

Tropical Cyclone Forecasts in the DIMOSIC Project – Medium-Range Forecast Models with Common Initial Conditions

Jan-Huey Chen¹, Linjiong Zhou², Linus Magnusson³, Ron McTaggart-Cowan⁴, and Martin Köhler⁵

¹NOAA GFDL/ UCAR

²Cooperative Institute for Modeling Earth Systems, and Program on Atmospheric and Oceanic Sciences, Princeton University

³European Centre for Medium-Range Weather Forecasts

⁴Environment and Climate Change Canada

⁵Deutscher Wetterdienst

December 16, 2022

Abstract

The Tropical cyclone (TC) forecast skill of the eight global medium-range forecast models which are participating in the DIMOSIC (Different Models, Same Initial Conditions) project is investigated in this study. Each model was used to generate 10-day forecasts from the same initial conditions provided by the European Centre for Medium-Range Weather Forecasts. There are a total of 123 initial dates spanning in one year from June 2018 to June 2019 with a 3-day interval. The TC track and intensity forecasts are evaluated against the best track dataset. TC-related precipitation and tropical cyclogenesis forecasts are also compared to explore the differences and similarities of TC forecasts across the models. This comparison of TC forecasts allows model developers in different centers to benchmark their model against other models, with the impact of the initial condition quality removed. The verifications reveal that most models show slow-moving and right-of-track biases in their TC track forecasts. Also, a common dry bias in TC-related precipitation indicates a general deficiency in TC intensity and convection in the models which should be related to insufficient model resolution. These findings provide important references for future model developments.

1
2 **Tropical Cyclone Forecasts in the DIMOSIC Project – Medium-Range Forecast**
3 **Models with Common Initial Conditions**

4 **Jan-Huey Chen^{1,2}, Linjiong Zhou³, Linus Magnusson⁴, Ron McTaggart-Cowan⁵, and**
5 **Martin Köhler⁶**

6 ¹National Oceanic and Atmospheric Administration/Geophysical Fluid Dynamics Laboratory,
7 Princeton, NJ, USA

8 ²University Corporation for Atmospheric Research, Boulder, CO, USA

9 ³Atmospheric and Oceanic Sciences Program, Princeton University, Princeton, NJ, USA

10 ⁴European Centre for Medium-Range Weather Forecasts, Reading, UK

11 ⁵Environment and Climate Change Canada, Montreal, Canada

12 ⁶Deutsche Wetterdienst, Offenbach, Germany

13
14 Corresponding author: Jan-Huey Chen (Jan-Huey.Chen@noaa.gov)
15

16 **Key Points:**

- 17
- 18 • Tropical cyclone forecasts are compared between global medium-range models from
19 leading modeling centers initialized with identical data.
 - 20 • Similarities and differences between the models set a benchmark of TC forecast with the
21 impact of the initial condition quality removed.
 - 22 • Common TC forecast biases indicate general deficiencies in the models and suggest a
direction for further model improvement.

23 **Abstract**

24 The Tropical cyclone (TC) forecast skill of the eight global medium-range forecast models
25 which are participating in the DIMOSIC (Different Models, Same Initial Conditions) project is
26 investigated in this study. Each model was used to generate 10-day forecasts from the same
27 initial conditions provided by the European Centre for Medium-Range Weather Forecasts. There
28 are a total of 123 initial dates spanning in one year from June 2018 to June 2019 with a 3-day
29 interval. The TC track and intensity forecasts are evaluated against the best track dataset. TC-
30 related precipitation and tropical cyclogenesis forecasts are also compared to explore the
31 differences and similarities of TC forecasts across the models. This comparison of TC forecasts
32 allows model developers in different centers to benchmark their model against other models,
33 with the impact of the initial condition quality removed. The verifications reveal that most
34 models show slow-moving and right-of-track biases in their TC track forecasts. Also, a common
35 dry bias in TC-related precipitation indicates a general deficiency in TC intensity and convection
36 in the models which should be related to insufficient model resolution. These findings provide
37 important references for future model developments.

38

39 **Plain Language Summary**

40 Despite recent improvements in our ability to predict the track and intensity of tropical cyclones,
41 these storms remain significant forecasting challenges. Forecasters rely heavily on the guidance
42 generated by numerical weather prediction systems, making the reliability of these systems
43 essential for accurate forecasts during these high-impact weather events. As a result,
44 improvement the quality of tropical cyclone guidance is an important numerical model
45 development objective. In this study, the TC forecast skills in the eight global medium-range
46 forecast models from the model development centers/institutes who participated in the
47 DIMOSIC (Different Models, Same Initial Conditions) project are examined. All models were
48 initialized from the same data provided by the ECMWF (European Centre for Medium-Range
49 Weather Forecasts) to investigate the differences and similarities among their TC forecasts
50 without the impact of the quality of initial conditions. Besides the general TC forecast
51 evaluation metrics including errors and biases of the track and intensity, the TC-related
52 precipitation and TC genesis skills are also evaluated to comprehensively explore the
53 performance of TC forecasts among all models. The comparison allows model developers in
54 different centers to benchmark their model against other participating models. Moreover, the
55 verification results provide important references for future model developments.

56

57 **1 Introduction**

58 Tropical cyclone (TC) prediction is an important mission for weather and climate
59 agencies in many countries. Over the past few decades, numerical models have become the most
60 important tools for operational centers to make TC forecasts on weather and sub-seasonal to
61 seasonal time scales. Therefore, improving the model performance of TC forecasts has been one
62 of the leading tasks in most operational centers or modeling research institutes working on model
63 development. In addition, the accurate depiction of physical processes that lead to a better TC

64 forecast in the model are also relevant to interesting scientific questions in the atmospheric
65 science research area more broadly.

66 The quality of initial conditions has a leading impact on short- to medium-range forecast
67 skill, including for TC forecasts. In Chen et al. (2019a), the fvGFS (finite volume Global
68 Forecasting System) model developed at the Geophysical Fluid Dynamics Laboratory (GFDL)
69 initialized with the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated
70 Forecasting System (IFS) data showed much-improved TC track forecasts for the 2017 Atlantic
71 hurricane season compared to its retrospective forecasts initialized with the data from the
72 National Centers for Environmental Prediction (NCEP) Global Forecast System (GFS) version
73 14. In Magnusson et al. (2019), the same approach was used, comparing the GFDL fvGFS model
74 forecasts to those from the IFS and GFS. The results showed that the choice of initial conditions
75 clearly dominated the forecast quality in the medium-range predictions, but that the model
76 formulation could also play a significant role.

77 Since major model development centers mostly develop their modeling systems
78 independently, the DIMOSIC project (Different Models, Same Initial Conditions; Magnusson et
79 al. 2022) was devised to investigate the relationship between the choice of model formulation
80 and forecast quality. Models developed by different world-leading modeling centers were
81 initialized from the same initial condition. In Magnusson et al. (2022), the differences and
82 similarities of the forecasts among the models were presented. The results found that some pairs
83 of models behaved more similarly than other pairs due to their sharing of partial physical
84 parameterizations, e.g. ECMWF IFS and DWD (Deutsche Wetterdienst) ICON (Icosahedral
85 Non-hydrostatic Model). On the other hand, ICON and GFDL SHiELD (System for High-
86 Resolution Prediction on Earth-to-Local Domains) showed relatively large forecast differences,
87 while both ranking among the best models. Regarding the influences from model formulations
88 on the forecasts, however, it was difficult to point out a single model component that had the
89 strongest impact on the forecast differences. Also, as pointed out by Magnusson et al. (2022) the
90 interaction between different model parameterizations and their respective configurations could
91 play a significant role as well.

92 In this study, the performance of TC forecasts from the DIMOSIC models is evaluated.
93 The TC track and intensity forecast skills among the models during the period of June 2018 to
94 June 2019 are compared. Since TC intensity in interpolated data does not reflect the actual TC
95 intensity at the native model resolution, the TC-related precipitation are also evaluated to provide
96 another perspective on forecasted TC activities for better exploring the differences and
97 similarities among the models. Also, the forecast skill of TC genesis was investigated by
98 comparing the hit/false alarm ratios among the models, as well as using the method based on the
99 lengths of TC genesis lead time introduced in Chen et al. (2019b) to examine the accuracy of TC
100 genesis timing in the model forecasts. These comparisons should be valuable for model
101 developers in different centers to benchmark their model's performance on TC forecasts against
102 that of other models, with the impact of the initial condition quality removed.

103 The models participating in the DIMOSIC project are introduced in section 2 which also
104 describes the observation data and methodology used in this study. The comparisons of track,
105 intensity, TC-related precipitation, and genesis forecasts among the models are contained in
106 section 3. Summary and discussion are presented in section 4.

107 2 DIMOSIC models, forecasts, and verification data

108 General information on the numerical models and their developing centers/institutes
 109 participating in the DIMOSIC project are listed in Table 1. The horizontal resolutions and the
 110 number of vertical levels of the models and their key references are included. Some
 111 centers/institutes submitted more than one model configurations to the project, but we only
 112 investigate one configuration of the model for each center/institute based on their suggestions.
 113 The only exception is to include two versions of IFS 45R1 and 47R3, to provide an example of
 114 the incremental change obtained for an upgrade of one model. For the sea surface temperature
 115 evolution in the models, the two IFSs used a partial coupling to the 3D ocean NEMO model
 116 (Mogensen et al. 2017), SHIELD is coupled with a 1D mixed layer ocean model (Pollard et al.
 117 1973), CMC used a thermodynamic mixed layer ocean model (Zeng and Beljaars 2005), and
 118 others used persistent anomalies from the analysis. Other detailed configurations of each model
 119 including dynamical cores and major physical parameterizations can be found in the sub-section
 120 of “Model descriptions” in the section of “Models and data” and in Table 2 in Magnusson et al.
 121 (2022), but is not repeated in this paper.

Acronyms	Models	Centers/Institutes	Resolution	Key references
ARPEGE	Action de Recherche Petite Echelle Grande Echelle (version: 46T1)	Meteofrance	5-25 km 105 levels	Roehrig et al. (2020)
CMC	Global Environmental Multiscale Model (GEM) (version: v5.0.2)	Canadian Meteorological Center (CMC)	15 km 80 levels	Girard et al. (2014) McTaggart-Cowan et al. (2019)
ICON	Icosahedral Non-hydrostatic Model (version: April 21)	Deutsche Wetterdienst (DWD)	13km 90 levels	DWD (2022)
IFS/ IFS-47R3	Integrated Forecasting System (versions: 45R1 and 47R3)	European Centre for Medium- range Weather Forecasts (ECMWF)	9 km 137 levels	ECMWF (2018, 2021)
JMA	Global Spectral Model (GSM) (version: GSM1705)	Japan Meteorological Agency (JMA)	20 km 100 levels	JMA (2019)
SHIELD	System for High-Resolution Prediction on Earth-to-Local Domains (version: rt2019)	Geophysical Fluid Dynamics Laboratory (GFDL)	13 km 91 levels	Harris et al. (2020)
UM	Unified Model	UK Met Office	10 km 70 levels	Walters et al. (2019)

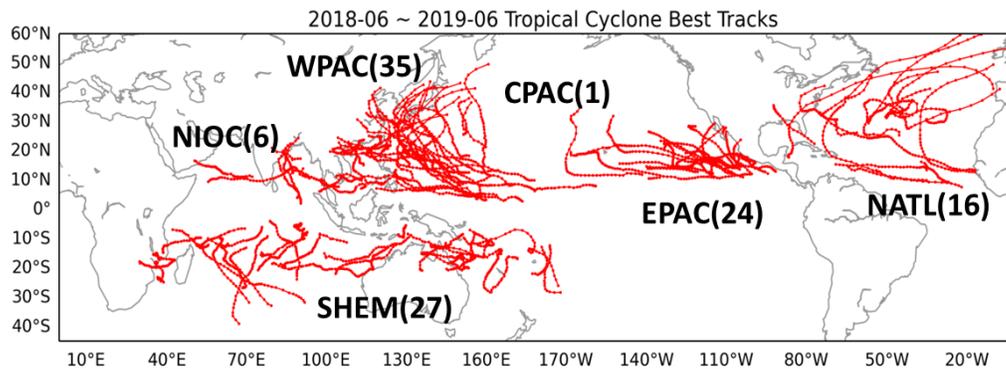
122 TABLE 1. The eight DIMOSIC models (the ECWMF contributed two IFS configurations). From left to right: model
 123 acronyms, model full names (and versions if applicable), centers/institutes that models belong to, horizontal
 124 resolutions and the number of vertical levels used in the models, and the key references of the models.

