

1 **A Tropical Cyclone Forecast Metric for Operations and Model Development**

2 Michael Fiorino

3 *Cooperative Institute for Research in Environmental Sciences, Univ. of Colorado Boulder and*
4 *NOAA Global Systems Laboratory, Boulder CO*

5
6
7
8

9

10

11

12

13

14

15

16

17

18

19 Submitted on 31 March 2021 as a Manuscript for Consideration in

20 *Weather and Forecasting*

21

22

23

24

25 Corresponding Author. Email: michael.fiorino@noaa.gov

26 ABSTRACT

27 Tropical Cyclone (TC) forecast metrics have two primary applications. For the human
28 forecaster, they suggest which numerical models and aids should be considered in the forecast
29 process, e.g., giving greater weight to a ‘good’ model. For the modeler, metrics define forecast
30 quality and help diagnose model error.

31 Position Error (PE; often called ‘track’ error) and Intensity Error (IE) are the current
32 standard TC forecast metrics but are only indirectly related to the actual TC forecast – the one-
33 minute-average, 10 m (surface) wind field. For both the Joint Typhoon Warning Center (JTWC)
34 and the National Hurricane Center (NHC), the official warnings that drive user response is the
35 onset of 34 kt winds. Thus, PE/IE statistics do not directly measure error in the TC forecast of
36 the surface wind, especially wind ≥ 34 kt.

37 A metric is proposed (Forecast Error or FE) that is more directly related to the actual TC
38 wind field forecast based on Integrated Kinetic Energy (IKE); FE is defined to be a function of
39 *both* PE and IE. We then demonstrate that the intensity forecast contributes at most 10-12% of
40 total FE, i.e., a perfect intensity forecast can only improve (lower) FE at most by 10-12%. Track
41 is still the most important part of the TC forecast and an 18-y model intercomparison shows that
42 the ECMWF global model makes the best PE forecasts in all basins and forecast times, despite
43 having the highest initial PE and IE.

44 Through a breakdown of PE/FE by storm we find that the not all large PE cases have large
45 FE and *vice versa*, i.e., the new metric can distinguish forecast errors with greater operational
46 impact.

47

48

49 **1. Introduction**

50 In 2006, a wise US Navy Captain told me:

51 *“You’re only as good as what you measure”¹*

52 This pithy statement captures the essence of metrics – quantification of quality or putting
53 numbers to ‘goodness.’ Among the many applications of metrics two stand out: 1) ‘guidance on
54 guidance’ for the human forecaster – what are the ‘good’ models and how ‘good’ is the current
55 ‘good’ model forecast?; and 2) diagnostics for model development. For mid-latitude weather
56 forecasting, the 5-day 500 hPa northern Hemisphere (NHEM) Anomaly Correlation (5DNAC), is
57 a *single* number found to be highly correlated with *general* numerical weather prediction (NWP)
58 skill and has had a long and strong influence on global model development most notably in the
59 early decades 1980-2000 (Benjamin et al. 2019 and Simmons and Hollingsworth (2002)).²

60 Tropical cyclone (TC) forecast metrics have similar history. In the 1960s, the annual
61 typhoon reports from the Joint Typhoon Warning Center (JTWC see JTWC 2020a) show a
62 struggle to define the forecast metric and its summary statistic, but by the late 1970s two TC
63 forecast metrics emerged that is now our current standard: position error (PE; great circle
64 distance between forecast and verifying, post-season ‘best track’ (BT) position) and intensity
65 error (IE; difference in maximum 1-min, 10-m wind speed between forecast and best track). The
66 summary statistics are mean PE (mPE) and mean absolute IE (maIE) at set forecast times,
67 typically 0, 12, 24, 36, 48, 72, 96 120 h.

¹ Victor Addison, Commanding Officer, Fleet Numerical Meteorology and Oceanography Center, Monterey CA

² In 1999 I developed a TC vitals assimilation scheme for ERA-40 while on secondment to ECMWF. The 5DNAC of runs with TCs was about 0.005 lower (worse) than without and the scheme was not used at ECMWF but was used at JMA for the JRA-25 reanalysis. At the time there was extreme sensitivity at ECMWF for any reduction in the 5DNAC of an model update/upgrade.

68 PE is often called ‘track error’ but this is both inaccurate in how PE is calculated (has
69 nothing to do with a path or line) and is only indirectly related to track. For example, a TC
70 forecast track can have a 0 n mi PE at 72 h but a very large PE at 24 or 120 h. This 0 n mi 72-h
71 PE would define the ‘track’ as ‘perfect’ when in fact it is not. To measure track and speed error,
72 PE is broken into parallel (along track) and orthogonal (cross track) components and while these
73 components are better related to track error they are still point-in-time measures. IE is similarly
74 point-in-time and a perfect intensity forecast at 72 h could be horribly wrong at 24 h and *vice*
75 *versa*.

76 A basic problem with the standard PE/IE TC forecast metrics is how they are indirectly
77 related to the actual TC forecast – the surface wind field defined at JTWC/NHC as the 1-min.
78 averaged, 10-m (above the surface) wind. Moreover, the areal extent and time-of-onset of
79 surface winds ≥ 34 kts (gale force) is the basis of NWS TC warnings/watch and the US Navy
80 ship avoidance zone in JTWC warnings.

81 Given the weak connection between PE/IE and the TC wind forecast, the primary
82 purpose of this paper is to propose a TC forecast error (FE) metric that is more directly related to
83 the operationally critical TC forecast of surface winds ≥ 34 kts. The metric is defined as a
84 function of both PE/IE and the BT mean radius of 34 kt winds (R34). When assuming a perfect
85 intensity forecast (IE=0) the FE metric, on a seasonal basis, is reduced by at most 12%. Thus,
86 approximately 90% of mean FE comes from PE, but not all large PE have large FE...

87 The record-setting 31-storm 2020 Atlantic (LANT) season is first analyzed to
88 demonstrate how FE compares to PE/IE and how individual storms contribute to the seasonal
89 means. We further show the how the FE v PE relationship depends on R34 and why large PE

90 can have small FE and vice versa. A challenge for model developers, using large-error cases for
91 testing, could be selecting the large-error metric – FE or PE?

92 We then broaden our examination to the three main NHEM basins: 1) the LANT; 2) the
93 central and Eastern north PACific (EPAC); and 3) the Western north PACific (WPAC) for the
94 14-y period 2007-2020. The same models used in the 2020 LANT are considered in the larger
95 intercomparison

96
97 iand the formulation; and additionally, how IE contributes at most 12% to FE. Second,
98 how mid-latitude forecast metrics that measure general NWP model skill are very different than
99 TC metrics and how PE is primary in measuring overall TC model forecast quality. Third, how
100 individual storms contribute to seasonal means and how large PE does not always imply large
101 FE.

102 The TC forecast is first defined in Section 2 and the nature of TC PE is compared to
103 standard Numerical Weather Prediction (NWP) metrics that are used by both forecasters (‘what
104 are the good models?’) and modeler developers (‘is my model getting better?’) in Section 3. We
105 then show that TC metrics cannot be used in the same way as NWP scores because they are
106 fundamentally different – an apples vs. oranges problem. More importantly, we will also see in
107 Section 3 that the standard TC metrics are indirectly related to the actual TC forecast of surface
108 wind. In section 4, the new FE metric is proposed based on error in Integrated Kinetic Energy
109 (IKE, (Powell and Reinhold 2007) and we demonstrate the impact of ‘perfect’ intensity forecasts
110 on FE. Application of TC metrics to model development is reviewed in Section 5 and how ‘bad’
111 PE forecasts are often different than big FE errors on an individual storm basis. In Section 6 we

112 conclude with a discussion how FE gives a different view of TC forecasts and how the metric
113 can be applied in both operations and modeling.

114

115 **2. A TC Forecast is...**

116 I have conducted informal surveys of both laymen (e.g., my wife) and professionals
117 asking the question ‘what is a TC forecast?’ The answer almost always include the phrase
118 ‘where the storm is going’ and certainly track and position are part of a TC forecast. In fact,
119 both the National Hurricane Center (NHC) and JTWC forecast the surface wind field (10-m, 1-
120 min average wind speed). The 2-D surface wind field is *parameterized* by position (latitude,
121 longitude), intensity, maximum surface wind speed (Vmax) at a radius of maximum wind
122 (Rmax), and the radii of gale force (34 kt, “R34”), storm force (50 kt, “R50”) and hurricane force
123 (64 kt, “R64”) winds by quadrant. Furthermore, the NHC hurricane and tropical storm
124 warnings/watches that call for action by emergency management end users typically depend on
125 the onset of gale forecast winds (NHC). Similarly, JTWC warnings show the areal extent of gale
126 force winds and defines, according to Commander Pacific (PACOM instruction ???), the ‘no-
127 sail’ zone where US Navy ships are to avoid. Thus, the forecast of 34 kt (and greater) winds is
128 perhaps the single most critical aspect of a TC forecast but is traditionally not verified directly.

129 Examples of NHC and JTWC TC forecasts are given in Fig. 1. Supertyphoon HAISEN of
130 2020 was the first supertyphoon (intensity ≥ 130 kt) of the western North Pacific (WPAC) season
131 and had large impacts on the US military, as well as Japan and the Korean peninsula. The JTWC
132 warning in Fig 1a shows the ‘no-sail’ zone (hatched) constructed from the R34 forecast plus the
133 mean PE at each forecast time and the R34/R50/R64 radii. Note how the no-sail region considers
134 both the extent of gale forecast winds and the JTWC mean PE. In Fig. 1b we see the NHC

135 equivalent for hurricane DELTA of 2020 which shows the track forecast and the forecast ‘cone’
136 the represents the region with the most probable track (NHC 2020b). The same NHC forecast in
137 Fig. 1b is converted to the JTWC form by Fleet Weather Center Norfolk that shows the no-sail
138 zone in Fig 1c comparable to Fig. 1a. Clearly, the largest contribution to FE comes from PE,
139 especially at the longer lead times.

140 Suppose the forecaster makes a perfect R34 forecasts (i.e., gets the size of the radius of
141 the 34 kt winds correct), but if the $PE > 2 * R34$ then the forecast winds will be outside the
142 observed extent of gale force winds and the forecast would be a (complete) ‘bust.’ In this case,
143 IE can be anything and the forecast of the onset of 34 kt winds will be always wrong. This is the
144 key concept behind the proposed FE.

145 Historically, FE was defined as *only* equal PE until the 1980s when IE was also assessed
146 in the JTWC Annual Typhoon/Tropical Cyclone Reports dating from 1959-present. In the early
147 era of TC forecasting (1960/70s) mean 24-h PE was order 150 n mi and the typical R34 in both
148 WPAC and the Atlantic (LANT) basins was around 100 n mi (see Fig. 7). Thus, forecasts were
149 essentially worthless when $PE > 200$ n mi. Of course, there is great variability in PE for
150 individual storms and/or seasons, but this scaling argument shows the importance of PE
151 (mathematically positive definite and unbounded) relative to R34 (physically bounded).

152

153 **3. Nature of the TC PE forecast metric**

154 In this section we compare a standard mid-latitude metric (anomaly correlation
155 coefficient or ACC) to TC PE to show important differences that can confound interpretation of
156 TC forecast quality. The essential problem is that the mean ACC has the same number of
157 forecasts or cases at each lead time over the same area (e.g., northern Hemisphere) and at the

158 same pressure level (e.g., 500 hPa). Further, the ACC score comes from the model solution (a
159 continuous process) whereas TC PE is calculated from a ‘cyclone tracker’ as a feature in the
160 model solution. Further, the number of cases at each tau will be different with the possibility of
161 two or more storms verifying at single time. From these considerations alone, statistics of mid-
162 latitude ACC are very different than statistics of TC PE – a classic ‘apples’ v ‘oranges’ problem.

163

164 *a. Mid-latitude die off curve*

165

166 As introduced in Section 1, the 5DNAC (5-d NHEM 500 hPa ACC) has been used to
167 measure progress in operational global NWP models over the last 40 y (Simmons and
168 Hollingsworth 2002) and is a headline score still used to by NWP centers for monitoring (e.g.,
169 detecting ‘dropouts’) and for model intercomparison. Synoptic meteorologists have also found
170 that when 5DNAC goes below 0.6, the model has little or no skill for deterministic weather
171 forecasts (Hollingsworth et al. 1980).

172 Fig. 2 gives a recent 30-d time series of 5DNAC, and the mean value, for the ‘big 5’ global
173 modeling systems. Figure 2 shows that for this period that the United Kingdom Meteorological
174 Office (UKMO) had the highest ACC score (0.919) followed by the European Centre for
175 Medium-Range Weather (ECMWF) (0.917) forecasts and the National Centers for
176 Environmental Prediction (NCEP) Global Forecast System (GFS) (0.909). The mean anomaly
177 correlation for lead times 0-192 h are then plotted to show how skill degrades in time. This so-
178 called ‘die off’ curve is given in Fig. 3 where the 0.6 line is added to indicate the ‘no-skill’
179 forecast time. For this 30-d period ECMWF has the longest skill lead time ~ 8.5 d. We also see
180 that the UKMO has a higher mean at day 5 (120 h from the time series in Fig.2) thus showing

181 how a single mean, while useful as a benchmark/standard, only partially represents model
182 forecast quality and error growth.

183

184 *b. TC mean PE curve*

185 For TC PE we analyze the historic 2020 LANT season featuring a record number of
186 storms (31), with many of them making landfall over the US (3 in Gulf of Mexico). Three
187 models are considered: 1) HWRF (Hurricane Weather Research Forecast model run at NCEP
188 using GFS initial and lateral boundary conditions see DTC 2018); 2) GFS (EMC 2020) and 3)
189 the deterministic run of ECMWF IFS called HRES (ECMWF 2020). Over the last 15 or so
190 years, ECMWF has been the ‘best’ model and is the gold standard for TC forecasting as will be
191 discussed below. We also note that we are comparing two versions of the GFS – the global
192 model itself and the GFS solution as dynamically downscaled by the HWRF limited-area model.

193 Fig. 4 gives two versions of the TC equivalent of the mid-latitude ‘die off’ curve; i.e.,
194 how the mean error grows as a function of forecast lead time. Fig. 4a shows the traditional line
195 plot as would be found in the yearly NHC verification reports e.g., Cangialosi 2020. The forecast
196 tracks for HWRF and the GFS come from NHC and the ECMWF tracks were calculated locally
197 using near-native resolution grids and the GFDL tracker (Marchok 2021). The post-season
198 reanalysis or ‘final’ best track is not yet available; we use instead the ‘working’ best track, as
199 analyzed by the NHC Hurricane Specialists during the forecast process. While the final post-
200 season, best track will differ from the operational working best track, there will be little change
201 in the mean PE mostly because of the many ‘fixes’ (e.g., observational position/intensity
202 estimates from aircraft, satellite, and radar) available in 2020 make the working best track more
203 accurate than in the best tracks for previous years, especially for position.

204 The basic result is that all three models have similar error growth, and that ECMWF has
205 generally lower mean PE, but this standard plot (Fig 4a) is not as informative as Fig. 4b which
206 shows the distribution of PE and the median. Here, the median is similar in the short range of 0-
207 36 h, but at the longer lead times (72-120 h) HWRF has a higher median at 72 and 96 h, but at
208 120 h HWRF has the lowest median. A basic difference in the 24 and 120 h mean PE is that the
209 number of cases is much lower (~60) at 120 h than at 24 h (~260). Another difference is that the
210 comparison in Fig. 4 is heterogenous (all cases) vice homogeneous (only cases in common are
211 included) for two reasons. The model tracker will not always find a TC vortex at all verifying
212 times because of communication and/or tracker failures but more commonly because the model
213 dissipated the storm. The other reason is that the more models in a homogeneous comparison,
214 the lower the number of cases so that a homogeneous mean is conditioned on *all* the models
215 being available simultaneously. The heterogeneous comparison better represents the *individual*
216 model statistics. In contrast to even 10 y ago, the probability of detection (forecasting the
217 observed storm) is around 90% for all the models at 120 h, i.e., the models almost always make a
218 forecast.

219 For statistical significance tests, homogeneous statistics should be used, but for
220 forecasting purposes even small (~10%) differences in mean PE that may not be statistically
221 significant are observable by the forecaster (based on personal experience at both NHC and
222 JTWC). Further, the change in number of cases with forecast time makes assessing general
223 model quality difficult as discussed below.

