



# AMERICAN METEOROLOGICAL SOCIETY

*Monthly Weather Review*

## **EARLY ONLINE RELEASE**

This is a preliminary PDF of the author-produced manuscript that has been peer-reviewed and accepted for publication. Since it is being posted so soon after acceptance, it has not yet been copyedited, formatted, or processed by AMS Publications. This preliminary version of the manuscript may be downloaded, distributed, and cited, but please be aware that there will be visual differences and possibly some content differences between this version and the final published version.

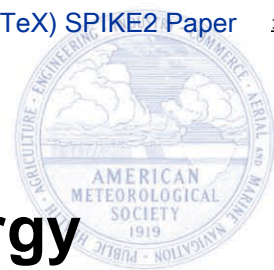
The DOI for this manuscript is doi: 10.1175/MWR-D-16-0030.1

The final published version of this manuscript will replace the preliminary version at the above DOI once it is available.

If you would like to cite this EOR in a separate work, please use the following full citation:

Kozar, M., V. Misra, and M. Powell, 2016: Hindcasts of Integrated Kinetic Energy in Atlantic Tropical Cyclones: A Neural Network Prediction Scheme. *Mon. Wea. Rev.* doi:10.1175/MWR-D-16-0030.1, in press.

© 2016 American Meteorological Society



# Hindcasts of Integrated Kinetic Energy in Atlantic Tropical Cyclones: A Neural Network Prediction Scheme

Michael E. Kozar<sup>1</sup>

Risk Management Solutions  
Tallahassee, FL 32304

Vasubandhu Misra

Center for Ocean-Atmospheric Prediction Studies,  
Department of Earth, Ocean and Atmospheric Science,  
Florida State University,  
Tallahassee, FL 32306

Mark D. Powell

Risk Management Solutions  
Tallahassee, FL 32304

---

<sup>1</sup> Corresponding author address: Michael Kozar; RMS Tallahassee; 612 Copeland St.; Tallahassee, FL 32304  
Email: Michael.Kozar@rms.com

Abstract

A new statistical-dynamical scheme is presented for predicting Integrated Kinetic Energy (IKE) in North Atlantic tropical cyclones from a series of environmental input parameters. Predicting IKE is desirable because the metric quantifies the energy across a storm's entire wind field, allowing it to respond to changes in storm structure and size. As such, IKE is especially useful for quantifying risks in large low-intensity high-impact storms such as Sandy in 2012. The prediction scheme, named the Statistical Prediction of Integrated Kinetic Energy Version 2 (SPIKE2), builds upon a previous statistical IKE scheme, by using a series of artificial neural networks instead of more basic linear regression models. By using a more complex statistical scheme, SPIKE2 is able to distinguish non-linear signals in the environment that could cause fluctuations in IKE. In an effort to evaluate SPIKE2's performance in a future operational setting, the model is calibrated using archived input parameters from Global Ensemble Forecast System (GEFS) control analyses, and is run in a hindcast mode from 1990 to 2011 using archived GEFS reforecasts. The hindcast results indicate that SPIKE2 performs significantly better than both persistence and climatological benchmarks.

52 **1. Introduction**

53 Integrated kinetic energy (IKE) is a recently developed metric that is designed to  
54 approximate the damage potential of landfalling tropical cyclones (Powell and Reinhold 2007).  
55 As its name suggests, IKE is defined as a summation of the kinetic energy within the near-  
56 surface wind field of a tropical cyclone (TC). By integrating energy across a large portion of a  
57 storm's wind field, IKE considers the overall structure of a TC. This is in stark contrast to many  
58 other existing hurricane metrics, which often quantify only a wind or pressure extreme at a  
59 single point within a TC. Intensity metrics such as maximum sustained wind speeds (VMAX) are  
60 undoubtedly useful for assessing the maximum potential damage caused by the winds in a TC  
61 (e.g. Emanuel 2005; Bell et al. 2000), but they do not paint a complete picture of storm damage  
62 potential.

63 In the decade following the landfall of Hurricane Wilma, no major hurricanes (VMAX >  
64 96kts) have made landfall in the United States. This drought is thought to be a rather rare event,  
65 (Hall and Hereid 2015), depending on the metric that is used to classify major hurricanes (Hart  
66 et al. 2015). Despite this perceived quiet period of significant United States hurricane activity,  
67 there has been no shortage of damaging storms that have made landfall in the past decade.  
68 According to initial estimates from the National Hurricane Center<sup>1</sup>, hurricanes Ike (AL092008),  
69 Irene (AL092011), and Sandy (AL182012) each caused more than \$15 million in losses across  
70 the United States during the major hurricane drought despite each storm's somewhat weak  
71 landfall intensity. This disconnect between VMAX and damage often occurs because storm size  
72 and structure must also be considered to properly evaluate storm surge potential (e.g. Irish et al.  
73 2008). Since Sandy, Irene, and Ike were such large storms, they were able to produce higher  
74 storm surge and damage totals than otherwise would be expected by storms of similar  
75 intensities. For this reason, it is likely that the IKE metric could add value to existing intensity

---

<sup>1</sup> <http://www.nhc.noaa.gov/data/tcr/>

76 metrics, by anticipating the higher damage potential of larger landfalling TCs (Powell and  
77 Reinhold 2007), especially considering that Ike, Sandy, and Irene all ranked very highly in terms  
78 of IKE relative to other storms in the historical record (Kozar and Misra 2014).

79         Despite the potential advantages of IKE, the concept of forecasting the energy metric in  
80 real-time is still in its infancy. Currently, operational forecasters have little to no guidance to  
81 predict IKE. Recently, Kozar and Misra (2014; hereafter KM14) explored whether or not it is  
82 feasible to fill that void with a simple statistical model in a proof-of-concept exercise. The  
83 resulting statistical model from that study was named the Statistical Prediction of Integrated  
84 Kinetic Energy (SPIKE), and it used linear regression to predict changes of IKE from a series of  
85 environmental predictors. Despite its simplicity, SPIKE was ultimately capable of outperforming  
86 a persistence forecast in a perfect-prognostic mode, indicating that statistical-dynamical  
87 forecasts of IKE might be possible in the future.

88         Building upon those results, the focus of this study is to further evaluate the operational  
89 potential of IKE forecasts using a more sophisticated statistical-dynamical scheme in a hindcast  
90 mode. Despite the successes of the proof-of-concept SPIKE model from KM14, linear  
91 regression is suboptimal for statistical weather prediction because the earth system is quite  
92 complex and contains several nonlinear signals. As such, the fixed linear regression coefficients  
93 in the SPIKE model will never be able to fully process the complex changing relationships  
94 between the environment and IKE variability within a TC. Therefore, a second-generation  
95 version of SPIKE is developed in this work by utilizing a more complex and nonlinear statistical  
96 framework in lieu of linear regression. More specifically, SPIKE version 2 (SPIKE2) utilizes a  
97 series of artificial neural networks (ANNs) to predict IKE tendency from a similar series of  
98 environmental input parameters. Ultimately, these networks are capable of learning and  
99 anticipating complex patterns in the environment, and as a result they are better suited to model  
100 a nonlinear system.

101           Furthermore, SPIKE2’s evaluation will be moved from a perfect-prognostic space  
102 previously used in the initial KM14 work to a hindcast space. Obviously, in an operational  
103 setting, a statistical-dynamical forecast scheme must contend with imperfect input parameters  
104 that contain forecast errors of increasing magnitude with increasing lead-time. Therefore, by  
105 running SPIKE2 in a hindcast mode with model data from the National Oceanic and  
106 Atmospheric Administration (NOAA)’s Second Generation Global Ensemble Reforecast Archive  
107 (Hamill et al. 2013), we are able to more comprehensively measure the potential performance of  
108 the IKE prediction models.

109           The next section discusses the historical and reforecast data that are used to calibrate  
110 and evaluate the SPIKE2 neural network system. In the subsequent sections, the discussion  
111 shifts towards the methodology and procedures for creating, calibrating, and evaluating the  
112 neural networks used in SPIKE2. Finally, the calibration and hindcast performance of SPIKE2 is  
113 compared against various persistence and climatology benchmarks for Atlantic TCs between  
114 1990 and 2011 in the penultimate section, preceding the concluding remarks.

115

## 116 **2. Historical and Model Reforecast Data**

117           Similar to KM14, a historical record of IKE in North Atlantic TCs is used to train and  
118 validate the SPIKE2 neural networks. This historical record covers the 1990 through 2011  
119 hurricane seasons, and includes over 5000 six-hourly fixes from nearly 300 individual storms  
120 (Misra et al. 2013). Since gridded wind fields are not available for all of these cases, the IKE  
121 values contained in this record are all estimated from operational wind radii and intensity metrics  
122 in the extended best track dataset (Demuth et al. 2006) using a series of equations from Powell  
123 and Reinhold (2007) and Misra et al. (2013). The mean value of IKE across all of the storm fixes  
124 included in our historical archive is 35 TJs, with a standard deviation of 43 TJ. The distribution of  
125 observed IKE values takes a somewhat lognormal shape with a long tail towards higher values

126 (KM14). As such, although most storms never reach 50 TJs of IKE, Hurricane Sandy likely had  
127 more than 400 TJ of IKE before it made landfall in New Jersey in 2012.

128           It should also be noted that past works have documented significant uncertainty within  
129 the historical record of wind radii that fluctuates depending on the data platforms that are  
130 available when analyzing each storm (e.g. Knaff et al. 2014; Landsea and Franklin, 2013).  
131 Therefore our historical IKE record, which again is based on the operation wind radii, likely  
132 inherits many of the same uncertainties found in the extended best track dataset.

133           Unlike KM14, SPIKE2's environmental input variables are drawn from a historical model  
134 reforecast database during the same 1990 to 2011 interval. The second generation Global  
135 Ensemble Reforecast Archive (Hamill et al. 2013) is selected as the source for this model data  
136 because it includes model runs dating back multiple decades using a static February 2012  
137 operational configuration (version 9.0.1) of National Centers for Environmental Prediction's  
138 (NCEP) Global Ensemble Forecast System (GEFS).

