

1     **Probabilistic Prediction of Tropical Cyclone Intensity from a**  
2                     **Multiple-Linear Regression Model**

3     CHIA-YING LEE, \* MICHAEL K. TIPPETT,<sup>2</sup> SUZANA J. CAMARGO,<sup>3</sup>  
  ADAM H. SOBEL<sup>2,3</sup>

                  \**International Research Institute of Climate and Society, Columbia University, Palisades, NY*

4                   <sup>2</sup>*Department of Applied Physics and Applied Mathematics, Columbia University, New York, NY*

5                           <sup>3</sup>*Lamont-Doherty Earth Observatory, Columbia University, Palisades, NY*

---

\**Corresponding author address:* Chia-Ying Lee, International Research Institute for Climate and Society,  
Columbia University, 61 Route 9W, 202 Monell, Palisades, NY 10964.

E-mail: cle@iri.columbia.edu

## ABSTRACT

7 The authors describe the development and verification of a statistical model relating tropical  
8 cyclone intensity to the local large-scale environment. A multiple linear regression framework  
9 is used to estimate the expected intensity of a tropical cyclone given the environmental  
10 and storm conditions. The uncertainty of the estimate is constructed from the empirical  
11 distribution of model errors. NCEP-NCAR reanalysis fields and historical hurricane data  
12 from 1981 to 1999 are used for model development, and data from 2000 to 2012 are used to  
13 evaluate model performance. Seven predictors are selected: initial storm intensity, the change  
14 of storm intensity in the past 12 hours, the storm translation speed, the difference between  
15 initial storm intensity and its corresponding potential intensity, deep-layer (850-200 hPa)  
16 vertical shear, atmospheric stability, and 200 hPa divergence. The system developed here  
17 predicts storm intensity changes in response to changes in the surrounding environment with  
18 skill comparable to existing operational tools. The probabilistic intensity predictions are  
19 shown to be reliable and skillful. Since one application of such a model is to predict changes  
20 in TC activities in response to natural or anthropogenic climate change, we examine the  
21 performance of the model using data that is readily available from global climate models,  
22 i.e., limited variables and monthly averages. We find that statistical models based on monthly  
23 data (as opposed to daily) with only a few essential predictors, e.g., the difference between  
24 storm intensity and PI, perform nearly as well at short leads as when daily predictors are  
25 used.

# 1. Introduction

The intensities of tropical cyclones (TCs) depend on their surrounding environments. The question of what distribution of storm intensities is consistent with a given environment is important both for short-range intensity forecasts and for long-term climate projections. As the global climate changes, environmental factors that influence TC intensity, such as large-scale circulation and atmospheric moisture, are also expected to change. The issue of how these changes could influence the climatology of TC occurrence and intensity has been the subject of a number of recent studies, either from observational (Emanuel 2005; Webster et al. 2005) or from global climate modeling (Zhao et al. 2009; Bender et al. 2010; Knutson et al. 2010; Camargo 2013; Chen and Lin 2013; Emanuel 2013) approaches. However, it remains difficult to confidently assess future changes in storm intensity. One of the primary reasons is that present climate models have horizontal resolutions that are too coarse to resolve the inner core convective structure in storms. This deficiency prevents these climate models from simulating major storms at their observed intensities and providing a reasonable distribution of storm intensity even in the current climate.

Downscaling is a strategy for connecting global scale predictions and quantities not resolved by the global model, in our case, TC intensity. Dynamical downscaling has been shown to be a very useful tool for understanding climate projections of TC activities (e.g., Emanuel 2006; Knutson et al. 2007, 2008; Bender et al. 2010, etc), but it is computationally expensive and still limited in the range of intensities that can be simulated. Statistical downscaling is much cheaper computationally. This method has been used for seasonal prediction of North Atlantic hurricane activity by Vecchi et al. (2011), who developed a basin-wide hurricane frequency forecast system based on a Poisson regression with two predictors: Atlantic main development region sea surface temperature (SST) and global tropical-mean SST from a suit of high-resolution global models. Their results from retrospective forecasts indicated that their statistical model can provide a skillful seasonal prediction.

As there is always an inherent uncertainty in TC predictions, in the absence of a per-

53 fect model (or perfect mathematical equations), a probabilistic forecast representing this  
54 uncertainty should be provided along with the deterministic forecast. There are two general  
55 approaches for constructing probabilistic forecasts: ensemble and statistical approaches. In  
56 the ensemble approach, a set of forecasts are conducted from either parallel multiple dynam-  
57 ical or statistical models (Krishnamurti et al. 1999; Puri et al. 2001; Weber 2005), or a single  
58 numerical model under various initial conditions (Zhang and Krishnamurti 1999), physical  
59 parametrization schemes, or stochastic perturbations (Chen et al. 2014). In the statistical  
60 approach, the uncertainties are based on the probability density function (PDF) of the past  
61 errors or on forecasts from a deterministic numerical or statistical model (DeMaria et al.  
62 2009). Vecchi et al. (2011) used both ensemble and statistical approaches to estimate the  
63 uncertainties.

64 From the perspective of short-range TC intensity forecasting, studies from DeMaria and  
65 Kaplan (1994, 1999); DeMaria et al. (2005, 2007); DeMaria (2009) show that statistical  
66 models (e.g., Statistical Hurricane Intensity Prediction Scheme, SHIPS; Logistic Growth  
67 Equation Model, LEGM) with the environmental parameters from global models, in addition  
68 to climatological and persistent predictors, are capable of predicting storm intensity with  
69 some skill. We can think of this as essentially a form of statistical downscaling applied to  
70 individual storms in the present climate. These studies' results suggest that the statistical  
71 downscaling method may provide us with an alternative option for climate projections of  
72 TC intensity distribution: utilizing a statistical model to relate the global climate model  
73 (GCM) environment to TC intensity. In other words, we could use a system that responds  
74 correctly to the changes in the surrounding (local) environment in the current climate, and  
75 apply the current relationship between TC intensity and environment to the future when  
76 environmental conditions are projected by climate models to change.

77 For short-range probabilistic predictions, the National Hurricane Center (NHC) started  
78 issuing wind speed probabilistic forecasts in 2006 based on a statistical approach (DeMaria  
79 et al. 2009, 2013). They use a Monte Carlo method to generate 1000 forecast tracks based on

80 random samples of track errors from the past 5 years of data. A similar method is applied  
81 to account for the uncertainties in storm intensity and structure. Emanuel et al. (2006)  
82 further used synthetic TC tracks generated from an algorithm with both deterministic and  
83 random components together with an idealized axisymmetric dynamical hurricane intensity  
84 model (run along the predicted track many times) to generate a set of intensity probabilistic  
85 forecasts for a specific geographic location. This approach has now been used to study  
86 tropical cyclone risk for a number of locations and from a number of perspectives (e.g. Lin  
87 et al. 2010).

88 Our ultimate goal is to understand the tropical cyclone intensity distribution in the  
89 projected climate scenarios from a statistical perspective. The present study represents an  
90 initial step of our approach: the development of a probabilistic TC intensity prediction  
91 model. Following the successful example of DeMaria et al.'s work on SHIPS, we develop  
92 a SHIPS-like multiple linear regression model from global reanalysis fields derived from  
93 observations of the current climate. However, different from SHIPS, we would like to select  
94 the smallest number of predictors possible, in order to minimize the required GCM fields.  
95 Because our goal is to make projections of TC intensity under future climate change scenarios  
96 (although such projections are not yet attempted in this study), we avoid predictors which  
97 we expect will make the statistical model overly tuned to the current climate. Sea surface  
98 temperature (SST) is the most obvious example; the absolute value of SST may be a useful  
99 predictor for individual storm intensity in the present climate, but we expect the relationship  
100 of SST to TC intensity to change in future climates. We also determine whether monthly  
101 data can be used instead of daily, as this greatly reduces the data requirement for climate  
102 studies.

103 Detailed description of data and the development of the deterministic multiple linear  
104 regression forecast model (called 'MLR' hereafter) will be illustrated in Section 2 and 3,  
105 respectively. Examples of forecasts from the MLR, and the verification of it will be shown  
106 in Section 4. The dependence of the MLR performance (errors) on predictors will also be

107 discussed. In Section 5, the possibility of using monthly model output is tested through  
108 examining the sensitivity of MLR performance to global model fields. In Section 6, we will  
109 demonstrate a skillful and reliable probabilistic forecast from the MLR which is conducted  
110 with the past error distribution. We summarize our findings in Section 7.

## 111 2. Data

112 The best-track dataset produced by NHC for North Atlantic TCs provides historical  
113 storm information, which includes storm location, maximum wind speed, and minimum  
114 sea level pressure (Jarvinen et al. 1984; Landsea and Franklin 2013). For synoptic condi-  
115 tions, we use the  $2.5^\circ \times 2.5^\circ$  daily and monthly National Centers for Environmental Predic-  
116 tion –National Center for Atmospheric Research (NCEP-NCAR) reanalysis dataset (Kalnay  
117 et al. 1996). In addition to atmospheric fields, we obtain upper ocean structure from the  
118  $\sim 1.8^\circ \times 1.8^\circ$  the NOAA NCEP Environmental Model Center (EMC) Climate Modeling  
119 Branch (CMB) global ocean data assimilation system (GODAS; Behringer and Xue 2004).  
120 We utilize the data from 1981 to 1999 for model training - both predictor selection and coef-  
121 ficient estimation. The data for the period 2000-2012 is used for testing model performance.  
122 We only consider storms that reach at least tropical storm strength ( $> 34$ kt) during their  
123 lifetime. As mentioned in DeMaria and Kaplan (1994), the statistical properties of storms  
124 over land are different from those over the ocean, so we only use the data for the period when  
125 storms are not over land. In the end, we have more than 4,000 cases of 12 hour forecasts in  
126 both the training and testing datasets. The number of forecasts decreases with increasing  
127 lead time. Here we consider 12-hourly lead times from 12 to 120 hours. There are only about  
128 2,000 cases left for both training and testing for the 120 hour lead forecast. Additionally, the  
129 official NHC and SHIPS TC forecasts from Automated Tropical Cyclone Forecasting System  
130 (ATCF; Sampson and Schrader 2000) are used for evaluating the MLR performance.

### 131 **3. Multiple-linear regression model development**

132 We begin by describing the initial pool of potential predictors for the MLR. We then  
133 choose a subset of those predictors for inclusion in the multiple linear regression model  
134 (MLR) using a forward selection (stepwise regression) procedure. The next step is to use  
135 a fitting procedure on the training data to derive the MLR, which consists of 10 linear  
136 equations for forecasts with 12-hourly lead times from 12 to 120 hours.

#### 137 *a. Initial pool of predictors*

138 Four variables are chosen that represent the storm itself: initial maximum wind speed  
139 ( $V_0$ ), minimum sea level pressure (MSLP), change of storm maximum wind speed in the  
140 previous 12 hours ( $dV/dt$ ), and the storm translation speed (trSpeed). The environmental  
141 (synoptic) conditions are represented by the vertical wind shear, moisture, environmental sta-  
142 bility, outflow temperature, upper tropospheric divergence, sea surface temperature (SST),  
143 and upper ocean thermal structures.

144 The most important synoptic predictor in the SHIPS includes the maximum potential  
145 intensity (PI), which is derived empirically (DeMaria and Kaplan 1994). Here, we use the  
146 Emanuel’s PI definition developed in Emanuel (1988) and modified in Emanuel (1995), and  
147 Bister and Emanuel (2002). The PI is a function of the ratio of the exchange coefficient for  
148 enthalpy to that for momentum, the ratio of outflow temperature to SST, and the vertical  
149 integral of air temperature along iso-entropy surface at the radius of maximum wind speed  
150 (Bister and Emanuel 2002). In Camargo et al. (2007) and Camargo et al. (2009), the  
151 importance of PI for predicting tropical cyclogenesis was shown using the PI computed  
152 both from daily and from monthly NCEP-NCAR reanalysis data. We use here the PI from  
153 Camargo’s dataset extended to 2012. The PI is first averaged over a disk extending 500 km  
154 from the storm center and then averaged over the forecast interval (e.g., 12, 24, ..., 120 hours)  
155 along the storm track. The PI enters the MLR as the difference of the mean PI and initial

156 intensity  $V_0$ , denoted  $dPI_{V_0}$ .

157 We do not use SST as a predictor, nor any function of SST which includes specific  
158 threshold values based on observations of the historical relationship between SST and TC  
159 intensity. The relationship of SST to TC intensity is expected to be change with the mean  
160 climate state (e.g., Vecchi and Soden 2007; Ramsay and Sobel 2011; Johnson and Xie 2010;  
161 Emanuel and Sobel 2013). The relative SST, or difference between the local SST and an  
162 appropriately defined tropical mean, may be a better predictor than absolute SST, but even  
163 that is limited in that it cannot (by definition) capture intensity changes that are related to  
164 global mean changes(e.g., Emanuel and Sobel 2013). PI is a theoretically derived quantity  
165 whose formulation does not include any explicit tuning to the present climate, and thus  
166 in principle should be able to capture intensity changes which result from global climate  
167 changes. Emanuel’s PI theory has been critiqued on theoretical grounds (e.g., Smith et al.  
168 2008) but appears at this point to be as good as any extant theory (e.g., Bryan and Rotunno  
169 2009). Our use of PI to the exclusion of any SST parameter thus appears to be a reasonable  
170 choice for our purpose. If advances in the theoretical understanding of TC intensity and its  
171 relationship to the environment are made in the future, our model can be easily re-derived  
172 to use any improved PI parameter that may be developed.