224

225 *c. R34 distribution as a function of forecast time*

226 To show how cases in a 72-h forecast are different that at 24 or 120 h consider the
227 histogram of R34 in Fig. 5. The histograms at 24, 72 and 120 h were constructed from best
228 track data in the LANT 2010-2020 by finding the verifying R34 assuming a forecast is made
229 from each time in the best track and following the standard verification rule that the storm must
230 be classified as a TC initially and at the verifying time. Initially there will be 100% cases, but
231 for short-lived storms (say 48 h) there would be no verifying R34 at 72 h, so the percentage of
232 ‘forecasts’ always decreases. For the LANT basin and in the 10-y period, the percentage of
233 verifying 24, 72, and 120 h is 85%, 51% and 31% respectively and while not shown, the number
234 of storms at these forecast times will also decrease to the point that one or two TCs could
235 dominate the seasonal mean PE at 120 h.

236 The fundamental problem is that the 120-h mean PE represents a different class of TCs
237 than at 24 h. Thus, the PE error growth curves cannot be interpreted in the same way as ACC die
238 off. From the histogram in Fig. 5 we see that the median R34 at tau0 is 85 n mi (not plotted), at
239 24 h 90 n mi, at 72 h 112 n mi and 120h 128 n mi. Not only are the number of cases and number
240 of storms different, but the storms are structurally different.

241

242 *d.. TC mean IE curve*

243 Intensity (Vmax) forecasts have many issues, but the main one for a model is that Vmax
244 is strongly dependent on spatial grid resolution. Lower resolution models will always have
245 lower intensity forecasts (Walsh et al. 2007). A solution used in operations is to start the model
246 intensity forecast from the same point through a bias correction procedure that applies an ‘offset’
247 or the initial model-observed intensity difference to subsequent forecasts. The full offset is
248 applied at tau 0 and then is decreases linearly to 0 at a later tau. For global models, the 0 offset

249 tau is 72 h, for limited-area models is 24 h. The procedure used here is identical to that in
250 NHC/JTWC operations.

251 The standard IE statistic is the *absolute* mean because unlike PE, IE can be both positive
252 (too strong) or negative (too weak). The mean IE (mIE) is thus a bias and for most models (bias
253 defined as model-observed) is negative. Fig. 6 gives the bias-corrected IE for the same cases as
254 in Fig. 4 except that we plot both the standard absolute mean (amIE; lines) and the distribution as
255 well (mIE; bars). All models have nearly 0 mean abs IE initially because of the bias correction
256 but by 72 h the amIE has nearly the same magnitude as the bias (thick horizontal black lines in
257 the mIE box-and-whiskers) for the global models, except for HWRF which is nearly unbiased.
258 This is truly remarkable. Unlike the global models, the HWRF amIE does *not* come from bias.

259 We also show in Fig. 6 the ECMWF intensity forecast from their tracker (dark gold) and
260 the one run locally using lower resolution grids (light gold). Somewhat surprisingly, the local
261 tracker, using 0.25 deg grids v the ECMWF tracker using ~0.10 deg grids, has very similar bias
262 although at 120 h the ECMWF tracker bias is somewhat less.

263 Another feature of the IE curves is that curves plateaus around 36-48 h unlike PE which
264 always grows in time. There must be fundamental difference in this metric as discussed by
265 (Kieu and Moon 2016) but in operational verification IE is treated in the same way as PE and is
266 seemingly given equal weight when assessing overall forecast skill or FE.

267

268 *e. PE relative to a baseline*

269 Another way to compare model PE is to calculate the percent improvement relative to a
270 baseline as in equation 1:

271
$$\%IMP_{PE} = \frac{(mPE_{baseline} - mPE_{model})}{mPE_{baseline}} \cdot 100 \quad (1)$$

272 If the model mPE is less than the baseline (lower PE) then $\%IMP_{PE}$ will be positive and bounded
273 to 100%. Conversely, a negative $\%IMP_{PE}$ means the model has a higher mPE and is unbounded.

274 NHC uses the ‘no-skill’ CLImatology and PERsistence model CLIPER (Abserson 1998)
275 as the baseline so that $\%IMP_{PE}$ represents ‘skill.’ In the 1970s few aids had ‘skill’ or were
276 unable to make predictions with mPE lower than CLIPER. Another reasons for using a CLIPER
277 baseline is in how the seasonal CLIPER mPE can represent the year-to-year variation in forecast
278 difficulty. For example, a season with mostly ‘straight runners’ (i.e., TCs that follow an
279 essentially linear path) will be more climatological (when averaging tracks in a 30-y period) and
280 CLIPER mPE will be lower. The problem with using CLIPER circa 2020 is that current the
281 official and model mPE is now almost 30% that of CLIPER whereas in the 1970s forecasters
282 could rarely ‘beat’ CLIPER. Consequently, $\%IMP_{PE}$ in the annual NHC reports is
283 approximately 50-60% for most models making inter-model comparison more difficult.

284 A better alternative is to use a global model as the baseline; to measure improvements
285 relative to the global model that drives all forecasts, especially limited-area models such as
286 HWRF. In Fig. 7 we show the $\%IMP_{PE}$ for HWRF and ECMWF relative to the GFS for the
287 same basin/year as in Fig. 4 (the 2020 LANT season). The effect of vortex initialization in the
288 HWRF is clearly seen in the much lower initial PE (HWRF mPE of 8 n mi v GFS mPE of 15 n
289 mi and a $\%IMP_{PE} > 50\%$ from Fig. 4 table). However, ECMWF is better than the GFS (~10%)
290 at *all* lead times which is remarkable given that ECMWF does no TC vortex initialization in
291 operations (and never has). Further, a mPE of ~100 n mi at 72 h in the LANT is incredible when
292 in the 1980s the mPE for any model was ~360 n mi!

293 To gain perspective on the operational significance of this massive reduction in mPE,
294 consider the area covered by a radius of mPE in Fig. 8. In the 1980s, a 72-h forecast of a storm

295 in the middle of the Gulf of Mexico could make landfall almost anywhere because its mPE was
296 360 n mi, whereas in 2020 the forecast would be much more accurate due to the reduction in
297 mPE to 100 n mi. In terms of area, this is a 92% reduction in the threat zone and must be
298 considered a singular and perhaps one of the greatest NWP success stories of the last several
299 decades.

300 To further put perspective on the 2020 mPE we plot the $\%IMP_{PE}$ from 2007-2020 and
301 for the main northern Hemisphere basins in Fig. 9a. ECMWF had better mPE than the GFS at all
302 forecast times not only in the LANT, but also in the other basins. Furthermore, ECMWF was
303 better than the GFS in all three basins and at all 5 lead times (i.e., ‘ran the board’) in 2010, 2011,
304 2015, 2017, 2018, and 2020.

305 Other notable features:

- 306 • Since 2012, HWRF rarely outperformed the global model providing its lateral
307 boundary conditions (GFS) and in 2020 only at tau 96 and 120 in the WPAC and
308 LANT. This is noteworthy as limited-area models are not expected to provide
309 long-term guidance when the TC forecast track is dominated by the lateral
310 boundary conditions, and thus suggests that the skill of the GFS is actually
311 decreasing at these longer forecast hours.
- 312 • The big change from 2007 to 2008 in the ECMWF mPE can be attributed to
313 change in physics as discussed in Fiorino 2009.
- 314 • The GFS improvement (ECMWF less green) from 2011 to 2012 can be mostly
315 attributed to improved data assimilation used in the GFS (hybrid 3DVAR-EnKF).
- 316 • The exception to ECMWF providing *the best short-term guidance* (tau 12 & 24)
317 was in the WPAC in 2019.

318 Applying the same analysis to the latest ECMWF reanalysis ERA5 (Hersbach et al. 2020)
319 and the operational run of the model (ECMWF HRES) is show in Fig. 9b. Note the vertical bar
320 to the right of the box plot and horizontal bar below. These give the mean for the 14-y period
321 (vertical) and the yearly mean (horizontal) where we can clearly see how ERA5 forecasts are
322 better than from the real-time run. The ERA5 model and data assimilation are representative of
323 the operational model circa 2016 with a resolution of ~ 31 km and more notably is a coupled
324 reanalysis using the latest 4DVAR data assimilation with a window that looks forward 12 h (not
325 possible in operations). The 2007 problem (Fiorino 2009) is gone in ERA5 and all years and tau-
326 basins are greener, especially before 2016 but even into 2020.

327

328 **4. TC forecast error metric**

329 We have seen that error growth curves of TC mean PE, while visually similar to AC die
330 off plots, represent very different measure of model forecast error because TC PE is for a unique
331 and discrete event (the TC) in the model solution whereas AC is calculated from the continuous
332 model solution itself. Also, the number of events in mean PE not only decreases with increasing
333 forecast time but come from a different class of storm. But perhaps the larger problem with using
334 PE and IE to measure forecast quality is the indirect relationship to the ‘actionable’ part of the
335 TC forecast – the onset and areal extent of 34+ kt winds that can range from near 0 n mi for
336 weak tropical storms or greater than 200 n mi for large (and generally intense) hurricanes.

337 To overcome these deficiencies we propose a metric FE that depends on R34 and is a
338 function of *both* PE and IE. First, we only require the model forecast position and intensity and
339 use the best track R34 as the forecast R34 (a perfect R34 forecast). We also use the best track
340 radius of maximum wind (Rmax) for the model forecast (a prefect Rmax forecast) when

341 available. When the best track Rmax is not available a statistical relationship between Rmax and
342 Vmax is used (Quiring et al. 2011). To restate, by using the verifying R34 (and Rmax) we only
343 require the model forecast position and intensity and the proposed FE metric is a lower bound.

344 We then construct a symmetric wind profile used by Fiorino and Elsberry 1989 and
345 shown in Eq. 2 below that only requires two (v, r) points to analytically specify the winds from
346 $r=0, R_{34}$:

$$347 \quad v(r) = V_{max} \left(\frac{r}{R_{max}} \right) \left\{ \exp \left(\frac{1}{b} \left(1 - \left(\frac{r}{R_{max}} \right)^b \right) \right\} \quad (2)$$

348
349 For the FE metric we use the ordered pairs (V_{max}, R_{max}) and the (34 kt, R34) to solve for b.

350 Given the wind profile, we integrate the square of $v(r)$ from $r=0, R_{34}$ and to calculate the
351 Integrated Kinetic Energy (IKE) (Powell and Reinhold 2007) for the model and best track.

352 Whereas IKE in Powell and Reinhold extends past R34, we only consider the kinetic energy of
353 the significant winds (≥ 34 kt).

354 FE is defined as the error IKE which is the IKE outside the overlap or ‘symmetric
355 difference’ as shown in Fig. 10. Python is used for the integration of the analytic wind profile
356 and to find the area outside the overlap. When $PE = 0$, $FE = 0$ even if $IE=0$ because there the
357 two R34 circles perfectly overlap – assuming a perfect R34 forecast. When $PE = 2R_{34}$, FE is
358 maximized and there is no overlap. For $PE > 2R_{34}$ we assume FE grows linearly from the $2R_{34}$
359 FE with increasing PE. The FE defined here is a lower bound for the best case of perfect R34
360 and Rmax forecasts.

361 The FE for the same models/season are shown in Fig. 11. Fortunately, the magnitude of
362 FE in TJ is the same order as PE in n mi; simplifies intercomparison.

363

364 **5. Application to Model Development**

365 Recently, advances in TC forecasts have largely been made via improvements to the data
366 assimilation (DA) methods and physics parameterizations; this section describes how TC metrics
367 can be used to aid in these developments. There are many ways to measure DA aspects and it is
368 naturally assumed that TC vortex initialization would be critical for good TC forecasts. Perhaps
369 the easiest way to measure the quality of TC vortex analysis and initialization is to quantify the
370 initial mPE and mIE. In Fig 13a. we find that over the 14-y period 2007-2020, HWRF mPE is
371 10 n mi, GFS is 14 n mi, ECMWF (HRES) and ERA5 are 19 and 18 n mi, respectively. For
372 initial mIE (Fig 13b), HWRF has a very slight negative bias of -1 kt, GFS is -14 kt, ECMWF
373 HRES is -18 kt and ERA5 is -20 kt.

374 Clearly, HWRF has a superior vortex initialization – the lowest initial PE and a nearly
375 unbiased initial intensity. However, even 12 h into the forecast both ECMWF(HRES) and ERA5
376 have lower PE than HWRF/GFS (Fig. 9b). If forecasting is the ‘acid test’ of an analysis, then
377 ECMWF must have the better vortex initialization, but they *do not assimilate* the TC ‘vitals’ or
378 TC position, motion and structure as analyzed in the real-time by the human forecasters (Trahan
379 and Sparling 2012). HWRF/GFS *do assimilate* TC vitals, although ‘vortex relocation’ was
380 recently turned off in the GFS (EMC 2020).

381 The inescapable conclusion is that for TC track forecasting model physics and DA are
382 critical, especially in the global models. Further, it is not obvious how TC metrics can be used in
383 model development, except possibly for identification of systematic physics errors for ‘large’ TC
384 forecast errors.

385

386 *a. 72-h PE and FE*

387 In Section 2 we discussed how the day 5 NHEM AC is a standard metric in global NWP.
388 For TCs, it is argued that a best single number metric is a day 3 or 72 h on two grounds. First,
389 the number of verifiable forecasts at 72 h is greater the 50% (i.e., 61% in the LANT and 67% in
390 the WPAC as shown in Fig. 5) and that by 72 h the forecast is less constrained by the initial
391 conditions and more determined by the model physics. The second reason is human forecast
392 process itself.

393 Typhoon Duty Officer (TDO equivalent of NHC hurricane specialist) training at JTWC
394 emphasizes that TC track forecast is a ‘connect-the-dots’ problem. Most of the time in a warning
395 cycle is devoted to analyzing TC location and structure – a process that involves the qualitative
396 assimilation of many observational and disparate data sources from ship reports to satellite
397 imagery that are often not assimilated in the models. Given a solid initial position (and motion),
398 the deterministic forecast problem is where will the storm be in 72 h (if we forecast the TC to
399 exist for 3 or more days)? Once that is set, typically constrained by the 72-h consensus forecast
400 and the previous warning, the issue is how to connect the initial and 72 h points and then to
401 extrapolate from 72 h to 120 h (day 5) in a way the makes both physical and synoptic sense.

402 As discussed in Section 4, when the PE is $\geq 2 \cdot R_{34}$, IE only changes the magnitude of FE
403 not the location component, i.e., the forecast high-wind area is outside the observed high-wind
404 area. This property of FE is a consequence of the perfect R_{34} forecast assumption and therefore
405 $PE \geq 2 \cdot R_{34}$ represent a bound on usefulness in terms of wind ≥ 34 kt. If a typical R_{34} is 100 n
406 mi and PE is 100 n mi and if the PE can be either left or right of the TC center then a PE 100 n
407 mi would be 200 n mi and equal to $2 \cdot R_{34}$. The forecast time when PE is around 100 n mi is
408 around 72 h as seen in Fig. 4a. For FE, note that the error growth in mean FE is linear between
409 0 and somewhere between 48 and 72 h in Fig. 11a. The similarity between the PE and FE

410 growth curves is not surprising given that the contribution of IE to FE is around 10% as shown in
411 Fig. 12. Thus, the 72 h lead time can be considered the ‘linchpin’ of the track forecast.

412

413 *b. large 72-h errors – PE v FE*

414 We identify problem storms/cases by plotting the tau72 mPE and mFE for each
415 individual storm in Fig. 14 where mFE by-storm is normalized by the ratio of LANT 2020
416 season 72-h mPE/mFE (1.5) so that the units and scale of FE is approximately the same. First
417 note that the number of TCs contributing to the 72-h mPE/mFE is 21 or a rate of 68% for the
418 total of 31 storms in the 2020 LANT season. This rate is substantially higher than the 10-y rate
419 of 51% (Fig. 5) indicating an active season. However, the number of cases per storm also varies
420 so that high PE/FE may make only a small contribution to the total (season in this paper) mean.
421 In Fig. 15, we first convert the PE/FE to $\%IMP_{PE}$ or contribution to the total mean following Eq.
422 1 and then normalize by the (number of cases per storm / total number of cases). This variable
423 more explicitly shows how each storm contributes to the seasonal mean.