139           These archived GEFS reforecasts include forecasts out to 16 days beyond the  
140 initialization time. The first eight days of the forecast are run at T254 horizontal resolution  
141 (~50km) with 42 vertical levels. The latter half of the forecast is run at a lower T190 resolution  
142 (~70km), with the same 42 vertical layers. Each of the reforecast runs is initialized once daily at  
143 00Z, as opposed to the six-hourly approach for generating operational GEFS forecasts. The  
144 initial conditions for the reforecast dataset are produced from the Climate Forecast System  
145 Reanalysis (CFSR; Saha et al. 2010) prior to February 2011 and operational Grid-Point  
146 Statistical Interpolation analyses after that time.

147           The reforecast archive includes ninety-eight output fields from initial time out to F+384  
148 hrs for each of the daily GEFS reforecasts<sup>2</sup>. The archived data is stored at 3-hourly intervals for  
149 the first 72 hours of the forecast and then at 6-hourly intervals after that time. Each  
150 meteorological field is bilinearly interpolated down to a somewhat coarse one-degree resolution

---

<sup>2</sup> <http://www.esrl.noaa.gov/psd/forecasts/reforecast2/>

151 global grid. In addition to the one-degree data sets, a smaller selection of twenty-eight fields is  
152 also stored in the GEFS's higher resolution native Gaussian grid ( $\sim 0.5^\circ$ ). However, the higher  
153 resolution fields are all single-level variables, primarily near the surface. As a result of this  
154 limitation, mid- and upper-atmospheric dynamic and thermodynamic fields (winds,  
155 temperatures, humidity, etc) are only available in the one-degree grids. Therefore, to maximize  
156 consistency, we use only the one-degree data to examine the dynamical and thermodynamical  
157 processes that relate to IKE variability for our SPIKE2 neural network system.

158         As its name suggests, the GEFS archive does not just include a single deterministic  
159 forecast. In fact, the reforecast dataset is comprised of eleven ensemble members (1 control  
160 run, and 10 perturbation runs) compared to 21 ensemble members in the operational GEFS. For  
161 the purposes of this work, only the control run in the GEFS reforecast set is considered, but  
162 future works can and should clearly expand upon these results to produce probabilistic  
163 forecasts that resolve the uncertainty in the model's initial environment.

164         This GEFS reforecast dataset includes some noteworthy biases with regards to  
165 resolving TCs that will be addressed here. Obviously, the GEFS reforecasts will include position  
166 and intensity errors, and the reforecasted environment is expected to be imperfect as well, with  
167 all errors increasing as lead time increases. For reference, Galarneau Jr. and Hamill (2015)  
168 analyzed track errors in the GEFS reforecast archive for TCs in the Gulf of Mexico between  
169 1985 and 2010 and found average positions errors to be 100km with a lead time of 24 hours,  
170 250km with a lead time of 72 hours, and 400km with a 120 hour forecast interval. Typically,  
171 these track reforecasts in the Gulf of Mexico were found to have a left and slow bias relative to  
172 the storms motion.

173         Furthermore, Galarneau Jr. and Hamill (2015) also indicated that the GEFS reforecasts  
174 had a significant and consistent low intensity bias. This does not come as a surprise and leads  
175 us to the potentially most concerning issue for using the GEFS reforecast database in this  
176 study. Simply put, the one-degree horizontal resolution data taken from the model will not be

177 sufficient to properly resolve the wind field of a TC. As a result, intensities will be  
178 underestimated, and wind fields may be too broad. In fact, the GEFS reforecasts may fail to  
179 generate a TC vortex altogether in some extreme scenarios.

180         However, since we are not trying to predict IKE directly from the model's wind field, but  
181 instead by relating environmental parameters to IKE variability, this low resolution data might  
182 still be sufficient, albeit less than ideal. By using the lower resolution GEFS reforecast data, we  
183 can estimate the lower bounds of skill for a real-time version of SPIKE2. Furthermore, the static  
184 model configuration provided by the GEFS reforecast dataset (Hamill et al. 2013) allows us to  
185 focus on the performance of the SPIKE2 predictive scheme independent of changes in the  
186 underlying dynamical model's configuration and performance.

187

### 188 **3. Selection of Input Parameters For Statistical-Dynamical Prediction**

189         Before the neural networks can be constructed, we must first establish which  
190 environmental and storm specific input parameters will be taken from the GEFS control  
191 reforecasts to produce predictions of IKE variability. The initial SPIKE model built in KM14  
192 utilized a series of 14 predictors that contained a significant linear relationship with IKE  
193 variability, many of which were taken directly from the Statistical Hurricane Intensity Prediction  
194 Scheme (SHIPS)'s developmental dataset (DeMaria and Kaplan 1999). These input parameters  
195 included various environmental predictors (both dynamical and thermodynamical), storm-  
196 specific parameters (e.g. position, minimum pressure), and persistence values of IKE based on  
197 known relationships between IKE and the environment (e.g. Maclay et al. 2008; Musgrave et al.  
198 2012). However, since these parameters were selected based on their linear relationships with  
199 IKE, it is necessary to reselect predictors to highlight the nonlinearities in the storm-environment  
200 system that hopefully can be captured by the more sophisticated neural network scheme utilized  
201 here for SPIKE2.

202 As was done in KM14, the goal in selecting these parameters should be to target  
203 physical processes that govern variability with a TC's structure and ultimately the IKE index.  
204 Therefore, we started with a large pool of predictors, including both predictors used in the linear  
205 model that had clear and justifiable relationships with IKE as well as control variables such as  
206 day of month. Properly tuning a nonlinear complex neural network is a bit more difficult than  
207 tuning a linear regression model as there are more weights and neurons than there are  
208 coefficients in a linear regression model. Nonetheless, as we constructed the neural network we  
209 removed predictors if network performance over the testing sample increased by subtracting the  
210 predictor. As such, each of the control parameters and a few of the other environmental  
211 predictors with weaker ties to IKE were not selected for the final version.

212 Ultimately, we settled on 18 input parameters for SPIKE2, each of which is related to  
213 targeted relationships between the environment and IKE, in order to maximize the neural  
214 networks potential predictive power. The specific predictors are listed in Table 1. From this point  
215 forward, each predictor will be referred to by its abbreviation in the table. This predictor list is  
216 very similar to those used in the linear SPIKE model, but does include a total of four additional  
217 predictors. As such, we acknowledge that a few of these predictors could be removed, and the  
218 performance of SPIKE2 would likely not change by an appreciable margin. However, removal of  
219 any of the predictors did not seem to improve validation performance, suggesting that the  
220 predictors were not setting the model back via overfitting. Therefore, we felt that by including  
221 some of these extra predictors, the neural network may have a better chance to resolve some of  
222 the nonlinear signals between the environment and TCs if we were careful to limit the number of  
223 neurons in the ANN, thus minimizing the chances of overfitting.

224 Nonetheless, we ran a series of perturbation tests and case studies to ensure that each  
225 individual variable had some physical relationship that could explain how it is affecting  
226 projections of IKE from SPIKE2. For brevity, the remaining discussion in this section is meant to  
227 highlight the physical relationships that can explain how each of the individual predictors affect

228 IKE variability, followed by a short explanation about how the predictors are directly calculated  
229 from the model fields.

230 Predictors such as D200, VORT, SHTD and SHRD are designed to represent the certain  
231 dynamical features (upper level divergence, low-level vorticity, weak easterly shear) that are  
232 favorable for TC development. These predictors were some of the more significant predictors in  
233 the linear regression SPIKE model, and their impact over SPIKE2's IKE projections remains  
234 strong. Meanwhile, SST, T150 and VMPI are meant to be tied to thermodynamical properties  
235 that govern the maximum intensity of the storm, the height of the tropopause, and how far a  
236 storm has to go before it reaches said maximum intensity (e.g. Emanuel 1988; Bister and  
237 Emanuel 1998). RHLO and RHMD capture well known relationships between moisture and TC  
238 development. MSLP, PENV, and VMAX are storm specific parameters that give some  
239 information about the TC's intensity and breadth at the validation time, wherein a more intense  
240 storm or a larger storm with all else being equal will have higher wind speeds and more IKE.  
241 LAT, LON, SDAY, and PDAY obviously give information about the storm's position and time.  
242 These can be useful for identifying climatological tendencies across the basin. Finally,  
243 predictors such as PIKE and dIKE12 give information about persistence (i.e. how much IKE the  
244 storm had previously, and was it gaining or losing IKE previously) that can be useful for  
245 predicting future trends in certain instances.

246 However, as alluded to in the opening section, the signals between IKE and these  
247 predictors are quite complex. Unlike traditional storm development, which has a somewhat  
248 straightforward relationship with some of these predictors (i.e. the combination of low shear and  
249 high SSTs typical translates to a stronger storm all else being equal), IKE is also tied to storm  
250 size and the many different processes that govern it. For instance, many storms tend to expand  
251 as they move poleward and interact with other baroclinic features or through extratropical  
252 transition (e.g. Evans and Hart 2008). As such, recurving TCs often gain IKE in mid-latitude

253 environments that would traditionally be considered non-favorable for development (Maclay et  
254 al. 2008).

255           Considering that extratropical transition occurs in just under half of all Atlantic TCs (Hart  
256 and Evans 2001), our prediction scheme must be calibrated to anticipate the correct IKE  
257 tendencies from these complex signals. As a result, the nonlinear equations within the ANNs will  
258 also use predictors such as LAT, SHRD, T150, RHLO, and SST to determine whether or not a  
259 storm is likely to expand in size (and also in IKE) from baroclinic forcings. Encouragingly, some  
260 simple case studies revealed that a hypothetical storm in the mid latitudes (high LAT), late in its  
261 lifecycle (high SDAY) will actually gain IKE as expected in a more baroclinic environment with  
262 lower SSTs and higher SHRD. However, if the storm is under a similar environment in the deep  
263 tropics or if shear and SSTs are too prohibitive in the mid latitudes, the neural networks will  
264 correctly identify that the storm is more likely to decay. Ultimately, by considering both baroclinic  
265 influences and traditional developmental mechanisms from this wide-ranging predictor base  
266 through a nonlinear system of equations, the ANNs should be able to improve upon the results  
267 of KM14.