173 Strong vertical wind shear tends to have a negative effect on TC intensification (DeMaria  
174 2006; Tang and Emanuel 2012), and is another key predictor in the SHIPS model (DeMaria  
175 et al. 2005). In this study, the deep-layer mean vertical wind shear (SHRD) is defined  
176 as the magnitude of the difference between the mean wind vectors at 200 and 850 hPa  
177 averaged over an annulus extending 200 km to 800 km around the storm center (Chen  
178 et al. 2006). Moisture is an another factor known to affect TC intensity. Kaplan and  
179 DeMaria (2003) found that rapidly intensifying TCs are often embedded in a high moisture  
180 environment. In the SHIPS model, four variables representing relative humidity at various  
181 levels (1000, 850-700, 700-500, and 500-300 hPa) have been used in the past. Until 2013, only

182 the 700-500 hPa mean relative humidity was used<sup>1</sup>. For the MLR, we initially consider two  
183 moisture variables: the mid level relative humidity (rhMid, 500-300 hPa) and the column-  
184 integrated relative humidity (rhCol; Bretherton et al. 2004). The latter accounts for the  
185 accumulated moisture from the top of the boundary layer to 200 hPa. For the environmental  
186 stability, we use the conditional instability ( $d\theta_{es}/dz$ ), defined as the vertical gradient of  
187 saturated equivalent potential temperature (dThetaEs) averaged from boundary layer to  
188 300 hPa. SHIPS considers convective instability (using equivalent potential temperature,  $\theta_e$ )  
189 instead of conditional instability. However, convective instability is usually related to the  
190 lifting of an entire layer, not just a parcel. The upper tropospheric divergence is usually  
191 used to estimate the large-scale forcing, and we use the divergence field at 200 hPa (div200)  
192 for this purpose. While there has been some suggestion that convective available potential  
193 energy (CAPE) is not a good predictor for deep tropical convection (Zipser 2003), CAPE is  
194 included in our initial predictor pool. The last atmospheric environmental predictor is the  
195 temperature at 200 hPa (T200), representing the outflow temperature.

196 Similar to the environmental predictors calculated in the SHIPS, rhMid, rhCol, dThetaEs,  
197 CAPE, and T200 are first averaged over an annulus centered at the storm location at forecast  
198 time with inner radius of 200 km and outer radius of 800 km. div200 is averaged over a circle  
199 with radius of 1000 km. To account for the time variation of these variables, they are also  
200 averaged over the forecast interval and along the storm track, as described for dPI-V<sub>0</sub>. As  
201 we are using the best-track data here, we do not apply a vortex removal scheme as used  
202 in SHIPS (DeMaria 2010) to account for differences in the instantaneous Global Forecast  
203 System (GFS) forecast fields and the NHC forecast of the storm center.

204 In 2004, SHIPS started incorporating predictors associated with the upper ocean (ocean  
205 heat content, ocean depth of 20 and 26 °C) from satellite altimetry data to account for the  
206 oceanic negative impact on the storm (DeMaria et al. 2005). Here, instead of upper ocean

---

<sup>1</sup>[http://rammb.cira.colostate.edu/research/tropical\\_cyclones/ships/docs/SHIPS\\_](http://rammb.cira.colostate.edu/research/tropical_cyclones/ships/docs/SHIPS_predictor_file_2013.doc)  
predictor\_file\_2013.doc

207 heat content, we use another air-sea interaction quantity, the ocean temperature averaged  
208 over the top 100 m (OT100; Price 2009), which represents the lowest SST that can be induced  
209 by a storm. As daily oceanic data are not readily available, monthly data from GODAS  
210 (Behringer and Xue 2004) is used here.

211 Our initial pool of predictors contains 13 variables: initial maximum wind speed ( $V_0$ ),  
212 initial minimum sea level pressure (MSLP), changes of storm intensity in the past 12 hours  
213 ( $dV/dt$ ), storm translation speed (trSpeed), difference between initial intensity and final PI:  
214  $dPI-V_0$ , vertical shear (SHRD), mid level and column-integrated relative humidity (rhMid  
215 and rhCol), conditional stability ( $d\Theta_{Es}$ ), CAPE, 200 mb divergence (div200) and tem-  
216 perature (T200), upper 100-m mean ocean temperature (OT100). The change of the max-  
217 imum wind speed over different forecast interval (e.g. 12, 24, 48, 72, 96, and 120 hours) is  
218 the predictand. The assumption of a linear relationship between intensity change and each  
219 predictor was tested as in Tippett et al. (2011, see their Fig. 12) and found (not shown) to  
220 be adequate.

#### 221 *b. Forward selection procedure*

222 In the forward selection procedure, the initial model contains no predictors other than the  
223 constant term. The reduction obtained in root-mean-square-error (RMSE) from adding a  
224 single variable is computed for each variable, and the variable that most reduces the RMSE is  
225 added. This procedure is repeated until all the variables are added. The RMSE is calculated  
226 with 10-fold cross validation; the training data are divided into 10 random subsets, 9 are  
227 used to estimate regression coefficients and the tenth is used to compute an estimate of the  
228 RMSE. Here, the random partitioning of the data is done 10 times, providing 100 estimates  
229 of the RMSE, similarly to Tippett et al. (2012). The mean RMSE as a function of the  
230 number of predictors for all lead times are shown in Fig. 1. The predictors are listed on the  
231 right of each panel in Fig. 1 in the order in which they enter the forward selection procedure.  
232 A reasonable rule for the number of predictors to include in the MLR is to stop including

233 predictors at the point where adding additional predictors does not significantly decrease  
234 the mean RMSE.

235 To estimate this stopping point robustly, we apply the forward selection procedure on  
236 the training data 10 times. While the order of the first few predictors is robust, the order  
237 of the remaining predictors is not and varies with the particular partition of the data. For  
238 example, for 12 hour forecast, the order of the 3rd - 13th predictors varies randomly each  
239 time we run the forward selection procedure, indicating little difference in the utility of the  
240 remaining predictors. Hence, we add the second stopping rule at the point (red lines in  
241 Fig. 1) when the order of selected predictors is stable.

242 For 12 hour forecasts (Fig. 1a), a substantial decrease in the RMSE occurs when the  
243 number of the predictors increases from 1 ( $dVdt$ ) to 2 ( $dVdt$  and  $dPI-V_0$ ). Additional  
244 predictors result in little reduction of the RMSE, indicating that the intensification trend  
245 and the difference between PI and the initial intensity are the most important predictors  
246 for short-lead forecasts. The failure of other environmental predictors to reduce RMSE may  
247 indicate that the predictable impact of these environmental parameters, such as shear, are  
248 implicitly included in the previous 12 hours intensity changes ( $dVdt$ ). For a 24 hour lead  
249 time,  $dPI-V_0$  is the most important predictor, followed by  $dVdt$  and SHRD. Therefore, for  
250 forecast times shorter than 1 day, three predictors,  $dVdt$ ,  $dPI-V_0$ , and SHRD, are adequate,  
251 and other predictors do not substantially reduce the RMSE.

252 The number of predictors required in the MLR (before RMSE stops decreasing) increases  
253 with the forecasting lead time. For a forecast with 72 hours interval, we need 6 predictors  
254 before the RMSE stops decreasing sharply, and the first 7 predictors selected have a fixed  
255 order. These predictors are  $dPI-V_0$ ,  $trSpeed$ , SHRD,  $V_0$ ,  $dThetaEs$ ,  $dVdt$ , and  $div200$ . For  
256 96 and 120 hour forecasts, however, the sharp decrease of RMSE happens at 4, and only  
257 5 predictors are selected in a fixed order each time we conduct the selection procedure. A  
258 possible reason could be that we need predictors other than those we have in order to further  
259 improve the 96 and 120 hours intensity forecasts. Another possible reason is that there is

260 a natural limitation of the linearity assumption for long-forecast intervals, i.e., we cannot  
261 continue to improve the forecast from a multiple-linear regression model with additional  
262 predictors, and a more sophisticated model is needed.

263 Overall we see that no more than 7 predictors ( $dPI_{V_0}$ ,  $dVdt$ ,  $trSpeed$ ,  $V_0$ ,  $SHRD$ ,  $div200$ ,  
264  $dThetaEs$ ) are needed to achieve most of the possible reduction in RMSE. The absence of  
265 any moisture parameters is noteworthy, and implies that the moisture variables do not add  
266 independent information. However, this result does not mean that moisture has no role. For  
267 instance, the impact of dry air is greatest when it reaches the core, which requires shear-dry  
268 air interaction. To capture such effect would require some other combination of shear and  
269 moisture parameters.  $OT100$ , the mean upper ocean temperature, is also not chosen. The  
270 analogous parameter in SHIPS, OHC, yields 4-8% improvement in the forecast error only  
271 when a threshold ( $50 \text{ KJcm}^{-2}$ ) is applied (DeMaria et al. 2005). It is possible that we need to  
272 use some kind of threshold when including  $OT100$ . Nevertheless, oceanic negative feedback  
273 also depends on the storm moving speed,  $trSpeed$ , which is selected here. Although  $T200$   
274 does not pass the selection test, the upper tropospheric temperature is included in the PI  
275 calculation. CAPE is eliminated as well, which is not too surprising since it has been noted  
276 that it might not be a useful predictor for tropical convection (Zipser 2003).

277 The signs of the coefficients (not shown) are consistent with the physical relationships  
278 between predictors and intensity change, though such a result is not guaranteed for corre-  
279 lated predictors. Storms are more likely to follow their historical development and therefore  
280 an intensifying storm would continue intensifying, yielding a positive coefficient of  $dVdt$ .  
281 Increasing storm speed reduces the negative oceanic feedback, and the coefficient of  $trSpeed$   
282 is positive. There is a greater potential for a storm far from its PI value to continue in-  
283 tensifying, and the coefficient of  $dPI_{V_0}$  is positive. Stronger shear prevents storms from  
284 strengthening, and the coefficient of  $SHRD$  is negative. The stronger the divergence and  
285 the more unstable atmosphere is, the better chance a system has to develop. Therefore the  
286 coefficients of  $div200$  and  $dThetaEs$  are positive and negative, respectively.

## 287 4. MLR intensity forecasts and errors

288 *a. Examples: Rita (2005), Earl (2010), and Irene (2011), Isaac (2012)*

289 After reducing the number of predictors from 13 to 7, we test the MLR with independent  
290 data from the period 2000 to 2012. As examples, we show the intensity predictions for four  
291 storms, Hurricanes Rita (2005), Earl (2010), Irene (2011), and Isaac (2012); their tracks are  
292 shown in Fig. 2. These storms encompassed a wide range of intensities, and have tracks  
293 that are typical for North Atlantic storms. Rita was a category 5 hurricane; Earl was strong  
294 category 4; Irene was category 3; and Isaac was weak category 1. Both Rita and Earl  
295 underwent rapid intensification. Rita and Isaac passed over the Gulf of Mexico while Earl  
296 and Irene recurved and moved toward the northeast US.

297 The MLR creates five day forecasts at 12 hourly intervals, and we apply it every 12  
298 hours during the lifetime of the storms (red lines in Fig. 3). For the purpose of comparison,  
299 we also plot forecasts from the SHIPS (blue line) in Fig. 3. The SHIPS predicts a faster  
300 intensification rate than the MLR does overall. Both models over-predict the intensity change  
301 (Isaac 18 to 21 Aug; Earl 25 to 28 Aug) while Isaac and Earl are in their early stage (tropical  
302 storm intensity). Aside from failing to keep a storm in its TS stage, neither the MLR nor  
303 the SHIPS is capable of capturing rapid intensification (RI). As addressed in DeMaria and  
304 Kaplan (1994), this is the limitation of a simple linear regression model. For a case where  
305 the storm strengthens gradually (Irene), both models do a fair job (Fig. 3d). In the decaying  
306 period, when the intensities drop rapidly, both models again have difficulty (after 22 Sep in  
307 Rita; after 03 Sep in Earl).

308 *b. Dependence of forecast errors on storm characteristics and synoptic conditions*

309 Errors in intensity predictions (from both dynamic and statistical models) are related  
310 to predictors particular to individual storms, such as initial intensity, translation speed,  
311 latitude, current intensity, PI, and wind shear (Bhatia and Nolan 2013). The relationship

312 between errors and predictors can differ among forecast times, and is not linear. Under-  
313 standing the relationship between errors and predictors can help us to better interpret and  
314 utilize forecast results. We stratify errors by predictor value and forecast time (not shown).  
315 Overall, the error increases with lead time till around 96 hour lead time, after which the  
316 error seems to be saturated. The weaker the initial storm intensity, the larger the error can  
317 be, especially for longer lead times. Usually, strong storms are closer to PI, therefore there  
318 is less chance for them to intensify or undergo RI. Strong storms are also more resilient to  
319 shear, and are unlikely to undergo rapid weakening. Therefore, there is an upper bound  
320 on errors from both intensifying and decaying storms. We can expect that the MLR might  
321 have larger errors while predicting weak storms. We also observe that the largest errors  
322 occur when  $dVdt$  is large (strong storm intensification fast in the past 12 hours), which is  
323 believed to be related to RI processes which are unresolved in the MLR. Larger errors also  
324 occur as storms move faster. Such situation usually happens when storms recurve to the  
325 mid-latitude. Another situation with possible larger errors occurs when the storms are in  
326 an environment which are favorable for RI (large PI and large positive  $dPI-V_0$ ), and when  
327 storms recurve (large negative  $dPI-V_0$ ). Forecast errors are also sensitive to SHRD and  
328 atmospheric stability, and the largest errors occur in unstable and moderate environmental  
329 shear. Moderate shear and more unstable atmosphere are typical environments for storms  
330 in the tropical Atlantic ( $10-20^\circ N$ ) and the Gulf of Mexico.