424 Comparing Fig. 15a (PE) to Fig. 15b (FE) we see a different distribution of problem
425 storms. For example, the biggest problem storm (longest bar below the 0 line) for the GFS was
426 **29L.2020** (the 29th TC in the atLantic in 2020 or hurricane **ETA**). For FE, **29L.2020** was
427 slightly below the 0 line, but TCs **17L.2020** (hurricane **PAULETTE**) and **20L.2020** (hurricane
428 **TEDDY**) were the larger problem storms. Table 1 summarizes these problem storms for PE v
429 FE and gives URLs for track plots that highlight the 72-h position. An example from the track
430 plot site is given in Fig. 16 for **14L.2020** (hurricane **MARCO**) at 2020082112 (12 UTC 21 Aug
431 2020) when the GFS had the largest 72-h PE error as the storm entered the Gulf of Mexico.

432 The storms with large 72-h FE in Table 1 are cases when the verifying TC was both large
433 (R34 > 100 n mi) and intense (Vmax > 65 kt). Table 2 shows how IKE varies as a function of
434 size (R34) and intensity. For a moderate hurricane (90 kt) the IKE for a large storm is over 10X
435 that of a small storm with an R34 of 50 n mi. Thus, FE will be large even for ‘good’ PE at 72 h
436 (~seasonal mPE) and furthermore, these storms would be more significant operationally because
437 the high IKE means greater storm surge and wind impacts.

438 To reiterate, examining error by TC using both the traditional PE and the new FE metric
439 may help to queue the modeler to find possible model problems and, perhaps more importantly,
440 identify TC forecasts errors that have the greatest operational significance.

441

442 **6. Summary and Conclusions**

443 At both JTWC and NHC, the extent and onset of 34 kt winds drives the warning process
444 and defines ‘actionable’ intelligence for impacted users. Thus, a TC forecast is a *field* of 1-min
445 (10-min average is the WMO standard) 10-m or *surface wind* and the new metric of Forecast
446 Error (FE) is based on Integrated Kinetic Energy (IKE; Powell and Reinhold 2007) within the
447 TC circulation of wind ≥ 34 kt.

448 The new metric for TC forecast quality is a function of *both* traditional metrics of PE and
449 IE. Current operational verification treats PE and IE separately (e.g., Cangialosi 2020) even
450 though TC intensity is physically related to track. Furthermore, science-based forecasting
451 demands a holistic consideration of TC structure, not just where the storm is going. Another
452 point of departure is how we compare models against a *model* baseline vice a no-skill CLIPER-
453 type baseline and how we breakdown the seasonal means by TC to identify how each individual
454 and unique storm contributes to the overall mean.

455 We considered two global models for the historic 2020 LANT hurricane season:

- 456 • NCEP GFS
 - 457 ○ the operational version of the GFS itself
 - 458 ○ the GFS dynamically downscaled by the limited-area HWRF model
- 459 • ECMWF global model
 - 460 ○ the operational version that makes the 10-d deterministic run (called HRES)
 - 461 ○ the circa 2016 version of HRES model cycle used in the latest ECMWF
 - 462 reanalysis ERA5.

463 The basic result is that all the models did well in 2020 regarding mean PE (e.g., relative
464 the 2015-2019 mean), but when using the GFS as a baseline, we found that ECMWF had
465 superior PE, expressed as a percentage improvement, at all forecast times except the model
466 initialization time (0 h). We analyzed a 14-y period (2007-2020), including the main NHEM
467 basins, and found that HWRF rarely improves the host model GFS TC track forecast. ECMWF
468 on the other hand was generally superior in its forecasted track relative to both GFS and HWRF.
469 We then compared the ECMWF HRES forecasts to ERA5 reanalysis forecasts, again using the
470 operational GFS as the baseline, and found substantially better and more consistent year-to-year
471 improvements in the ERA5 over the real-time HRES model. The conclusion is that the model
472 very much matters for TC track prediction and more so than low initial position and intensity
473 error (vortex initialization aspects).

474 The new metric of FE was defined as the IKE outside the union of two circles spanned by
475 the forecasted R34 and the observed R34 but offset by the PE as illustrated in Fig. 10. We
476 assumed the forecast R34 was the same as observed to make FE only a function of PE and IE.
477 The mFE error growth curves are similar to mPE growth, except mFE grows faster when the PE

478 is greater than 2R34 after tau48. However, the more important finding is that if the model makes
479 a perfect intensity forecast, the forecast error is only reduced 10-12% and mostly for the early
480 forecasts (0-48 h). In other words, track is still 90% of the TC forecast problem.

481 Application of the standard PE and new FE metric to model diagnosis and verification
482 was then considered. If the quality of the TC vortex analysis is measured by initial PE and initial
483 IE, the HWRF forecast would be considered the best. However, ECMWF had the highest initial
484 PE/IE but the best 12- and 24-h PE which implies that the HWRF vortex initialization degrades
485 the short-term forecasts and hence the initial motion and development of the TC. What is
486 particularly impressive is that the mPE of the ECMWF forecasts at 12 h is about 25 n mi, which
487 is close to observational uncertainty in TC location itself (Landsea and Franklin 2013)!

488 We then considered how each storm contributed to the seasonal mean at 72 h – a critical
489 forecast time for storms predicted to last more the 3 d. In 2020, 68% out of 31 TCs in the LANT
490 had one or more verifying positions at 72 h. This is greater than the previous 10-y mean of 51%
491 as shown in Fig. 5. The most important conclusion regarding FE is that large ($R_{34} > 120$ n mi)
492 and intense ($V_{max} > 90$ kt) storms will have much higher IKE error (FE) for a PE that is close to
493 the mPE. All storms are not equal in how they affect the seasonal mean error.

494 Finally, the FE analysis presented here gives a different view of TC forecast error.
495 Application to other basins and years should help identify the big and operationally significant
496 TC forecast error to guide model and forecaster development.

497

498 *Acknowledgements.*

499 I first want to acknowledge and thank Vic Addison, CAPT USN (ret), for his very
500 thought-provoking comment during Captain's call while I was on active duty at FNMOC prior to
501 joining NHC. There are few days when I do not reflect on CAPT Addison's "you're only as
502 good as what you measure" statement. The late Charlie Neumann, the father of American digital
503 TC data, suggested to me in 1975 that 'forecast error = position error' and that intensity error
504 does not contribute to forecast error. Finally, my long association with JTWC as their models
505 officer and as a qualified Typhoon Duty Officer greatly influenced and contributed to my
506 thinking on TC forecasts and metrics, as did my 3 years as the NHC lead techniques
507 development (Charlie Neuman's billet). The experience of forecasting at both JTWC and NHC
508 has been personally invaluable. The ERA5 forecasts were kindly provided by Hans Hersbach,
509 the ECMWF reanalysis project lead.

510 *Data Availability Statement.*

511 All TC forecast and best track data available upon request.

512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534

REFERENCES

Aberson, S. D., 1998: Five-Day Tropical Cyclone Track Forecasts in the North Atlantic Basin. *Weather Forecast*, **13**, 1005-1015.

Benjamin, S. G., Brown, J. M., Brunet, G., Lynch, P., Saito, K., & Schlatter, T. W., 2019: 100 Years of Progress in Forecasting and NWP Applications, *Meteorological Monographs*, **59**, 13.1-13.67. Retrieved Jan 16, 2021, from <https://journals.ametsoc.org/view/journals/amsm/59/1/amsmmonographs-d-18-0020.1.xml>

Cangialosi, J.P. 2020: National Hurricane Center Forecast Verification Report. 2019 Hurricane Season. NHC: https://www.nhc.noaa.gov/verification/pdfs/Verification_2019.pdf

Development Testbed Center, 2018: Hurricane WRF (HWRF) Documentation. Available online at: <https://www.dtcenter.org/community-code/hurricane-wrf-hwrf/documentation>

European Centre for Medium-range Weather Forecasts, 2020: IFS documentation. Available online at: <https://www.ecmwf.int/en/publications/ifs-documentation>.

Fiorino, M., 2009: Record-setting performance of the ECMWF IFS in medium-range tropical cyclone track prediction. ECMWF Newsletter No. 118, 20-27 <http://www.ecmwf.int/sites/default/files/elibrary/2008/14607-newsletter-no118-winter-200809.pdf>

535 Fiorino, M., and R. L. Elsberry, 1989: Some Aspects of Vortex Structure Related to Tropical
536 Cyclone Motion. *J Atmos Sci*, **46**, 975-990.

537

538 Hersbach, H., and Coauthors, 2020: The ERA5 global reanalysis. *Q J Roy Meteor Soc*, **146**,
539 1999-2049.

540

541 Hollingsworth, A., K. Arpe, M. Tiedtke, M. Capaldo, and H. Savijarvi, 1980: The Performance
542 of a Medium-Range Forecast Model in Winter-Impact of Physical Parameterizations. *Mon*
543 *Weather Rev*, **108**, 1736-1773.

544 JTWC, 2020a: Annual Typhoon/Tropical Cyclone Reports:
545 <https://www.metoc.navy.mil/jtwc/jtwc.html?cyclone>

546 JTWC, 2020b: Indo-Pacific Command Instruction governing JTWC operations:
547 <https://www.usno.navy.mil/NOOC/nmfc-ph/RSS/jtwc/pubref/3140.html>

548

549 Landsea, C. W., and J. L. Franklin, 2013: Atlantic Hurricane Database Uncertainty and
550 Presentation of a New Database Format. *Mon Weather Rev*, **141**, 3576-3592.

551

552 Liu, Q., T. Marchok, H.-L. Pan, M. Bender, and S. J. Lord, 2000: Improvements in hurricane
553 initialization and forecasting at NCEP with global and regional (GFDL) models. NWS Tech.
554 Procedures Bull. 472, 7 pp. [Available online at:
555 [https://www.researchgate.net/publication/237459506_Improvements_in_Hurricane_Initialization](https://www.researchgate.net/publication/237459506_Improvements_in_Hurricane_Initialization_and_Forecasting_at_NCEP_With_Global_and_Regional_GFDL_Models)
556 [_and_Forecasting_at_NCEP_With_Global_and_Regional_GFDL_Models](https://www.researchgate.net/publication/237459506_Improvements_in_Hurricane_Initialization_and_Forecasting_at_NCEP_With_Global_and_Regional_GFDL_Models)]

557

558 Marchok, T.P., 2021: Important factors in the tracking of tropical cyclones in operational
559 models. *J Appl. Meteo Climo.*, submitted for publication.
560

561 National Hurricane Center, 2020a: National Hurricane Operations Plan [available on online at:
562 <https://www.icams-portal.gov/publications/nhop/nhop2.htm>
563

564 National Hurricane Center, 2020b: NHC Forecast Cone[available online at:
565 <https://www.nhc.noaa.gov/aboutcone.shtml>
566

567 Environmental Modeling Center, 2020: Global Spectral Model Reference Page available online
568 at:
569 https://www.emc.ncep.noaa.gov/emc/pages/numerical_forecast_systems/gfs/documentation.php
570

571 Powell, M. D., and T. A. Reinhold, 2007: Tropical cyclone destructive potential by integrated
572 kinetic energy. *B Am Meteorol Soc*, **88**, 513-+.
573

574 Simmons, A. J., and A. Hollingsworth, 2002: Some aspects of the improvement in skill of
575 numerical weather prediction. *Q J Roy Meteor Soc*, **128**, 647-677.
576

577 Trahan, S., and L. Sparling, 2012: An analysis of NCEP tropical cyclone vitals and potential
578 effects on forecasting models. *Wea. Fore.*, **27**, 744-756, [https://doi.org/10.1175/WAF-D-11-](https://doi.org/10.1175/WAF-D-11-00063.1)
579 [00063.1](https://doi.org/10.1175/WAF-D-11-00063.1).
580

581 Quiring, S., A. Schumacher, C. Labosier, and L. Y. Zhu, 2011: Variations in mean annual
582 tropical cyclone size in the Atlantic. *J Geophys Res-Atmos*, **116**.

583

584 Walsh, K. J. E., M. Fiorino, C. W. Landsea, and K. L. McInnes, 2007: Objectively determined
585 resolution-dependent threshold criteria for the detection of tropical cyclones in climate models
586 and reanalyses. *J Climate*, **20**, 2307-2314.

587

588

589

590

591

592

TABLES AND FIGURES

593 Table 1. Problem storms for the GFS as defined by 72-h PE and FE for 2020 LANT season

TC	Error Type	Comments	URL for track plots
14L (MARCO)	PE	GFS too far West and South as storm entered the Gulf of Mexico	http://tctrkveri.wxmap2.com/trk-14L-2020-MOD-CUR.htm
17L (PAULETTE)	FE	Large R34 increases makes FE larger than PE	http://tctrkveri.wxmap2.com/trk-17L-2020-MOD-CUR.htm
20L (TEDDY)	FE	Large R34 as with 17L, particularly before extratropical transition	http://tctrkveri.wxmap2.com/trk-20L-2020-MOD-CUR.htm
29L (ETA)	PE FE	Two cases of PE > 700 n mi	http://tctrkveri.wxmap2.com/trk-29L-2020-MOD-CUR.htm
31L (IOTA)	FE	R34 > 100 n mi	http://tctrkveri.wxmap2.com/trk-31L-2020-MOD-CUR.htm

594

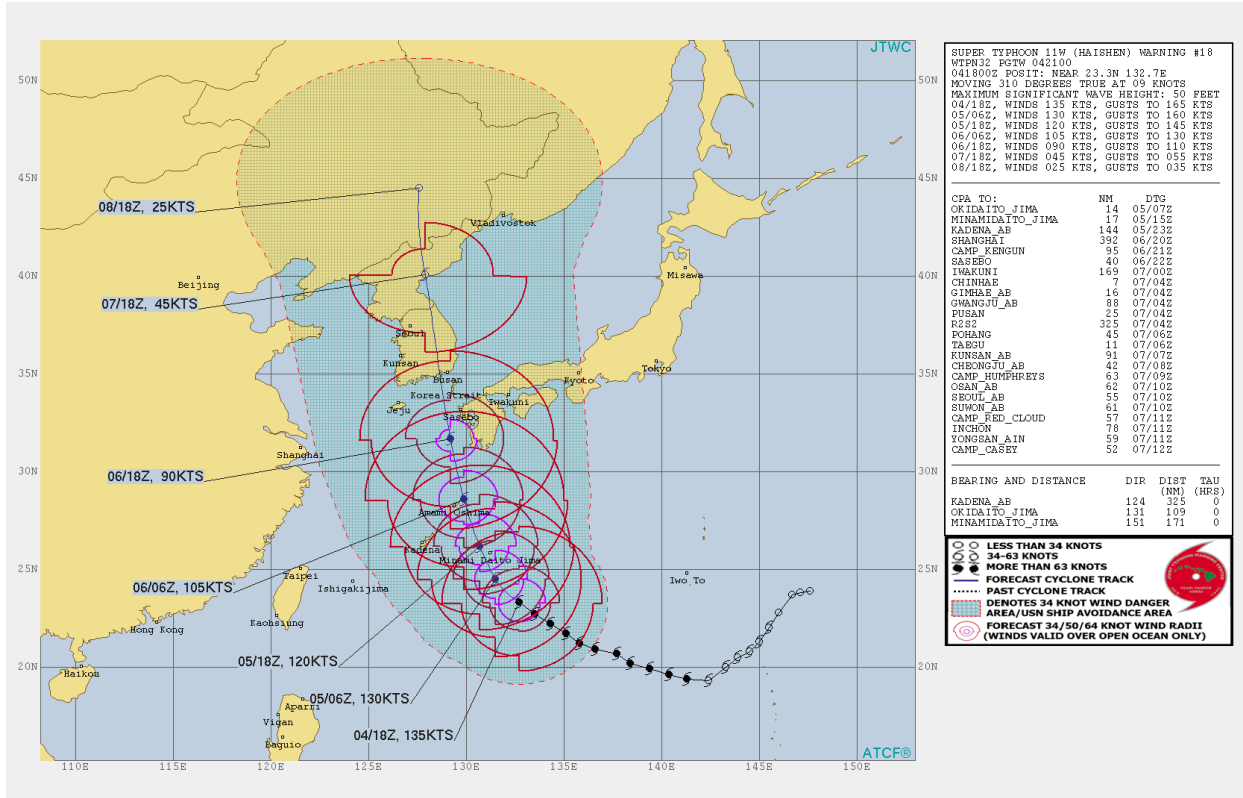
595

596 Table 2. IKE (TJ) for various Vmax and R34 using the Quirling et al. (2011) Vmax minus Rmax
597 relationship and the wind profile from Fiorino and Elsberry (1989)
598