125 All models conducted 10-day forecasts from the same initial conditions: ECMWF
 126 operational analyses based on the IFS model cycle 45R1 (ECMWF 2018). The 9-km analyses on
 127 137 vertical levels are generated from a 4DVar data assimilation system (Rabier et al. 2000). All
 128 participating institutes received the interpolated data at 0.1 degree for their model initialization.
 129 Detailed procedures for handling the initial conditions in each model are described in section 2c
 130 and Table 3 of Magnusson et al. (2022).

131 The total 123 forecasts are conducted with initialization dates spanning one year from
 132 June 2018 to June 2019 at 3-day intervals. The 10-day forecast outputs from each model were
 133 interpolated to a common 0.5 degree grid using an average interpolation method available in the
 134 EcCodes/MIR package (<https://confluence.ecmwf.int/display/ECC/ecCodes+Home>). The GFDL
 135 simpler tracker (Harris et al. 2016) was used with a warm-core criterion to track TCs in the

136 forecasts of the eight models based on the fields of sea-level pressure, 10-m wind speed, 850-hPa
137 vorticity, and mean temperature between 500-300 hPa.

138 There were 109 observed TCs in the DIMOSIC period. Their storm tracks based on the
139 best track data (b-deck) in Automated Tropical Cyclone Forecast (ATCF) dataset (Miller et al.
140 1990; Sampson and Schrader 2000) are shown on the map in Fig 1. There were 35 TCs in the
141 northwest Pacific basin (WPAC), 24 in the northeast Pacific basin (EPAC), and 16 in the North
142 Atlantic basin (NATL). In the Southern Hemisphere (SHEM), there were 27 TCs in the
143 combined South Indian Ocean and South Pacific Ocean. Besides the overview of global
144 analyses, the individual TC forecast skill in the above four major sub-regions will be investigated
145 and compared with each other in the following sections.



146
147 FIGURE 1. All TCs in the DIMOSIC period. The best tracks are from the ATCF dataset. The numbers of TCs are
148 indicated in the brackets next to the acronyms of the six regions: WPAC: northwest Pacific basin; EPAC: northeast
149 Pacific basin; CPAC: north-central Pacific basin; NATL: North Atlantic basin; NIOC: north Indian Ocean; SHEM:
150 Southern Hemisphere.

151 The ATCF dataset was used to evaluate the forecast errors of TC track and intensity, and
152 the skill of TC genesis forecasts in the models at six-hour intervals. For the TC-related
153 precipitation, the NASA Global Precipitation Measurement (GPM) Integrated Multi-satellitE
154 Retrievals for GPM (IMERG) dataset was used (Hong et al. 2004) for verification. To equally
155 compare to the model output field of total precipitation, this high-resolution (0.1 degree) satellite
156 observational dataset was interpolated into 0.5 degree.

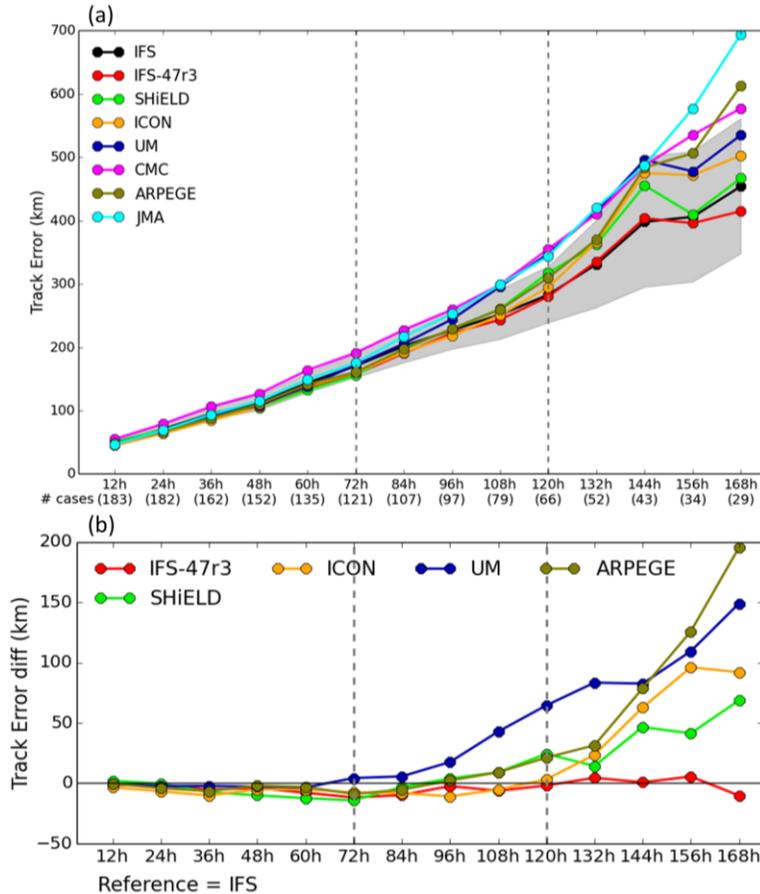
157

158 3. Results

159 3.1 TC track forecasts

160 The prediction of TC path, or track, is the one of the most important factors for taking
161 necessary precautions against possible impacts from hurricanes or typhoons. The homogeneous
162 comparisons of global mean TC track forecast errors along with the forecast lead time at 12-h
163 intervals are shown in Fig. 2a. The differences among the models are small through the 72-hour
164 lead time with the exception of CMC which shows a relatively higher error than others. During
165 the 72 to 120-hour head time, TC track errors diverge into two groups. The two IFSs, ICON,
166 SHIELD, and ARPEGE show lower errors than UM, CMC, and JMA. Both versions of IFS show
167 the lowest TC track errors after the 120-hour lead time all the way to the 168-hour lead time,
168 followed by ICON and SHIELD. Note that most models are within the 95% confidence levels of
169 IFS, which means the differences in forecast skills are not statistically significant. To highlight

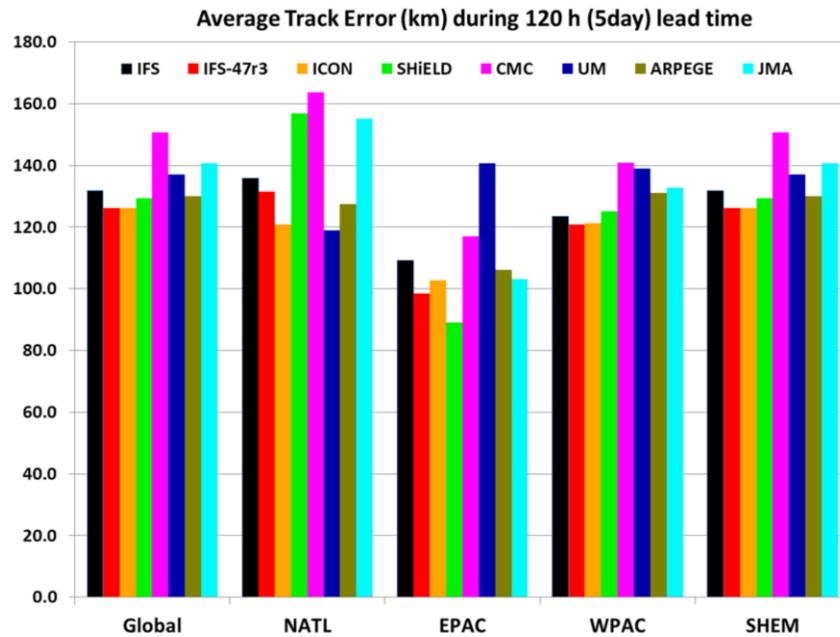
170 the difference between the leading models, Fig. 2b shows the differences in TC track errors of
 171 five of the models compared to the IFS. The newer IFS-47R43 performs slightly better (negative
 172 values in track error differences) or equivalently during the entire 7 days. In the early lead times
 173 (36-84 hours), most of the five models also show slightly better forecasts than IFS. ICON
 174 displays similar (or slightly better) skill to the two IFS versions until the 120-hour lead time.
 175 Both SHiELD and ARPEGE perform very well before the 96-hour lead time. After the 120-hour
 176 lead time, SHiELD shows lower forecast errors than other models except for the two IFSs.



177
 178 FIGURE 2. (a) Global mean TC track forecast errors (km) at every 12 hour forecast lead time for IFS (black), IFS-
 179 47r3 (red), SHiELD (green), ICON (yellow), UM (blue), CMC (magenta), ARPEGE (grass green), and JMA (light
 180 blue). The 95% confidence levels for IFS are indicated by the gray color shading. Numbers of homogeneous cases
 181 for individual lead times are listed in the brackets at the bottom of each abscissa. Vertical gray dotted lines indicate
 182 72 and 120 hour forecast lead times. (b) Global mean TC track forecast error differences of IFS-47r3 (red), SHiELD
 183 (green), ICON (yellow), UM (blue), and ARPEGE (grass green) comparing to IFS.

184 Figure 3 shows the 5-day average TC track errors for all models in the entire globe and in
 185 the four major sub-regions individually. For the two IFSs, IFS-47R3 shows lower track errors
 186 than IFS globally and in all major sub-regions. ICON shows competitive low track errors to IFS-
 187 47R3 in the WPAC and the SHEM, and a very low track error in the NATL. SHiELD performs
 188 the lowest track error in the EPAC, and low track errors besides the two IFSs and ICON in the
 189 SHEM and the WPAC. However, SHiELD has a much larger track error in the NATL, which
 190 results from the slow moving bias shown in the forecasts of Hurricane Florence and the bias of
 191 direction of motion shown in the forecasts of Hurricane Leslie. Detailed investigations can be
 192 found in Text S1 and Figs. S1-S3. The performance of the TC track forecast in UM is notable.

193 The model shows the lowest track error in the NATL but the largest track error in the EPAC,
 194 while its track errors in the WPAC and the SHEM are also at the high end compared to other
 195 models. JMA performs competitively to IFS-47R3 and ICON in EPAC, but not in other sub-
 196 regions. With regard to CMC, it shows larger track errors than other models in most sub-regions
 197 except for in the EPAC during this targeted year. From Fig. 3, we can also see that the globally-
 198 averaged track forecast errors among the models is dominated by the errors in the WPAC and the
 199 SHEM since the majority of the TCs were found in these two sub-regions.

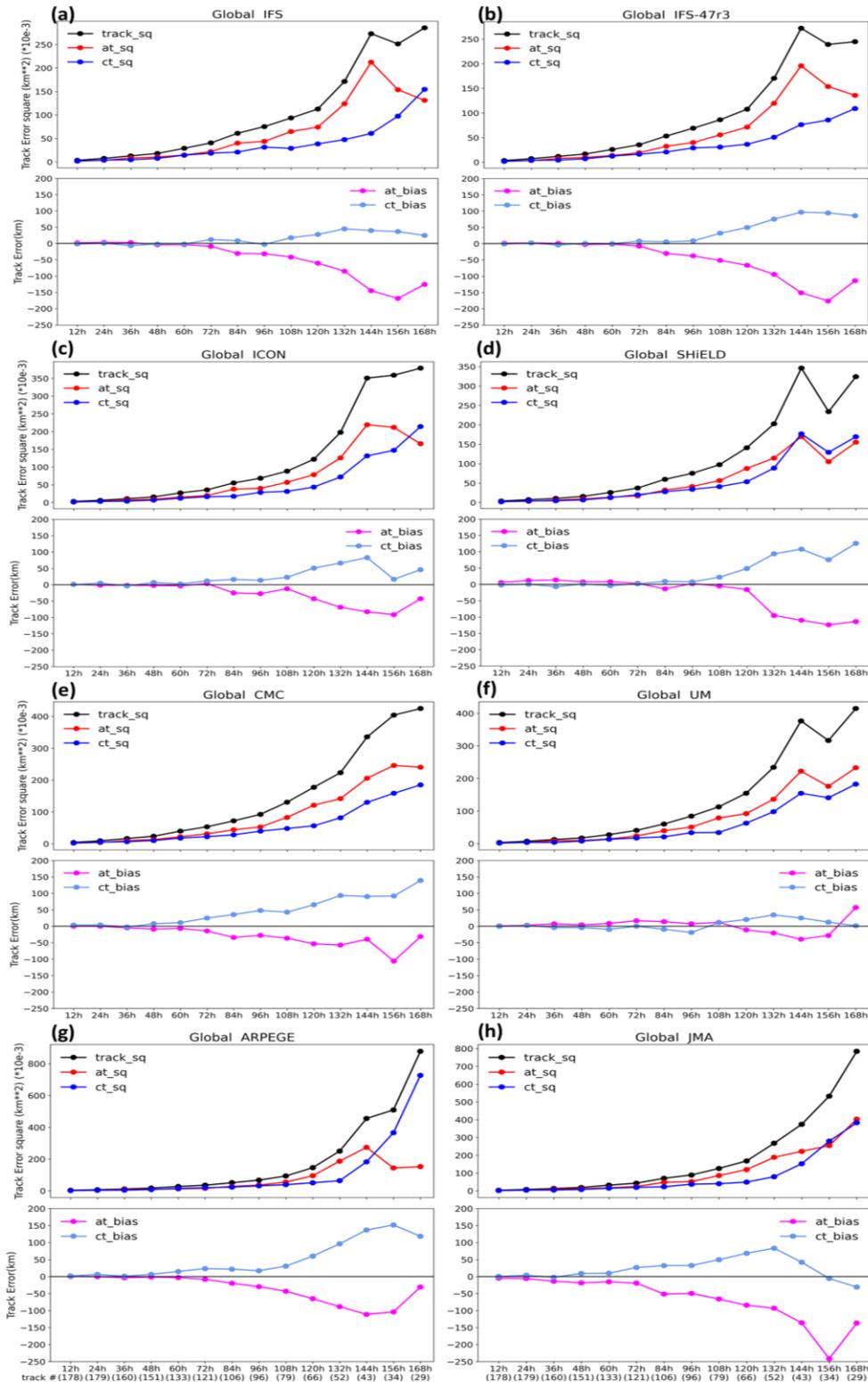


200
 201 FIGURE 3 Averaged Track errors (km) in globe and 4 sub-regions during the 120-hour lead time for the 8 models.
 202 Abbreviations and colors used for the models are the same as in Fig. 2a. Abbreviations used for the sub-regions on
 203 the abscissa are the same as in Fig. 1.