331 To understand the relationship between the predictors and the forecast error, in Fig. 4 we  
332 show the sparkline chart of the standard deviation ( $\sigma$ ) of the forecast error for the individual  
333 predictors for 24 hours lead time. The larger the standard deviation's variability over the  
334 individual predictor's magnitude range, the more sensitive is the error to the predictor.  
335 Among all predictors, the forecast error seems to be most sensitive to the initial storm  
336 intensity, somewhat less to PI and  $dPI-V_0$ , and not very sensitive to the remaining predictors.

337 The dependence of the MLR error on storm location is shown in Fig. 5. In general, the  
338 error seems to show no systematic dependence on location, especially for shorter lead times.

339 However, there are negative biases over the Caribbean sea and western Gulf of Mexico. Over  
340 these regions, there is usually a very warm upper ocean, which is thought to be essential for RI  
341 (Lin et al. 2009; Hui and Wang 2014). Therefore such bias could be due to lacking predictors  
342 related to warm upper ocean structure in the MLR, which leads to an underestimation of  
343 storm intensity prediction. However, because OT100 (the ocean temperature averaged over  
344 the upper 100 m) does not passed the predictor selection procedure, considering upper ocean  
345 structure might only be necessary for the storms passing through the warm upper ocean area.  
346 For 96 hour forecasts, there is a positive bias (the MLR over-predicts the storm intensity)  
347 over the west Atlantic, especially along the east coast of United States.

348 *c. Verification of the MLR*

349 To evaluate the overall performance of the MLR, we calculate the mean absolute error  
350 (MAE) and the RMSE against the best track data from all the testing cases (Fig. 6a).  
351 Errors from the NHC official forecast (OFCL), SHIPS, and the Statistical Hurricane Intensity  
352 Forecast (SHIFOR; Knaff et al. 2003) from 2000 to 2012 are also shown. It is important  
353 to note that the specific data used by the MLR (e.g., reanalysis rather than forecast fields)  
354 gives an advantage to MLR over the operational forecasts, a point we will discuss in more  
355 detail later. Both the MAE and RMSE of the MLR increase with lead time from 12 to  
356 96 hours forecasts. The error does not continue to grow from 96 to 120 hours in the MLR  
357 and OFCL. In theory, we expect the error to increase with forecast time, but the error can  
358 not grow indefinitely. In other words, there should be an upper bound on the forecast error,  
359 and we expect the error to level off after a certain lead time, as does the error from OFCL.  
360 However, although not significant, there is small decline of the error from 96 to 120 hours in  
361 the MLR, which could be due to the small sample size.

362 Conventionally, for any TC intensity prediction model to be considered skillful, the model  
363 has to provide forecasts that are better than SHIFOR. As expected, the MLR, the SHIPS,  
364 and the OFCL all have skill by this criterion at almost all lead times. An additional baseline

365 model is the persistence model, whose forecast intensity at all leads is simply the initial inten-  
366 sity (black-dashed lines in Fig. 6) here. The MLR is skillful for all forecasting times against  
367 this baseline as well (red line has smaller errors than the black dashed line in Fig. 6a, b).

368 Compared to the SHIPS, the MLR has smaller MAE and RMSE for the 12-hour forecasts,  
369 and the 12-hour SHIPS MAE is about 0.5 kt larger than that of the baseline model. However,  
370 this behavior is due to the fact that the baseline model (as well as the MLR) is using the  
371 storm track and initial intensity from best-track data rather the NHC intensity forecast  
372 (provided by ATCF) which is the best information available for SHIPS to use. Since all the  
373 models are verified against the best-track data, it is reasonable that our baseline model and  
374 the MLR have smaller errors than does SHIPS.

375 The MAE of the MLR 12-hour forecasts best-track data is 5.8 kt. Replacing the best-  
376 track data for current and past storm intensity with those from ATCF, the 12-hour MLR  
377 MAE increases to 6.5 kt (not shown), which is close to that of SHIPS. The MAE of the MLR  
378 for longer lead times (96-120 hours) increases by 1 kt (not shown) when ATCF data is used.  
379 To assess the impact of track forecast error on the performance of the MLR, we replace  
380 best-track location data with NHC forecast tracks (MLR\_FstTrack, blue lines in Fig. 6c,  
381 and d). After the 12-hour forecasts, the MAE in the MLR\_FstTrack forecasts is on average  
382 2-3 kts higher than the MLR. To estimate the impact of using reanalysis data rather than  
383 GFS forecast data for environmental fields, we conduct an experiment using predictors taken  
384 from the NCEP reanalysis fields at the initial time and then keeping these values fixed (i.e.  
385 persistence) for all lead times (MLR\_initNCEP, orange line in Fig. 6c, and d). MAE in the  
386 MLR\_initNCEP experiment increases by about 1-2 kt in average.

387 Another important difference between MLR and SHIPS is the choice of predictors. SHIPS  
388 includes many more environmental parameters than the MLR does. Till 2013, SHIPS had  
389 24 predictors (Schumacher et al. 2013). We conducted experiments to see how varying  
390 the numbers of predictors impacts the errors. Figure 6e, and f show the errors from three  
391 of these experiments. When using  $dPI-V_0$  ( $dPI-V_0$ : blue line) as the only predictor, we

392 surprisingly find out that the intensity errors, both MAE and RMSE, are not too far from  
393 those obtained with the full MLR. Adding  $dVdt$  (dPI-V<sub>0</sub>+dVdt: solid sky-blue line) as  
394 the second predictor results in a small improvement. These results suggest that with only  
395 one or two predictors, the regression model still has some degree of skill. If we continue  
396 adding predictors, we continue to see an improvement in the MLR performance, such as  
397 adding SHRD (...+SHRD: dashed-sky-blue line). Using all 13 predictors from initial pool  
398 of predictors (MLR\_13: orange line), however, results in larger errors for 24 to 60 hours  
399 forecasts, an indication of over-fitting.

400 The difference between the performance of the MLR and SHIPS results from using best-  
401 track or NHC forecasting TC data, using global reanalysis or forecast fields, and the choice  
402 of predictors. Another difference is the data used to estimate the coefficients. SHIPS co-  
403 efficients are re-estimated each year; MLR coefficients are not. The use of best-track data  
404 and reanalysis fields are the main reason why the MLR performs better than the SHIPS.  
405 It is not clear whether the use of more predictors causes better performance in the SHIPS  
406 than in the MLR for 24 to 60-hour lead times. Nevertheless, the experiments with different  
407 combination of predictors suggest that a multiple-linear regression model with even just a  
408 few of the most essential predictors has substantial skill.

## 409 **5. Monthly v.s. daily reanalysis data**

410 Given our interest in the variations of environment and TC intensity on long time-scales,  
411 an interesting question is to what extent the MLR can predict intensity changes given  
412 monthly rather than daily data. In an application where climate model output is used,  
413 for example, being able to use monthly data is a great advantage. With the same choice of  
414 predictors, we use monthly NCEP-NCAR reanalysis data to train and test the MLR (called  
415 MLR\_Monthly). The forecasting errors of the MLR\_Monthly are comparable to those of the  
416 MLR (Fig. 7). There is almost no discrepancy between the forecast errors of the MLR and

417 the MLR\_Monthly for lead times less than 48 hours. For longer lead times ( $> 72$  hours), on  
418 the other hand, the advantage of using daily data is clear.

419 The good performance of the MLR\_Monthly suggests that there is some usable relation  
420 between the daily and monthly predictors. Comparing the predictors from monthly data to  
421 those from daily data, we find that the monthly  $dPI-V_0$  is larger than the daily one (Fig. 8),  
422 and that this difference is a result of larger monthly PI. There is almost no relationship  
423 between monthly and daily  $div200$ , and a forward selection with monthly predictors (not  
424 shown) does not choose  $div200$ . This is because the monthly  $div200$  is too small (and  
425 uncorrelated) compared to the daily one. There are also some positive biases for smaller  
426 values of SHRD and  $dThetaEs$ , but negative biases for larger values in monthly data. To  
427 assess which one is the primary cause of the larger error in MLR\_Monthly, we replace each  
428 of the four predictors with daily data in the MLR\_Monthly separately (not shown). Our  
429 results suggest that the difference in  $dPI-V_0$  in monthly and daily data is the main cause for  
430 the differences between MLR and MLR\_Monthly.

## 431 **6. Probabilistic intensity prediction**

### 432 *a. Forecast distribution*

433 As described in the introduction, accounting for the uncertainty of the MLR forecast  
434 increases its value and utility. The simplest and most direct estimate of forecast uncertainty  
435 is given by the empirical distribution of errors presented in the training process (using all  
436 the training data). This approach is feasible because of the size of the training dataset used  
437 here. Although we recognize that rare events may be poorly sampled and may benefit from  
438 parametric descriptions, we only consider using the empirical distribution of errors here. We  
439 condition the error distribution on forecast lead time, and show in Fig. 9a the increasing  
440 spread of the error PDF as the forecast lead time increases from 12 to 72 hours. There  
441 is no significant change in the error PDF between 96 and 120 forecast hours; the errors

442 appear to have saturated at that point. The errors are close to being normally distributed  
443 for all lead times. Figure 9b is the cumulative distribution function (CDF) and Fig. 9c  
444 compares the distribution of the MLR errors for 24 hour forecast to the standard normal  
445 distribution. Except for the extreme points, most of the points fit well within the theoretical  
446 normal distribution profile. All the large negative errors ( $< -35$  kt) are from the cases with  
447 rapid intensification. For instance, the largest negative error in Fig. 9b is about -86 kt,  
448 and corresponds to a forecast of Hurricane Emily in 1987 when Emily underwent rapid  
449 intensification. Our probabilistic intensity forecasts consist of the PDF with mean value  
450 given by the MLR and spread given by the lead-time dependent error distribution (Fig. 10).

451 The analyses in Section 3b suggest that the forecast error and its distribution are sensitive  
452 to the predictors, and are most sensitive to  $V_0$  (initial storm intensity). Strictly speaking,  
453 such dependence is beyond the linear regression framework, in which the error variance  
454 is assumed independent of the predictors. However, regression based-models with varying  
455 spread are common in weather forecasting (Gneiting et al. 2005; Wilks and Hamill 2007).  
456 Here, we examine the sensitivity of forecast error distribution to  $V_0$ . We illustrate the  
457 dependence of the error distribution on initial wind speed in the case of the 24 hour forecast  
458 error (Fig. 11). The standard deviation of all 24-hour forecast errors from all training data  
459 is 12 kts (dashed line in Fig. 11a, which is calculated from the blue line in Fig. 9a). However,  
460 binning 24-hour forecast errors based on the corresponding values of  $V_0$  shows that the  $\sigma$  of  
461 errors increases, roughly linearly, as  $V_0$  increases (black diamond-line in Fig. 11c). Now, if  
462 we do not account for the fact that the distribution of forecast error is a function of  $V_0$ , the  
463 joint CDF of errors between lead times and  $V_0$  would look like Fig. 11b. If we do, the joint  
464 CDF would look like Fig. 11c. To understand how accounting for the dependence of error  
465 distribution to  $V_0$  changes the results, we develop two probabilistic models. We name the  
466 probabilistic model with error distribution as a function of lead time only MLR\_1, while the  
467 one utilizing error distribution as a function of both lead time and  $V_0$  is named MLR\_2.

468 In addition to MLR\_1 and MLR\_2, we construct another probabilistic model as a baseline

469 reference model from the persistence model (Fig. 6a) and its error distribution (Fig. 10b).  
470 Fig. 10b indicates that the mean of the distribution is near zero at any given lead time  
471 and that the spread increases with lead time. The spread here is larger than the spread in  
472 Fig. 10a. This is consistent with results in Fig. 6b: the  $\sigma$  (which can also be represented  
473 as RMSE) in persistence model is larger than that in the MLR. We call this model the  
474 probabilistic persistence model. When making a probabilistic forecast with the persistence  
475 model, we can simply shift the distribution in Fig. 10b to be centered on the initial intensity.

476 The MLR\_1 and MLR\_2 are considered skillful to the extent to which their probabilistic  
477 performance is better than that of the probabilistic persistence model. The fundamental  
478 difference is that the probabilistic persistence distribution is conditioned only on the current  
479 intensity and historical changes in storm intensity, while the MLR\_1 and MLR\_2 distributions  
480 are additionally conditioned on the current and forecast environment.

#### 481 *b. Probabilistic forecasts and verification*

482 Probabilistic forecasts (persistence, MLR\_1, and MLR\_2) are made for all the cases in the  
483 testing data set. Figure 12 shows probabilistic forecasts for Hurricane Earl initialized on 0000  
484 UTC 28 August, 2010, a period during which Earl underwent RI (36 - 60 forecast hours).  
485 For the 12 - 36 hour lead times, the observed storm intensity falls into the 25-75 percentile  
486 region in MLR\_1 and MLR\_2 forecast, while in the persistence model it does not. Once  
487 Earl started to undergo RI, the best-track data falls beyond the highest 50 percentile region  
488 in all models, although it is located at higher probabilities in MLR\_1 and MLR\_2 than in  
489 the persistence model. Qualitatively, the case in Fig. 12a-c exhibits the advantage in using  
490 MLR\_1 and MLR\_2 instead of the persistence model for probabilistic forecast. Knowing  
491 current and forecast environmental conditions helps probabilistic intensity prediction.