268           The majority of the predictors discussed above (LAT, LON, MSLP, VORT, D200, etc) are  
269 calculated directly from the corresponding TC signature within 3-D atmospheric fields from the  
270 GEFS's control run. However, it should be noted that the GEFS dataset by itself is insufficient to  
271 calculate all eighteen of the input parameters. For instance, some of the input parameters  
272 require information about the ocean surface (VMPI, SST), time and date of year (SDAY, PDAY)  
273 and past values of IKE (PIKE, dIKE12). Therefore, to obtain hindcasts for each of the input  
274 parameters the GEFS reforecast dataset will be supplemented with a number of other datasets.  
275 Daily one-degree NOAA Optimum Interpolation SST ("OI SST"; Reynolds et al. 2007) is used to  
276 estimate observed ocean surface conditions. The historical IKE record (derived from the  
277 extended best track dataset) is used to produce the persistence parameters, and finally the  
278 NHC best track dataset is used to get the time information for each storm fix.

279           Once the input parameters are calculated from the GEFS control run for all forecast  
280 hours between initial time and T+72 hr, each parameter is normalized by its sample within the  
281 GEFS control run for all storms between 1990 and 2011. Normalizing the input parameters  
282 offers the benefit of filtering out some of the systematic biases in the GEFS, which in turn should  
283 enhance the performance of the operational IKE prediction schemes.

284

## 285 **4. Setup of Artificial Neural Network for SPIKE2**

286           With the predictors and data sources now established, this section details how the  
287 artificial neural networks are constructed, calibrated, and then run in a hindcast mode. As  
288 highlighted earlier, ANNs are chosen for SPIKE2 because of their ability to resolve and adapt to  
289 changing nonlinear signals in a certain system (e.g. Kriesel 2007). Thanks in part to their  
290 versatility, ANNs have been used in meteorology over the past several years to complete a wide  
291 array of tasks. A non-exhaustive list of tasks that ANNs have been used for includes evaluating  
292 uncertainty in hurricane wind analyses (DiNapoli et al. 2012), processing remotely sensed data  
293 (e.g. Atkinson et al. 2010), classifying circulation patterns (e.g. Cawley and Dorling 2005),  
294 predicting troposphere ozone levels (e.g. Abdul-Wahab and Al-Alawi 2002), forecasting wind  
295 speeds (Cao et al. 2012), forecasting precipitation and flooding (e.g. Hapuarachchi et al. 2011),  
296 and predicting the strength of the Indian monsoon on a seasonal scale (Shukla et al. 2011). A  
297 more detailed summary of earlier ANN applications in meteorology can be found in a review by  
298 Gardner and Dorling (1998).

299

### 300 *4.1 Network Hierarchy and Algorithms*

301           The SPIKE2 prediction scheme will be built using a system of multiple two-layer feed-  
302 forward ANNs. Our two-layer feed-forward networks' hierarchy includes a hidden layer with  
303 twenty artificial neurons and an output layer with a single neuron that will ultimately produce the

304 desired results from the input parameters. Twenty neurons were chosen for the hidden layer to  
305 maximize predictive skill based on the results of an exhaustive search test, in which we found  
306 that this number of neurons corresponded to the best validation performance over the test  
307 subsample. By testing model performance with a wide varying number of neurons in this  
308 exhaustive search, we were able to find the approximate point at which ANN complexity is small  
309 enough to minimize the chance of overfitting, without compromising its ability to recognize and  
310 generalize the nonlinear signals in the TC-environment system.

311 In our case, the output of the neural networks will be IKE tendency for a given forecast  
312 hour, or in other words the difference between IKE at validation time and IKE at initialization  
313 time. Meanwhile, the eighteen normalized parameters discussed in Section 3 are selected as  
314 the input parameters of the neural network. As such, the goal of each ANN is to produce an  
315 estimate of IKE tendency from environmental and storm specific values within a model solution.

316 Ultimately, each of these ANNs within the SPIKE2 scheme are trained using a shared  
317 learning algorithm, wherein the networks are calibrated using a set of input parameters and  
318 known target (IKE tendency). The weights of the network's neurons are designed to adapt from  
319 a somewhat random initial value to a more optimal value, as the error function reaches a  
320 minimum. More specifically, the learning algorithm uses a Levenberg-Marquardt  
321 backpropagation algorithm (Marquardt 1963) to find this error minimum. This specific algorithm  
322 is designed to solve non-linear least squared problems and is typically thought to be an efficient  
323 and stable method for converging at an optimal solution in neural network learning (e.g. Hagan  
324 and Menhaj 1994).

325

## 326 *4.2 Training, Validation, and Test Samples for Calibration*

327 To avoid overfitting and to promote generalization in the above supervised learning  
328 algorithm, the historical input and target output data series that are used to construct the ANNs

329 will be randomly split into three subsets. The first subset of data, named the training sample, is  
330 comprised of 70% of the input and target series. As its name suggests, the training sample is  
331 used to train the network by establishing the optimal weights within the neurons. The validation  
332 sample is a smaller subset, comprised of 15% of the historical input and target series. This  
333 subset is ultimately used to determine when the neural network can stop learning based on the  
334 network's ability to generalize effectively. As such, the learning algorithm searches for the point  
335 at which the neural network has the least amount of error over the validation subset during  
336 calibration. Finally, the third subset of input and target data is called the testing sample. This  
337 test sample is not used in the calibration of the model in any way. Instead it simply provides a  
338 more accurate measurement of out-of-sample network performance during calibration.

339           It should be noted that the three subset samples used in calibration are not entirely  
340 independent from one another because of storm-based serial correlation. The general population  
341 of calibration data for any given forecast hour contains multiple target IKE tendency values from  
342 long-lived storms, but will not ever contain multiple sets of predictors from the same model run.  
343 Furthermore, each GEFS run that predictors are taken from is separated by at least 24hrs from  
344 the next closest analysis, as the GEFS is only initialized once daily in the NCEP reforecast data  
345 set. As KM14 showed, past IKE change did not have significant ties to future IKE tendency  
346 beyond the first 24 hours. Therefore, storm-based serial correlation between subsequent target  
347 IKE tendency values for each forecast hour should be somewhat limited. Nonetheless, these  
348 three subsets are only used in the calibration of individual neural networks. Once the weights  
349 are established with analyses as detailed in Section 5, the evaluation exercises done in Section  
350 6, will use out of sample hindcast data to drive the neural network in an effort to best simulate  
351 how the models may perform in real time.

352  
353  
354

### 355 *4.3 Neural Network Random Variability*

356  
357 Inevitably, the methodology used to construct the neural networks introduces random  
358 variability into each individual ANN. Specifically, random variability is first introduced when the  
359 general population of input and target parameters from 1990 to 2011 is randomly split into the  
360 three separate subset samples. Additional random variability is introduced to the neural  
361 networks because the weights within the neurons are initialized somewhat randomly before  
362 arriving at their optimal weights. Ultimately, the random variability makes it all but impossible for  
363 two ANNs to be exactly identical to one another, even if they are calibrated on the exact same  
364 input and target output datasets. Each neural network weighs connections in the nonlinear  
365 system somewhat differently, and as a result, some of the networks will seem more accurate in  
366 certain situations but less accurate in other situations. Therefore, it is insufficient to base  
367 SPIKE2 off a forecast from just a single neural network.

368 Instead, SPIKE2 will utilize a system of one hundred individual neural networks to make its  
369 prediction of IKE tendency for each of the forecast intervals (i.e. a separate system of neural  
370 networks for each forecast interval). As shown by the schematic of SPIKE2 (Figure 1), the  
371 system of neural networks will produce 100 separate independent predictions of IKE tendency  
372 from a single set of input parameters. A deterministic forecast of IKE tendency from SPIKE2 will  
373 be taken from the median of these 100 individual predictions. Using the median from a large  
374 sample of ANNs helps to minimize the random variability present in a single neural network's  
375 forecast, thus allowing SPIKE2 to focus on the true skill of the neural networks. The overall skill  
376 of this deterministic forecast will be discussed at length in Section 6.

377 Not covered in this paper for brevity, but in development, are a series of probabilistic  
378 products that take each of the 100 individual ANNs within SPIKE2 into account, rather than just  
379 the median estimation. These SPIKE2 probabilistic products are used to evaluate the  
380 uncertainty within the ANN statistical scheme from a single set of input parameters, as each

381 ANN within SPIKE2 has a slightly different set of weights. Probability of exceedance, uncertainty  
382 ranges, and error bars, are just a small sample of some of the probabilistic utilities that are  
383 possible with SPIKE2 before even considering the idea of forcing the model with a wide array of  
384 input parameters from multiple forecast models or ensembles.

385

#### 386 *4.4 Comparison of Neural Network and Linear Model Performance*

387       Once again, the primary goal for developing SPIKE2 is to create a statistical dynamical  
388 model that is capable of predicting IKE in a hindcast mode, which would mark a significant step  
389 toward moving to real time operational forecasts of IKE. This objective differs from the goals of  
390 the linear regression version of SPIKE in KM14. That previous linear regression model was  
391 designed in a perfect prognostic mode to prove that IKE prediction is possible when given  
392 accurate environmental predictors. As a result of the different objectives and the different  
393 running environments (perfect prognostic vs hindcast), SPIKE and SPIKE2 are trained on two  
394 completely different sets of predictors, even if the storms in the calibration and evaluation  
395 sample are identical.