204 The sources of track errors can be due to biases either in forecasts of the TC translational
 205 speed or the TC direction of motion. The lower sub-panels in Fig. 4 show the globally averaged
 206 along-track (AT) errors and cross-track (CT) errors (perpendicular to the track) for all models.
 207 Both AT and CT errors are calculated as great circle distances. Most of the models start to show
 208 negative AT biases and positive CT biases during 72 to 120-hour lead time, which indicates that
 209 the TC track errors during the later lead times are mostly due to the slow and northward (for an
 210 easterly moving TC) moving biases. In general, UM shows the smallest AT and CT biases
 211 among all models, and the biases of SHIELD take place at longer forecast lead times than other
 212 models.

213 By the Pythagorean Theorem the square of the total error equals the squares of the AT
 214 and CT errors (Chen et al. 2019a). The squares of total track errors, CT errors, and AT errors are
 215 plotted in upper sub-panels in Fig. 4 to illustrate the proportion of contributions from the AT and
 216 CT errors respectively to the total error. From Figs. 4a,b, we can see that the AT error
 217 contributes more than the CT error to the total track error in the two IFSs, which indicates that
 218 the track errors in the IFSs are dominated by the slow-moving bias (negative AT biases). The
 219 characteristic of the consistently larger AT error square than the CT error square can also be
 220 found in ICON, CMC, and UM, but the differences between their AT and CT error squares are

221 smaller than those in the IFSs. SHIELD, ARPEGE, and JMA show relatively closer AT and CT
222 error squares, especially SHIELD. This indicates that the contributions of the slow and
223 northward moving biases to the total track errors are similar in these models. For ARPEGE, the
224 CT error is the dominant track error in late lead times, while the AT error contributes more to the
225 JMA's total TC track errors during the 120 to 144-hour lead time.



226

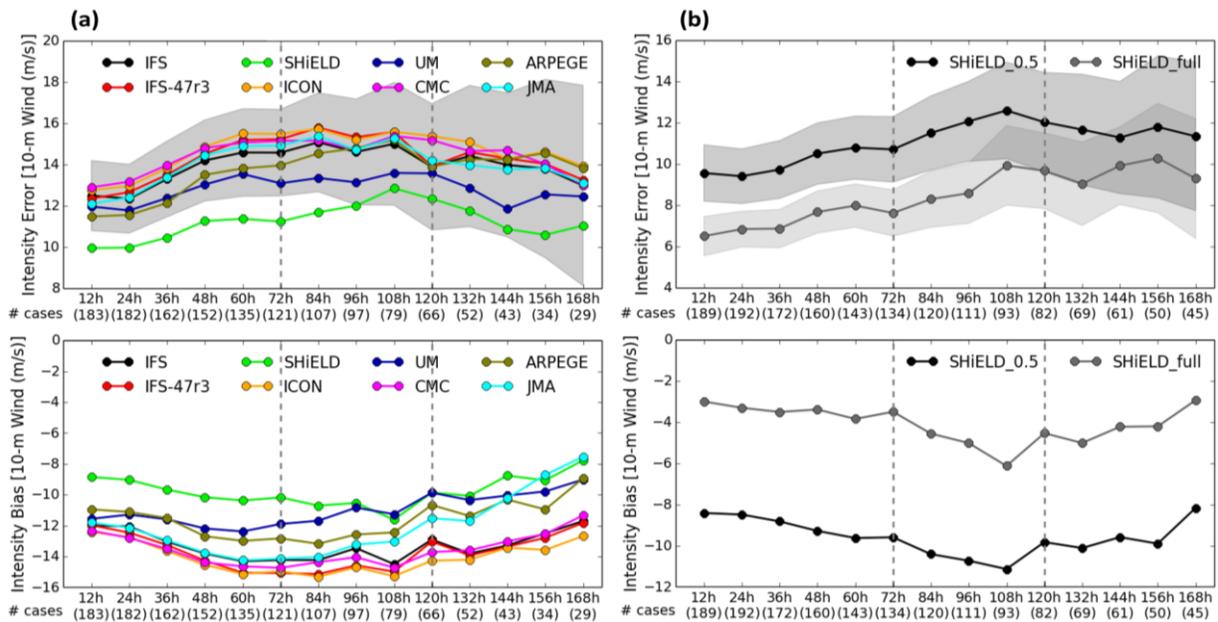
227 FIGURE 4 Global analyses of along-track (AT) error and cross-track (CT) error for (a) IFS, (b) IFS-47R3, (c)
 228 ICON, (d) SHiELD, (e) CMC, (f) UM, (g) ARPEGE, and (h) JMA. The squares of total track errors (black), along-
 229 track errors (red), and cross-track errors (blue) are in the upper panels for each model. The biases of along-track
 230 (magenta) and cross-track (light blue) errors are in the lower panels. Numbers of homogeneous cases for each lead
 231 time are listed at the bottom of lower panels.

232 It is also found that models have different AT and CT errors in different sub-regions. All
 233 models showed a slow-moving bias and a poleward bias in the NATL and the WPAC, but in the
 234 EPAC, except for JMA, most models show a fast-moving bias. In contrast, there are no
 235 consistently slow or fast moving biases among models in the SHEM. Detailed analyses of AT
 236 and CT errors for all eight models in the four major sub-regions can be found in Text S2 and
 237 Figs. S4-S7.

238

239 **3.2 TC intensity forecasts**

240 It has been more challenging to predict TC intensity than track, especially for global
 241 models which usually cannot resolve fine scale interactions between thermal dynamics and
 242 dynamics due to insufficient resolutions. As outlined in Section 2, the 10-day forecast outputs
 243 were interpolated to a common 0.5 degree for each model. Therefore, the model-predicted TC
 244 intensities found by the tracker are underestimated due to the low data resolution. However, it is
 245 still of interest to compare the relative differences of TC intensities among the models for the
 246 interpolated output data. The global mean TC intensity errors and biases based on the maximum
 247 10-m wind speed are presented in Fig. 5a. SHiELD predicts a much stronger TC intensity than
 248 other models, followed by UM and ARPEGE. Figure 5b uses SHiELD as an example to
 249 demonstrate the differences between the TC intensities obtained from the native resolution (13
 250 km grid) outputs and in the interpolated 0.5 degree resolution data. The average differences of
 251 total error and bias are between 3 to 5 ms^{-1} .

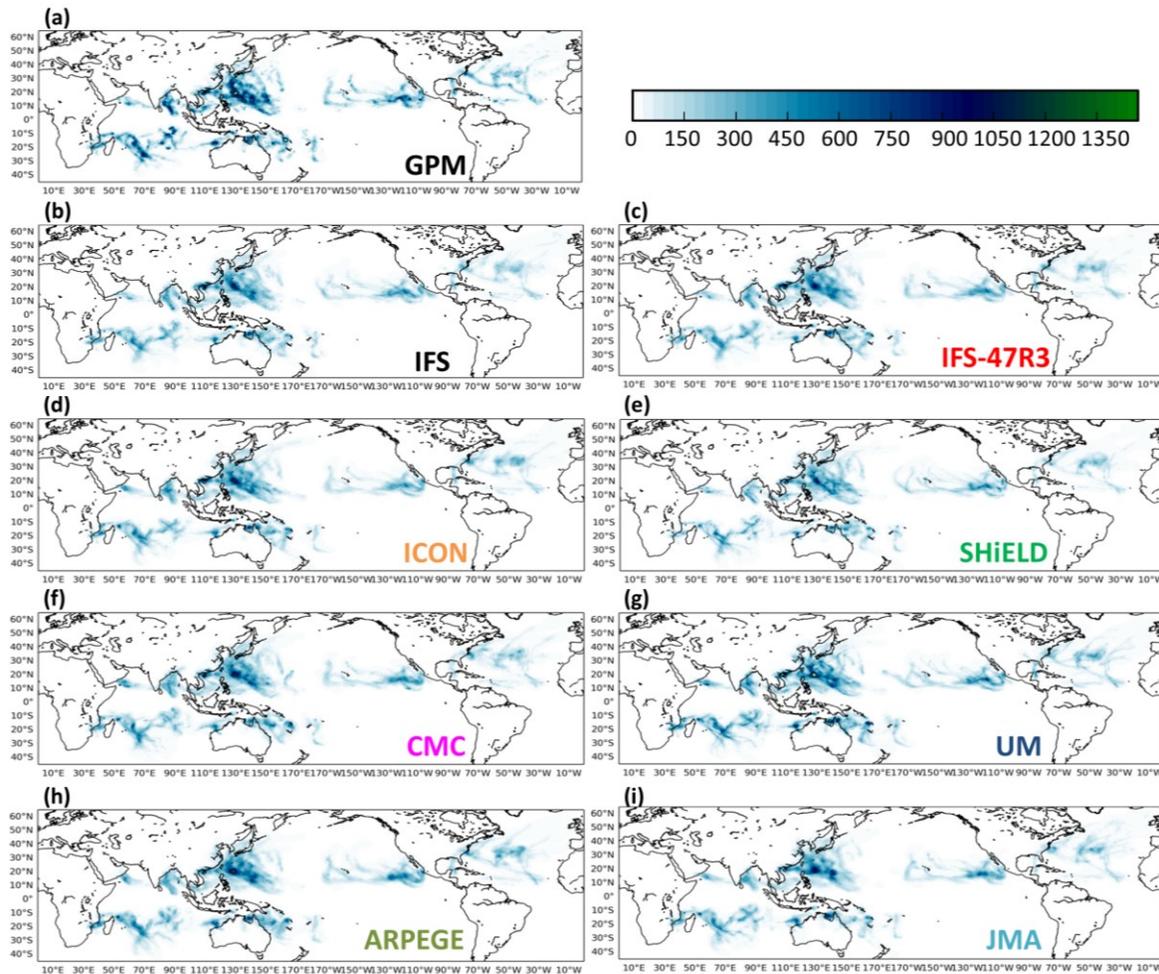


252

253 FIGURE 5. (a) Global mean TC intensity errors and biases. Upper panel: Absolute error of the maximum 10-m wind
 254 speed (m s^{-1}) along with the model forecast lead time for 8 models. Abbreviations and colors used for the models are
 255 the same as in Fig. 2a. The 95% confidence levels for IFS are indicated by the gray color shading. Numbers of
 256 homogeneous cases for individual lead times are listed in the brackets at the bottom. Vertical grey dotted lines
 257 indicate 72 hour and 120 hour lead times. Lower panel: As in the upper panel, but for the bias of the maximum 10-m
 258 wind speed (m s^{-1}). (b) As in (a), but for SHiELD native resolution data (black; SHiELD_full) and SHiELD 0.5
 259 degree interpolated data (gray; SHiELD_0.5). The 95% confidence levels for each resolution data are indicated by
 260 the same medium and light transparent grey shading areas, with their overlapping region denoted by dark grey
 261 shading.