Vmax (kt)	R34=50 n mi	R34 =100 n mi	R34=150 n mi
65	5	35	76
90	8	56	108
130	15	79	148

599
600

601
602




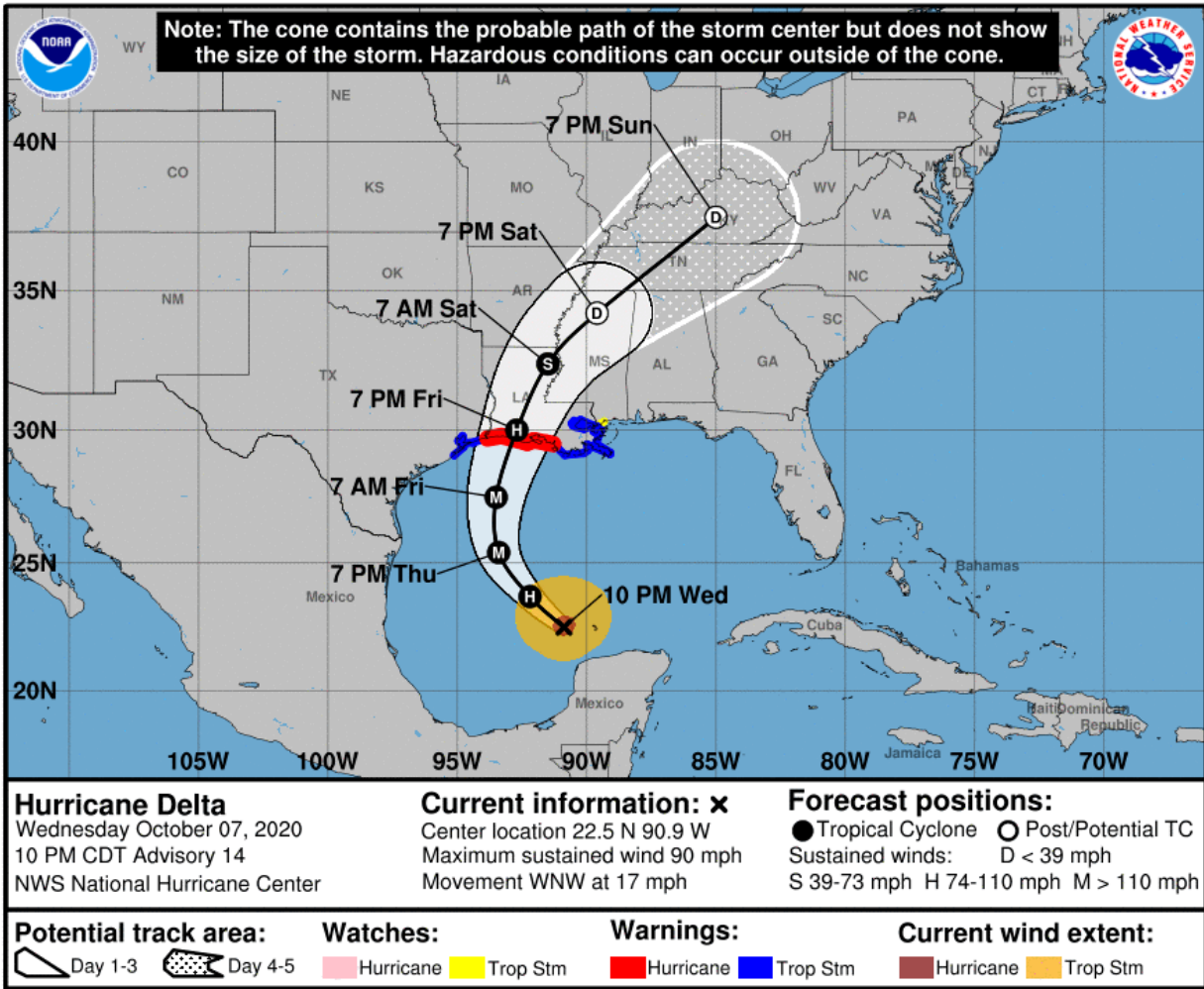
SUPER TYPHOON 11W (HAISHEN) WARNING #18
 WTPN32 PGTW 042100
 041800Z POSIT. NEAR 23.3N 132.7E
 MOVING 310 DEGREES TRUE AT 09 KNOTS
 MAXIMUM SIGNIFICANT WAVE HEIGHT: 50 FEET
 04/18Z, WINDS 135 KTS, GUSTS TO 165 KTS
 05/06Z, WINDS 130 KTS, GUSTS TO 160 KTS
 05/18Z, WINDS 120 KTS, GUSTS TO 145 KTS
 06/06Z, WINDS 105 KTS, GUSTS TO 130 KTS
 06/18Z, WINDS 090 KTS, GUSTS TO 110 KTS
 07/18Z, WINDS 045 KTS, GUSTS TO 055 KTS
 08/18Z, WINDS 025 KTS, GUSTS TO 035 KTS

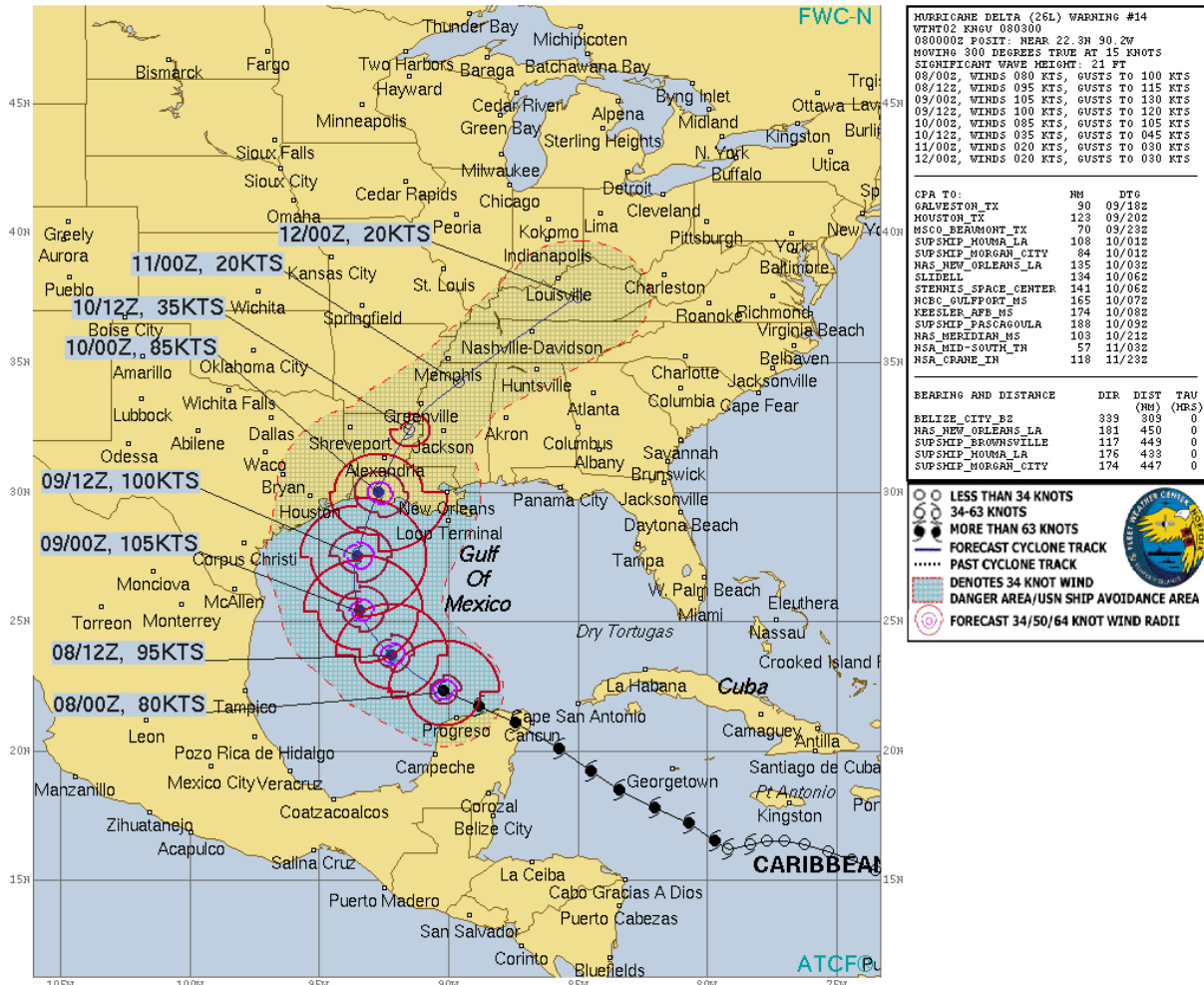
CPA TO:	NM	DTG
OKIDAITO JIMA	14	05/07Z
MINAMIDAITO JIMA	17	05/15Z
KADENA AB	144	05/23Z
SHANGHAI	392	06/20Z
CAMP KENGUN	95	06/21Z
SASEO	40	06/22Z
IMARUNI	169	07/00Z
CHINHAE	7	07/04Z
GIMHAE AB	16	07/04Z
GHWANGJU AB	88	07/04Z
PUSAN	25	07/04Z
P2B2	325	07/04Z
POHANG	45	07/06Z
ZASU	11	07/06Z
KUNSAN AB	91	07/07Z
CHEONGJU AB	42	07/08Z
CAMP HUMPHREYS	63	07/09Z
OSAN AB	62	07/10Z
SEOUL AB	55	07/10Z
SUNON AB	61	07/10Z
CAMP RED CLOUD	57	07/11Z
INCHON	78	07/11Z
YONGSAN AIN	59	07/11Z
CAMP CASEY	52	07/12Z

BEARING AND DISTANCE	DIR	DIST	TAU
		(NM)	(HRS)
KADENA AB	124	325	0
OKIDAITO JIMA	131	109	0
MINAMIDAITO JIMA	151	171	0

○ LESS THAN 34 KNOTS
 ⊙ 34-63 KNOTS
 ● MORE THAN 63 KNOTS
 PAST CYCLONE TRACK
 [Red Dashed Circle] DENOTES 34 KNOT WIND DANGER AREA/USN SHIP AVOIDANCE AREA
 [Red Solid Circle] FORECAST 34/50/64 KNOT WIND RADII (WINDS VALID OVER OPEN OCEAN ONLY)



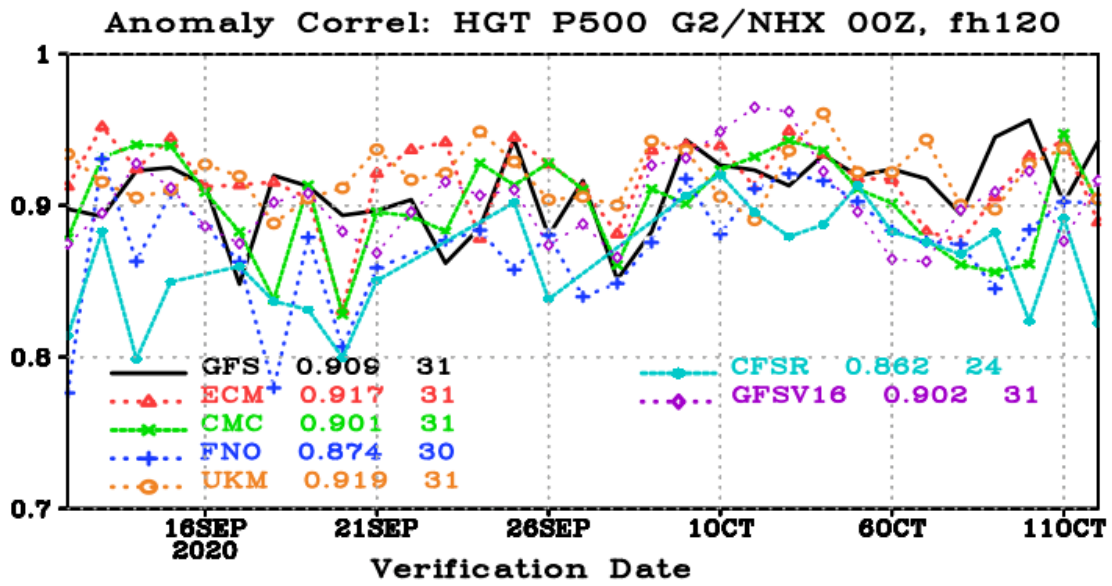




609

610 Figure 1: a) 18UTC 04 September 2020 (2020090418) JTWC warning (equivalent to NHC
 611 advisory) for Typhoon HAISEN; b) 00 UTC 08 October 2020 (2020100800) NHC advisory for
 612 hurricane DELTA; and c) Navy equivalent of same DELTA advisory which shows the no-sail
 613 zone and the wind radii.