396       Therefore, to explicitly show the advantages of the new neural network methodology over  
397 the previously used linear regression model, we also created a perfect prognostic version of our  
398 neural networks using the same data source (developmental SHIPS; DeMaria and Kaplan 1999)  
399 that was used for calibration and evaluation of the SPIKE model in KM14. Unsurprisingly the  
400 added flexibility provided by the nonlinear equations allowed the neural network to outperform  
401 the linear regression model across a 1995-2011 comparison period. For instance, SPIKE2 has a  
402 mean absolute error of 12.6 TJ over its training sample at its longest 72-hour forecast interval.  
403 In contrast, the linear regression model has a comparable mean absolute error of 14 TJ at a  
404 much shorter 24-hour forecast interval. As such, we will progress onward out of the perfect  
405 prognostic space, and begin to calibrate the neural networks with predictors from numerical

406 analyses in Section 5, in an effort to prepare the neural networks for evaluation with hindcast  
407 predictors in Section 6.

408

## 409 **5. Calibration of Neural Networks Using GEFS Analyses**

410           Ultimately, to establish the weights of the ANNs, we calibrate the entire system with  
411 targets of IKE tendency taken from our historical record and normalized input parameters taken  
412 from the control zero-hour analyses (F00) at validation time within the GEFS reforecast archive.  
413 These analyses represent the best estimation for observed environmental conditions within the  
414 model data's one-degree resolution. As such, the SPIKE2 system will be calibrated to accept  
415 reforecast input parameters from the same coarse resolution when it is ultimately evaluated in a  
416 hindcast mode. It is important to note that the F00 analyses do not include forecast errors.  
417 Therefore, only the persistence IKE predictor will change with advancing forecast hour as the  
418 persistence IKE value becomes further removed from the validation time.

419           Since the GEFS reforecast runs are only initialized at 00Z, the maximum sample size for  
420 the F00 calibration dataset is the 1377 storm fixes that occur at 00Z between 1990 and 2011.  
421 However, SPIKE2 requires persistence parameters of varying forecast lead times. Therefore,  
422 the sample size of the calibration dataset will decrease with increasing forecast hour because  
423 short-lived storms will not have longer-term persistence values. For comparison purposes, there  
424 are 1097 fixes at a forecast hour of 12 hours and 614 fixes for the 72-hour forecast interval.

425           The performance of SPIKE2's deterministic forecast for the analysis-based calibration  
426 dataset is shown in Table 2. Similar to SHIPS (DeMaria and Kaplan 1994) and SPIKE (KM14),  
427 the explained variance for the targeted tendency value increases with increasing forecast hour.  
428 The correlation between IKE tendency predictions from SPIKE2 with GEFS F00 data and the  
429 observed historical dataset was  $r=0.73$  at a 12-hour forecast window compared to  $r=0.91$  at 72  
430 hours. This seemingly counterintuitive result can be explained by considering that the

431 magnitude of IKE tendencies increases with growing forecast hour, such that random  
432 fluctuations and observational biases are less impactful at longer forecast hours. Furthermore,  
433 forecast errors are not present in any of the GEFS F00 input parameters such that the input  
434 parameters are no less accurate at 72 hours than they are at 12-hours.

435 In addition to predicting IKE tendency, SPIKE2 can also predict the actual value of IKE  
436 at the validation time by adding its IKE tendency prediction to the persistence IKE value from  
437 the model's initialization time. KM14 found that the IKE metric was somewhat resistant to  
438 change because it considers the energy across a storm's entire wind field. As a result, it is  
439 unsurprising that SPIKE2 performs better at predicting IKE than it does at predicting IKE  
440 tendency because it can use the decent performance of a persistence IKE forecast to its  
441 advantage, especially in short forecast windows. At a 12-hour forecast window, SPIKE2  
442 explains 91% of the observed variance ( $r=.95$ ) when using the GEFS F00 input parameters.  
443 That performance does not degrade sharply, as the explained variance remains near 80%  
444 ( $r=0.90$ ) at a longer 72-hour window.

445 While these high correlations are promising, they are somewhat meaningless if similar  
446 performance can be achieved by simply using a persistence forecast. Encouragingly, the  
447 SPIKE2 calibration model has a lower 72-hour forecast error than does a much shorter 24-hour  
448 persistence forecast. To provide another metric for comparison, we have evaluated the mean  
449 squared error (MSE) from SPIKE2 over its calibration dataset against a persistence forecast at  
450 each corresponding forecast hour in Figure 2. Overall, SPIKE2 has lower MSE than does  
451 persistence by a fair margin (45% at a 12 hour forecast window, climbing up to 82% by 72  
452 hours). The improvements over persistence are statistically significant at a  $p=0.025$  level for all  
453 forecast intervals based on a two-sample bootstrapping test.

454 Also shown in Figure 2 are the reproduced results from the original linear version of  
455 SPIKE detailed in KM14. These results are also calculated over the model's calibration interval,  
456 1990-2011, but as noted earlier, the two models used predictors from entirely different datasets

457 making this comparison uneven. Nonetheless, the calibration statistics indicate that the linear  
458 SPIKE model simply cannot measure up to the neural networks in SPIKE2, although both  
459 models offer substantial improvement over persistence. Like the results of Section 4.4, this  
460 evidence continues to support our hypothesis that the neural networks will be superior to simple  
461 linear regression because it can account for the nonlinearities in the TC-environment system.

462 Although these initial calibration results appear to be encouraging, it should once again  
463 be noted that the hindcast version of SPIKE2 discussed in the following section will use  
464 imperfectly reforecasted input parameters from the GEFS control runs. As such, it would be  
465 unfair to expect SPIKE2's hindcasts in the following section to achieve these high performance  
466 benchmarks. Instead, the performance metrics shown in Table 2 can be viewed as the  
467 maximum potential skill that can be obtained by SPIKE2. The intent of these performance  
468 benchmarks is to determine how the model will degrade when forecast errors are introduced to  
469 the model input fields. Nonetheless, the exercise proved useful by identifying a set of weights  
470 within the artificial neurons that can be used to produce hindcasts of IKE from the GEFS  
471 reforecasts.

472

## 473 **6. Performance of SPIKE2 Hindcasts Using GEFS Reforecasts**

474 In this section, we will adapt the SPIKE2 ANN system to run in a hindcast mode with the  
475 GEFS reforecast control run from 1990 to 2011. As just discussed, the network will retain the  
476 same neuron weights that were calibrated in the previous exercise with GEFS control analyses.  
477 However, unlike the calibration exercises the neural networks will be given imperfect input  
478 parameters from the GEFS reforecast control run at various lead times out to 72 hours. This will  
479 enable us to determine how forecast errors affect SPIKE2's ability to predict IKE. We can  
480 understand from this analysis of predictive skill whether or not SPIKE2 might offer skillful  
481 operational support in a real time environment.

482           Much like the last section, we will evaluate the deterministic forecast from SPIKE2 using  
483 the target IKE tendency and IKE values as the historical baseline. Statistics such as correlations  
484 and mean absolute errors will be used to detect the magnitude of performance deterioration  
485 relative to the maximum potential performance levels obtained in the calibration exercise. As  
486 was done in the earlier calibration exercises and in KM14, SPIKE2's deterministic skill will be  
487 evaluated relative to simple persistence forecasts. However, in addition, a new more  
488 challenging benchmark will also be introduced by way of a simple statistical model that  
489 considers climatology and other non-forecast parameters

490           Such a benchmark model would follow in the footsteps of the Statistical Hurricane  
491 Intensity Forecast model (SHIFOR), which uses seven known parameters at initialization time to  
492 set the baseline performance for operational intensity forecasts (Jarvinen and Neumann 1979;  
493 Knaff et al. 2003). The exact parameters of SHIFOR include: Julian day, initial storm intensity,  
494 previous 12-hour intensity change, initial latitude, initial longitude, initial zonal component of  
495 storm motion, and initial meridional component of storm motion. These SHIFOR climatology and  
496 persistence predictors are somewhat relevant to IKE tendency as well. Therefore, an IKE  
497 statistical persistence model named the "Benchmark of Integrated Kinetic Energy (BIKE)" is  
498 created to predict IKE tendency in a simple linear regression model using the same seven input  
499 parameters, with two exceptions. First, the 12-hr intensity change parameter will be switched  
500 out for a 12-hr IKE change parameter. Second, the initial or persistence value of IKE will be  
501 added as an eighth predictor. This BIKE regression model is trained using all 00Z storm fixes  
502 from 1990-2011, such that its calibration fit will be compared to the GEFS-SPIKE2 hindcasts at  
503 lead times of 24hrs, 48hrs, and 72hrs for the same 1990-2011 interval.

504           Case studies act as a good first step to evaluate the SPIKE2 hindcasts relative to their  
505 assortment of benchmarks in an effort to see how IKE forecasts might perform during significant  
506 landfalling events. To that end, Figure 3 contains a plot of SPIKE2 hindcasts shown against  
507 historical values of IKE just prior to landfall for Hurricanes Floyd (AL081999), Katrina

508 (AL122005), Ike (AL112008), and Irene (AL122011). Each of these four storms gained  
509 considerable IKE as they approached land, and as a result, a persistence forecast would have  
510 greatly underestimated the storm's destructive potential at landfall. BIKE proves to be a more  
511 challenging benchmark for SPIKE2 in these four case studies, as it arguably outperforms  
512 SPIKE2 for Hurricane Floyd. Nonetheless, the SPIKE2 hindcasts outperform BIKE in most other  
513 cases. The SPIKE2 hindcasts for Irene were particularly impressive as the green curve  
514 representing the hindcast remains very close to the black line representing the observations  
515 throughout the 72 hour forecast period. One final item of note, in nearly each case, the SPIKE2  
516 hindcast using reforecasted predictors performs worse than the SPIKE2 calibration model using  
517 predictors from analyses. This result is expected, as it suggests that the performance of the  
518 ANNs will degrade with the introduction of forecast errors in the series of input predictors.

519         Moving to a more general perspective, mean error and correlation statistics are shown  
520 on a line plot in Figure 4 for all of the storm fixes within the 1990-2011 evaluation sample. The  
521 SPIKE2 hindcasts are capable of explaining more than 80% of the variance in the historical IKE  
522 record with a day of lead time, and mean absolute errors are approximately 12 TJs in the same  
523 24-hour forecast window. As lead time increases, hindcast performance expectedly decays, but  
524 the model is still capable of explaining 70% of the historical IKE variance at 48hrs and 62% at  
525 72-hours, with errors of 16.6 and 20.7 TJs at those times respectively.