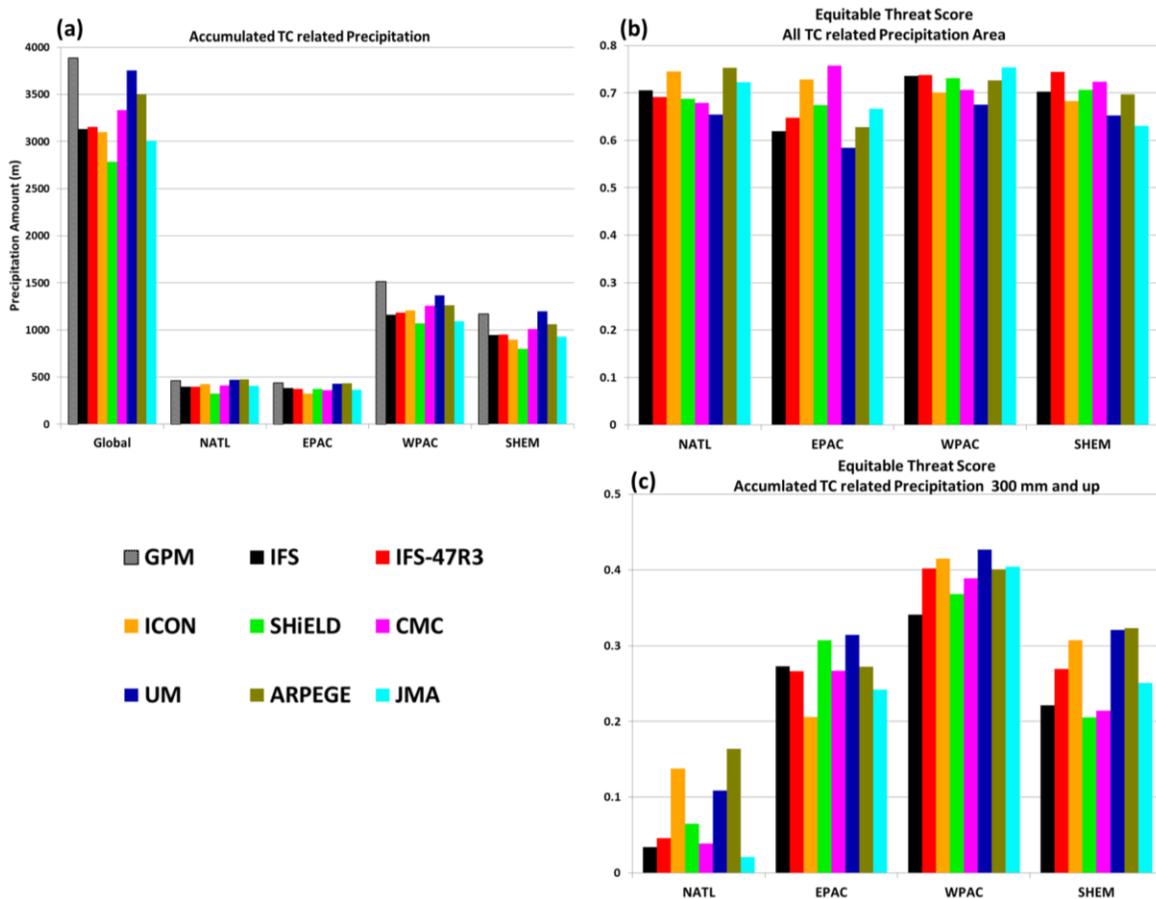
262 3.3 Forecasts of TC-related precipitation

263 Since the performance of TC intensity forecasts cannot be fully represented by the
 264 interpolated data, here, the TC-related precipitation is evaluated to provide another perspective
 265 on the forecasted TC characteristics in the models. Using the TC track information, the
 266 precipitation within 350 km of each TC center is used to investigate the TC-related precipitation
 267 for each model. Figure 6 shows the accumulated total precipitation for all TCs during the
 268 DIMOSIC period in each model compared to the Global Precipitation Measurement (GPM)
 269 observational data (Fig. 6a). The comparison shows that all models under-predict the amount of
 270 precipitation, especially in the most active areas of the WPAC and the EPAC. From a broad
 271 visual comparison, UM and SHIELD appear to have produced the highest and lowest amounts of
 272 precipitation among all the models, respectively. This can be confirmed by comparing the
 273 accumulated precipitation of models to the GPM data presented in Fig. 7a. The UM shows a
 274 much larger amount of precipitation than other models, followed by ARPEGE and CMC. In
 275 contrast, SHIELD shows the least TC-related precipitation among all models, except in the
 276 EPAC. The ranks of models are similar in different sub-regions, while the global precipitation
 277 amounts are dominated by those in the WPAC and the SHEM, as expected.

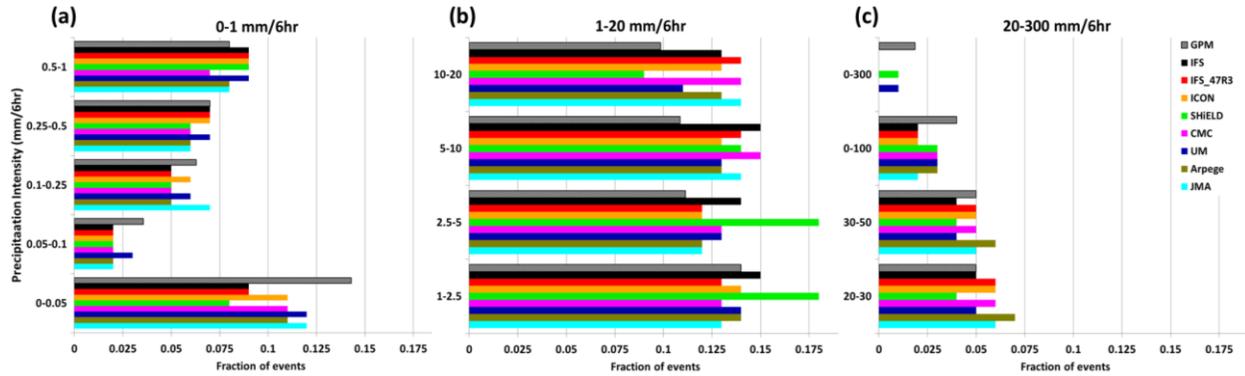


278
 279 FIGURE 6. Accumulated TC-related precipitation (unit: mm) for all TCs in the DIMOSIC period in (a) Global
 280 Precipitation Measurement (GPM) observations, (b) IFS, (c) IFS-47R3, (d) ICON, (e) SHIELD, (f) CMC, (g) UM,
 281 (h) ARPEGE, and (i) JMA.

282 To more objectively compare the forecasted locations of TC-related precipitation in each
 283 model to the GPM observations, the equitable threat scores (ETSs; Schaefer 1990) are computed.
 284 The ETSs for all TC-related precipitation areas in the four sub-regions for all eight models are
 285 compared in Fig. 7b. Note that although UM shows the closest precipitation amount to the GPM
 286 observation data (Fig. 7a), its ETS (skill) is lower than other models when considering all TC-
 287 related precipitation areas. This could be related to its relatively larger track errors (Fig.3) that
 288 cause the displacement of precipitation locations. However, for the areas with at least 300 mm of
 289 accumulated TC-related precipitation, the ETSs of UM are generally higher than those of other
 290 models (Fig. 7c). This is likely due to its relatively better prediction of precipitation amounts
 291 (Fig. 7a). In contrast, SHIELD under-predicts the precipitation amounts, but owing to its better
 292 track forecast in the EPAC (Fig. 3), it is able to achieve relatively higher ETSs in this sub-region
 293 (Figs. 7b,c). When comparing the two IFSs, Fig. 7 shows that their accumulated precipitation
 294 amounts are similar, but the newer IFS-47R3 generally had higher ETS scores (Figs. 7b,c).
 295 Finally, JMA shows relatively higher ETSs in the WPAC in both categories, while in the NATL,
 296 the highest ETSs in the two categories are achieved by ICON and ARPEGE (Figs. 7b,c).



297
 298 FIGURE 7. (a) Global accumulated TC-related precipitation (unit: m) and for the four sub-regions for all TCs in the
 299 DIMOSIC period. The equitable threat scores (ETSs) for (b) all TC-related precipitation areas and (c) the areas with
 300 300 mm and up accumulated TC-related precipitation. The GPM analysis data in (a) is shown in the bars with the
 301 grey checkerboard pattern. Abbreviations and colors used for the models and abbreviations used for the sub-regions
 302 on the abscissa are the same as in Fig. 3.



303

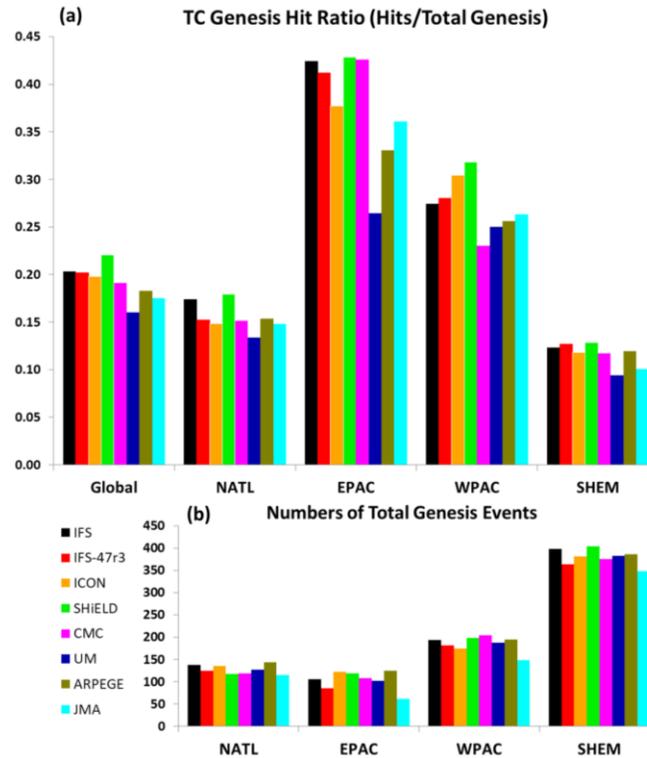
304 FIGURE 8. Fractions of precipitation events (GPM observations and model forecasts) in each precipitation intensity
 305 bin. (a) 5 bins for precipitation intensities from 0 to 1 mm(6 h)⁻¹. (b) four bins for precipitation intensities from 1 to
 306 20 mm(6 h)⁻¹. (c) four bins for precipitation intensities from 20 to 300 mm(6 h)⁻¹. Abbreviations and colors used for
 307 the model are the same as in Fig. 7.

308 TC-related precipitation based on different precipitation intensities was also analyzed.
 309 Figure 8 shows the fractions of precipitation events in different precipitation intensity bins. Most
 310 of the models under-predicted light (weaker than 0.5 mm (6 h)⁻¹; Fig. 8a) and heavy (stronger
 311 than 50 mm (6 h)⁻¹; Fig. 8c) precipitations, but over-predicted medium precipitation events (Fig.
 312 8b). Although SHIELD and UM were able to predicts some heavy precipitation events in the bin
 313 of 100-300 mm (6 h)⁻¹, SHIELD significantly over-predicted the events between 1-5 mm (6 h)⁻¹.
 314 The SHIELD development team at GFDL will closely examine the precipitation forecasts in the
 315 model in the near future, particularly to better isolate the possible reasons for these excessive
 316 precipitation amounts.

317

318 3.4 Forecasts of TC genesis

319 When a timeline contains an observed TC genesis, the track and intensity forecasts of the
 320 TC are verified based on the model forecasts initialized at or after the observed TC genesis time.
 321 In contrast, to investigate the models' performance for TC genesis, the 10-day forecast runs
 322 initialized before the observed TC genesis time which is based on the first "TD (tropical
 323 depression)" recorded in the ATCF best track data are considered. All TCs found by the GFDL
 324 simple tracker in these forecasts but not existing as TCs in the initial conditions for the forecast
 325 are counted as genesis events in the models. If a TC genesis has a track that "matches" an
 326 observed TC track, the genesis case is categorized as a "hit event". Otherwise, it is a "false
 327 alarm". The same criteria is used as in Chen et al. (2019b) to judge a model storm was a
 328 successful prediction of an observed genesis event.