492 We use the rank probability skill score (RPSS) to verify quantitatively whether the  
493 MLR\_1 and MLR\_2 have more skill than the persistence model. Rank probability score  
494 (RPS) is a squared-error score with respect to the observational probability, which is 1 if

495 the forecast event occurs, and 0 if the event does not occur. The cumulative forecasts and  
 496 observations, denoted as  $P_m$  and  $O_m$ , are:

$$P_m = \sum_{j=1}^m p_j, \quad m = 1, \dots, J, \quad (1)$$

497 and

$$O_m = \sum_{j=1}^m o_j, \quad m = 1, \dots, J. \quad (2)$$

498 where  $p_j$  is the forecast probability of the storm intensity falling in the  $j$ -th bin, the observed  
 499 probability  $o_j$  is 1 if the observations falls in the  $j$ -th bin and zero otherwise, and  $J$  is the  
 500 total number of bins. Here, we stratify the probabilities by the corresponding  $V_{max}$  every  
 501 5 kt from 0 to 200 kt. The RPS is the sum of the squared differences between the cumulative  
 502 probabilities  $P_m$  and  $O_m$ :

$$RPS = \sum_{m=1}^J (P_m - O_m)^2. \quad (3)$$

503 RPS is oriented so that smaller values indicate better forecasts. A correct forecast with no  
 504 uncertainty has a RPS of 0. The RPSS compares the average RPS to that of the persistence  
 505 forecast:

$$RPSS = 1 - \frac{\overline{RPS}}{RPS_{persistence}}. \quad (4)$$

506 If the probabilistic model is skillful, RPSS is positive, and the larger the value is, the  
 507 more skill the model has. The RPS and RPSS can be computed for each forecast. Figure 12d  
 508 shows the RPSS for all of the forecasts of Earl. Compared to the persistence model, the  
 509 skill of MLR\_1 and MLR\_2 increases with the length of forecast interval, reflecting the fact  
 510 that the value of knowing the environment increases as lead time increases. As the forecast  
 511 interval becomes longer, the environmental conditions become more crucial for intensity  
 512 prediction. Hence, the MLR\_1 and MLR\_2 forecasts show greater skill compared to the  
 513 persistence model with increasing lead time. For all the testing cases (Fig. 13), RPSS for  
 514 the MLR\_1 and MLR\_2 again increases with time. The modest advantage in RPSS that the  
 515 MLR\_2 shows over the MLR\_1 indicates that considering the sensitivity of forecast errors to  
 516 initial storm intensity does improve probabilistic intensity predictions.

517 The performance of probabilistic forecasts of particular intensity categories (TS, cate-  
518 gory 1-2, and 3-5 hurricane) can also be verified. Figure 14 shows the reliability diagrams of  
519 MLR\_1, MLR\_2, the persistence model, and the NHC official probabilistic forecast (OFCL)  
520 for categorized cases from 2008 to 2012 at all lead times. Forecast probabilities at each  
521 forecast time for each intensity category are grouped into bins at 0.1 (10%) interval. The  
522 probabilistic forecasts of the occurrence of TS and hurricane intensity Saffir-Simpson cate-  
523 gories 1-2 by the two MLR models, persistence model and OFCL all show a fair reliability as  
524 the observed frequency is close to the forecast probability (Fig. 14a, c). The MLR\_1, MLR\_2  
525 and persistence model forecasts of TS occurrence show almost perfect reliability for high  
526 probabilities ( $> 0.6$  bin) and under-forecasting for low probabilities. OFCL over-forecasts  
527 the occurrence of TS in the 0.1 - 0.6 (10% - 60%) range. MLR\_1, MLR\_2, and OFCL pre-  
528 dict occurrence of hurricane intensity categories 1-2 reliably for forecasts probabilities of  
529 0 - 0.3 (0% - 30%), while the persistence model under-forecasts it. For higher probabilities  
530 ( $> 0.3$  bin), they all over-forecast the occurrence.

531 For major hurricanes (category 3-5), the reliability diagram curves are very different  
532 among these four forecast models. MLR\_1 and MLR\_2 are still fairly reliable. The persis-  
533 tence model and OFCL, however, are substantially less reliable. They both overforecast the  
534 occurrence of major hurricanes. For the persistence model, the probabilistic forecast always  
535 maintains relatively higher probabilities close to the initial storm intensity, as shown in Fig.  
536 12a. Hence, the persistence model could easily overforecast decaying major storms. Com-  
537 pared to the MLR\_1, MLR\_2 is more reliable (the red line in Fig. 14 is closer to the one-to-one  
538 line than the pink line). Note that OFCL overestimates all categories. This implies that  
539 OFCL must underestimate tropical depression events, which could be from overforecasting  
540 the development from tropical depressions to TS or underforecasting the decaying rate of  
541 tropical storms.

## 7. Conclusions

This study describes the development of a probabilistic tropical cyclone (TC) intensity prediction system based on a multiple linear regression model (MLR) and the past error distribution. While such models are important for operational intensity forecasts, our interest is more directed at the capability of such model to connect TC intensity to the large-scale environment. The ability of the MLR to relate variations in the large-scale environment to TC intensity has applications beyond short-term forecasting and can be a tool towards understanding the distribution of TC intensities in a climate undergoing natural or anthropogenic changes. The MLR contains 7 predictors selected based on their physical importance to TC development and their ability to reduce forecast errors. The predictors are the initial maximum wind speed ( $V_0$ ), change of storm intensity in the past 12 hours ( $dVdt$ ), storm translation speed ( $trSpeed$ ), the difference between PI and  $V_0$  ( $dPI-V_0$ ), 850-200 mb vertical wind shear magnitude (SHRD), conditional stability ( $dThetaEs$ ), 200 mb divergence ( $div200$ ). Among these predictors,  $dVdt$  is the most important predictors for forecasts with short lead times ( $< 24$  hours).  $dPI-V_0$  and SHRD becomes more important with increasing lead times. We try to avoid predictors whose relationship to TC intensity is likely to change in a changed global mean climate, such as sea surface temperature (SST).

We find that the magnitudes of forecast errors depends on the synoptic situation. In particular, larger forecast errors occur when the storm is weaker and when there is a higher chance of rapid intensification. The MLR does not represent RI well, similarly to other linear models. Additionally, the MLR has relatively larger errors under moderate shear, smaller upper tropospheric divergence, and a conditionally more unstable atmosphere.

A relatively simple description of the synoptic environment is able to characterize much of the predictable component of the variability of TC intensity. The errors of a single predictor ( $dPI-V_0$ ) model are not too far from those of the full 7 predictors model. Likewise, a model based on monthly averaged environment performs quite well, with its primary weakness being related to the lack of daily PI information. Given these results, monthly data from

569 global climate models with only the most essential predictors might be sufficient for some  
570 climate applications.

571 In addition to the deterministic forecasts, we demonstrate the ability of the MLR to  
572 produce probabilistic predictions. We construct a forecast PDF with mean given by the  
573 MLR and uncertainty given by the empirical error distribution from the training data. Two  
574 probabilistic models are developed here. In MLR\_1, the error distribution is only a func-  
575 tion of lead time while in MLR\_2, it is also a function of initial storm intensity, reflecting  
576 the dependence of forecast errors on the environment. A reference probabilistic forecast is  
577 constructed from historical storm intensity changes (independent of environment) and used  
578 to compute the rank probability score skill (RPSS). Both MLR\_1 and MLR\_2 have good  
579 skill, especially for a long lead forecasts when their advantage over the persistence forecast  
580 is greatest. The analysis of the reliability diagram further shows that the reliability among  
581 two MLR probabilistic models, the persistence model, and official NHC forecasts are similar  
582 to each other for TS, and hurricanes of categories 1 and 2. Nevertheless, MLR\_1 and MLR\_2  
583 are much more reliable than OFCL and persistence model for major hurricanes. Further-  
584 more, both RPSS and reliability analyses suggest that using an environment-dependent error  
585 distribution (MLR\_2) improves the forecast skill.

586 In short, in this study we have developed a system that reliably predicts the distribution  
587 of storm intensity conditional on the surrounding local environment. We have developed  
588 and verified the system in a forecast framework using regression methods and probabilistic  
589 verification to assess resolution and reliability of the intensity distribution. This reliable  
590 probabilistic model could be a component of a system for future risk assessment.

591 *Acknowledgments.*

592 We thank Dr. Mark DeMaria for the useful comment. The research was supported by  
593 Office of Naval Research under the research grant of MURI (N00014-12-1-0911).

## REFERENCES

- 596 Behringer, D. W. and Y. Xue, 2004: Evaluation of the global ocean data assimilation system  
597 at NCEP: The Pacific Ocean. *Eighth Symposium on Integrated Observing and Assimila-*  
598 *tion Systems for Atmosphere, Oceans, and Land Surface*, AMS 84th Annual Meeting,  
599 Washington State Convention and Trade Center, Seattle, WA, 11–15.
- 600 Bender, M. A., T. R. Knutson, R. E. Tuleya, J. J. Sirutis, G. A. Vecchi, S. T. Garner, and  
601 I. M. Held, 2010: Modeled impact of anthropogenic warming on the frequency of intense  
602 Atlantic hurricanes. *Science*, **327**, 454–458.
- 603 Bhatia, K. T. and D. S. Nolan, 2013: Relating the skill of tropical cyclone intensity forecasts  
604 to the synoptic environment. *Wea. Forecasting*, **28**, 961–980.
- 605 Bister, M. and K. A. Emanuel, 2002: Low frequency variability of tropical cyclone potential  
606 intensity 1. Interannual to interdecadal variability. *J. Geophys. Res.: Atmospheres*, **107**,  
607 4801.
- 608 Bretherton, C. S., M. E. Peters, and L. E. Back, 2004: Relationships between water vapor  
609 path and precipitation over the tropical oceans. *J. Climate*, **17**, 1517–1528.
- 610 Bryan, G. H. and R. Rotunno, 2009: Evaluation of an analytical model for the maximum  
611 intensity of tropical cyclones. *J. Atmos. Sci.*, **66 (10)**, 3042–3060.
- 612 Camargo, S. J., 2013: Global and regional aspects of tropical cyclone activity in the CMIP5  
613 models. *J. Climate*, **26**, 9880–9902.
- 614 Camargo, S. J., K. A. Emanuel, and A. H. Sobel, 2007: Use of a genesis potential index to  
615 diagnose ENSO effects on tropical cyclone genesis. *J. Climate*, **20**, 4819–4834.

- 616 Camargo, S. J., M. C. Wheller, and A. H. Sobel, 2009: Diagnosis of the MJO modulation of  
617 tropical cyclogenesis using an empirical index. *J. Atmos. Sci.*, **66**, 3061–3074.
- 618 Chen, J.-H. and S.-J. Lin, 2013: Seasonal predictions of tropical cyclones using a 25-km-  
619 resolution general circulation model. *J. Climate*, **26**, 380–398.
- 620 Chen, S. S., F. Judt, J. Berner, C.-Y. Lee, M. Curcic, C. Snyder, and R. Rotunno, 2014:  
621 Distinct characteristics of hurricane ensemble forecasts using physical parameterizations  
622 vs. stochastic perturbations. AMS 94th Annual Meeting, Atlanta, GA.
- 623 Chen, S. S., J. A. Knaff, and F. D. Marks, 2006: Effects of vertical wind shear and storm  
624 motion on tropical cyclone rainfall asymmetries deduced from TRMM. *Mon. Wea. Rev.*,  
625 **134**, 3190–3208.
- 626 DeMaria, M., 2006: The effect of vertical shear on tropical cyclone intensity change. *J.*  
627 *Atmos. Sci.*, **53**, 2076–2088.
- 628 DeMaria, M., 2009: A simplified dynamical system for tropical cyclone intensity prediction.  
629 *Mon. Wea. Rev.*, **137**, 68–82.
- 630 DeMaria, M., 2010: Tropical cyclone intensity change predictability estimates using a  
631 statistical-dynamical model. 29th AMS Conference on Hurricanes and Tropical Meteorology,  
632 Tucson, AZ.
- 633 DeMaria, M. and J. Kaplan, 1994: A Statistical Hurricane Intensity Prediction Scheme  
634 (SHIPS) for the Atlantic basin. *Wea. Forecasting*, **9**, 209–220.
- 635 DeMaria, M. and J. Kaplan, 1999: An updated Statistical Hurricane Intensity Prediction  
636 Scheme (SHIPS) for the Atlantic and Eastern North Pacific basins. *Wea. Forecasting*, **14**,  
637 326–337.
- 638 DeMaria, M., J. A. Knaff, R. Knabb, C. Lauer, C. R. Sampson, and R. T. DeMaria, 2009:

639 A new method for estimating tropical cyclone wind speed probabilities. *Wea. Forecasting*,  
640 **24**, 1573–1591.

641 DeMaria, M., J. A. Knaff, and C. Sampson, 2007: Evaluation of long-term trends in tropical  
642 cyclone intensity forecasts. *Meteorology and Atmospheric Physics*, **97 (1-4)**, 19–28.

643 DeMaria, M., M. Mainelli, L. K. Shay, J. A. Knaff, and J. Kaplan, 2005: Further improve-  
644 ments to the Statistical Hurricane Intensity Prediction Scheme (SHIPS). *Wea. Forecasting*,  
645 **20**, 531–543.

646 DeMaria, M., et al., 2013: Improvements to the operational tropical cyclone wind speed  
647 probability model. *Wea. Forecasting*, **28**, 586–602.

648 Emanuel, K. and A. H. Sobel, 2013: Response of tropical sea surface temperature, precipita-  
649 tion, and tropical cyclone-related variables to changes in global and local forcing. *J. Adv.*  
650 *Model. Earth Syst.*, **5**, 1942–2466.

651 Emanuel, K. A., 1988: The maximum intensity of hurricanes. *J. Atmos. Sci.*, **45**, 1143–1155.

652 Emanuel, K. A., 1995: Sensitivity of tropical cyclones to surface exchange coefficients and a  
653 revised steady-state model incorporating eye dynamics. *J. Atmos. Sci.*, **52**, 3969–3976.