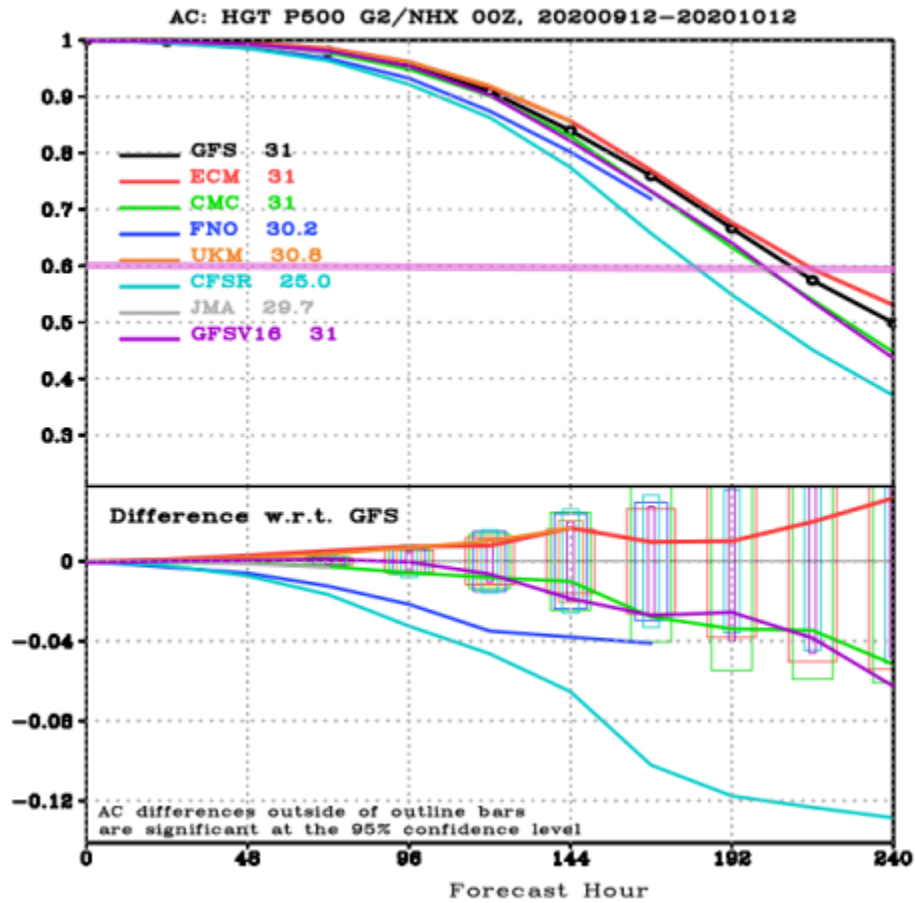
614



615

616 Figure 2: The 30-d time series of 5DNAC for 7 modeling systems over the period 20200912-

617 20201012 courtesy of https://www.emc.ncep.noaa.gov/gmb/SATS_vsdb/



619

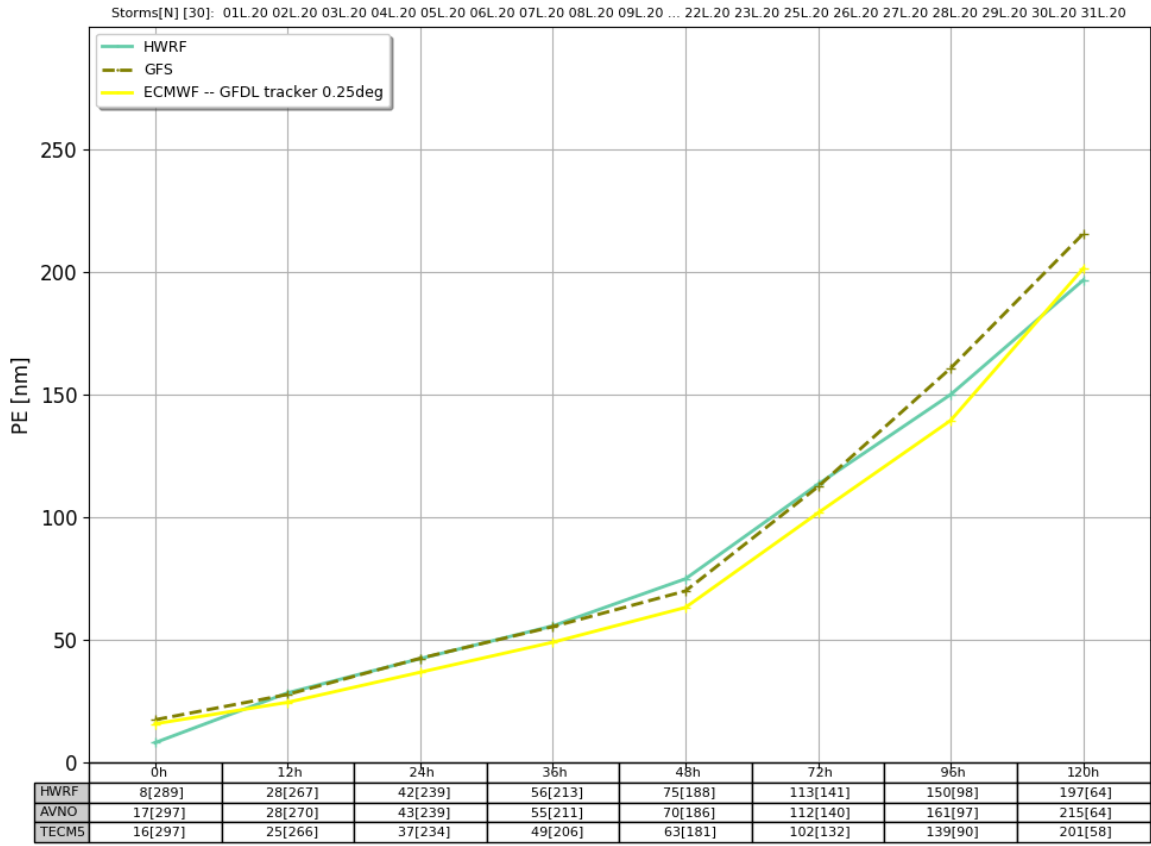
620

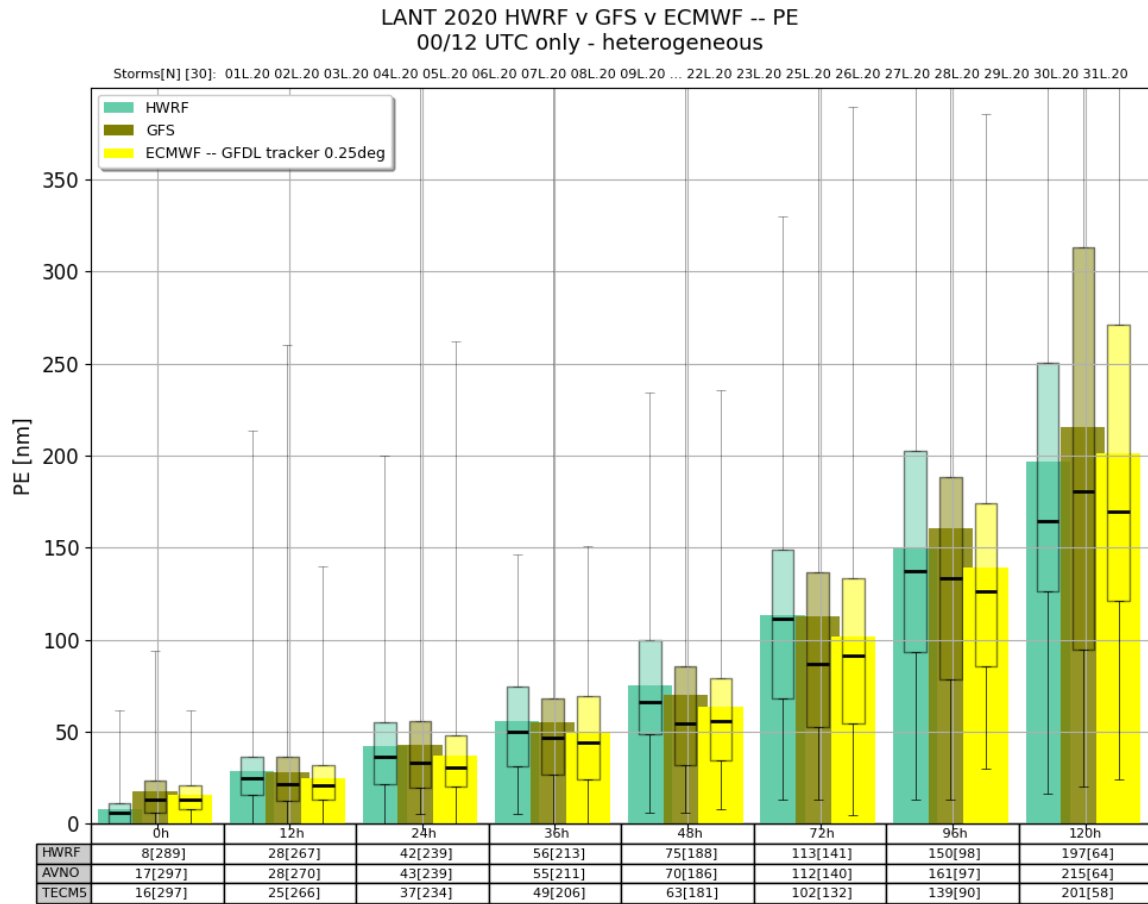
621 Figure 3: 500 hPa ‘die off’ curve from the mean of the time series in Fig. 2. The 0.6 line is
 622 drawn to show the point the point in time where the model forecast cannot be distinguished from
 623 climatology, i.e., no value in deterministic weather forecasting. Available from:

624 https://www.emc.ncep.noaa.gov/gmb/STATS_vsdb/

625

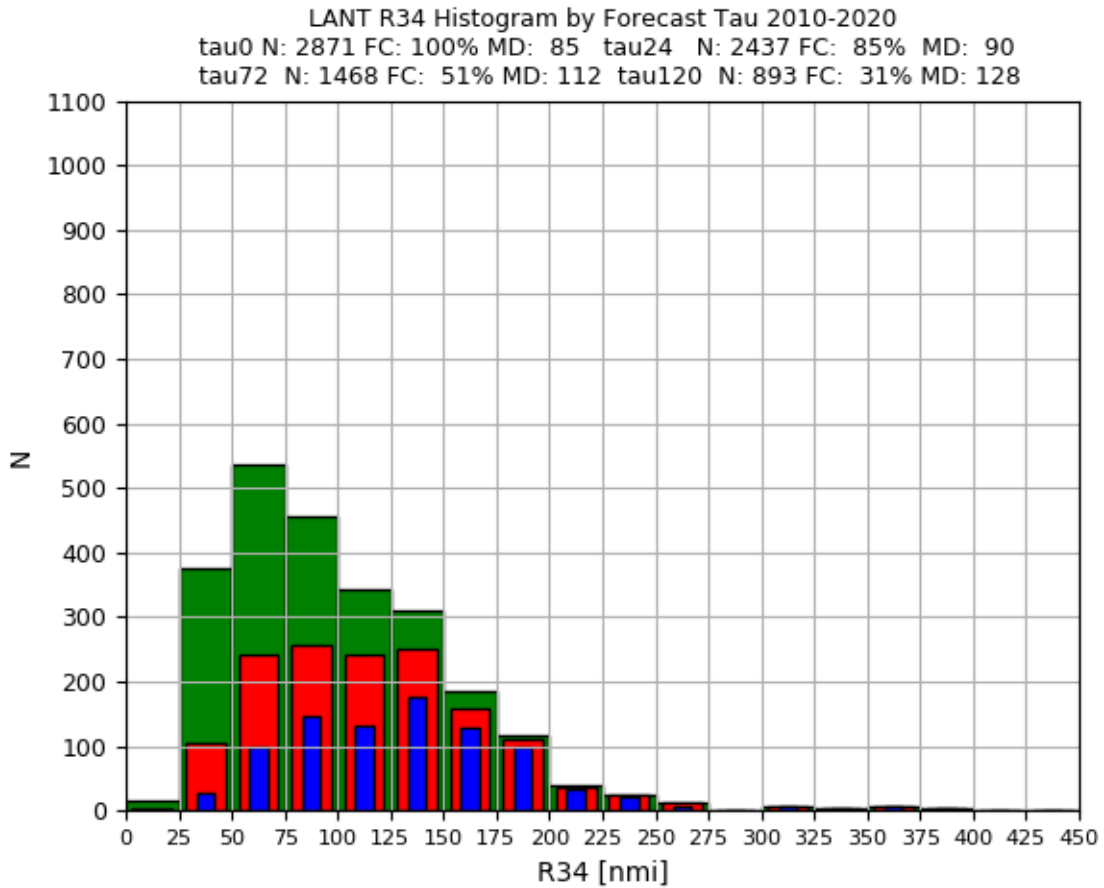
LANT 2020 HWRf v GFS v ECMWF -- PE -- traditional line
00/12 UTC only - heterogeneous





630

631 Figure 4. 2020 atANTic mean position error (PE) at standard forecast times for the HWRf, GFS
 632 and ECMWF models for 00/12 UTC only. The comparison is heterogeneous. Panel A is the
 633 traditional line plot similar to the 500 hPa dieoff curve. Panel B gives a box-whiskers version of
 634 the data in 4A showing the distribution of PE. In both panels, the mean value and [number of
 635 cases] is shown in the table below the plot. The thick bar in 4B is the mean and the thinner bar
 636 extends from the 25th to 75th percentiles with the median indicated by the solid black line.

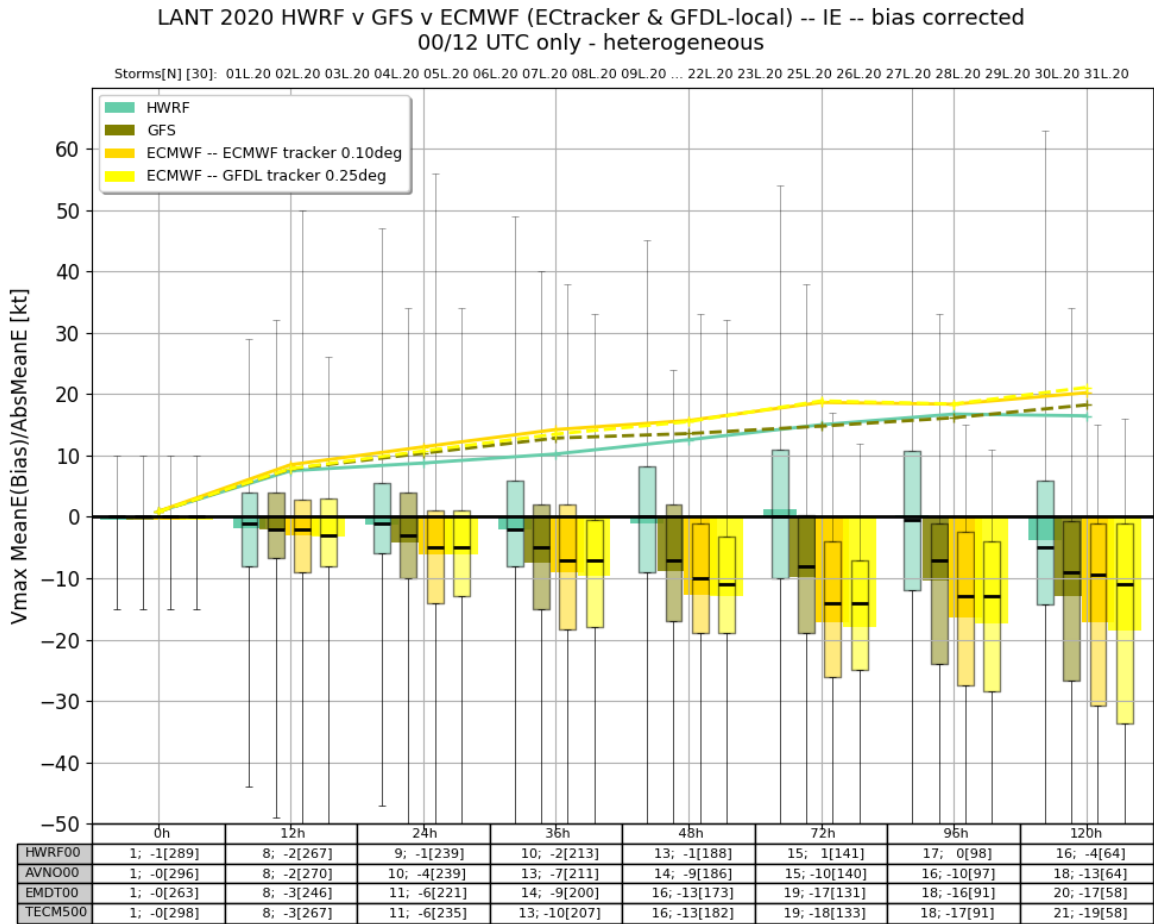


637
638

639

640 Figure 5. Histogram of the radius of 34 kt (R34) for the atLANTic 2010-2020 at three forecast
641 times 24, 72, and 120 h (green, red, and blue bars respectively). The % of forecasts and median
642 R34: 0 h) 100% and 85 n mi; 24 h) is 85% and 90 n mi; 72 h) is 51% and 112 nmi; and 120 h)
643 31% and 128 n mi.

644



645

646 Figure 6. 2020 atANTic mean absolute IE (lines) at standard forecast times for the HWRf, GFS

647 and ECMWF models for 00/12 UTC only. The comparison is heterogeneous as in Fig. 4. The

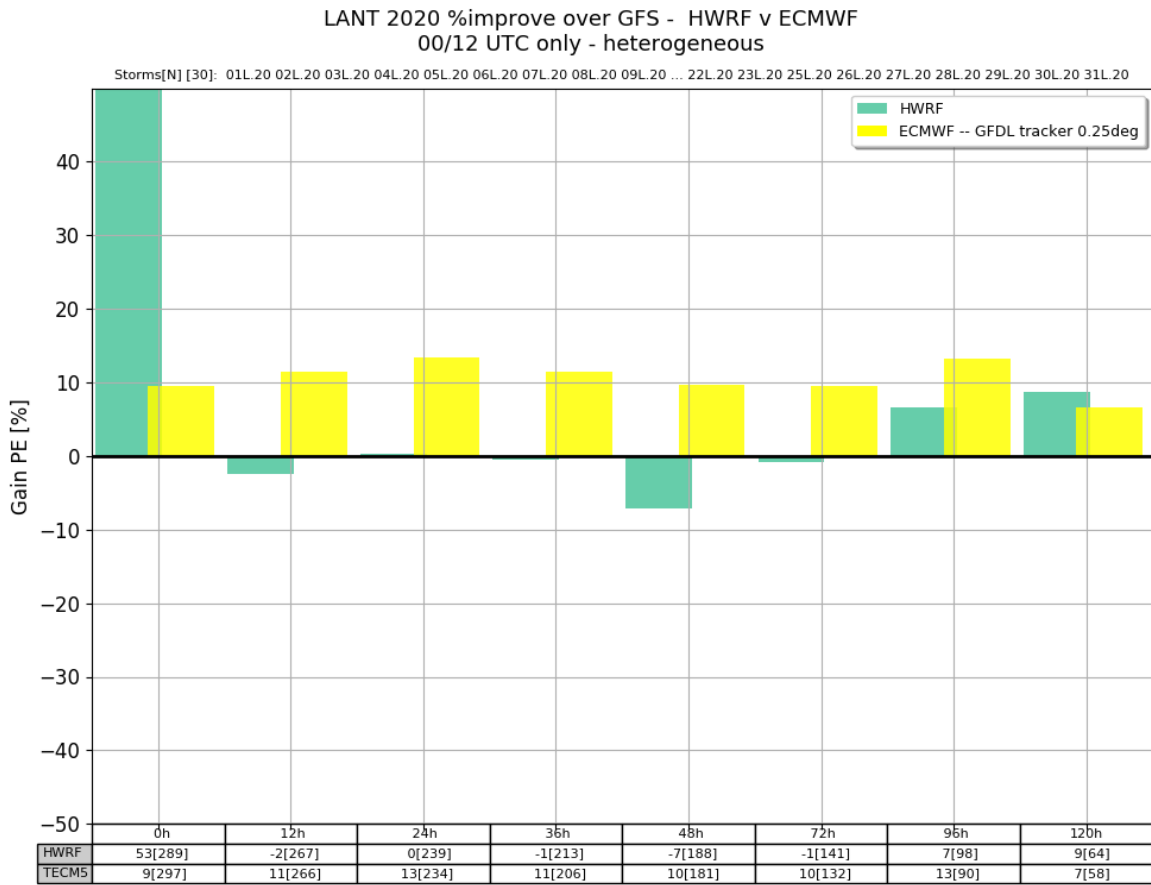
648 bars are the mean IE (bias, defined as model-observed). The thick bar is the mean and the