526         The performance of the hindcast easily exceeds the performance benchmark set forth by  
527 a persistence forecast. For instance, a 72-hour SPIKE2 hindcast has comparable error on  
528 average to a half-as-long 36-hour persistence forecast. Mean squared error statistics paint a  
529 similar picture (Figure 5), as the SPIKE2 model offers a 60% reduction of MSE at 24-hours  
530 relative to persistence. This reduction in MSE relative to persistence holds steady as forecast  
531 hour increases, fluctuating between 50% and 70% between 24hrs and 72hrs of lead time. The  
532 lack of a trend with advancing forecast hour in this MSE reduction metric (outside of the first 12-  
533 24hrs, where persistence forecasts excel) is likely attributed to a balance between rapidly

534 increasing persistence error (lowering the benchmark), and increasing forecast errors in the  
535 input data holding back the SPIKE2 scheme (decreased hindcast performance). The  
536 significance of these improvements is once again tested with a two-sample bootstrapping  
537 exercise. Results, indicate that the SPIKE2 hindcasts are significantly better than persistence at  
538 the  $p=0.05$  level for all forecast windows greater than 12hrs, and at the  $p=0.01$  level for all  
539 forecast windows greater than or equal to 48hrs.

540           Unsurprisingly, the BIKE model is indeed a tougher benchmark than just a simple  
541 persistence forecasts as noted by both the correlation and mean error metrics. For instance,  
542 BIKE has a 12% lower mean absolute error at 72-hours than does persistence. Nonetheless,  
543 the SPIKE2 hindcasts still clearly outperform BIKE. For instance, BIKE's mean absolute errors  
544 are more than 30% higher than the SPIKE2 hindcast errors at all three of the shown forecast  
545 windows.

546           On the other hand, the performance of the hindcasts falls short of the higher  
547 performance levels found during the calibration exercises. Again, this result was expected  
548 because statistical-dynamical prediction schemes are only as accurate as the input data going  
549 into the statistical model. In this case, the GEFS reforecasts include forecast errors that were  
550 not present in the analyses, which results in this degradation of performance. Furthermore, a  
551 lesser decrease in performance should also be expected just by running the ANNs on a dataset  
552 that they were not calibrated with (i.e. the GEFS F00 analyses).

553           Nonetheless, the drop in performance from the calibration tests to the hindcast tests is  
554 not a hindrance. Mean errors only increased by less than 15% and correlations only decreased  
555 by less than 7% inside the shorter 12 and 24-hour windows. Growing inaccuracies in the GEFS  
556 input variables, led to a more dramatic decrease in performance at larger longer forecast  
557 windows relative to the maximum potential performance level in the calibration exercises.  
558 However, once again, these hindcasts are still convincingly skillful relative to a persistence  
559 forecast. In fact, the hindcast performance metrics (green curve) are much closer to the

560 potential performance metrics in the calibration runs (blue curve) than they are to the  
561 persistence performance benchmarks (red curves).

562

## 563 **7. Conclusions and Outlook for Future Operational Development**

564 Despite the promise of the hindcast results presented above, there is still some work left  
565 to be done to adapt this model for operational use. For example, the neural networks would  
566 likely need be recalibrated to receive operationally predicted input parameters from a desired  
567 model in real-time unless the targeted model is similar to the GEFS reforecast data used here  
568 (such as the operational GEFS). If recalibration is needed, a sufficiently long historical database  
569 will once again be needed to normalize the predictors and set the neuron weights. This  
570 limitation is one of the primary reasons for using the control run in the available GEFS  
571 reforecast database. Despite its coarse one-degree resolution, the GEFS archive contained a  
572 long record of data from a static version of the same model. Unfortunately, few operational  
573 models have long archives of forecasts or hindcasts that are readily available. Therefore,  
574 adapting SPIKE2 to be used with a higher resolution operational model or model ensembles is  
575 dependent upon securing an archive for the desired model. As such, adapting SPIKE2 to the  
576 rest of the GEFS ensemble members at a similar resolution or to archived model data stored in  
577 the Observing System Research and Predictability Experiment (THORPEX)'s Interactive Grand  
578 Global Ensemble (TIGGE) archive would be easier to accomplish than would be adapting  
579 SPIKE2 to work with predictors from the Hurricane Weather Research and Forecasting Model  
580 (HWRF).

581 In addition to calibration, future work must focus on determining whether or not SPIKE2  
582 forecasts can be made in a timely manner. SPIKE2's products almost certainly cannot be issued  
583 instantaneously at initialization time. Although the neural networks themselves can be run fairly  
584 quickly, an operational version of SPIKE2 still requires dynamically forecasted input parameters,

585 and unfortunately, the output from most modern dynamical models is not available until a few  
586 hours after their initialization time. Therefore, statistical-dynamical models dependent upon  
587 forecast model data, such as Model Output Statistics (MOS) and in the future SPIKE2, cannot  
588 come out until the dynamical models' run time concludes. As a result, a 72-hour SPIKE2  
589 forecast using GFS or GEFS data would be already a few hours into its forecast period by the  
590 time it was issued, thus shortening it to a 66-70 hour forecast depending on the actual issuance  
591 time.

592         A large delay between issuance and initialization would be detrimental to the usefulness  
593 of SPIKE2 because most operational forecasters are required to issue their forecasts at regular  
594 intervals. To alleviate this concern some operational statistical-dynamical models, are run in a  
595 so called "early cycle" mode, wherein each product uses environmental predictors from the  
596 previous dynamical model run, which is typically initialized six hours earlier (i.e. the 00Z forecast  
597 uses dynamical predictors from an 18Z model). Adapting this early cycle approach to SPIKE2  
598 will ensure that its IKE forecasts are in advance of each forecast advisory. Consequently,  
599 SPIKE2's dynamical predictors in an early cycle mode would be several hours old before  
600 SPIKE2 is even issued. As such, the need to forecast the input parameters an additional few  
601 hours into the future would likely result in a slight degradation to model skill. Considering that  
602 SPIKE2 hindcasts outperform much shorter persistence forecasts, this may not have a  
603 substantial effect on performance. Nonetheless, it will be necessary to compare the pros and  
604 cons of running SPIKE2 in this early cycle mode against attempting to develop an interpolator to  
605 issue SPIKE2 forecasts in a "late cycle" mode as real time development of SPIKE2 continues.

606         Nonetheless, the results presented in the earlier sections serve as a proof of concept,  
607 suggesting that SPIKE2 could be a viable product in an operational setting once these hurdles  
608 are cleared. In calibration exercises, the deterministic scheme is capable of explaining the  
609 majority of variance in the historical IKE archive, and offers a significant improvement over  
610 persistence. Importantly, the addition of imperfectly predicted input parameters from a coarse

611 GEFS control run archive did not cause SPIKE2's performance to drop off severely. Instead,  
612 SPIKE2 hindcasts from 1990 to 2011 still exhibit significant skill over any known IKE persistence  
613 or climatology metrics, despite the inclusion of forecast errors from a rather coarse resolution  
614 GEFS dataset. Not to mention, the briefly discussed SPIKE2 probabilistic products add value to  
615 the deterministic IKE forecasts by offering a quantitative estimate of uncertainty in the statistical  
616 neural network scheme.

617         With the inclusion of input parameters from a higher resolution dataset that is capable of  
618 better resolving some of the storm specific predictors, it may be possible to improve SPIKE2's  
619 skill even further. Nonetheless, if even the level of performance by SPIKE2 with the GEFS  
620 reforecast data can be maintained when adapting SPIKE2 for operations, it would surpass the  
621 ability of any known guidance specifically targeted for deterministic IKE prediction.

622

## 623 **Acknowledgements**

624 Thanks to Robert Hart, Phillip Sura, Allan Clarke, Ming Ye, and Mark Bourassa for their helpful  
625 comments and feedback. This work was graciously supported by grants from NOAA  
626 (NA12OAR4310078, NA10OAR4310215, NA10OAR4320143). Finally, we greatly appreciate  
627 the constructive comments and suggestions given to us by two anonymous reviewers during the  
628 submission process.

629

630

631

632

633

634

635

636 **References**

637 Abdul-Wahab, S.A., and S.M Al-Alawi, 2002: Assessment and prediction of tropospheric ozone  
638 concentration levels using artificial neural networks. *Environmental Modelling & Software*,  
639 **12**, 219-228.

640

641 Atkinson, P.M., and A.R.L. Tatnall, 1997: Introduction Neural networks in remote sensing.  
642 *International J. of Remote Sensing*, **18**, 699-709.

643

644 Bell, G. D., and Coauthors, 2000: Climate assessment for 1999. *Bull. Amer. Meteor. Soc.*,**81**,  
645 1328–1378.

646

647 Bister, M. and K.A. Emanuel, 1998: Dissipative heating and hurricane intensity. *Meteor. Atm.*  
648 *Phys.*, **52**, 233-240.

649

650 Cao Q., B.T. Ewing, M.A. Thompson, 2012: Forecasting wind speed with recurrent neural  
651 networks. *European J. of Operational Res.*, **221**, 148-154.

652

653 Cawley G.C. and S.R. Dorling, 1996: Reproducing a subjective classification scheme for  
654 atmospheric circulation patterns over the United Kingdom using a neural network.  
655 *Proceedings International Conference on Neural Networks*, 281-286.

656

657 Demuth, J.L., M. DeMaria., and J. A. Knaff, 2006: Improvement of advanced microwave  
658 sounder unit tropical cyclone intensity and size estimation algorithms. *J. Appl. Meteorol.*  
659 *Climatol.*, **45**, 1573-1581.

