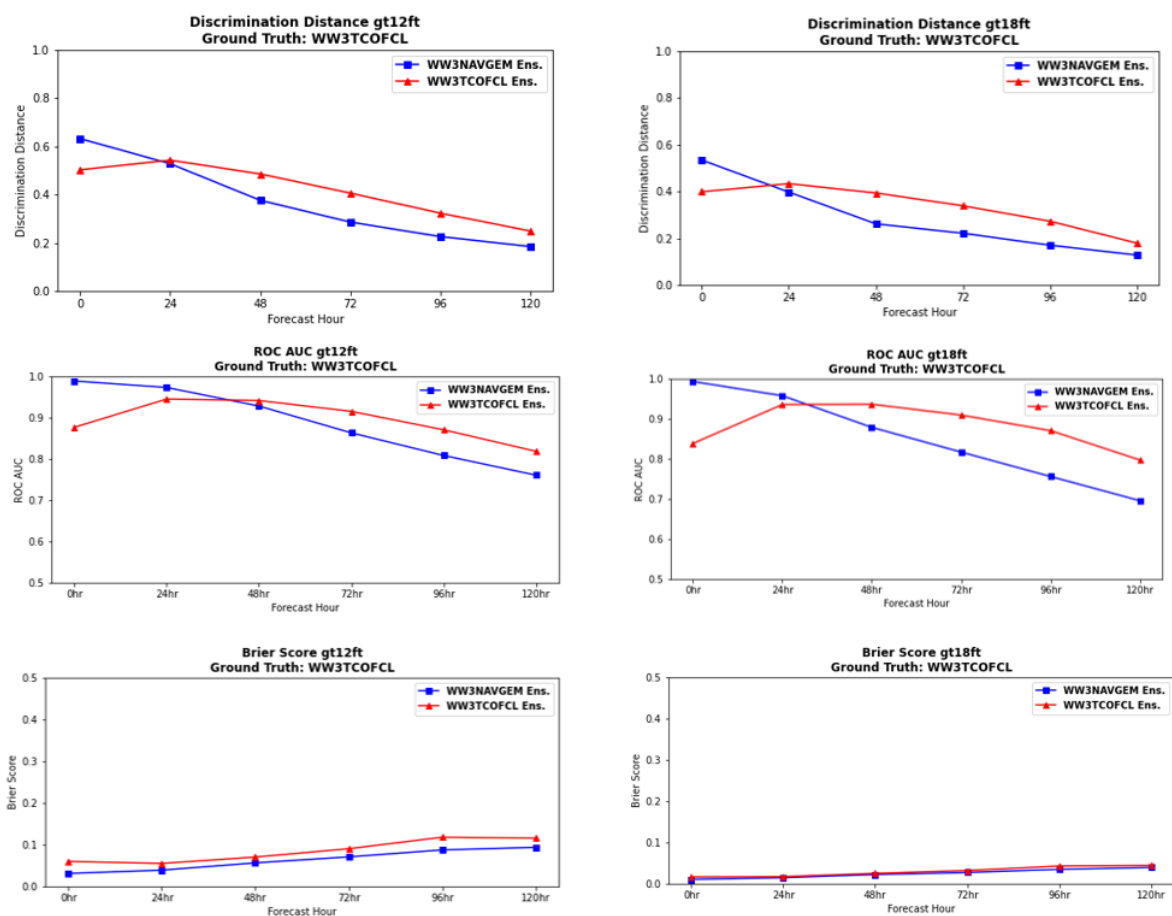
594

33

595

596   Figure 6.  ROC Diagrams for WW3NAVGEM Ensemble (left) and WW3TCOFCL Ensemble (right)

597   with WW3TCOFCL deterministic model analysis employed as ground truth.   Sequence

598   progresses from (top) 24-h to (middle) 72-h to (bottom) 120-h forecast.  Dashed line indicates

599   no skill.  See Table 3 for numbers of head-to-head cases.

600

601



602

603     Figure 7.  Discrimination Distances (top), ROC AUC (middle), and Brier Scores (bottom) for

604     WW3TCOFCL Ensemble and WW3NAVGEM Ensemble.  12-ft (left) and 18-ft significant wave

605     height (right) shown with WW3TCOFCL deterministic model analysis employed as ground truth.

606     Sequence progresses from (top) 24-h to (middle) 72-h to (bottom) 120-h forecast.  See Table 3

607     for numbers of head-to-head cases.

608