654 Emanuel, K. A., 2005: Increasing destructiveness of tropical cyclones over the past 30 years.  
655 *Natural*, **436**, 686–688.

656 Emanuel, K. A., 2006: Climate and tropical cyclone activity: A new model downscaling  
657 approach. *J. Climate*, **19**, 4797–4802.

658 Emanuel, K. A., 2013: Downscaling CMIP5 climate models shows increased tropical cyclone  
659 activity over the 21st century. *PNAS*, published online.

660 Emanuel, K. A., S. Ravela, E. Vivant, and C. Risi, 2006: A statistical deterministic approach  
661 to hurricane risk assessment. *Bull. Amer. Meteor. Soc.*, 299–313.

662 Gneiting, T., A. E. Raftery, A. H. Westveld, and T. Goldman, 2005: Calibrated probabilistic  
663 forecasting using ensemble model output statistics and minimum CRPS estimation. *Mon.*  
664 *Wea. Rev.*, **133**, 1098–1118.

665 Hui, W. and Y. Wang, 2014: A numerical study of Typhoon Megi (2010). Part I: Rapid  
666 Intensification. *Mon. Wea. Rev.*, **142**, 29–48.

667 Jarvinen, B. R., C. J. Neumann, and M. A. S. Davis, 1984: A tropical cyclone data tape for  
668 the North Atlantic basin, 1886-1983: Contents, limitations, and uses. NOAA Technical  
669 Memorandum NWS NHC 22, NOAA, Coral Gables, Florida, 21pp.

670 Johnson, N. C. and S.-P. Xie, 2010: Changes in the sea surface temperature threshold for  
671 tropical convection. *Natural Geosci.*, **3**, 842–845.

672 Kalnay, E., et al., 1996: The NCEP/NCAR 40-year reanalysis project. *Bull. Amer. Meteor.*  
673 *Soc.*, **77**, 437–471.

674 Kaplan, J. and M. DeMaria, 2003: Large-scale characteristics of rapidly intensifying tropical  
675 cyclones in the North Atlantic basin. *Wea. Forecasting*, **18**, 1093–1108.

676 Knaff, J. A., M. DeMaria, C. R. Sampson, and J. M. Gross, 2003: Statistical, 5-day tropical  
677 cyclone intensity forecasts derived from climatology and persistence. *Wea. Forecasting*,  
678 **18**, 80–92.

679 Knutson, T. R., J. J. Sirutis, S. T. Garner, I. M. Held, and R. E. Tuleya, 2007: Simulation of  
680 the recent multidecadal increase of Atlantic hurricane activity using a 18-km-grid regional  
681 model. *Bull. Amer. Meteor. Soc.*, **88**, 1549–1565.

682 Knutson, T. R., J. J. Sirutis, S. T. Garner, G. A. Vecchi, and I. M. Held, 2008: Simulated  
683 reduction in Atlantic hurricane frequency under twenty-first-century warming condition.  
684 *Natural Geosci.*, **1**, 359–364.

685 Knutson, T. R., et al., 2010: Tropical cyclones and climate change. *Natural Geosci.*, **3**,  
686 157–163.

687 Krishnamurti, T. N., C. M. Kishtawal, T. E. LaRow, D. R. Bachiochi, Z. Zhang, C. E.  
688 Williford, S. Gadgil, and S. Surendran, 1999: Improved weather and seasonal climate  
689 forecast from multimodel superensemble. *Science*, **285**, 1548–1550.

690 Landsea, C. W. and J. L. Franklin, 2013: Atlantic hurricane database uncertainty and  
691 presentation of a new database format. *Mon. Wea. Rev.*, **141**, 3576–3592.

692 Lin, I.-I., C.-H. Chen, I.-F. Pun, W. T. Liu, and C.-C. Wu, 2009: Warm ocean anomaly, air  
693 sea fluxes, and the rapid intensification of Tropical Cyclone Nargis (2008). *Geophys. Res.*  
694 *Lett.*, **36**, L03 817.

695 Lin, N., K. A. Emanuel, J. A. Smith, and E. Vanmarcke, 2010: Risk assessment of hurricane  
696 storm surge for New York City. *J. Geophys. Res.: Atmospheres*, **115**, D18 121.

697 Price, J. F., 2009: Metrics of hurricane-ocean interaction: vertically-integrated or vertically-  
698 averaged ocean temperature? *Ocean Science*, **5 (3)**, 351–368.

699 Puri, K., J. Barkmeijer, and T. Palmer, 2001: Ensemble prediction of tropical cyclones using  
700 targeted diabatic singular vectors. *Q. J. R. Meteor. Soc.*, **127**, 709–731.

701 Ramsay, H. A. and A. H. Sobel, 2011: Effects of relative and absolute sea surface temperature  
702 on tropical cyclone potential intensity using a single-column model. *J. Climate*, **24**.

703 Sampson, C. R. and A. J. Schrader, 2000: The automated tropical cyclone forecasting system  
704 (version 3.2). *Bull. Amer. Meteor. Soc.*, **81**, 1231–1240.

705 Schumacher, A., M. DeMaria, and J. Knaff, 2013: Summary of the new statistical-dynamical  
706 intensity forecast models for the Indian Ocean and Southern Hemisphere and resulting  
707 performance. Tech. rep., JTWC.

708 Smith, R. K., M. T. Montgomery, and S. Vogl, 2008: A critique of Emanuel’s hurricane  
709 model and potential intensity theory. *Quart. J. Roy. Meteor. Soc.*, **134** (632), 551–561,  
710 doi:10.1002/qj.241.

711 Tang, B. and K. A. Emanuel, 2012: A ventilation index for tropical cyclones. *Bull. Amer.*  
712 *Meteor. Soc.*, **93**, 1901–1912.

713 Tippett, M., S. J. Camargo, and A. H. Sobel, 2011: A poisson regression indec for tropical  
714 cyclone genesis and the role of large-scale vorticity in genesis. *J. Climate*, **21**, 2335–2357.

715 Tippett, M. K., A. H. Sobel, and S. J. Camargo, 2012: Association of monthly U.S. tornado  
716 occurrence with large-scale atmospheric parameters. *Geophys. Res. Lett.*, **39**, L02 801.

717 Vecchi, G. A. and B. J. Soden, 2007: Effect of remote sea surface temperature change on  
718 tropical cyclone potential intensity. *Natural*, **450**, 1066–1070.

719 Vecchi, G. A., M. Zhao, G. Villarini, R. Anthony, A. Kumar, I. M. Held, and R. Gudgel,  
720 2011: Statisitcal-dynamical predictions of seasonal north atlantic hurricane activity. *Mon.*  
721 *Wea. Rev.*, **139**, 1070–1082.

722 Weber, H. C., 2005: Probabilistic prediction of tropical cyclones. Part II: Intensity. *Mon.*  
723 *Wea. Rev.*, **133**, 1853–1864.

724 Webster, P. J., G. Holland, J. A. Curry, and H.-R. Chang, 2005: Changes in tropical cyclone  
725 number, duration, and intensity in a warming environment. *Science*, **309**, 1844–1846.

726 Wilks, D. S. and T. M. Hamill, 2007: Comparison of ensemble-MOS methods using GFS  
727 reforecasts. *Mon. Wea. Rev.*, **135**, 2379–2390.

728 Zhang, Z. and T. N. Krishnamurti, 1999: A perturbation method for hurricane ensemble  
729 predictions. *Mon. Wea. Rev.*, **127**, 447–469.

730 Zhao, M., I. M. Held, S.-J. Lin, and G. A. Vecchi, 2009: Simulations of global hurricane cli-  
731 matology, interannual variability, and response to global warming using a 50-km resolution  
732 GCM. *J. Climate*, **22**, 6653–6678.

733 Zipser, E., 2003: *Some views on “Hot Towers” after 50 years of tropical field programs and*  
734 *two years of TRMM data*. 29, Amer. Meteor. Soc., 49-57 pp.

735 **List of Tables**

736 1 Descriptions of experiments.

32

TABLE 1. Descriptions of experiments.

experiment	description
MLR	a multiple linear regression model with all 7 selected predictors, NHC best-track data, daily NCEP-NCAR reanalysis fields
MLR_FstTrack	same as MLR, but with environmental predictors calculated along the NHC official forecast track
MLR_InitNCEP	same as MLR, but with environmental predictors calculated from NCEP-NCAR reanalysis fields at the initial time
dPI_V <sub>0</sub>	a linear regression model with only dPI_V <sub>0</sub> as a predictor, NHC best-track data, daily NCEP-NCAR reanalysis fields
dPI_V <sub>0</sub> +dVdt	same as dPI_V <sub>0</sub> , but with one more predictor: dVdt
... + SHRD	same as dPI_V <sub>0</sub> +dVdt, but with one more predictor: SHRD
MLR_Monthly	same as MLR, but with environmental predictors calculated from monthly NCEP-NCAR reanalysis fields

## 737 List of Figures

- 738 1 Root-mean-square-error (RMSE) as a function of number of predictors used  
739 in MLR for forecast lead time: (a) 12, (b) 24, (c) 48, (d) 72, (e) 96, and (f)  
740 120 hours. The predictors are listed on the right of each panel in the order  
741 that they are selected by the forward selection procedure. The order of the  
742 predictor selection is stable above the red line. 36
- 743 2 Tracks for hurricanes Rita (2005), Earl (2010), Irene (2011) and Isaac (2012).  
744 The intensity predictions for each storm are shown in Fig. 3. 37
- 745 3 Five-days intensity forecasts made every 12 hours by MLR (red) and SHIPS  
746 (blue) of (a) Rita, (b) Isaac, (c) Earl, and (d) Irene up to landfall. The black  
747 line shows the best-track intensity of each storm. 38
- 748 4 A sparkline chart of the standard deviation ( $\sigma$ ) of forecast errors for 24 hours  
749 forecasts. The  $\sigma$  is calculated with data that have been binned based on  
750 the magnitude of each individual predictor, and therefore the  $\sigma$  of errors is a  
751 function of the predictor here. The unit of the x-axis is one standard deviation  
752 of the individual predictor. 39
- 753 5 (a) Horizontal map of storm locations of all 12 hours forecasts, where the  
754 color indicates the intensity prediction errors from the MLR. (b) Horizontal  
755 map of the MLR errors from (a). Each tile represents a  $4^\circ \times 4^\circ$  box; the color  
756 indicates the mean error averaged over each tile. We show only the boxes  
757 with more than 5 data points. (c) and (d) are similar to (a) and (b), but for  
758 the 96 hours intensity prediction errors. 40
- 759 6 (a) Mean absolute error (MAE) of intensity prediction from the persistence  
760 (baseline, black-dotted line), SHIFOR (black-dashed line), NHC (OFCL, gray-  
761 solid line), SHIPS (black-solid line), and MLR forecasts. (c) and (e) are similar  
762 to (a), but for different sets of experiments, which are described in Table 1.  
763 (b), (d) and (f) are the same as (a), (c), and (e), but for RMSE. 41

764	7	(a) MAE of intensity prediction from the persistence (baseline, black-dotted line), SHIPS (black-solid line), MLR (MLR with daily data, red), and MLR_Monthly (MLR with monthly data, blue)	42
765			
766			
767	8	Scatter plots overlaid on two-dimensional histograms of daily (x-axis) and monthly (y-axis) environmental predictors used in MLR: (a) the difference between PI and initial maximum wind speed: $dPI_{V_0}$ , (b) the vertical wind shear: SHRD, (c) the 200 hPa divergence: $div_{200}$ , and (d) the conditional stability: $d\Theta_{Es}$ . The black-solid lines indicate the best-fit lines while the black-dashed lines are the one-to-one lines.	43
768			
769			
770			
771			
772			
773	9	(a) Probability density function (PDF) of the MLR forecast errors from all the training data. (b) Similar to (a) but for cumulative density function (CDF). In (a) and (b), colors indicate various forecast times as listed on the legend. (c) Quantile-Quantile plot showing the fit of the MLR errors for 24 hours forecast (from training data) to the normal distribution. The quantiles of the MLR errors are shown on y-axis while the quantiles of the normal distribution are on the x-axis. The unit of x-axis is one standard deviation of MLR errors.	44
774			
775			
776			
777			
778			
779			
780	10	Joint CDF (%) of the magnitude of errors and lead time from (a) MLR_1, and (b) the probabilistic persistence model based on training data. In both (b) and (a), the white-dashed lines indicate where the mean error is, and the black solid and dashed lines show where the highest 10% and 50% probabilities are.	45
781			
782			
783			
784	11	(a) Standard deviation ( $\sigma$ ) of error distribution for 24 hours intensity predictions from the training data. Black-dashed line is calculated with all the errors (MLR_1), while black-diamond line is calculated with errors that have binned based on the initial storm intensity ( $V_0$ ) (MLR_2). (b) Joint CDF (%) from MLR_1. Note that the CDF is not a function of $V_0$ . (c) Similar to (b), but now the CDF is is calculated from MLR_2, and is also a function of $V_0$ .	46
785			
786			
787			
788			
789			