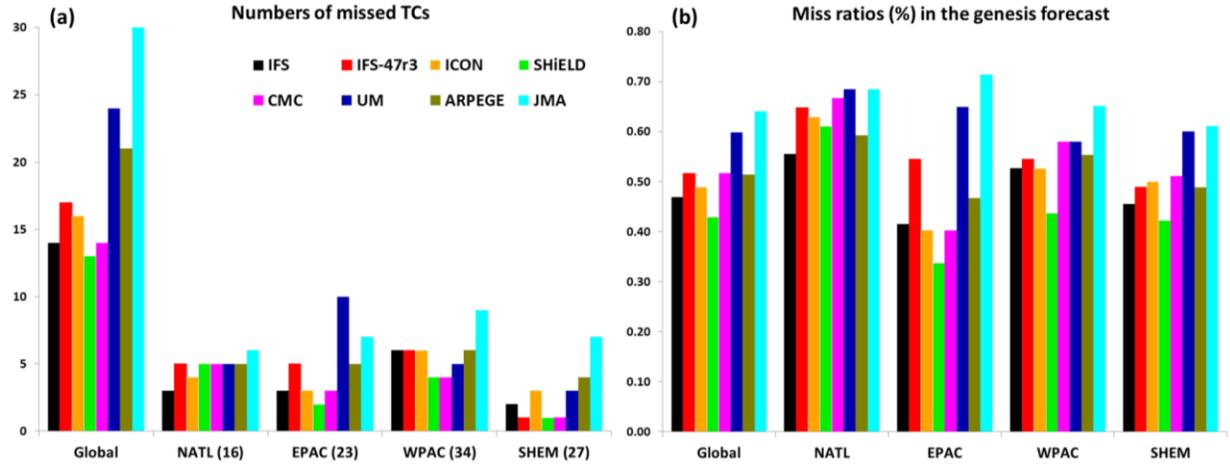


329

330 FIGURE 9. (a) Ratios of hit events to the total number of genesis events and (b) Numbers of total genesis events
 331 (sum of hit events and false alarms) for all models. Abbreviations and colors used for the models and abbreviations
 332 used for the sub-regions on the abscissa are the same as in Fig. 7.

333 Figure 9 shows the TC genesis ratios (hits to total predicted genesis events) and the
 334 number of total genesis events for each of the eight models in different regions. The sum of hit
 335 events and false alarms is equal to the number of total forecasted genesis events. We first find
 336 that all models show the highest hit ratios (Fig. 9a) with the fewest total genesis events in the
 337 EPAC. This indicates that models can predict TC genesis more skillfully in the EPAC than in
 338 other sub-regions. In contrast, models show the lowest hit ratios but the largest numbers of total
 339 genesis events in the SHEM, which indicates that models generate more false alarms in this sub-
 340 region than in the others. In general, SHIELD demonstrates the highest hit ratios both globally
 341 and in all sub-regions, followed by the two IFSs and ICON. CMC also shows high hit ratios in
 342 the EPAC. UM shows the lowest hit ratios in most regions with the exception of the WPAC.

343 During the DIMOSIC period, there were 16, 23, 34, and 27 TCs generated in the NATL,
 344 the EPAC, the WPAC, and the SHEM, respectively. However, not all observed TC geneses were
 345 predicted by the models. Figure 10a lists the number of TCs which were completely missed by
 346 models in each sub-region. It shows that JMA missed the genesis of total 30 TCs globally which
 347 is the most among all models. Most TC missed by JMA were in the WPAC and the SHEM. UM
 348 also missed many TC geneses in the EPAC. SHIELD had the least number of missed TCs
 349 globally followed by IFS and CMC. The newer IFS-47R3 missed more TCs than IFS in the
 350 NATL and the EPAC.



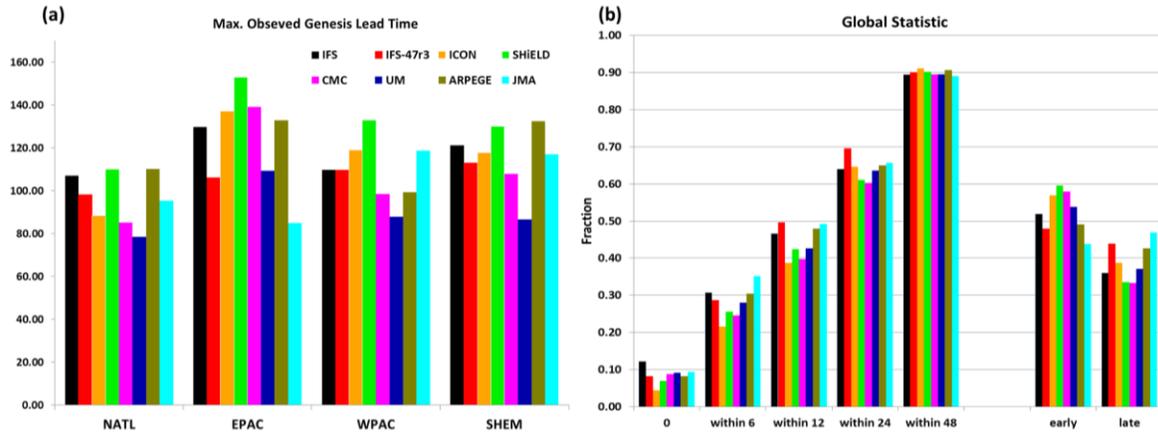
351
 352 FIGURE 10. (a) Numbers of missed TCs and (b) miss ratios (%) in the genesis forecasts from the eight models.
 353 Abbreviations and colors used for the models and abbreviations used for the sub-regions on the abscissa are the
 354 same as in Fig. 7. The numbers on the abscissa in (a) indicate the observed TC numbers in each sub-regions.

355 In the DIMOSIC period, models were initialized every three days. Hence, during the 10
 356 days before an observed TC genesis event, a model could have three or four 10-day forecasts
 357 initialized and these runs are expected to predict this genesis event. Therefore, besides counting
 358 the number of completely missed TCs, the “miss ratio” can be computed as the number of
 359 missing cases compared to the number of expected genesis hit events (Chen et al. 2019b). The
 360 miss ratios for each model in the different sub-regions are shown in Fig. 10b to better reveal the
 361 differences among the models. It shows that SHIELD shows the lowest miss ratios generally,
 362 except for in the NATL, while JMA and UM still struggle with relatively high miss ratios
 363 globally. We note that IFS-47R3 shows higher miss ratios than IFS in all sub-regions (Fig. 10b),
 364 including in the WPAC and the SHEM where IFS-47R3 shows the same or fewer numbers of
 365 completely missed TCs than IFS (Fig. 10a).

366 Following Chen et al. (2019b), beyond the scores of hit events, false alarms, and missing
 367 cases, we also investigate how precisely a model could predict the timing of TC genesis by
 368 comparing the “length of lead time” (see Fig. 8 in Chen et al. 2019b). The observed genesis lead
 369 time (OLT) is defined as the difference in time between the model initial time and the time at
 370 which observed TC genesis occurred. On the other hand, the time span from the model initial
 371 time to the model-predicted TC genesis lead time is referred as the model genesis lead time
 372 (MLT). The differences between the MLT and OLT (DMO) can indicate how accurate a model
 373 is in generating storms at the observed genesis time. If a model-predicted TC genesis occurred
 374 exactly at the observed TC genesis time, the DMO of this hit event is “zero”. A positive DMO
 375 means that the model hit event occurs later than the observed TC genesis time, while negative
 376 DMO values are associated with early initiation of the TC in the model

377 For each observed TC, it is expected that more than one hit event will happen in the set of
 378 10-day forecasts that cover the observed genesis time. To assess the predictive skill of each
 379 model, we only consider the maximum OLT, corresponding to the integration that identified the
 380 observed TC at the longest lead time. Figure 11a shows the mean values of the maximum OLT
 381 of all observed TCs in the four major sub-regions. In the NATL, both SHIELD and ARPEGE
 382 show a 110-hour OLT which is longer than the OLTs of other models, e.g. 78-hour OLT of UM.
 383 This indicates that SHIELD and ARPEGE could, on average, predict a hit TC genesis event 32

384 hours earlier than UM in the NATL. SHIELD also shows the earliest hit events in the EPAC and
 385 the WPAC, while ARPEGE shows the earliest hit events in the SHEM. The models in this study
 386 generally predict hit events earlier in the EPAC than in other sub-regions, except for JMA, which
 387 performs better in the WPAC and the SHEM than in other sub-regions. It is also interesting to
 388 see that IFS shows earlier hit events than the newer IFS-47R3 in most sub-regions.



389

390 FIGURE 11. (a) Mean values of maximum observed genesis lead time (in hours) of all storms in each sub-region for
 391 the eight models. Abbreviations used for the sub-regions on the abscissa are the same as in Fig. 7. (b) Fractions of
 392 global total hit events in each model that occurred within a certain DMO length. On the abscissa, “0” is for hit
 393 events which happened at the observed genesis time. “Within 6 (12, 24, or 48)” is for hit events with DMO lengths
 394 in 6 (12, 24, or 48) hours. “Early” is for all hit events with negative DMOs and “late” is for all hit events with
 395 positive DMOs. Abbreviations and colors used for the models in both (a) and (b) are the same as in Fig. 7.

396 Figure 11b shows the fraction of global total hit events in each model which occurred
 397 within a certain length of DMO (indicated on the abscissa). From the definition of DMO, a
 398 model showing more genesis cases with short DMOs indicates that the model could predict more
 399 accurate genesis timings of its hit events. It can be found that IFS shows the highest fraction
 400 among all of the eight models in the “zero” DMO length category. The results indicate that IFS
 401 shows the highest ratio of its hit events forecasted at the observed TC genesis time among the
 402 models. Besides the two IFSSs, JMA and ARPEGE also accurately predict TC genesis timings
 403 within the first three categories (“zero”, “within 6”, and “within 12”) which is the 24-hour
 404 window (12 hours before or after) centered on the observed TC genesis time. In contrast, ICON
 405 shows the smallest fractions in the first three DMO length categories, which implies that the
 406 accuracy of TC genesis timing of ICON is relatively low compared to the other models.

407 From the result of the “within 48” DMO in Fig. 11b, we can see that more than 89% of
 408 hit events in each of the models occur within the 48 hours before or after the observed TC
 409 genesis time. When comparing the results in the “early” and “late” categories, it is seen that most
 410 models forecast their hit events before the observed TC genesis time, except for JMA. The ratios
 411 of “early” to “late” cases are larger in SHIELD and CMC compared to other models, while the
 412 IFS-47R3 shows relatively even number of cases of hit events generated before or after the
 413 observed TC genesis time.