649 thinner bar extends from the 25th to 75th percentiles with the median indicated by the solid black

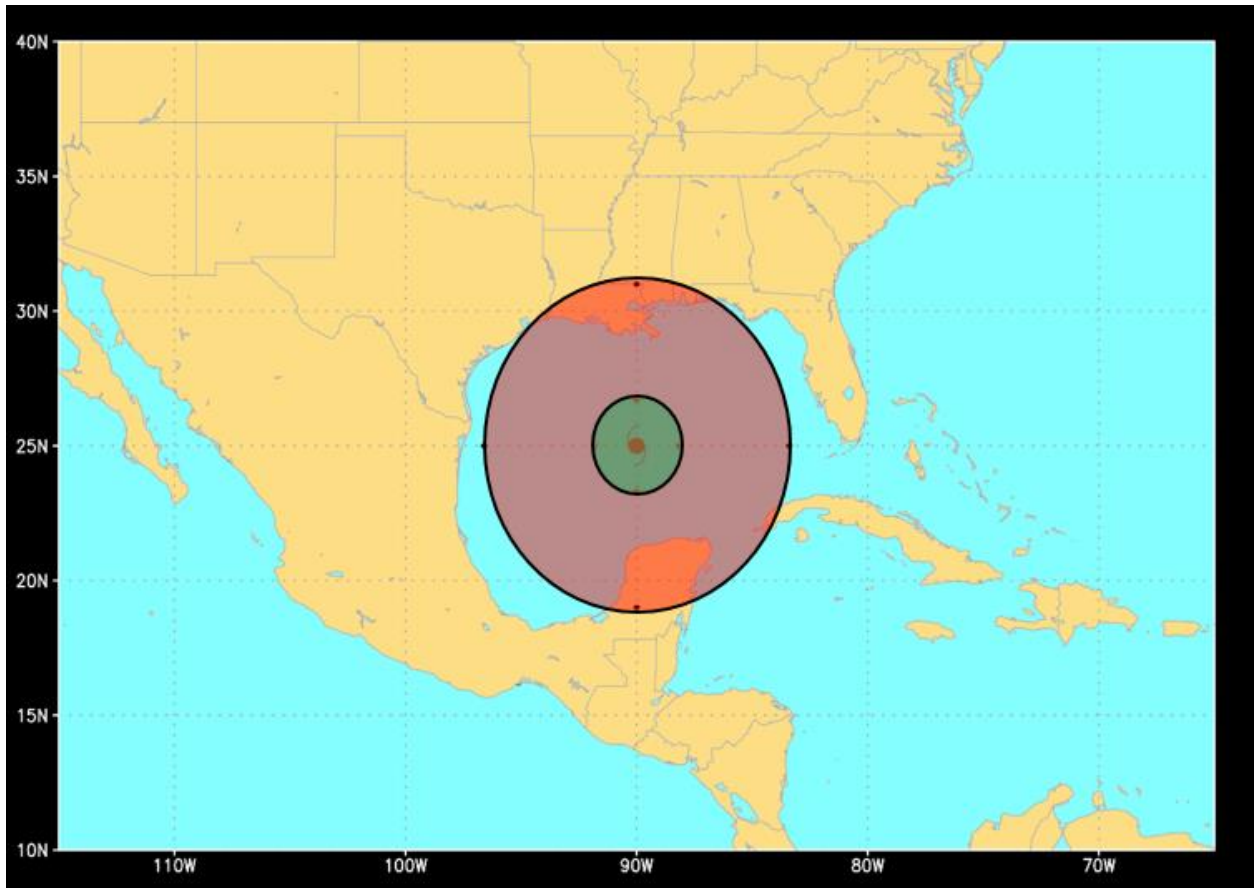
650 line (box and whisker). The mean value and [number of cases] is shown in the table below the

651 plot.

652 . The bias correction is the same as used in operations. The line is the absolute mean and the
 653 bars are the mean or intensity bias.



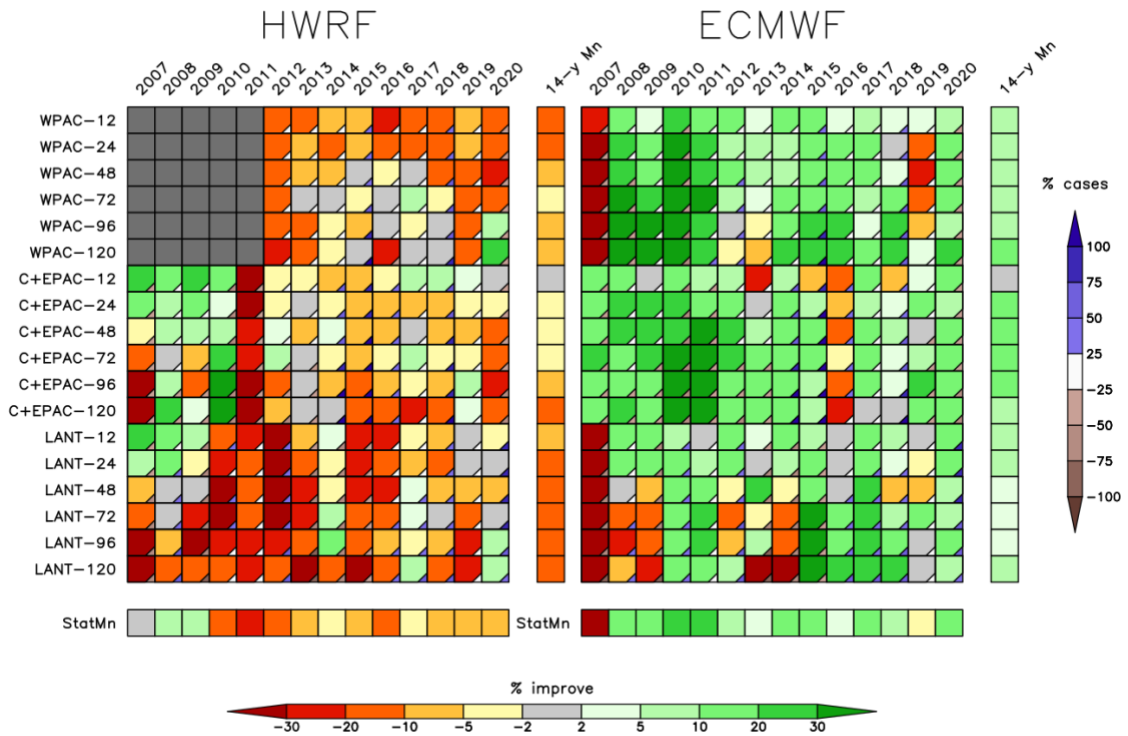
654
 655 Figure 7. Percent improvement in PE (lower PE is a positive improvement see Eq. 1) of HWRF
 656 v ECMWF using the GFS as the baseline for the 2020 LANT season (same data as in Fig. 4).
 657 This plot is similar to NHC ‘skill’ where climatology and persistence is the baseline (Cangialosi
 658 2020 Figure 5).



659
660
661

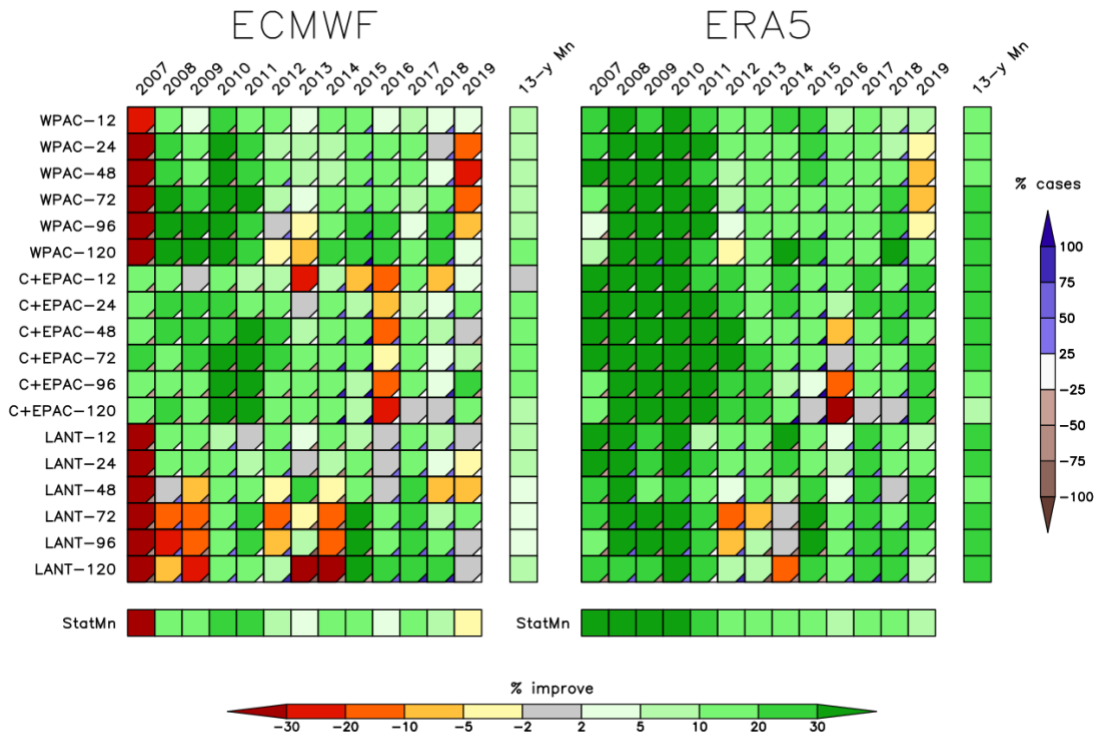
Figure 8. Area of a 360 n mi mPE (red) vs area of 100 n mi mPE (green).

HWRF/ECMWF Mean Position Error %improve over GFS [%]



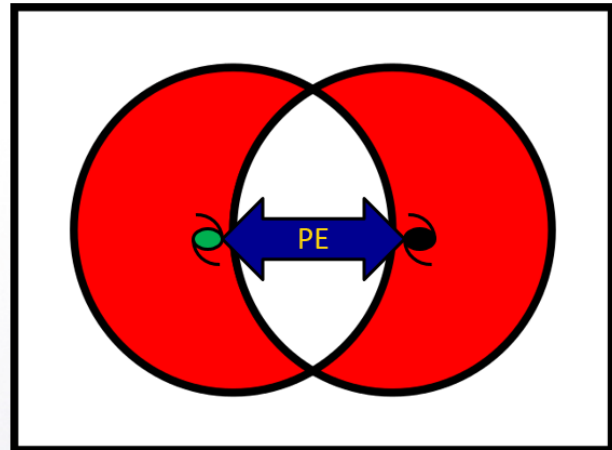
663 Figure 9. $\%IMP_{PE}$ for the NHEM basins 2007-2020. The $\%IMP_{PE}$ is colored so that green is
 664 ‘good’ (lower mPE) than the GFS baseline and red is ‘worse’ (higher mPE). The % cases
 665 contributing to the mean is shown in the triangle in the lower right-hand corner of each box.
 666 Panel A compares HWRF to ECMWF operational deterministic run. Panel B compares
 667 ECMWF (HRES) to forecasts from the latest ECMWF reanalysis (ERA5).
 668

ECMWF/ERA5 Mean Position Error %improve over GFS [%]



672

- $FE \equiv \text{IKE}_{fc} \Delta \text{IKE}_{bt}$ assuming:
 - $R_{34}^{fc} = R_{34}^{bt}$
 - $r_{max}^{fc} = r_{max}^{bt}$
- $FE = f(PE, IE)$; $IE = 0 @ 2R_{34}$
- FE ranges from 0 to a max at $PE = 2R_{34}$
for $PE > 2R_{34}$ linear increase