660

661 DiNapoli S.M, M.A. Bourassa, and M.D. Powell, 2012: Uncertainty and Intercalibration Analysis  
662 of H\*Wind. *J. Atmos. and Ocean. Tech.*, **29**, 822-833  
663  
664 Emanuel K., 1988: The Maximum Intensity of Hurricanes. *J. Atmos. Sci.*, **45**, 1143-1155.  
665  
666 Evans, C., and R.E. Hart, 2008: Analysis of the Wind Field Evolution Associated with the  
667 Extratropical Transition of Bonnie (1998). *Mon. Wea. Rev.*, **136**, 2047-2065.  
668  
669 —, 2005: Increasing destructiveness of tropical cyclones over the past 30 years. *Nature*, **436**,  
670 686–688.  
671  
672 Galarneau, T., and T. Hamill, 2015: Diagnosis of Track Forecast Errors for Tropical Cyclone  
673 Rita (2005) Using GEFS Reforecasts. *Wea. Forecasting*, in press.  
674  
675 Gardner M.W., and S.R. Dorling, 1998: Artificial neural networks (the multilayer perceptron)-a  
676 review of applications in the atmospheric sciences. *Atmos. Environment*, **32**, 2627-2636.  
677  
678 Hagan, M.T. and Menhaj, M.B., 1994: Training feedforward networks with the Marquardt  
679 algorithm. *IEEE Transactions on Neural Networks*, **5**, 989-993.  
680  
681 Hall T., and K. Hereid, 2015: The frequency and duration of U.S. hurricane droughts. *Geo. Res.*  
682 *Let*, **42**, 3482-3485.  
683  
684 Hamill, T. M., G. T. Bates, J. S. Whitaker, D. R. Murray, M.I Fiorino, T. J. Galarneau, Y. Zhu,  
685 and W. Lapenta, 2013: NOAA's Second-Generation Global Medium-Range Ensemble  
686 Reforecast Dataset. *Bull. Amer. Meteor. Soc.*, **94**, 1553–1565.

687 Hapuarachchi, H.A.P., Q.J. Wang, and T.C. Pagano, 2011: A review of advances in flash flood  
688 forecasting. *Hydrol. Process.*, **25**, 2771-2784.

689

690 Hart R.E., and J.L. Evans, 2001: A Climatology of the Extratropical Transition of Atlantic  
691 Tropical Cyclones. *J. Climate*, **14**, 546-564.

692

693 —, D.R. Chavas, and M.P. Guishard, 2015: The arbitrary definition of the current Atlantic  
694 major hurricane landfall drought. *Bull. Amer. Soc.*, in press.

695

696 Irish, J.L., D T. Resio, and J.J. Ratcliff, 2008: The influence of Storm Size on Hurricane Surge.  
697 *J. Phys. Oceanogr.*, **38**, 2003-2013.

698

699 Jarvinen, B. R., and C. J. Neumann, 1979: Statistical forecasts of tropical cyclone intensity for  
700 the North Atlantic basin. NOAA Tech. Memo. NWS NHC-10, 22 pp.

701

702 Knaff J.A., M. DeMaria, C.R. Sampson, and J.M. Gross, 2003: Statistical, 5-Day Tropical  
703 Cyclone Intensity Forecasts Derived from Climatology and Persistence. *Wea. Forecasting*,  
704 **18**, 80–92.

705

706 —, S.P. Longmore, D.A. Molenaar, 2014: An Objective Satellite-Based Tropical Cyclone Size  
707 Climatology, *J. Cli.*, **27**, 455-476.

708

709 Kozar, M.E. and V. Misra, 2014: Statistical Prediction of Integrated Kinetic Energy in Atlantic  
710 Tropical Cyclones, *Mon. Wea. Rev.*, e-view.

711

712 Kriesel D., 2007: A Brief Introduction to Neural Networks. [available online:  
713 [http://www.dkriesel.com/en/science/neural\\_networks](http://www.dkriesel.com/en/science/neural_networks)]  
714

715 Landsea, C.W., and J.L. Franklin, 2013: Atlantic Hurricane Database Uncertainty and  
716 Presentation of a New Database Format. *Mon. Wea. Rev.*, **141**, 3576-3592.  
717

718 Maclay K.S., M. DeMaria, and T.H. Vonder Haar, 2008: Tropical Cyclone Inner-Core Kinetic  
719 Energy Evolution. *Monthly Weather Review* **136**:12, 4882-4898.  
720

721 Marquardt D., 1963: An algorithm for least squares estimation of non-linear parameters. *J. Soc.*  
722 *Ind. Appl. Math*,**11**, pp.431 -441  
723

724 Musgrave K.D., R. K.Taft, J. L.Vigh, B. D.McNoldy, and W. H.Schubert, 2012: Time evolution of  
725 the intensity and size of tropical cyclones, *J. Adv. Model. Earth Syst.*, **4**, M08001.  
726

727 Powell, M.D., and T.A. Reinhold, 2007: Tropical Cyclone Destructive Potential by Integrated  
728 Kinetic Energy. *Bull. Amer. Meteor. Soc.*, **88**, 513–526.  
729

730 Reynolds, R. W., T. M. Smith, C. Liu, D. B. Chelton, K. S. Casey, and M. G. Schlax, 2007: Daily  
731 high-resolution blended analyses for sea surface temperature. *J. Climate*, **20**, 5473-5496.  
732

733 Saha S., and coauthors, 2010: The NCEP Climate Forecast System Reanalysis. *Bull. Amer.*  
734 *Meteor. Soc.*, **91**, 1015-1057.  
735

736 Shukla R.P., K.C. Tripathi, A.C. Pandley, and I.M.L. Das, 2011: Prediction of Indian summer  
737 monsoon rainfall using Niño indices: A neural network approach. *Atmos. Res.*, **102**, 99-109.

738 **Figure Captions**

739

740 **Figure 1:** Schematic of the SPIKE2 neural network system. A single set of input parameters is  
741 passed into each of the one-hundred independent artificial neural networks (ANN1, ANN2, ...  
742 ANN100) that make up SPIKE2. Each network produces its own separate prediction of IKE  
743 tendency based on the same input parameters. The median of these predictions is used as  
744 SPIKE2's deterministic prediction, but each individual prediction can be used for probabilistic  
745 forecasting.

746

747 **Figure 2:** Evaluation of SPIKE2 skill in a calibration mode with GEFS analyses relative to a  
748 persistence forecast. Calibration skill is measured as a percent reduction of mean squared error  
749 (MSE) for the model's deterministic predictions from 1990 to 2011, with respect similar MSE  
750 calculations for a persistence forecast at various forecast hours. A reduction of MSE is plotted  
751 as a positive percentage, indicating that the model has outperformed persistence. SPIKE2 has  
752 significantly lower MSE than persistence at the  $p=0.05$  level for all forecast hours. For  
753 reference, also shown is the reproduced results of the linear regression version of SPIKE as  
754 detailed in KM14.

755

756 **Figure 3:** 72-hour runs of SPIKE2 plotted against benchmarks and the observed IKE values  
757 (black lines) for notable hurricanes immediately prior to their landfalls. SPIKE2 hindcasts utilize  
758 GEFS reforecasted predictors from a run initialized at the time specified in each legend. SPIKE2  
759 calibration runs utilize analyzed predictors from the GEFS archive valid at each forecast time.  
760 The benchmark model is initialized at the same time as the SPIKE2 hindcast run for direct  
761 comparison purposes. Not explicitly shown is a persistence forecast which would be

762 represented by a horizontal line stretching from the first observed IKE value on the left through  
763 the entire 72hr period.

764  
765 **Figure 4:** Performance statistics for SPIKE2. Correlation (panel A) and mean absolute error  
766 values, in units of TJ, (panel B) are shown between the historical record and SPIKE2 in various  
767 modes or a persistence forecast. The correlation value that is shown in these plots is for IKE,  
768 not IKE tendency. Calibration statistics are identical to those in Table 2, and are used as a  
769 maximum potential reference point to determine the degradation of skill when forecast error is  
770 introduced to the model in the hindcast runs via the input parameters from the GEFS reforecast.  
771 Also shown for reference is the performance of a persistence forecast and the statistical  
772 climatological and persistence model, BIKE.

773  
774 **Figure 5:** Evaluation of SPIKE2 skill relative to a persistence forecast. Performance is once  
775 again measured as a percent reduction of mean squared error (MSE) for the model's  
776 deterministic predictions from 1990 to 2011, with respect similar MSE calculations for a  
777 persistence forecast at various forecast hours. The calibration skill is reproduced from Figure 2  
778 and is shown alongside the skill of the SPIKE2 hindcasts with reforecasted input parameters  
779 relative to persistence. The hindcast model is significantly better than persistence at the  $p=0.05$   
780 level for all forecast hours greater than or equal to 24-hours in length.

781

782

783

784

785

786

787 **Figures and Tables**

788 **Table 1:** Variables used in the SPIKE2 neural networks. These input parameters are obtained  
 789 from GEFS reforecasts and analyses, NOAA OI SSTs, and the historical record. Many of these  
 790 predictors are inspired by the parameters contained in the SHIPS developmental dataset.

791  
792

Variable	Definition	Units
PIKE	persistence of IKE	TJ
dIKE12	previous 12hr change of IKE	TJ
VMAX	maximum sustained wind speed	kts
VMPI	Difference between maximum potential intensity and VMAX	kts
LAT	latitude of storm's center	°N
LON	longitude of storm's center	-°W
MSLP	minimum sea level pressure	hPa
PENV	average surface pressure ( <i>averaged from r=0-800km</i> )	hPa
VORT	850 hPa vorticity ( <i>r=0-1000km</i> )	10 <sup>-7</sup> s <sup>-1</sup>
D200	200 hPa divergence ( <i>r=0-1000km</i> )	10 <sup>-7</sup> s <sup>-1</sup>
SHRD	850-200 hPa shear magnitude ( <i>r=0-800km</i> )	kts
SHTD	850-200 hPa shear direction ( <i>r=0-800km</i> )	°
RHLO	850-700 hPa relative humidity ( <i>r=0-800km</i> )	%
RHMD	700-500 hPa relative humidity ( <i>r=0-800km</i> )	%
T150	150 hPa temperatures ( <i>r=0-800km</i> )	°C
SST	sea surface temperatures	°C
SDAY	time after tropical storm genesis	days
PDAY	time from peak of season (Sept. 10)	days

793

794

795

796

797

798

799 **Table 2:** Performance of SPIKE2's deterministic forecast when evaluated with the calibration  
800 input set from the GEFS F00 analyses. Sample size indicates the amount of storm fixes that  
801 were included at each forecast hour.  $R_{\text{tendency}}$  measures the correlation between the observed  
802 IKE tendency value and the predicted IKE tendency value from SPIKE2's output.  $R_{\text{IKE}}$  measures  
803 the correlation between the observed IKE value at validation time and the predicted IKE value  
804 calculated by adding SPIKE's tendency prediction to the existing persistence value from  
805 initialization time. Mean error is simply the mean absolute difference between the predictions  
806 from SPIKE2 and the observed IKE values.