- 790 12 (a) Intensity probabilistic prediction in CDF for Hurricane Earl (2010) based  
791 on the probabilistic persistence model (Fig. 10b). (b) Same as (a), but the  
792 CDF is based on MLR\_1 (Fig. 10a and Fig. 11b ). (c) Similar to (b), but the  
793 CDF is based on MLR\_2 (Fig. 11c). The initial time for (a-c) is 0000 UTC  
794 28 August, 2010. The white-solid lines indicate the observed maximum wind  
795 speed, while the white-dashed lines are the MLR and the persistence model  
796 predicted intensity, which is also where the mean error is in Fig. 10. Black  
797 solid and dashed lines show where the highest 10% and 50% probabilities  
798 are. (d) Ranked probability skill score (RPSS) of the two MLR probabilistic  
799 predictions relative to those from the probabilistic persistence model for all  
800 forecasts in Earl's life cycle. 47
- 801 13 Rank probability skill score (RPSS) of the two MLR probabilistic predictions  
802 relative to those from probabilistic persistence model for all testing cases. 48
- 803 14 Reliability diagrams(a, c, e) showing observed frequency as a function of  
804 forecast probabilities for (a) tropical storm (TS), (b) category 1-2 hurricane,  
805 and (c) category 3-5 hurricane, respectively. Colors indicates forecasts from  
806 MLR\_1(pink), MLR\_2 (red), persistence (blue) and NHC (OFCL, gray) from  
807 2008 - 2012. (b), (d), and (f) show the sample size used in (a), (c), (e) for  
808 each forecast. 49

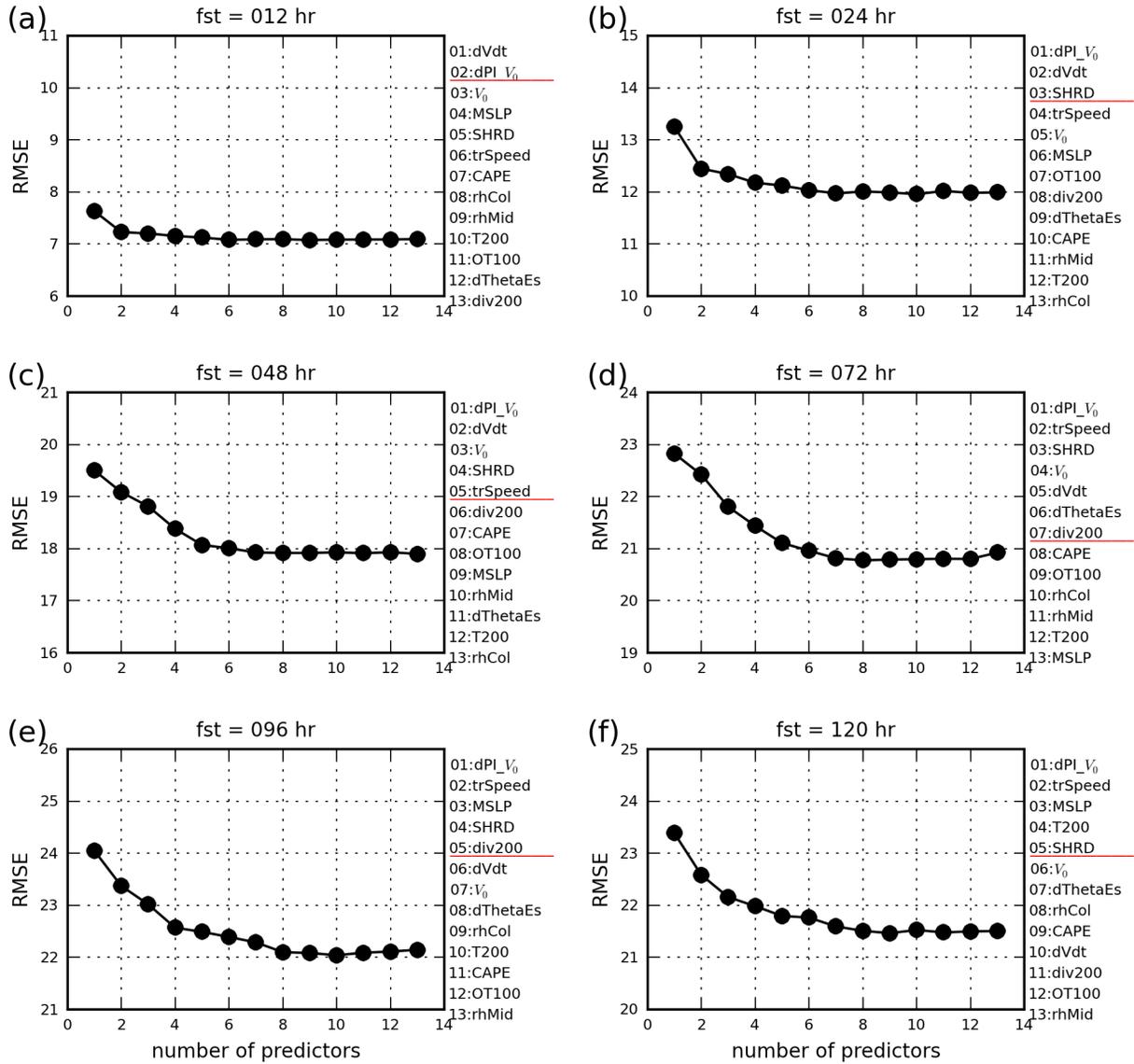


FIG. 1. Root-mean-square-error (RMSE) as a function of number of predictors used in MLR for forecast lead time: (a) 12, (b) 24, (c) 48, (d) 72, (e) 96, and (f) 120 hours. The predictors are listed on the right of each panel in the order that they are selected by the forward selection procedure. The order of the predictor selection is stable above the red line.

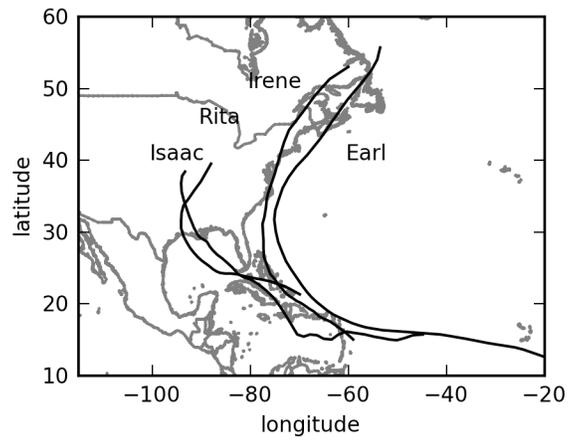


FIG. 2. Tracks for hurricanes Rita (2005), Earl (2010), Irene (2011) and Isaac (2012). The intensity predictions for each storm are shown in Fig. 3.

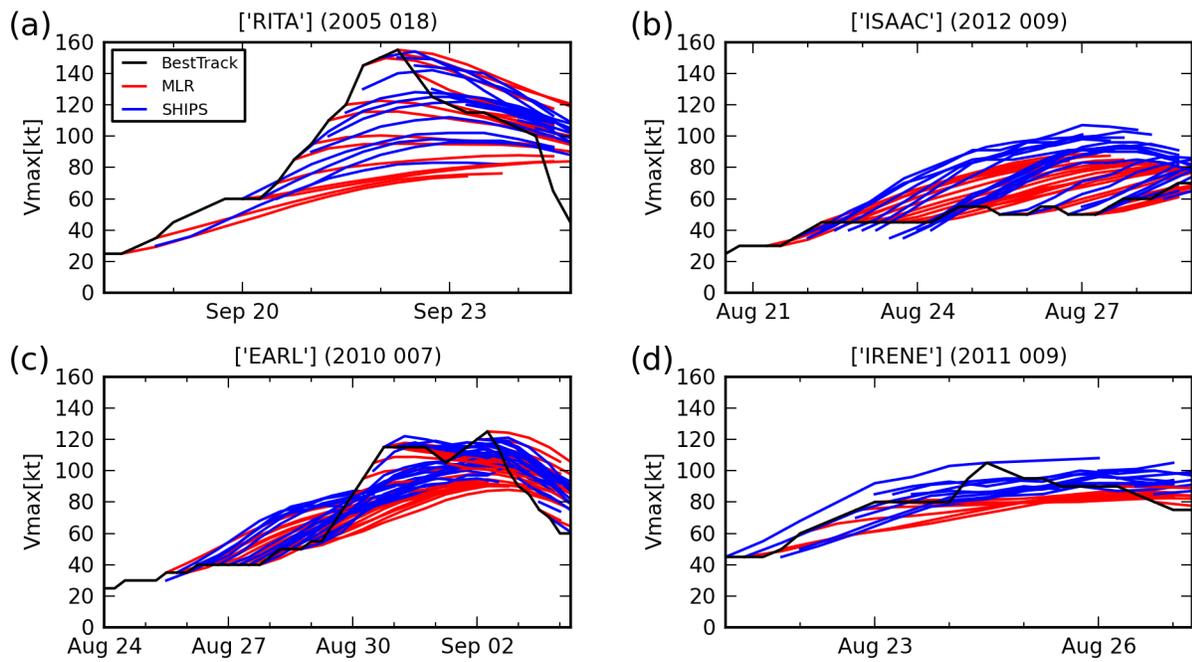


FIG. 3. Five-days intensity forecasts made every 12 hours by MLR (red) and SHIPS (blue) of (a) Rita, (b) Isaac, (c) Earl, and (d) Irene up to landfall. The black line shows the best-track intensity of each storm.

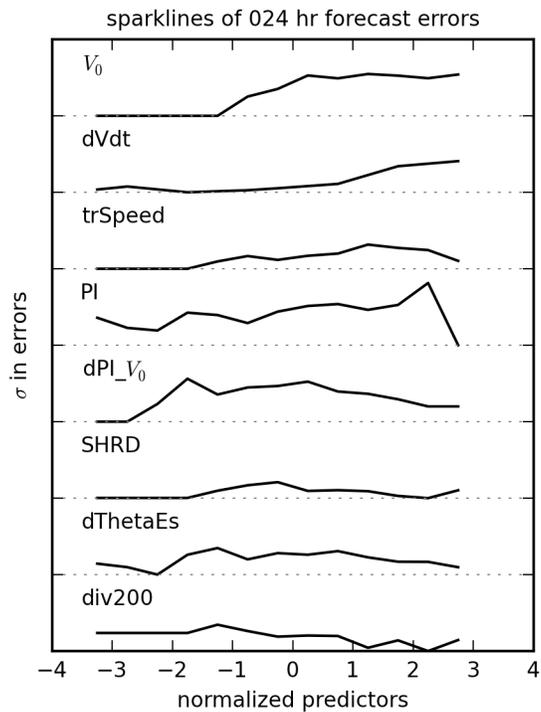


FIG. 4. A sparkline chart of the standard deviation ( $\sigma$ ) of forecast errors for 24 hours forecasts. The  $\sigma$  is calculated with data that have been binned based on the magnitude of each individual predictor, and therefore the  $\sigma$  of errors is a function of the predictor here. The unit of the x-axis is one standard deviation of the individual predictor.

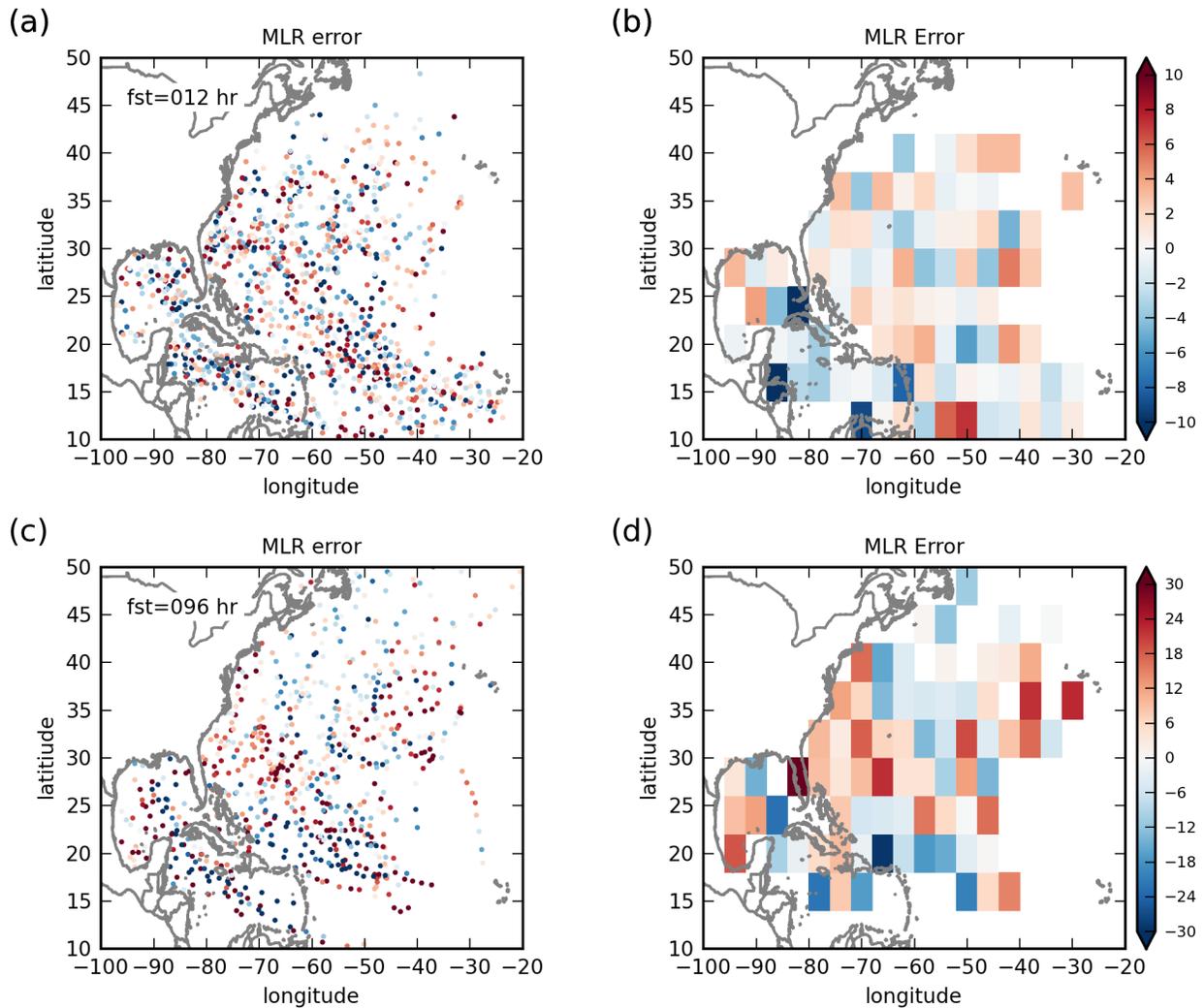


FIG. 5. (a) Horizontal map of storm locations of all 12 hours forecasts, where the color indicates the intensity prediction errors from the MLR. (b) Horizontal map of the MLR errors from (a). Each tile represents a  $4^\circ \times 4^\circ$  box; the color indicates the mean error averaged over each tile. We show only the boxes with more than 5 data points. (c) and (d) are similar to (a) and (b), but for the 96 hours intensity prediction errors.