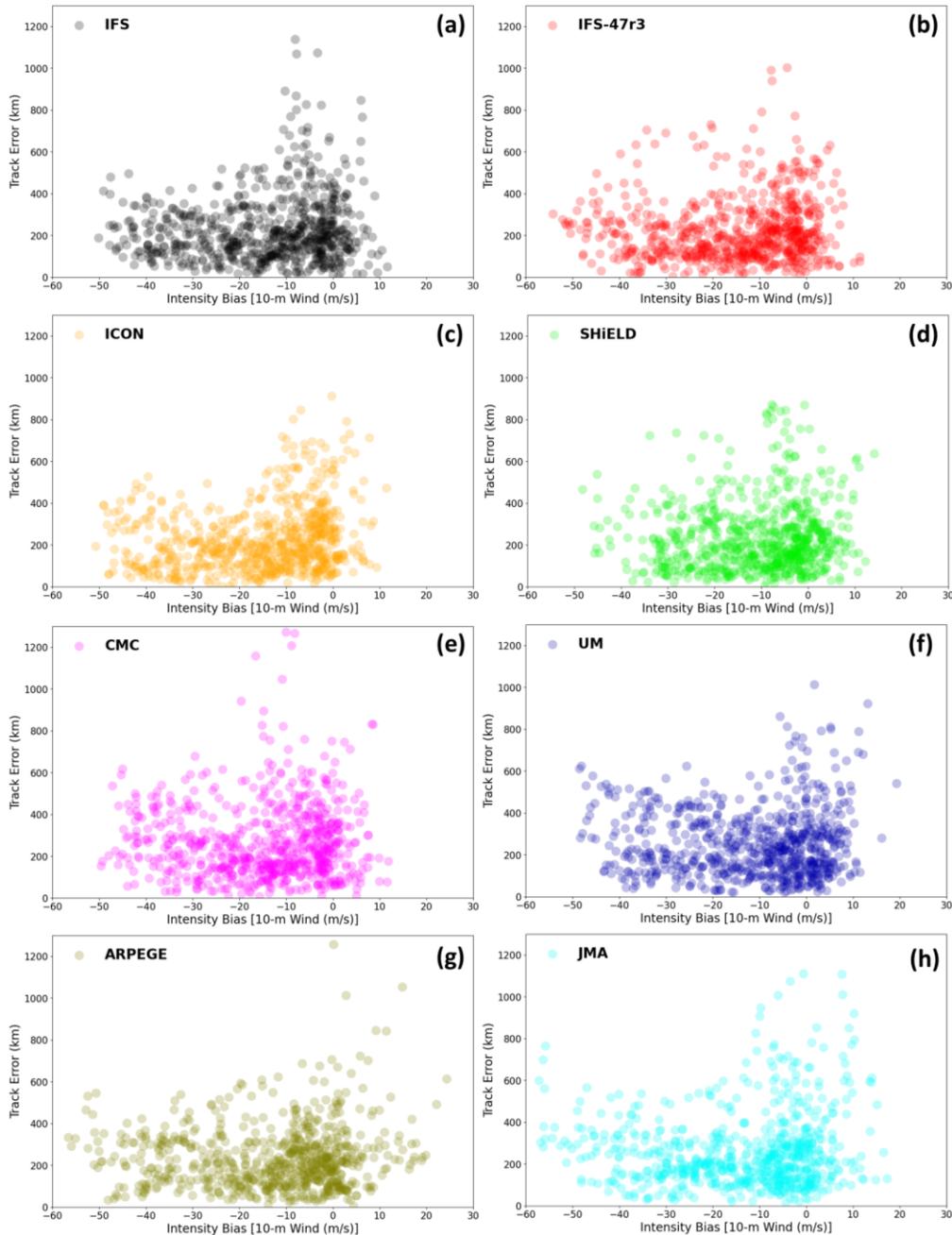
414 **4 Summary and Discussion**

415 The DIMOSIC project provides a great opportunity to engage the worldwide community
416 of medium-range modeling centers on cooperative model research and development. This study
417 investigated TC forecast skills in the eight participating global medium-range forecast models
418 during the year-long DIMOSIC period (June 2018 to June 2019). All models conducted 10-day
419 forecasts from the same initial conditions based on the ECMWF IFS model cycle 45R1. The
420 horizontal resolutions of the eight models ranged from 5 to 25km, and there were different
421 choices of dynamical cores and physics parameterizations across the models (Magnusson et al.
422 2022). The forecast skills of TC track and intensity have been presented for the eight models.
423 The TC-related precipitation and the performance of TC genesis forecasts have also been
424 evaluated.

425 Comparing the model forecasts to the observations for the 109 TCs in the DIMOSIC
426 period, IFS (45R1) and the updated version IFS-47R3 shows the best global averaged TC track
427 forecasts, followed by ICON and SHIELD. CMC shows a relatively higher error than others
428 before the 72-hour lead time. Based on our preliminary investigations, it could be related to the
429 initializing moisture shock given that the CMC has much moister analyses than IFS which may
430 induce convection collapses after the initialization with IFS ICs. For the TC track forecasts in
431 different sub-regions, UM and ICON show the lowest track errors in the NATL, while SHIELD
432 had the best track forecasts in the EPAC. In the WPAC and the SHEM, both IFS-47R3 and
433 ICON show the lowest TC track errors, followed by SHIELD. From the analyses of along-track
434 (AT) and cross-track (CT) errors, the models behave differently in different sub-regions. All
435 models showed a slow-moving bias and a poleward bias in the NATL and the WPAC, but in the
436 EPAC, except for JMA, most models show a fast-moving bias. In contrast, there are no
437 consistently slow- or fast-moving biases among models in the SHEM.

438 For TC intensity forecasts, based on the TC tracker results using the interpolated 0.5
439 degree resolution data, SHIELD performs relatively better than other models, followed by UM
440 and ARPEGE. From Table 2, we can see that the resolutions of the models range between 5 and
441 25 km, and the resolution of SHIELD is in the middle of that range. Therefore, the outperforming
442 TC intensity by SHIELD may imply that the resolution is not the only major factor limiting TC
443 intensity in global models. The use of dynamics and physics in the model also plays important
444 role. The performance of TC track and intensity forecasts could reveal some of the
445 characteristics of a model especially related to its dynamics and physics interactions. In Chen et
446 al. (2019b), it has been demonstrated that updating the GFS dynamical core to the nonhydrostatic
447 FV3 (Lin 2004; Putman and Lin 2007; Harris et al. 2020) can largely improve TC intensity
448 forecasts, and additional improvements in TC intensity and genesis forecasts were seen when
449 replacing the Zhao-Carr cloud microphysics scheme with the advanced GFDL cloud
450 microphysics scheme (Zhou et al. 2019).

451 Here, we attempt to probe into the characteristics of the models based on their biases of
452 TC track and intensity. Figure 12 shows the scatter plots of TC track and intensity errors for all
453 of the forecasts during the 72-120-hour lead time in each model. Some similarities can be found
454 in the scattered distributions of the two IFSs and ICON, including both ranges of intensity bias
455 and track error. This is consistent with the findings in Magnusson et al. (2022) that IFS and
456 ICON behave relatively similarly due to the sharing of partial physical parameterizations sharing
457 between ECWMF and DWD. In contrast, SHIELD, ARPEGE, and JMA show rather unique
458 patterns their own.



459

460 FIGURE 12. Scatter plot distribution of track errors (unit: km; ordinate) and intensity biases (the maximum 10-m
 461 wind speed; unit: m s^{-1} ; abscissa) of all forecasts during the lead time of 72-120 hour for (a) IFS, (b) IFS-47r3, (c)
 462 ICON, (d) SHIELD, (e) CMC, (f) UM, (g) ARPEGE, and (h) JMA.

463 Figure 12 also shows that in most models, forecasts with larger track error (>700 km) are
 464 usually accompanied by smaller intensity biases ($<10 \text{ m s}^{-1}$). In contrast, for those forecasts with
 465 larger intensity biases, their track errors are not consistently larger. At GFDL, it has been noticed
 466 that when the performance of TC track forecasts was improved by using an advection scheme
 467 with a stronger damping in the dynamics, a degradation of TC intensity was observed. The two-
 468 delta filter in the non-monotonic advection scheme and the monotonicity constraint in the tracer
 469 advection affect the model diffusivity which can also impact the diabatic heating and the location

470 of the TC deep convection relative to the eye (Gao et al. 2021). This was attributed to the impact
471 of stronger damping, which suppresses finer-scale features and activities, e.g. grid-scale
472 convection in the TCs, which further suppresses the TC intensities. An in-depth study of the
473 impact of grid-scale convection activity on TC track forecasts in SHIELD is in preparation.

474 Since the interpolated data cannot fully represent the performance of TC intensity
475 forecasts in the models, the TC-related precipitation was also evaluated to provide another
476 perspective on forecasted TC characteristics. Compared to the GPM observation data, all models
477 under-predict the amount of TC-related precipitation, especially in the Pacific Ocean. UM better
478 captures the regions with annual accumulated TC-related precipitation of more than 300mm
479 compared to the other models. However, when considering all TC-related precipitation areas, the
480 ETS of UM is generally lower than that of other models. This could be related to the relatively
481 large track errors of UM, especially in the EPAC. In contrast, SHIELD shows the largest dry bias
482 in TC-related precipitation among all models, but it still achieves relatively high ETSs in the
483 EPAC due to its better track forecasts in this sub-region. As to the intensity of TC-related
484 precipitation, most models over-predict medium precipitation events but under-predict light and
485 heavy precipitation events. Among all models, SHIELD noticeably over-predicts precipitation
486 events at the intensity of $1\text{-}5\text{ mm (6 h)}^{-1}$. The SHIELD development team at GFDL will take a
487 close look at its low precipitation amount and over-predicted medium intensity precipitation
488 events in the future.

489 The assessment of TC genesis forecast skill was based here on hits and misses, measures
490 that showed significant inter-model variability across the different sub-regions of interest. All
491 models show the highest hit ratios with the fewest total genesis events in the EPAC, which
492 indicates that models can better predict TC genesis in the EPAC. In contrast, models generate
493 more false alarms in the SHEM. SHIELD shows the highest hit ratios globally, followed by the
494 two IFSs and ICON. CMC also achieves high hit ratios in the EPAC. In contrast, UM shows
495 lower hit ratios than other models. As for the missed TC genesis cases, JMA missed 30 of the
496 100 observed TC geneses during the target period, which is the most among all models. Note that
497 JMA uses the coarsest model resolution among all participating models, which could impair its
498 TC genesis performance. In contrast, SHIELD shows the least number of missed TCs globally,
499 which may be benefiting from its better TC intensity forecast.

500 We also investigate how well the participating models predict the timing of TC genesis
501 by comparing the “length of lead time” proposed by Chen et al. (2019b). The results show that
502 models can generally accurately predict TC formation earlier in the EPAC than in other sub-
503 regions, except for JMA which predicts the WPAC and the SHEM hit events earlier than in other
504 sub-regions. SHIELD generally predicts the earliest hit events globally, while ARPEGE also
505 predicts TC genesis events earlier than other models in the NATL and the SHEM. Based on the
506 differences between the model genesis lead time and the observed genesis lead time, IFS shows
507 the most accurate timing of the TC genesis forecast, followed by JMA and ARPEGE. In contrast,
508 the accuracy of genesis timing in ICON is relatively lower than that of other models. We also
509 found that most models develop TCs earlier than observed in the best track, except for JMA. In
510 addition, more than 89% of hit events in each of the models occur within 48 hours of the
511 observed genesis time.

512 The comparison between IFS (version 45R1) and IFS-47R3 provides an opportunity to
513 examine the incremental change obtained for an upgrade of one model. The upgrade from
514 version 45R1 to 47R3 includes many changes in data assimilation and model physics. The

515 changes and meteorological impacts have been documented by ECMWF on the website:
516 <https://www.ecmwf.int/en/forecasts/documentation-and-support/changes-ecmwf-model>. One of
517 the listed impacts from the upgrade is the improvement of TC position errors. From our analysis,
518 the average track errors during the first 120 hours of IFS-47R3 are 2.5 to 10.8 km less than IFS
519 (Fig. 3) in the four major sub-regions, which is consistent with the ECMWF implementation
520 report. However, from the analyses of the along-track and cross-track errors, the biases of slow
521 and poleward movement are similar in these two model versions. Possibly associated with the
522 major upgrade to moist physics (Bechtold et al. 2020), IFS-47R3 shows a slightly larger negative
523 TC intensity biases than the older IFS version which was also found in Magnusson et al. (2021).
524 Although total precipitation predictions remain similar across the two IFS versions, the newer
525 IFS-47R3 achieves higher ETS scores for large accumulations, especially in the WPAC and the
526 SHEM. We also found that IFS-47R3 has more precipitation events with stronger precipitation
527 intensity (10-50 mm (6 h)⁻¹) than IFS. However, as to the TC genesis forecast, IFS-47R3 shows
528 some degradation from IFS, including more missed TC genesis event, higher miss ratios, shorter
529 genesis lead times, and less accuracy on TC genesis timing. These degradations in TC genesis
530 performance may be related to the weaker TC intensities in the newer version.

531 To summarize, in this study, extensive evaluation was made of the performance of TC
532 forecasts in the DIMOSIC models based on one year of model predictions initialized with the
533 same initial conditions. Although it is hard to precisely isolate the influence of individual
534 components in different model formulations on TC forecast skill in such an overview, the
535 comparisons based on different evaluation metrics highlight important similarities and
536 differences between the models. The results will be valuable for model developers in
537 participating centers as a benchmark of TC forecast skill with the impact of the initial condition
538 quality removed. Also, common forecast biases of the TC movement and TC-related
539 precipitations indicate general deficiencies in DIMOSIC models and point out a direction for
540 model developers for further model improvement.

541

542

543 **Acknowledgments**

544 The authors thank Kun Gao, Jie Chen, Morris Bender, and Tom Knutson for GFDL internal
545 review, and Lucas Harris, James Doyle, and Simon Lang for their comments helped to improve
546 this article. Authors also would like to thank other DIMOSIC participants, Duncan Ackerley,
547 Yves Bouteloup, K. C. Kwon, Yoonjin Lim, Mio Mastueda, Takumi Matsunobu, and Yamaguchi
548 Munehiko for their contribution to the DIMOSIC project.

549

550 **Open Research**

551 All DIMOSIC model interpolated data can be requested from Linus Magnusson. All TC
552 analyses are archived in the GFDL Tape Archive System at /archive/jhc/DIMOSIC/Analysis/TC
553 and can be requested from Jan-Huey Chen.