Forecast Hr	Sample Size	$R_{\text{tendency}}$	$R_{\text{IKE}}$	Mean Error
12	1097	0.73	0.95	7.8 TJ
24	974	0.83	0.92	10.7 TJ
36	859	0.83	0.90	12.5 TJ
48	773	0.86	0.89	13.4 TJ
60	679	0.89	0.91	13.2 TJ
72	614	0.91	0.90	14.1 TJ

807

808

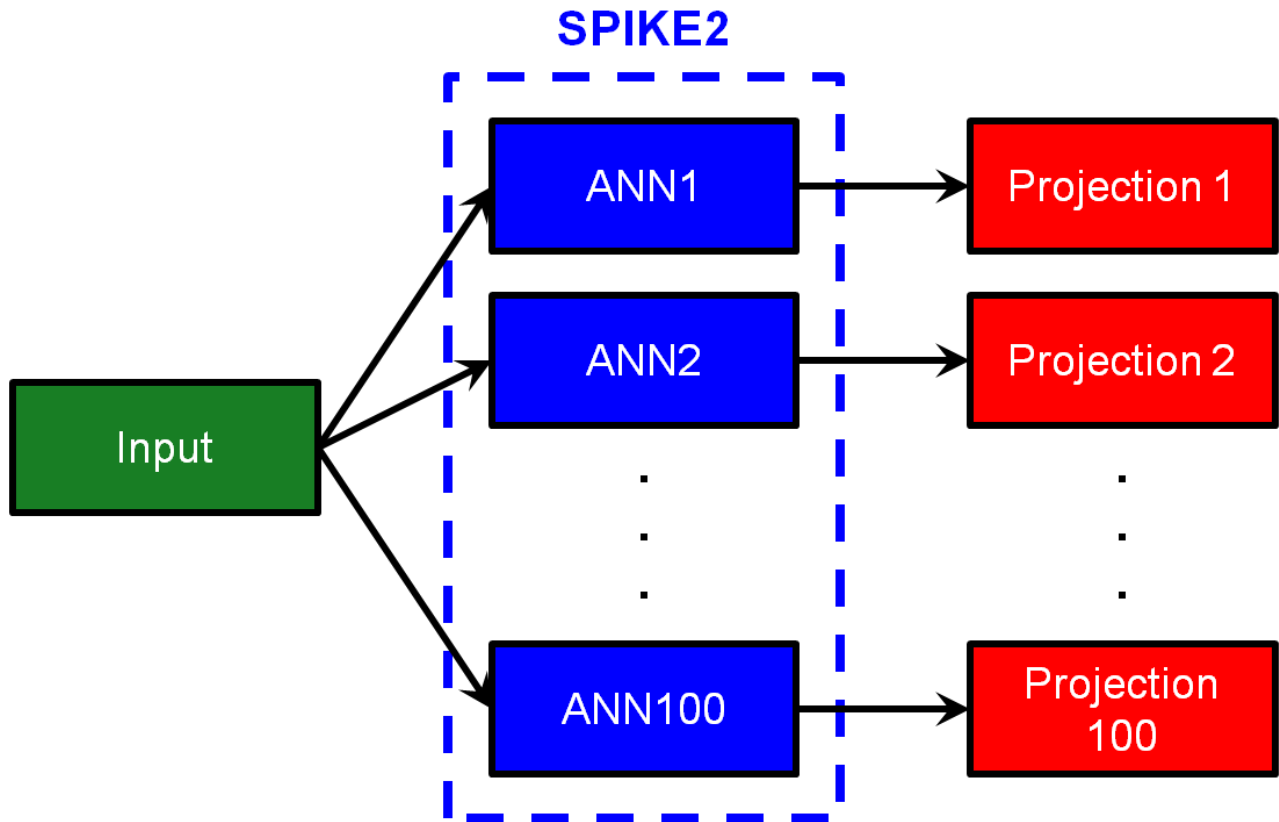
809

810

811

812

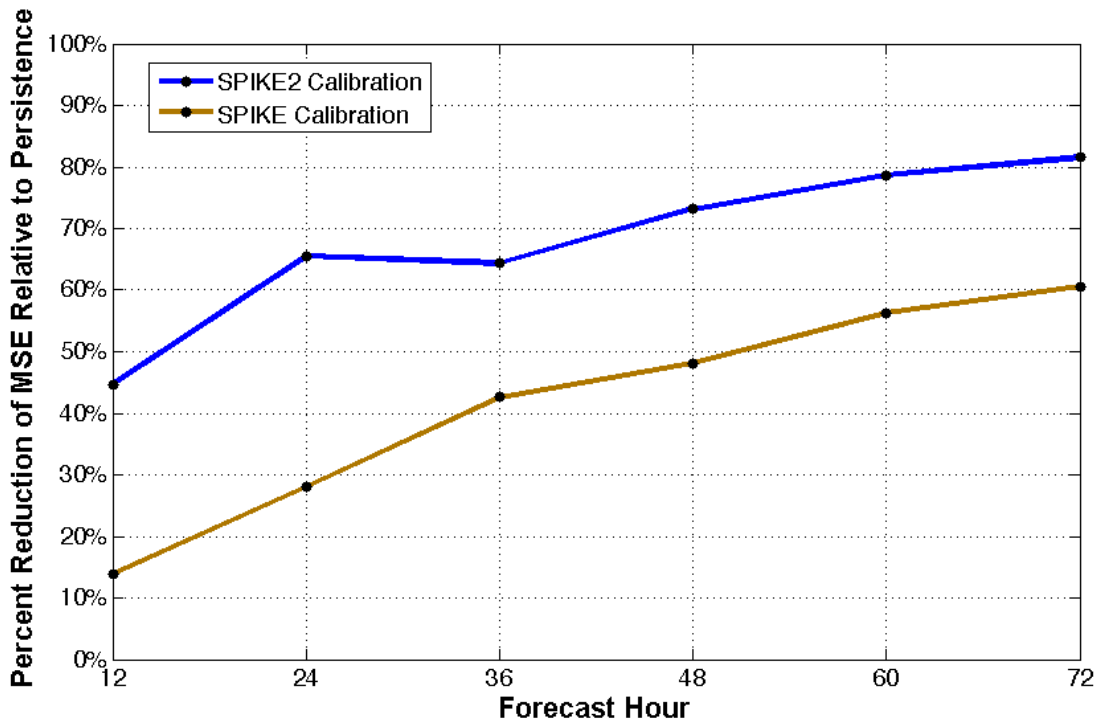
813



814

815 **Figure 1:** Schematic of the SPIKE2 neural network system. A single set of input parameters is  
 816 passed into each of the one-hundred independent artificial neural networks (ANN1, ANN2, ...  
 817 ANN100) that make up SPIKE2. Each network produce its own separate prediction of IKE  
 818 tendency based on the same input parameters. The median of these predictions is used as a  
 819 SPIKE2's best deterministic prediction, but each individual prediction can be used for  
 820 probabilistic forecasting.

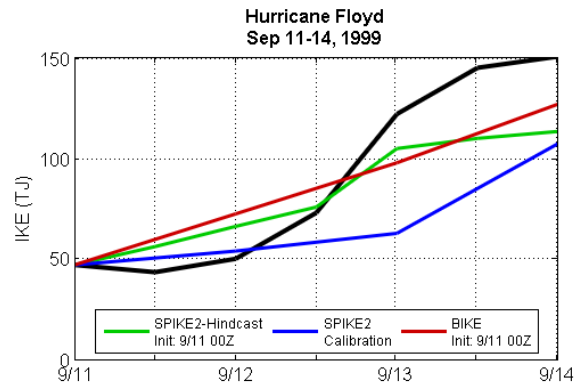
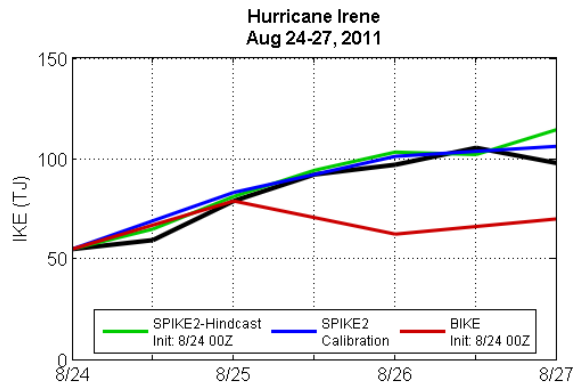
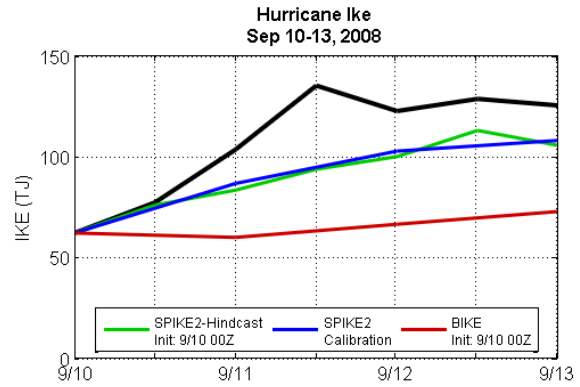
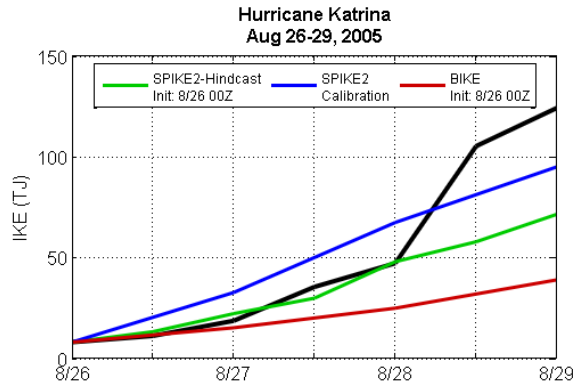
821  
 822  
 823  
 824  
 825  
 826  
 827  
 828