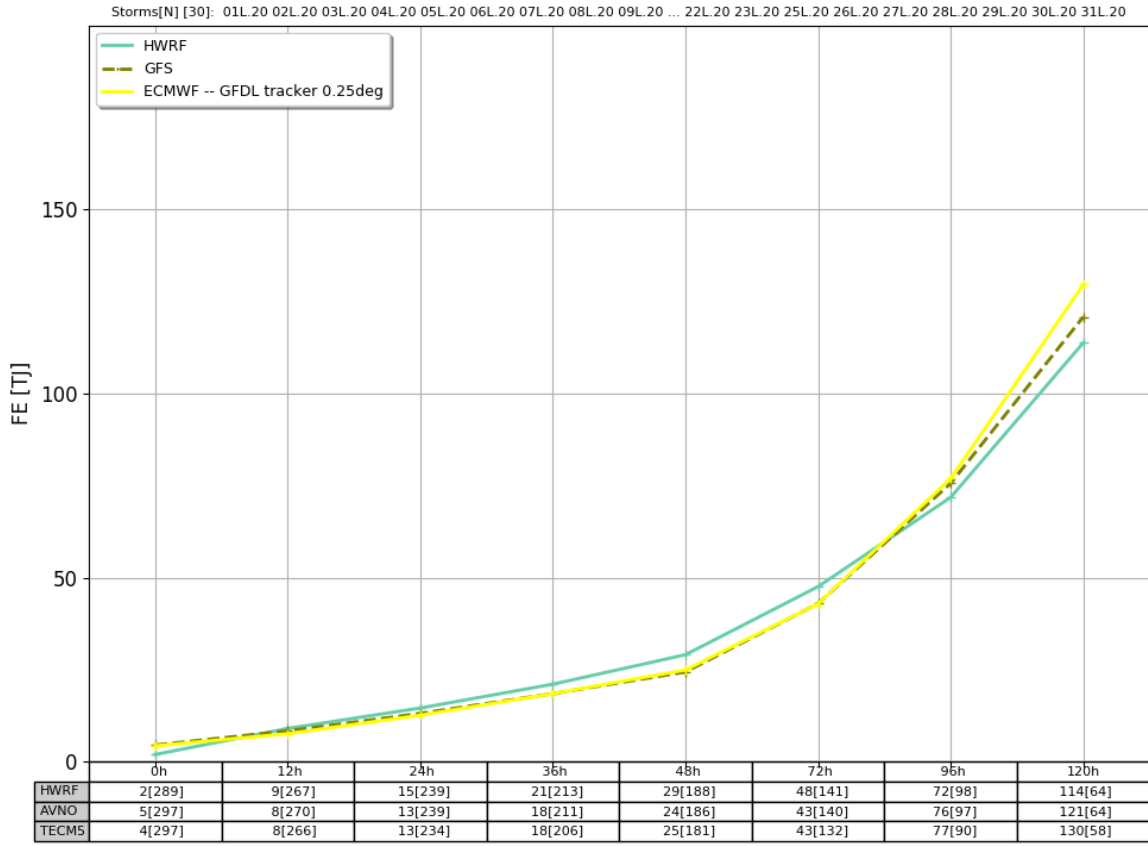
673

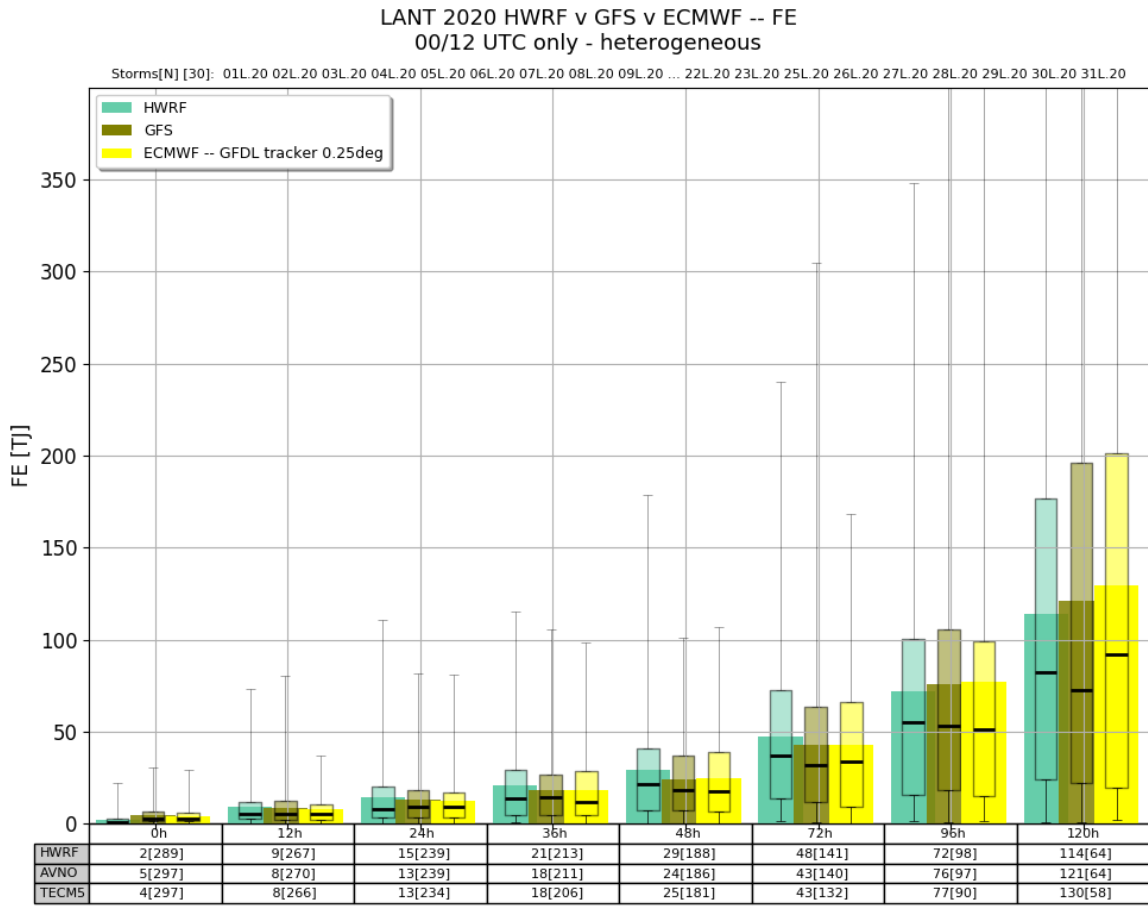
674 Figure 10. FE is the IKE in the 'symmetric difference' (red area) and how FE is related to PE.

675 Details of the calculation are also given. The *fc* superscript indicates a forecast and the *bt*

676 superscript the verifying best track.

LANT 2020 HWRf v GFS v ECMWF -- FE traditional line
00/12 UTC only - heterogeneous



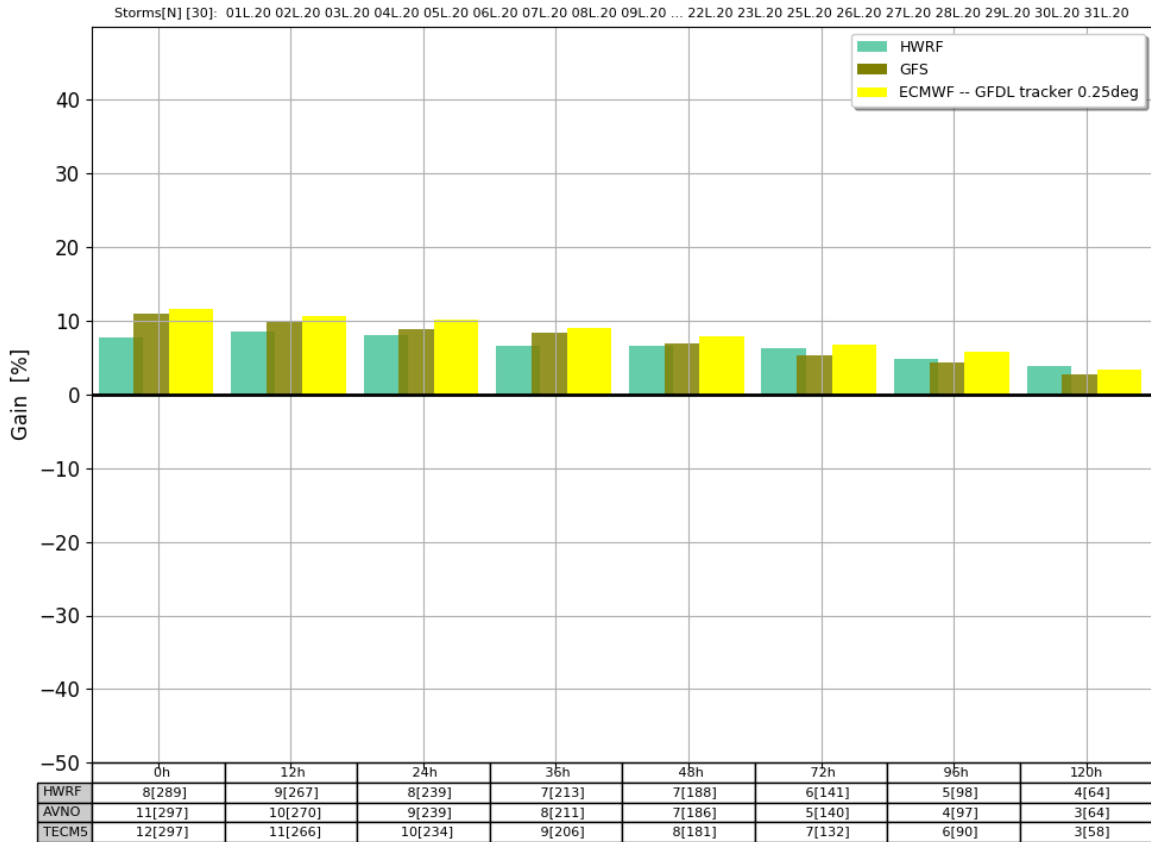


680

681 Figure 11. as in Fig. 4 except the metric is FE units: TJ. Panel A is the traditional line and B is

682 the box-whiskers version.

LANT 2020 HWRf v GFS v ECMWF %improve FE with IE=0
00/12 UTC only - heterogeneous



683

684 Figure 12. as in Fig. 4 except metric is percent improve in FE assuming a perfect intensity

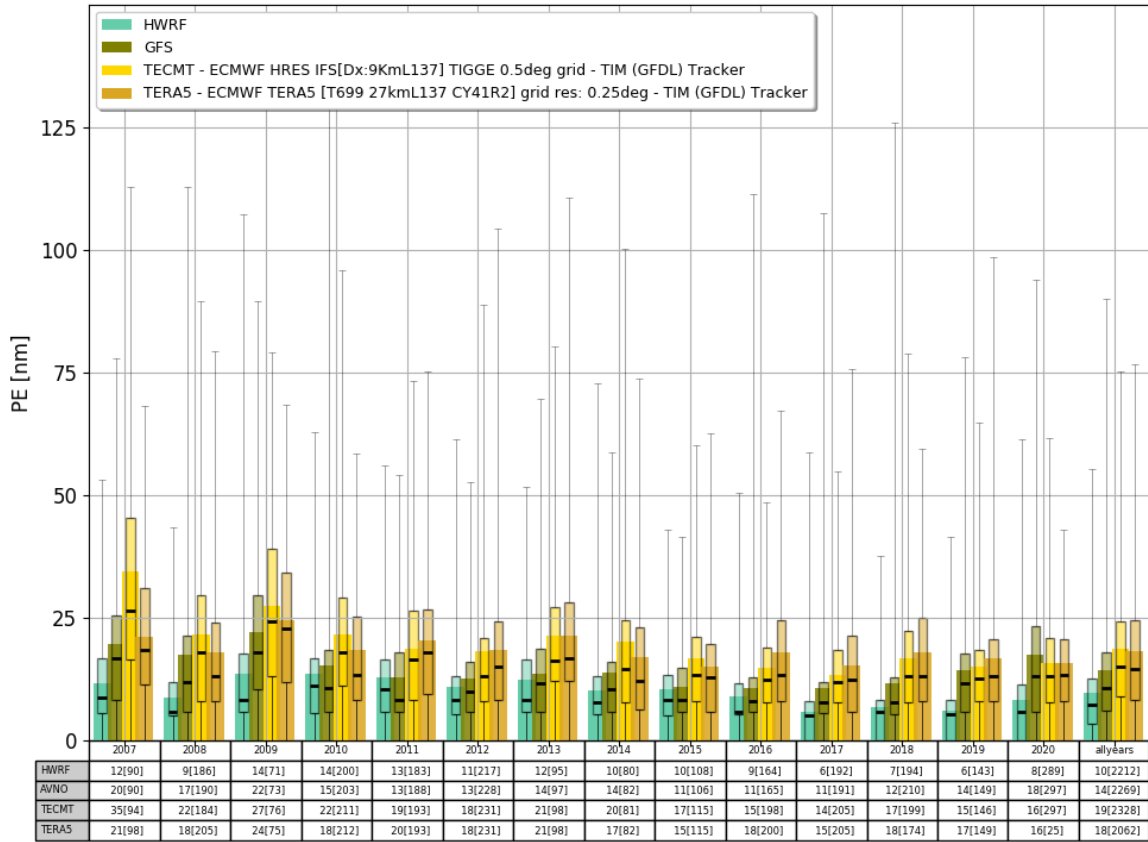
685 forecast or IE=0.

686

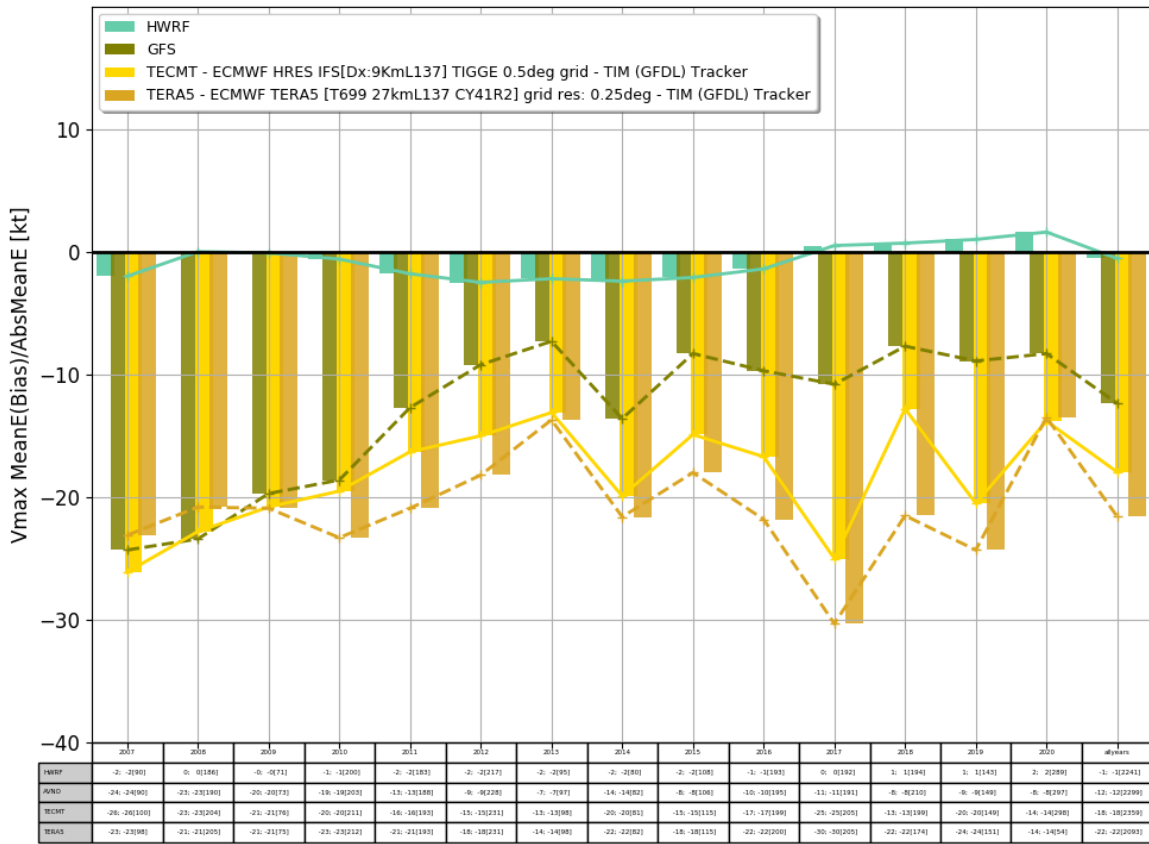
687

688

LANT 2007-2020 Initial PE - HWRF v GFS v ECMWF(HRES) v ERA5
00/12 UTC Heterogeneous



LANT 2007-2020 Initial IE - HWRf v GFS v ECMWF(HRES) v ERA5
00/12 UTC Heterogeneous