554

555 **References**

- 556 Bechtold, P., R. Forbes, I. Sandu, S. Lang, & M. Ahlgrimm (2020), A major moist physics
557 upgrade for the IFS. *ECMWF Newsletter*, 164, URL:
558 <https://www.ecmwf.int/en/newsletter/164/meteorology/major-moist-physics-upgrade-ifs>
- 559 Chen, J.-H., S.-J. Lin, L. Magnusson, M. A. Bender, X. Chen, L. Zhou, B. Xiang, S. L. Rees, M.
560 J. Morin, & L. Harris (2019a), Advancements in hurricane prediction with NOAA's next
561 generation forecast system. *Geophysical Research Letters*, 46 (8), 44954501,
562 doi:10.1029/2019GL082410.
- 563 Chen, J.-H., Lin, S., Zhou, L., Chen, X., Rees, S., Bender, M., & Morin, M. (2019b), Evaluation
564 of Tropical Cyclone Forecasts in the Next Generation Global Prediction System, *Monthly*
565 *Weather Review*, 147(9), 3409-3428, doi: 10.1175/MWR-D-18-0227.1
- 566 DWD (2022), ICON : Icosahedral Nonhydrostatic Weather and Climate Model. *Technical*
567 *Report, DWD*. URL: <https://code.mpimet.mpg.de/projects/iconpublic/wiki/Documentation>
- 568 ECMWF (2018), IFS Documentation CY45R1. *Technical Report, ECMWF*. URL:
569 <https://www.ecmwf.int/en/publications/ifs-documentation>
- 570 ECMWF (2021), IFS Documentation CY47R3. *Technical Report, ECMWF*. URL:
571 <https://www.ecmwf.int/en/publications/ifs-documentation>

- 572 Gao, K., L. Harris, L. Zhou, M. A. Bender, and M. Morin (2021), On the sensitivity of hurricane
573 intensity and structure to horizontal tracer advection schemes in FV3. *Journal of the*
574 *Atmospheric Sciences*, 78(9), doi:10.1175/JAS-D-20-0331.13007-3021.
- 575 Girard, C., & Coauthors (2014), Staggered vertical discretization of the Canadian Environmental
576 Multiscale (GEM) model using a coordinate of the log-hydrostatic-pressure type. *Monthly*
577 *Weather Review*, 142 (3), 1183.
- 578 Harris, L. M., S.-J. Lin, & C. Tu (2016), High-resolution climate simulations using GFDL
579 HiRAM with a stretched global grid. *Journal of Climate*, 29, 4293–4314, doi:10.1175/JCLI-
580 D-15-0389.11196, doi:10.1175/MWR-D-13-00255.1.
- 581 Harris, L., & Coauthors (2020), GFDL SHIELD: A unified system for weather to seasonal
582 prediction. *Journal of Advances in Modeling Earth Systems*, 12 (10),
583 doi:10.1029/2020MS002223.
- 584 Hong, Y., K.-L. Hsu, S. Sorooshian, & X. Gao (2004), Precipitation estimation from remotely
585 sensed imagery using an artificial neural network cloud classification system. *Journal of*
586 *Applied Meteorology*, 43 (12), 1834-1853, doi:10.1175/JAM2173.1.
- 587 JMA (2019), Outline of the operational numerical weather prediction at the Japan
588 Meteorological Agency, *Technical Report, Japan Meteorological Agency*. URL
589 <http://www.jma.go.jp/jma/jma-eng/jma-center/nwp/outline2019-nwp/index.htm>.
- 590 Lin, S.-J. (2004), A “vertically Lagrangian” finite-volume dynamical core for global models,
591 *Monthly Weather Review*, 132, 2293-2307.
- 592 Magnusson, L., J.-H. Chen, S.-J. Lin, L. Zhou, & X. Chen (2019), Dependence on initial
593 conditions vs. model formulations for medium-range forecast error variations. *Quarterly*
594 *Journal of the Royal Meteorological Society*, 145, doi:10.1002/qj.3545

- 595 Magnusson, L., & co-authors (2021), Tropical cyclone activities at ECMWF. *ECMWF Technical*
596 *Memorandum 888*
- 597 Magnusson, L., D. Ackerley, Y. Bouteloup, J.-H. Chen, J. Doyle, P. Earnshaw, Y. C. Kwon, M.
598 Köhler, S. T. K Lang, Y.-J. Lim, M. Matsueda, T. Matsunobu, R. McTaggart-Cowan, A.
599 Reinecke, M. Yamaguchi1, & L. Zhou (2022), Skill of medium-range forecast models using
600 the same initial conditions. *Bulletin of the American Meteorological Society*, 103(9), E2050-
601 E2068, doi:10.1175/BAMS-D-21-0234.1
- 602 McTaggartCowan, R., & Coauthors (2019), Modernization of atmospheric physics
603 parameterization in Canadian NWP. *Journal of Advances in Modeling Earth Systems*, 11
604 (11), 35933635, doi:10.1029/2019MS001781.
- 605 Miller, R. J., A. J. Scrader, C. R. Sampson, & T. L. Tsui (1990), The Automated Tropical
606 Cyclone Forecast System (ATCF). *Weather and Forecasting*, 5, 653–660.
- 607 Mogensen, K. S. , L. Magnusson , & J.-R. Bidlot (2017), Tropical cyclone sensitivity to ocean
608 coupling in the ECMWF coupled model: Tropical cyclone sensitivity. *J. Geophys. Res.*
609 *Oceans*, 122, 4392–4412, doi:10.1002/2017JC012753.
- 610 Pollard, R. T. , P. B. Rhines , & R. O. R. Y. Thompson (1973), The deepening of the wind-mixed
611 layer. *Geophys. Fluid Dyn.* , 4, 381–404, doi:10.1080/03091927208236105.
- 612 Putman, W. M., & S.-J. Lin (2007) Finite-volume transport on various cubed-sphere grids.
613 *Journal of computational physics*, 227, 55-78, doi:10.1016/j.jcp.2007.07.022.
- 614 Rabier, F., H. Jaärvinen, E. Klinker, J.-F. Mahfouf, & A. Simmons (2000), The ECMWF
615 operational implementation of four-dimensional variational assimilation. I: Experimental
616 results with simplified physics. *Quarterly Journal of the Royal Meteorological Society*, 126
617 (564), 24 1143–1170, doi:10.1002/qj.49712656415

- 618 Roehrig, R., & Coauthors (2020), The cnrm global atmosphere model arpegeclimat 6.3:
619 Description and evaluation. *Journal of Advances in Modeling Earth Systems*, 12 (7),
620 doi:10.1029/2020MS002075.
- 621 Schaefer, J. T. (1990) The critical success index as an indicator of warning skill. *Weather and*
622 *Forecasting*, 5 , 570–575.
- 623 Sampson, C. R., & A. J. Schrader (2000), The Automated Tropical Cyclone Forecasting System
624 (version 3.2). *Bulletin of the American Meteorological Society*, 81, 1231–1240.
- 625 Walters, D., & Coauthors (2019), The met office unified model global atmosphere 7.0/7.1 and
626 jules global land 7.0 configurations. *Geoscientific Model Development*, 12(5), 1909-1963,
627 doi:10.5194/gmd-12-1909-2019.
- 628 Zeng, X. , & A. Beljaars (2005), A prognostic scheme of sea surface skin temperature for
629 modeling and data assimilation: Sea surface skin temperature scheme. *Geophysical*
630 *Research Letters*, 32, L14605, doi:10.1029/2005GL023030.
- 631 Zhou, L., S.-J. Lin, J.-H. Chen, L. M. Harris, X. Chen, & S. Rees (2019), Toward convective
632 scale prediction within the Next Generation Global Prediction System. *Bulletin of the*
633 *American Meteorological Society*. 100 (7): 1225-43, doi:10.1175/BAMS-D-17-0246.1

1
2 **Tropical Cyclone Forecasts in the DIMOSIC Project – Medium-Range Forecast**
3 **Models with Common Initial Conditions**

4 **Jan-Huey Chen^{1,2}, Linjiong Zhou³, Linus Magnusson⁴, Ron McTaggart-Cowan⁵, and**
5 **Martin Köhler⁶**

6 ¹National Oceanic and Atmospheric Administration/Geophysical Fluid Dynamics Laboratory,
7 Princeton, NJ, USA

8 ²University Corporation for Atmospheric Research, Boulder, CO, USA

9 ³Atmospheric and Oceanic Sciences Program, Princeton University, Princeton, NJ, USA

10 ⁴European Centre for Medium-Range Weather Forecasts, Reading, UK

11 ⁵Environment and Climate Change Canada, Montreal, Canada

12 ⁶Deutsche Wetterdienst, Offenbach, Germany

13
14 Corresponding author: Jan-Huey Chen (Jan-Huey.Chen@noaa.gov)
15

16 **Key Points:**

- 17
- 18 • Tropical cyclone forecasts are compared between global medium-range models from
19 leading modeling centers initialized with identical data.
 - 20 • Similarities and differences between the models set a benchmark of TC forecast with the
21 impact of the initial condition quality removed.
 - 22 • Common TC forecast biases indicate general deficiencies in the models and suggest a
direction for further model improvement.

23 **Abstract**

24 The Tropical cyclone (TC) forecast skill of the eight global medium-range forecast models
25 which are participating in the DIMOSIC (Different Models, Same Initial Conditions) project is
26 investigated in this study. Each model was used to generate 10-day forecasts from the same
27 initial conditions provided by the European Centre for Medium-Range Weather Forecasts. There
28 are a total of 123 initial dates spanning in one year from June 2018 to June 2019 with a 3-day
29 interval. The TC track and intensity forecasts are evaluated against the best track dataset. TC-
30 related precipitation and tropical cyclogenesis forecasts are also compared to explore the
31 differences and similarities of TC forecasts across the models. This comparison of TC forecasts
32 allows model developers in different centers to benchmark their model against other models,
33 with the impact of the initial condition quality removed. The verifications reveal that most
34 models show slow-moving and right-of-track biases in their TC track forecasts. Also, a common
35 dry bias in TC-related precipitation indicates a general deficiency in TC intensity and convection
36 in the models which should be related to insufficient model resolution. These findings provide
37 important references for future model developments.

38

39 **Plain Language Summary**

40 Despite recent improvements in our ability to predict the track and intensity of tropical cyclones,
41 these storms remain significant forecasting challenges. Forecasters rely heavily on the guidance
42 generated by numerical weather prediction systems, making the reliability of these systems
43 essential for accurate forecasts during these high-impact weather events. As a result,
44 improvement the quality of tropical cyclone guidance is an important numerical model
45 development objective. In this study, the TC forecast skills in the eight global medium-range
46 forecast models from the model development centers/institutes who participated in the
47 DIMOSIC (Different Models, Same Initial Conditions) project are examined. All models were
48 initialized from the same data provided by the ECMWF (European Centre for Medium-Range
49 Weather Forecasts) to investigate the differences and similarities among their TC forecasts
50 without the impact of the quality of initial conditions. Besides the general TC forecast
51 evaluation metrics including errors and biases of the track and intensity, the TC-related
52 precipitation and TC genesis skills are also evaluated to comprehensively explore the
53 performance of TC forecasts among all models. The comparison allows model developers in
54 different centers to benchmark their model against other participating models. Moreover, the
55 verification results provide important references for future model developments.

56

57 **1 Introduction**

58 Tropical cyclone (TC) prediction is an important mission for weather and climate
59 agencies in many countries. Over the past few decades, numerical models have become the most
60 important tools for operational centers to make TC forecasts on weather and sub-seasonal to
61 seasonal time scales. Therefore, improving the model performance of TC forecasts has been one
62 of the leading tasks in most operational centers or modeling research institutes working on model
63 development. In addition, the accurate depiction of physical processes that lead to a better TC

64 forecast in the model are also relevant to interesting scientific questions in the atmospheric
65 science research area more broadly.

66 The quality of initial conditions has a leading impact on short- to medium-range forecast
67 skill, including for TC forecasts. In Chen et al. (2019a), the fvGFS (finite volume Global
68 Forecasting System) model developed at the Geophysical Fluid Dynamics Laboratory (GFDL)
69 initialized with the European Centre for Medium-Range Weather Forecasts (ECMWF) Integrated
70 Forecasting System (IFS) data showed much-improved TC track forecasts for the 2017 Atlantic
71 hurricane season compared to its retrospective forecasts initialized with the data from the
72 National Centers for Environmental Prediction (NCEP) Global Forecast System (GFS) version
73 14. In Magnusson et al. (2019), the same approach was used, comparing the GFDL fvGFS model
74 forecasts to those from the IFS and GFS. The results showed that the choice of initial conditions
75 clearly dominated the forecast quality in the medium-range predictions, but that the model
76 formulation could also play a significant role.