829  
 830  
 831  
 832  
 833  
 834  
 835  
 836  
 837  
 838  
 839

**Figure 2:** Evaluation of SPIKE2 skill in a calibration mode with GEFS analyses relative to a persistence forecast. Calibration skill is measured as a percent reduction of mean squared error (MSE) for the model's deterministic predictions from 1990 to 2011, with respect similar MSE calculations for a persistence forecast at various forecast hours. A reduction of MSE is plotted as a positive percentage, indicating that the model has outperformed persistence. SPIKE2 has significantly lower MSE than persistence at the  $p=0.05$  level for all forecast hours. For reference, also shown is the reproduced results of the linear regression version of SPIKE as detailed in KM14.

840  
 841  
 842  
 843  
 844  
 845



846

847 **Figure 3:** 72-hour runs of SPIKE2 plotted against benchmarks and the observed IKE values  
 848 (black lines) for notable hurricanes immediately prior to their landfalls. SPIKE2 hindcasts utilize  
 849 GEFS reforecasted predictors from a run initialized at the time specified in each legend. SPIKE2  
 850 calibration runs utilize analyzed predictors from the GEFS archive valid at each forecast time.  
 851 The benchmark model is initialized at the same time as the SPIKE2 hindcast run for direct  
 852 comparison purposes. Not explicitly shown is a persistence forecast which would be  
 853 represented by a horizontal line stretching from the first observed IKE value on the left through  
 854 the entire 72hr period.

855

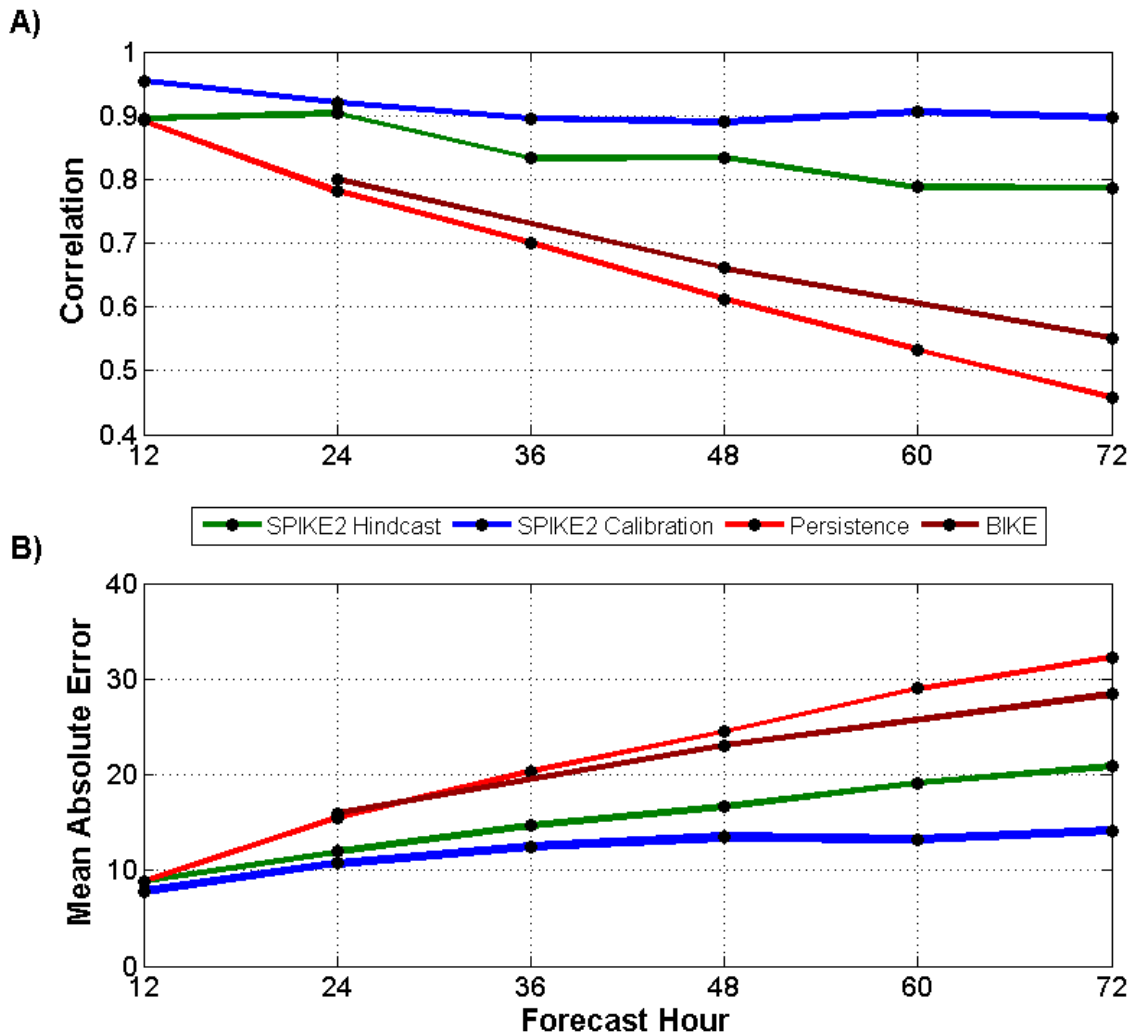
856

857

858

859

860



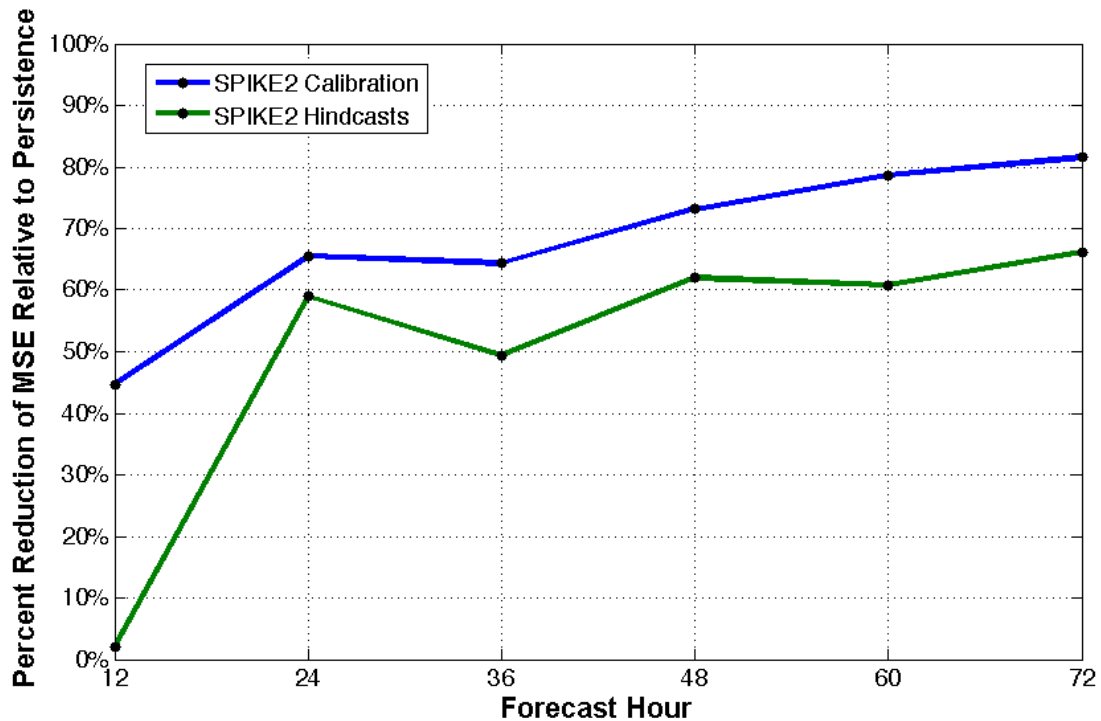
861

862 **Figure 4:** Performance statistics for SPIKE2. Correlation (panel A) and mean absolute error  
 863 values, in units of TJ, (panel B) are shown between the historical record and SPIKE2 in various  
 864 modes or a persistence forecast. The correlation value that is shown in these plots is for IKE,  
 865 not IKE tendency. Calibration statistics are identical to those in Table 2, and are used as a  
 866 maximum potential reference point to determine the degradation of skill when forecast error is  
 867 introduced to the model in the hindcast runs via the input parameters from the GEFS reforecast.  
 868 Also shown for reference is the performance of a persistence forecast and the statistical  
 869 climatological and persistence model, BIKE.

870

871

872



873

874 **Figure 5:** Evaluation of SPIKE2 skill relative to a persistence forecast. Performance is once  
 875 again measured as a percent reduction of mean squared error (MSE) for the model's  
 876 deterministic predictions from 1990 to 2011, with respect similar MSE calculations for a  
 877 persistence forecast at various forecast hours. The calibration skill is reproduced from Figure 2  
 878 and is shown alongside the skill of the SPIKE2 hindcasts with reforecasted input parameters  
 879 relative to persistence. The hindcast model is significantly better than persistence at the  $p=0.05$   
 880 level for all forecast hours greater than or equal to 24-hours in length.

881

882

883

884

885

886

887

888

889

890

891