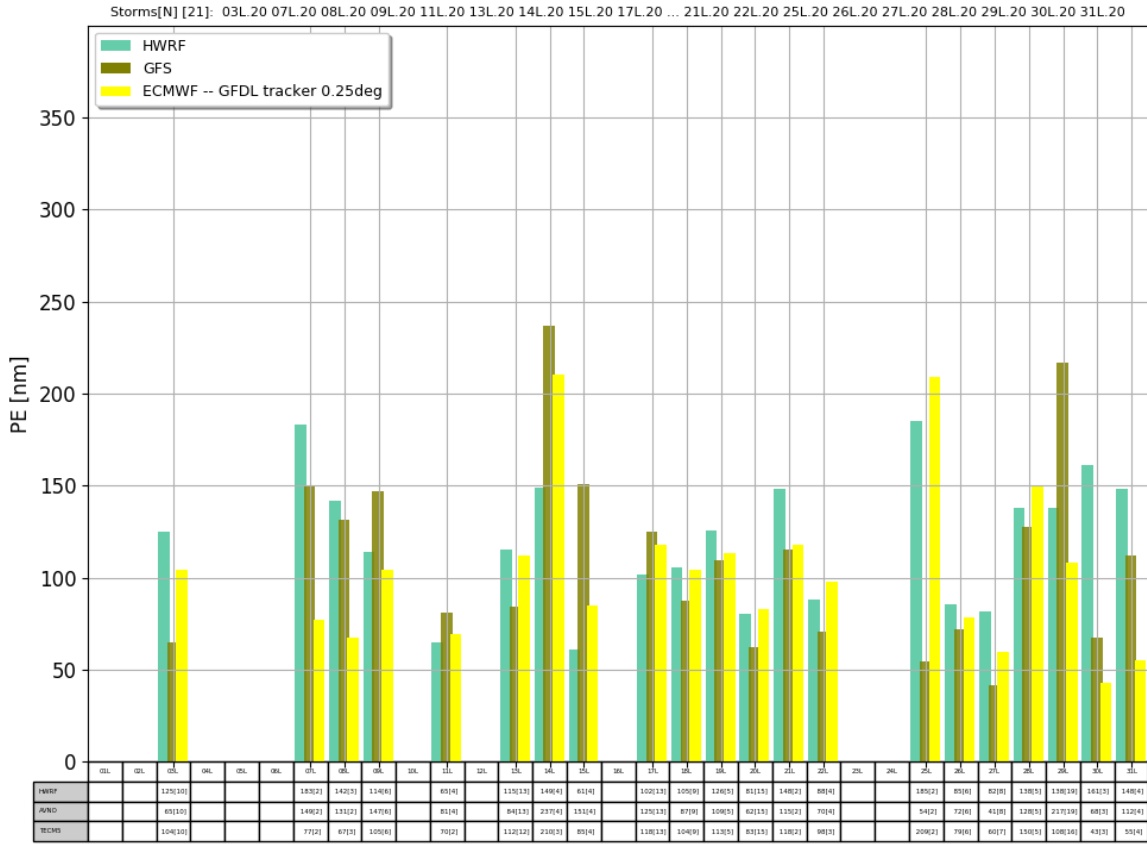


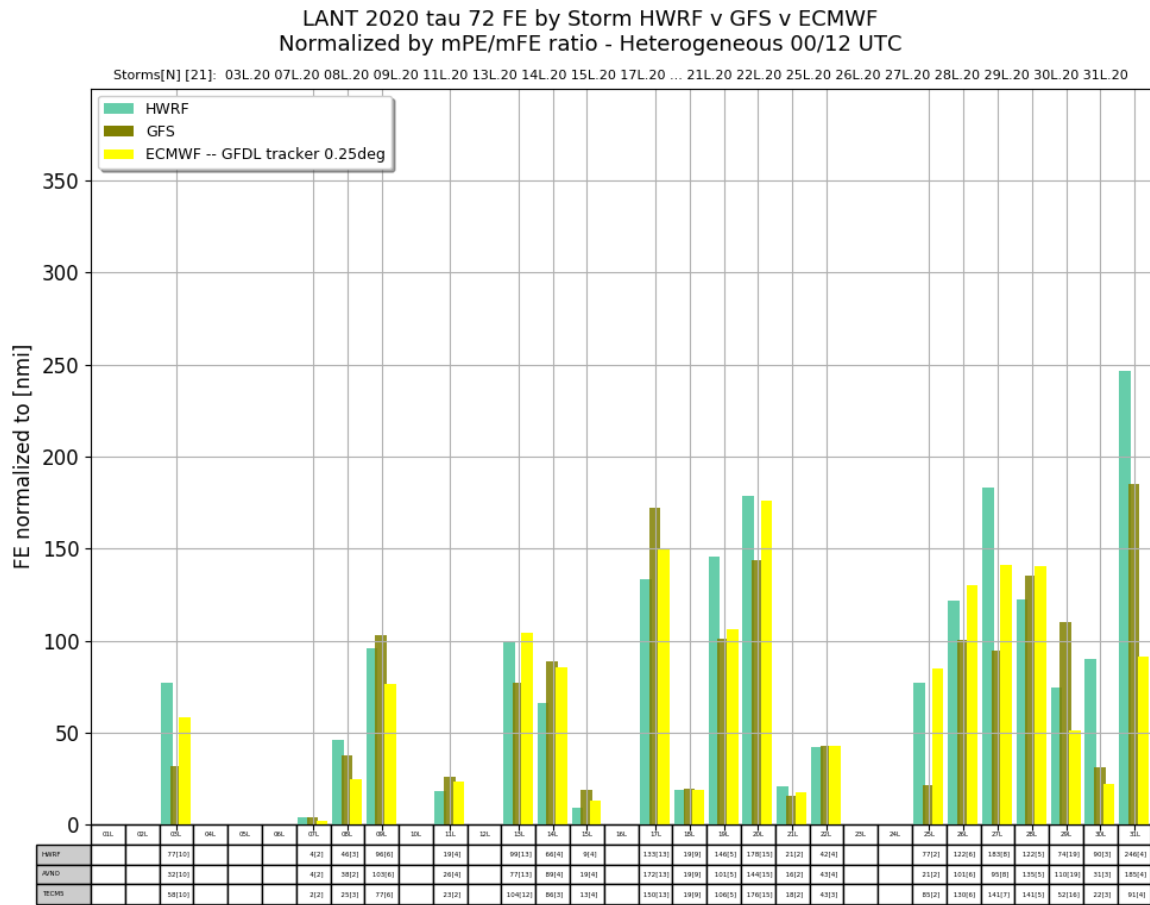
693

694 Figure 13. initial PE (A) and initial IE (B) for HWRf v GFS v ECMWF HRES v ERA5 in the
695 LANT 2007-2020.

696

LANT 2020 tau 72 PE by Storm HWR v GFS v ECMWF
Heterogeneous 00/12 UTC





702

703 Figure 14. 72-h mPE (A) and 72-h mFE (B) by storm normalized by mean ratio of mPE/mFE at

704 72 h to put FE in same units and scale as PE. The total number of storms verifying at 72 h is 21

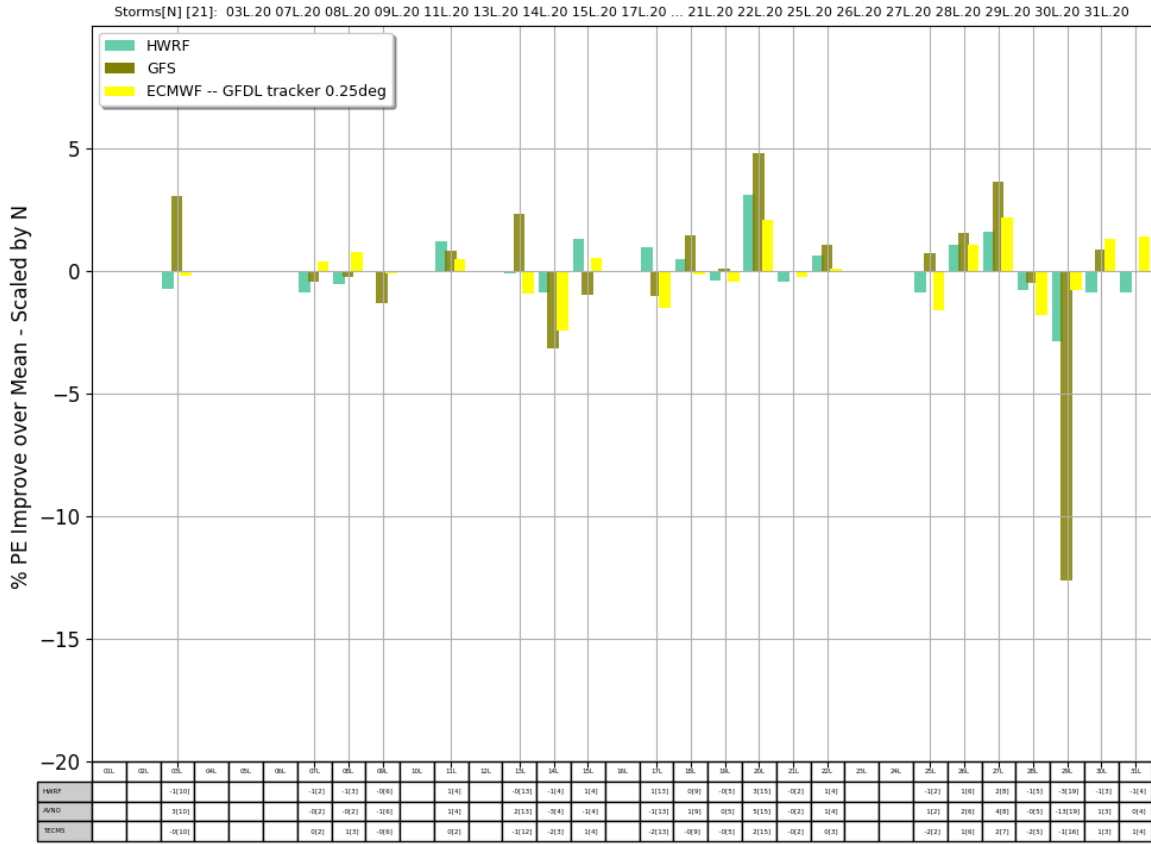
705 out of 31 for the season. The number of verifying 72-h cases is 68% which is higher than 10-y

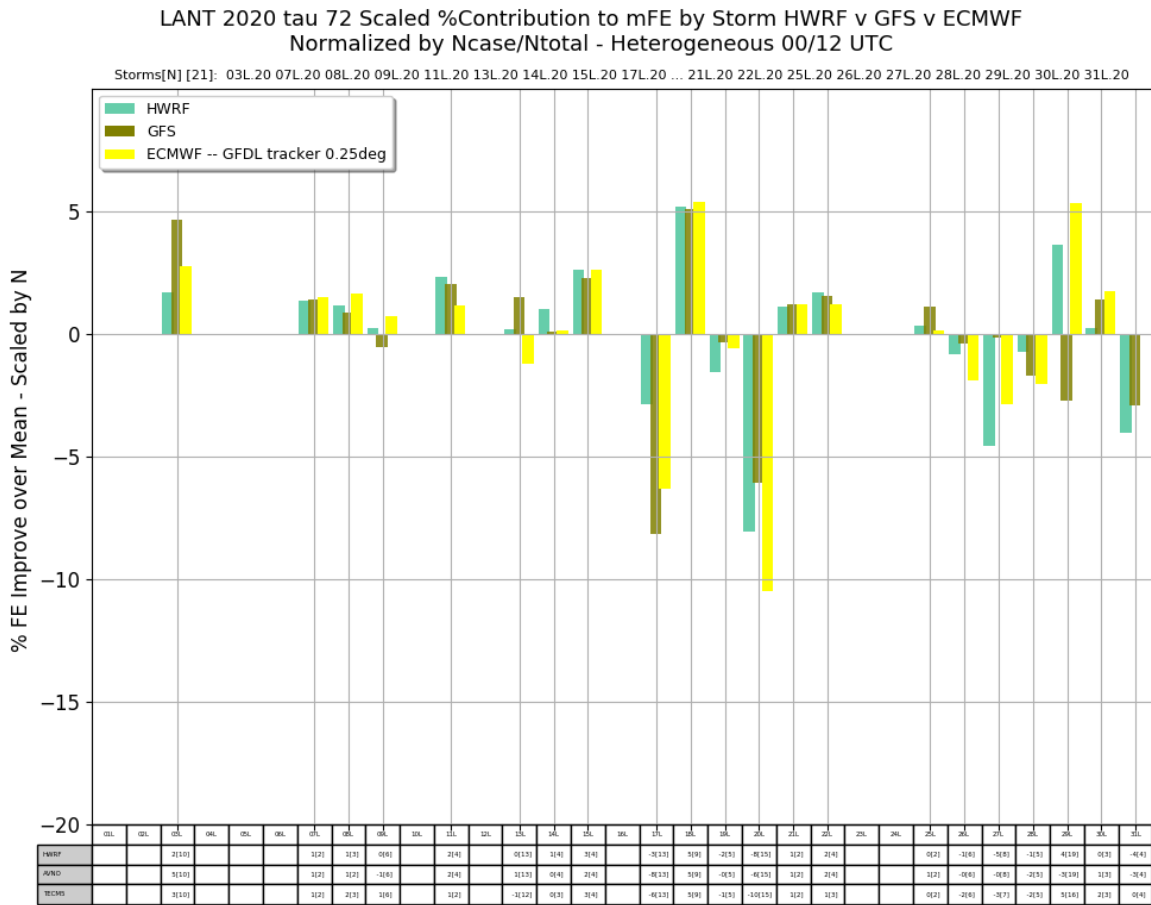
706 average of 51% as shown in Fig. 5.

707

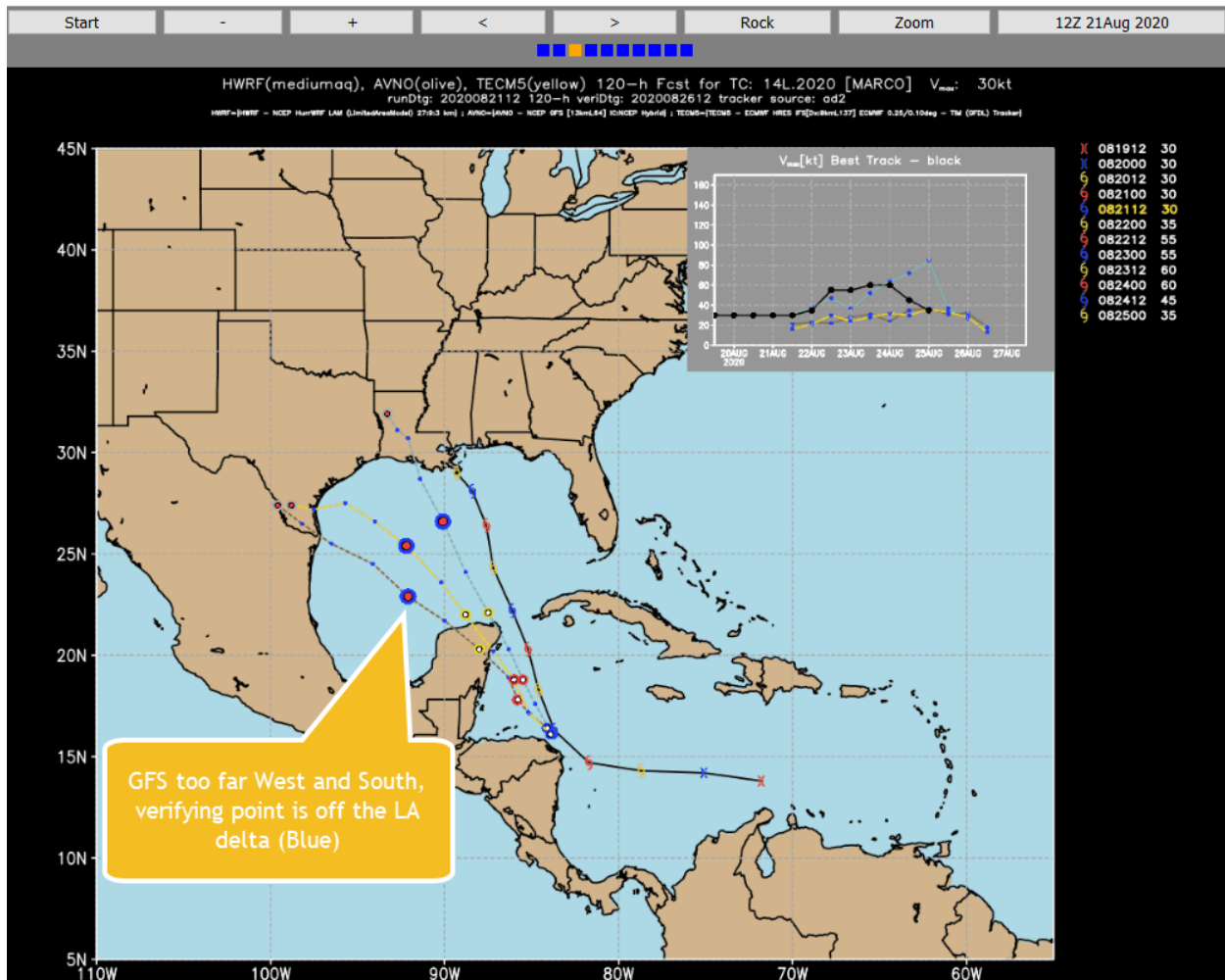
708

LANT 2020 tau 72 Scaled %Contribution to mPE by Storm HWRf v GFS v ECMWF
 Normalized by Ncase/Ntotal - Heterogeneous 00/12 UTC





713
 714 Figure 15. Same cases/storms as in Fig. 14 but the metric is $\%IMP_{P/FE}$ of the contribution by
 715 each storm to the 2020 season to the 72-h mPE (A) and mFE (B) normalized by ratio of the
 716 number of cases by storm and divided by the total number of cases. The mean across all storms
 717 for this scaled percent improvement is 0.



718

719 Figure 16. HWRf (light blue) / ECWWF (gold) / GFS (olive green) forecast for 14L (MARCO)
 720 on 2020082112 (12UTC 21 Aug 2020). The intensity forecast is show in the upper-right hand
 721 insert using the same color code as the tracks with the thick black line the working best track
 722 intensity. See: <http://ctrkveri.wxmap2.com/trk-14L-2020-MOD-CUR.htm> for all forecasts used
 723 in the verification.

724

725

726 Kieu, C. Q., and Z. Moon, 2016: Hurricane Intensity Predictability. *B Am Meteorol Soc*, **97**,
 727 1847-1857.

728 Powell, M. D., and T. A. Reinhold, 2007: Tropical cyclone destructive potential by integrated
 729 kinetic energy. *B Am Meteorol Soc*, **88**, 513-+.

730 Simmons, A. J., and A. Hollingsworth, 2002: Some aspects of the improvement in skill of
731 numerical weather prediction. *Q J Roy Meteor Soc*, **128**, 647-677.
732 Walsh, K. J. E., M. Fiorino, C. W. Landsea, and K. L. McInnes, 2007: Objectively determined
733 resolution-dependent threshold criteria for the detection of tropical cyclones in climate models
734 and reanalyses. *J Climate*, **20**, 2307-2314.
735