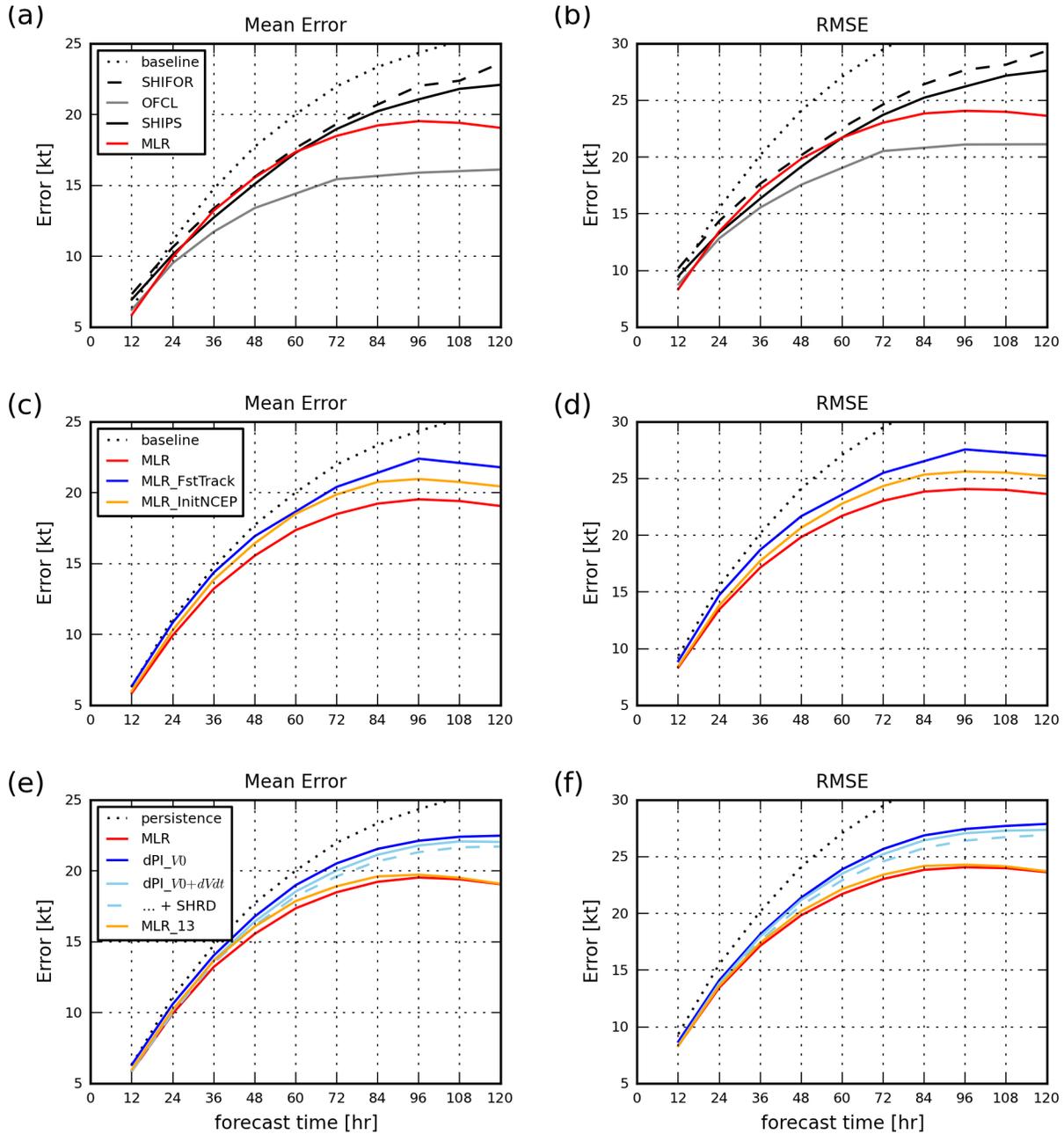


FIG. 6. (a) Mean absolute error (MAE) of intensity prediction from the persistence (baseline, black-dotted line), SHIFOR (black-dashed line), NHC (OFCL, gray-solid line), SHIPS (black-solid line), and MLR forecasts. (c) and (e) are similar to (a), but for different sets of experiments, which are described in Table 1. (b), (d) and (f) are the same as (a), (c), and (e), but for RMSE.

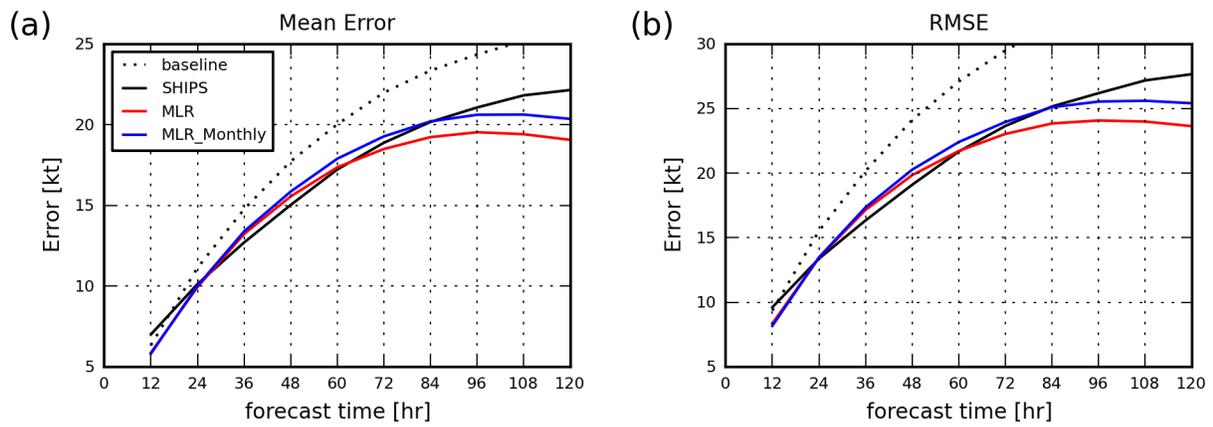


FIG. 7. (a) MAE of intensity prediction from the persistence (baseline, black-dotted line), SHIPS (black-solid line), MLR (MLR with daily data, red), and MLR\_Monthly (MLR with monthly data, blue)

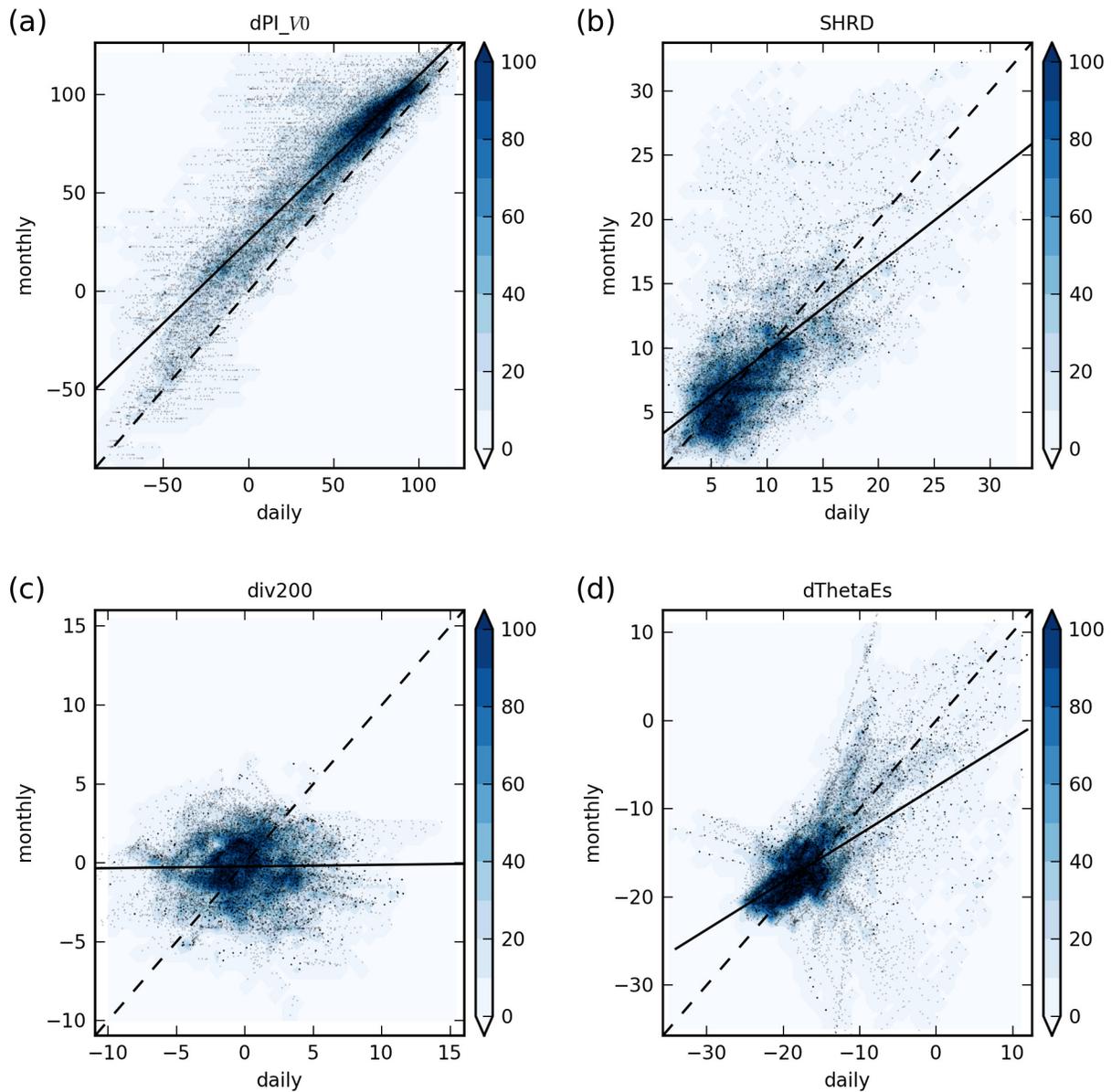


FIG. 8. Scatter plots overlaid on two-dimensional histograms of daily (x-axis) and monthly (y-axis) environmental predictors used in MLR: (a) the difference between PI and initial maximum wind speed:  $dPI_{V_0}$ , (b) the vertical wind shear: SHRD, (c) the 200 hPa divergence:  $div200$ , and (d) the conditional stability:  $d\Theta_{Es}$ . The black-solid lines indicate the best-fit lines while the black-dashed lines are the one-to-one lines.

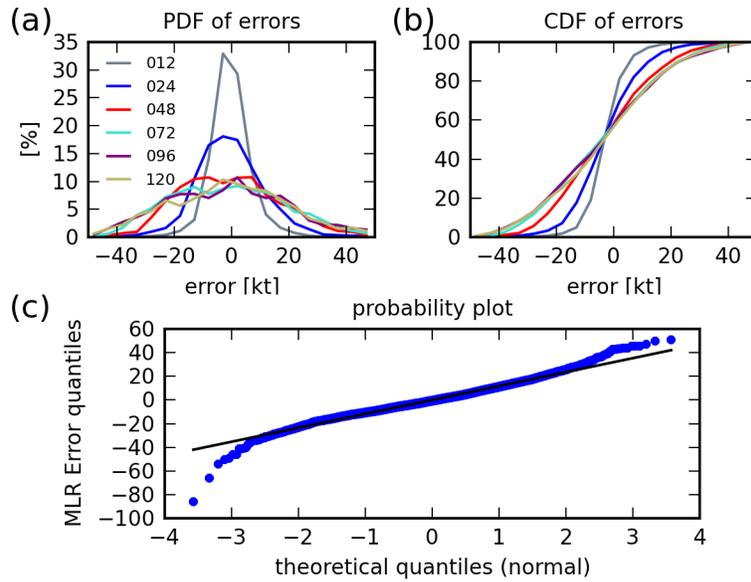


FIG. 9. (a) Probability density function (PDF) of the MLR forecast errors from all the training data. (b) Similar to (a) but for cumulative density function (CDF). In (a) and (b), colors indicate various forecast times as listed on the legend. (c) Quantile-Quantile plot showing the fit of the MLR errors for 24 hours forecast (from training data) to the normal distribution. The quantiles of the MLR errors are shown on y-axis while the quantiles of the normal distribution are on the x-axis. The unit of x-axis is one standard deviation of MLR errors.

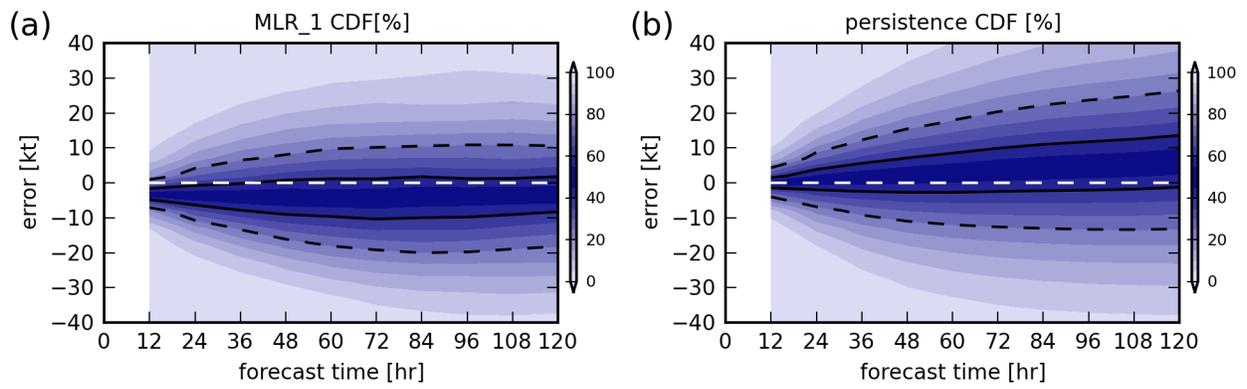


FIG. 10. Joint CDF (%) of the magnitude of errors and lead time from (a) MLR\_1, and (b) the probabilistic persistence model based on training data. In both (b) and (a), the white-dashed lines indicate where the mean error is, and the black solid and dashed lines show where the highest 10% and 50% probabilities are.

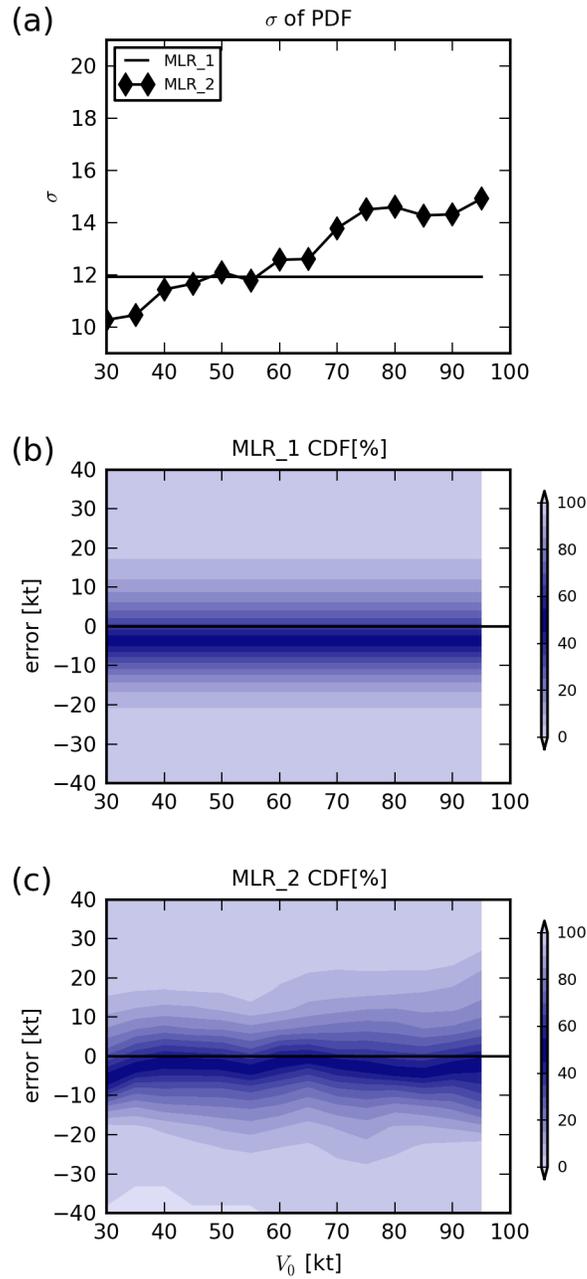


FIG. 11. (a) Standard deviation ( $\sigma$ ) of error distribution for 24 hours intensity predictions from the training data. Black-dashed line is calculated with all the errors (MLR\_1), while black-diamond line is calculated with errors that have binned based on the initial storm intensity ( $V_0$ ) (MLR\_2). (b) Joint CDF (%) from MLR\_1. Note that the CDF is not a function of  $V_0$ . (c) Similar to (b), but now the CDF is calculated from MLR\_2, and is also a function of  $V_0$ .

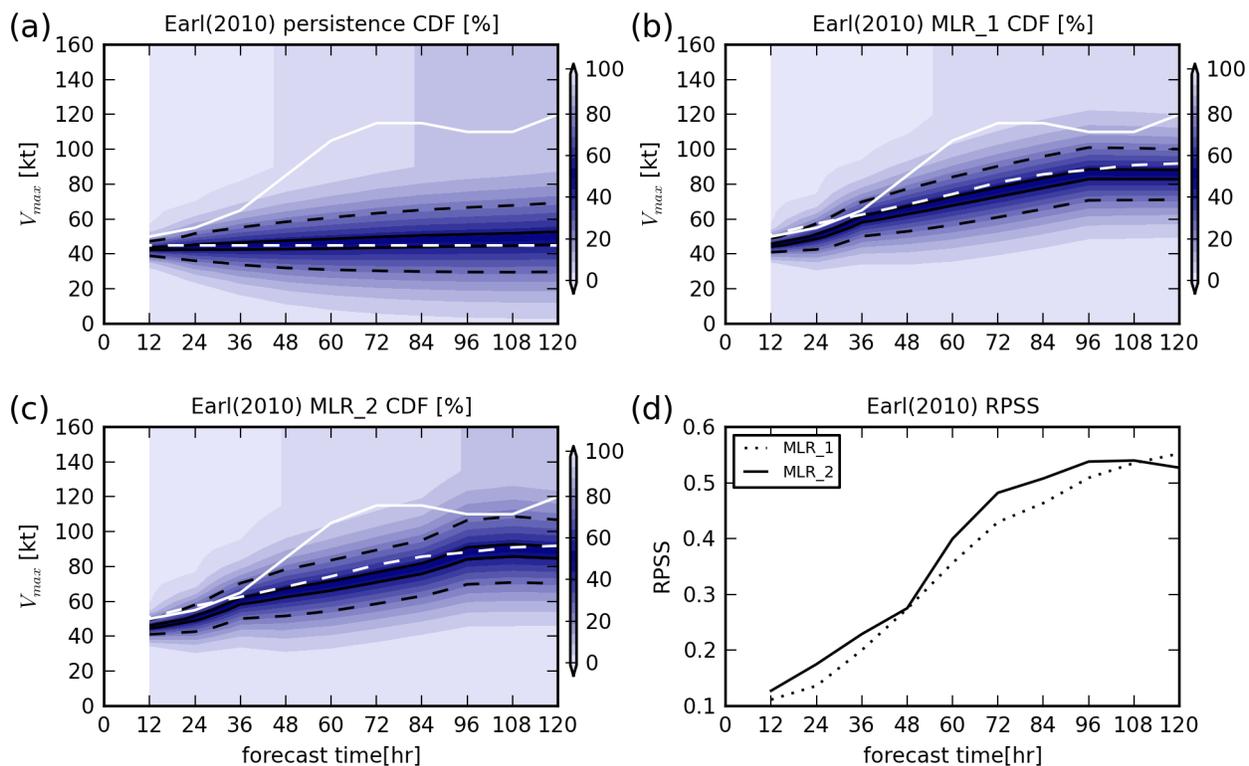


FIG. 12. (a) Intensity probabilistic prediction in CDF for Hurricane Earl (2010) based on the probabilistic persistence model (Fig. 10b). (b) Same as (a), but the CDF is based on MLR\_1 (Fig. 10a and Fig. 11b). (c) Similar to (b), but the CDF is based on MLR\_2 (Fig. 11c). The initial time for (a-c) is 0000 UTC 28 August, 2010. The white-solid lines indicate the observed maximum wind speed, while the white-dashed lines are the MLR and the persistence model predicted intensity, which is also where the mean error is in Fig. 10. Black solid and dashed lines show where the highest 10% and 50% probabilities are. (d) Ranked probability skill score (RPSS) of the two MLR probabilistic predictions relative to those from the probabilistic persistence model for all forecasts in Earl's life cycle.

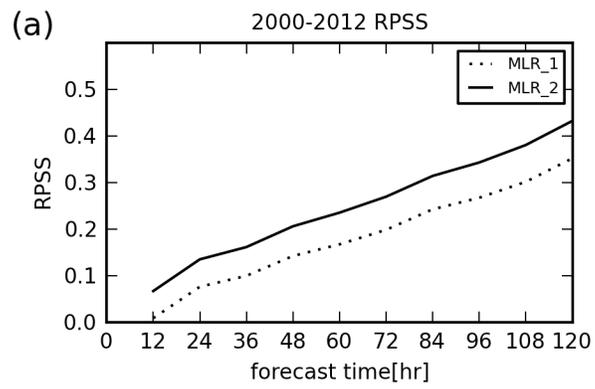


FIG. 13. Rank probability skill score (RPSS) of the two MLR probabilistic predictions relative to those from probabilistic persistence model for all testing cases.

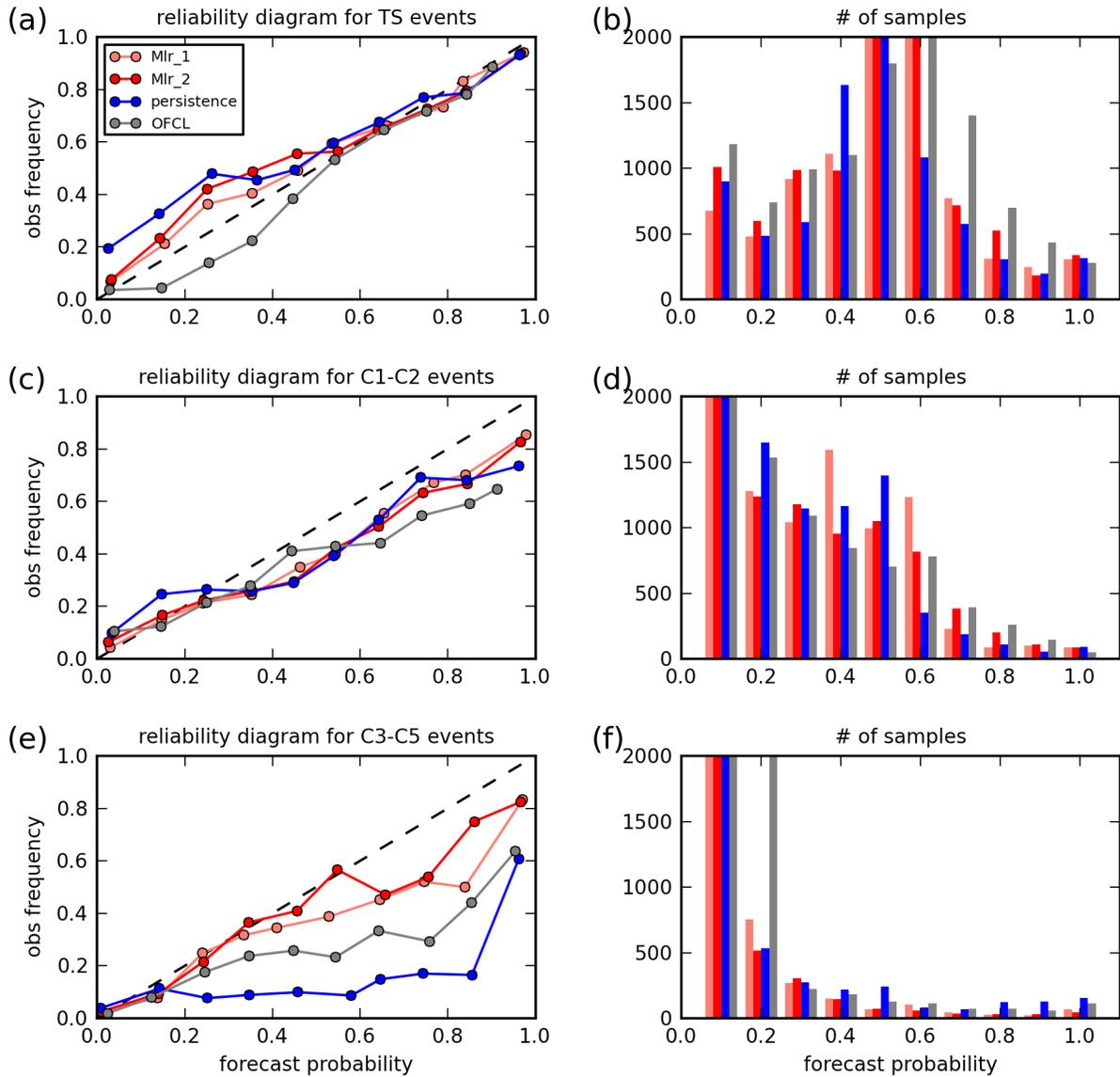


FIG. 14. Reliability diagrams(a, c, e) showing observed frequency as a function of forecast probabilities for (a) tropical storm (TS), (b) category 1-2 hurricane, and (c) category 3-5 hurricane, respectively. Colors indicates forecasts from MLR\_1(pink), MLR\_2 (red), persistence (blue) and NHC (OFCL, gray) from 2008 - 2012. (b), (d), and (f) show the sample size used in (a), (c), (e) for each forecast.