77 Since major model development centers mostly develop their modeling systems
78 independently, the DIMOSIC project (Different Models, Same Initial Conditions; Magnusson et
79 al. 2022) was devised to investigate the relationship between the choice of model formulation
80 and forecast quality. Models developed by different world-leading modeling centers were
81 initialized from the same initial condition. In Magnusson et al. (2022), the differences and
82 similarities of the forecasts among the models were presented. The results found that some pairs
83 of models behaved more similarly than other pairs due to their sharing of partial physical
84 parameterizations, e.g. ECWMF IFS and DWD (Deutsche Wetterdienst) ICON (Icosahedral
85 Non-hydrostatic Model). On the other hand, ICON and GFDL SHiELD (System for High-
86 Resolution Prediction on Earth-to-Local Domains) showed relatively large forecast differences,
87 while both ranking among the best models. Regarding the influences from model formulations
88 on the forecasts, however, it was difficult to point out a single model component that had the
89 strongest impact on the forecast differences. Also, as pointed out by Magnusson et al. (2022) the
90 interaction between different model parameterizations and their respective configurations could
91 play a significant role as well.

92 In this study, the performance of TC forecasts from the DIMOSIC models is evaluated.
93 The TC track and intensity forecast skills among the models during the period of June 2018 to
94 June 2019 are compared. Since TC intensity in interpolated data does not reflect the actual TC
95 intensity at the native model resolution, the TC-related precipitation are also evaluated to provide
96 another perspective on forecasted TC activities for better exploring the differences and
97 similarities among the models. Also, the forecast skill of TC genesis was investigated by
98 comparing the hit/false alarm ratios among the models, as well as using the method based on the
99 lengths of TC genesis lead time introduced in Chen et al. (2019b) to examine the accuracy of TC
100 genesis timing in the model forecasts. These comparisons should be valuable for model
101 developers in different centers to benchmark their model's performance on TC forecasts against
102 that of other models, with the impact of the initial condition quality removed.

103 The models participating in the DIMOSIC project are introduced in section 2 which also
104 describes the observation data and methodology used in this study. The comparisons of track,
105 intensity, TC-related precipitation, and genesis forecasts among the models are contained in
106 section 3. Summary and discussion are presented in section 4.

107 2 DIMOSIC models, forecasts, and verification data

108 General information on the numerical models and their developing centers/institutes
 109 participating in the DIMOSIC project are listed in Table 1. The horizontal resolutions and the
 110 number of vertical levels of the models and their key references are included. Some
 111 centers/institutes submitted more than one model configurations to the project, but we only
 112 investigate one configuration of the model for each center/institute based on their suggestions.
 113 The only exception is to include two versions of IFS 45R1 and 47R3, to provide an example of
 114 the incremental change obtained for an upgrade of one model. For the sea surface temperature
 115 evolution in the models, the two IFSs used a partial coupling to the 3D ocean NEMO model
 116 (Mogensen et al. 2017), SHIELD is coupled with a 1D mixed layer ocean model (Pollard et al.
 117 1973), CMC used a thermodynamic mixed layer ocean model (Zeng and Beljaars 2005), and
 118 others used persistent anomalies from the analysis. Other detailed configurations of each model
 119 including dynamical cores and major physical parameterizations can be found in the sub-section
 120 of “Model descriptions” in the section of “Models and data” and in Table 2 in Magnusson et al.
 121 (2022), but is not repeated in this paper.

Acronyms	Models	Centers/Institutes	Resolution	Key references
ARPEGE	Action de Recherche Petite Echelle Grande Echelle (version: 46T1)	Meteofrance	5-25 km 105 levels	Roehrig et al. (2020)
CMC	Global Environmental Multiscale Model (GEM) (version: v5.0.2)	Canadian Meteorological Center (CMC)	15 km 80 levels	Girard et al. (2014) McTaggart-Cowan et al. (2019)
ICON	Icosahedral Non-hydrostatic Model (version: April 21)	Deutsche Wetterdienst (DWD)	13km 90 levels	DWD (2022)
IFS/ IFS-47R3	Integrated Forecasting System (versions: 45R1 and 47R3)	European Centre for Medium- range Weather Forecasts (ECMWF)	9 km 137 levels	ECMWF (2018, 2021)
JMA	Global Spectral Model (GSM) (version: GSM1705)	Japan Meteorological Agency (JMA)	20 km 100 levels	JMA (2019)
SHIELD	System for High-Resolution Prediction on Earth-to-Local Domains (version: rt2019)	Geophysical Fluid Dynamics Laboratory (GFDL)	13 km 91 levels	Harris et al. (2020)
UM	Unified Model	UK Met Office	10 km 70 levels	Walters et al. (2019)

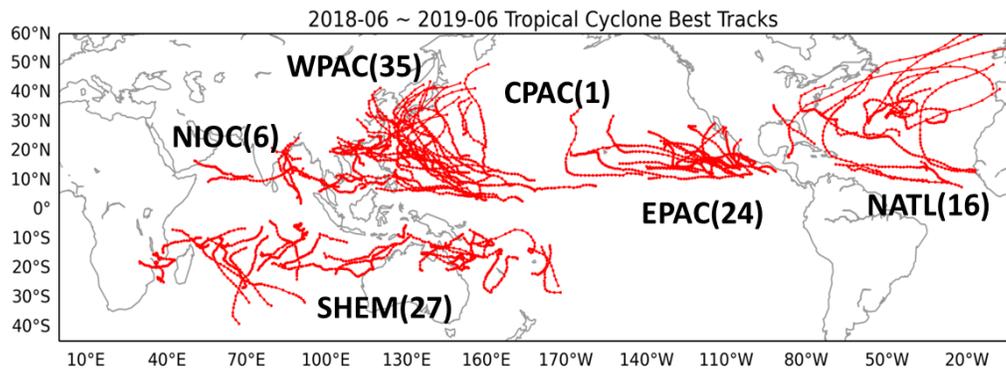
122 TABLE 1. The eight DIMOSIC models (the ECWMF contributed two IFS configurations). From left to right: model
 123 acronyms, model full names (and versions if applicable), centers/institutes that models belong to, horizontal
 124 resolutions and the number of vertical levels used in the models, and the key references of the models.

125 All models conducted 10-day forecasts from the same initial conditions: ECMWF
 126 operational analyses based on the IFS model cycle 45R1 (ECMWF 2018). The 9-km analyses on
 127 137 vertical levels are generated from a 4DVar data assimilation system (Rabier et al. 2000). All
 128 participating institutes received the interpolated data at 0.1 degree for their model initialization.
 129 Detailed procedures for handling the initial conditions in each model are described in section 2c
 130 and Table 3 of Magnusson et al. (2022).

131 The total 123 forecasts are conducted with initialization dates spanning one year from
 132 June 2018 to June 2019 at 3-day intervals. The 10-day forecast outputs from each model were
 133 interpolated to a common 0.5 degree grid using an average interpolation method available in the
 134 EcCodes/MIR package (<https://confluence.ecmwf.int/display/ECC/ecCodes+Home>). The GFDL
 135 simpler tracker (Harris et al. 2016) was used with a warm-core criterion to track TCs in the

136 forecasts of the eight models based on the fields of sea-level pressure, 10-m wind speed, 850-hPa
137 vorticity, and mean temperature between 500-300 hPa.

138 There were 109 observed TCs in the DIMOSIC period. Their storm tracks based on the
139 best track data (b-deck) in Automated Tropical Cyclone Forecast (ATCF) dataset (Miller et al.
140 1990; Sampson and Schrader 2000) are shown on the map in Fig 1. There were 35 TCs in the
141 northwest Pacific basin (WPAC), 24 in the northeast Pacific basin (EPAC), and 16 in the North
142 Atlantic basin (NATL). In the Southern Hemisphere (SHEM), there were 27 TCs in the
143 combined South Indian Ocean and South Pacific Ocean. Besides the overview of global
144 analyses, the individual TC forecast skill in the above four major sub-regions will be investigated
145 and compared with each other in the following sections.



146
147 FIGURE 1. All TCs in the DIMOSIC period. The best tracks are from the ATCF dataset. The numbers of TCs are
148 indicated in the brackets next to the acronyms of the six regions: WPAC: northwest Pacific basin; EPAC: northeast
149 Pacific basin; CPAC: north-central Pacific basin; NATL: North Atlantic basin; NIOC: north Indian Ocean; SHEM:
150 Southern Hemisphere.

151 The ATCF dataset was used to evaluate the forecast errors of TC track and intensity, and
152 the skill of TC genesis forecasts in the models at six-hour intervals. For the TC-related
153 precipitation, the NASA Global Precipitation Measurement (GPM) Integrated Multi-satellitE
154 Retrievals for GPM (IMERG) dataset was used (Hong et al. 2004) for verification. To equally
155 compare to the model output field of total precipitation, this high-resolution (0.1 degree) satellite
156 observational dataset was interpolated into 0.5 degree.

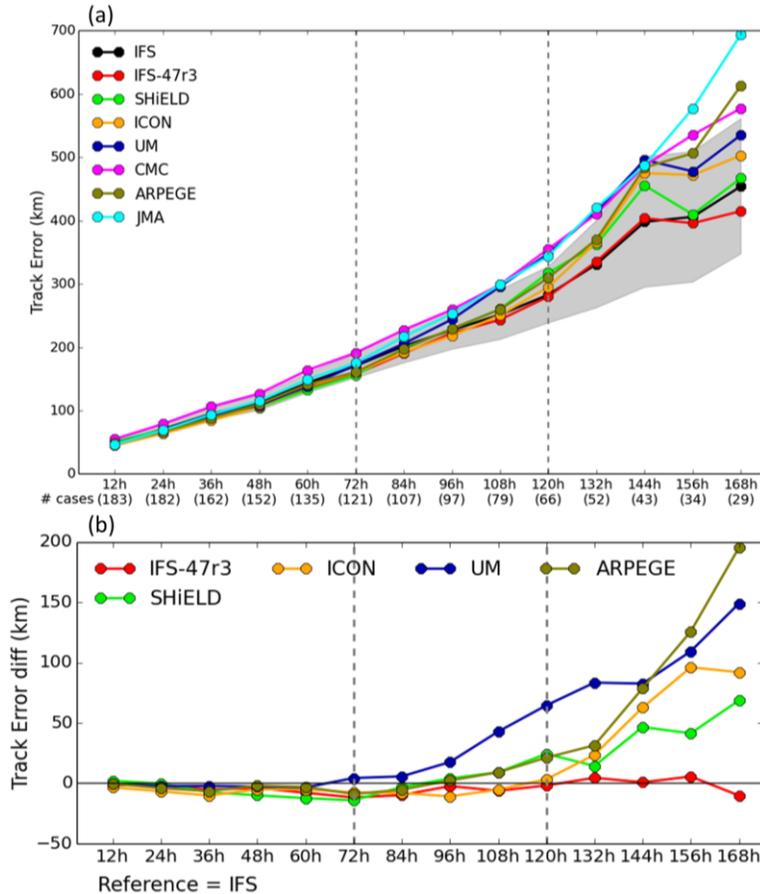
157

158 3. Results

159 3.1 TC track forecasts

160 The prediction of TC path, or track, is the one of the most important factors for taking
161 necessary precautions against possible impacts from hurricanes or typhoons. The homogeneous
162 comparisons of global mean TC track forecast errors along with the forecast lead time at 12-h
163 intervals are shown in Fig. 2a. The differences among the models are small through the 72-hour
164 lead time with the exception of CMC which shows a relatively higher error than others. During
165 the 72 to 120-hour head time, TC track errors diverge into two groups. The two IFSs, ICON,
166 SHIELD, and ARPEGE show lower errors than UM, CMC, and JMA. Both versions of IFS show
167 the lowest TC track errors after the 120-hour lead time all the way to the 168-hour lead time,
168 followed by ICON and SHIELD. Note that most models are within the 95% confidence levels of
169 IFS, which means the differences in forecast skills are not statistically significant. To highlight

170 the difference between the leading models, Fig. 2b shows the differences in TC track errors of
 171 five of the models compared to the IFS. The newer IFS-47R43 performs slightly better (negative
 172 values in track error differences) or equivalently during the entire 7 days. In the early lead times
 173 (36-84 hours), most of the five models also show slightly better forecasts than IFS. ICON
 174 displays similar (or slightly better) skill to the two IFS versions until the 120-hour lead time.
 175 Both SHiELD and ARPEGE perform very well before the 96-hour lead time. After the 120-hour
 176 lead time, SHiELD shows lower forecast errors than other models except for the two IFSs.

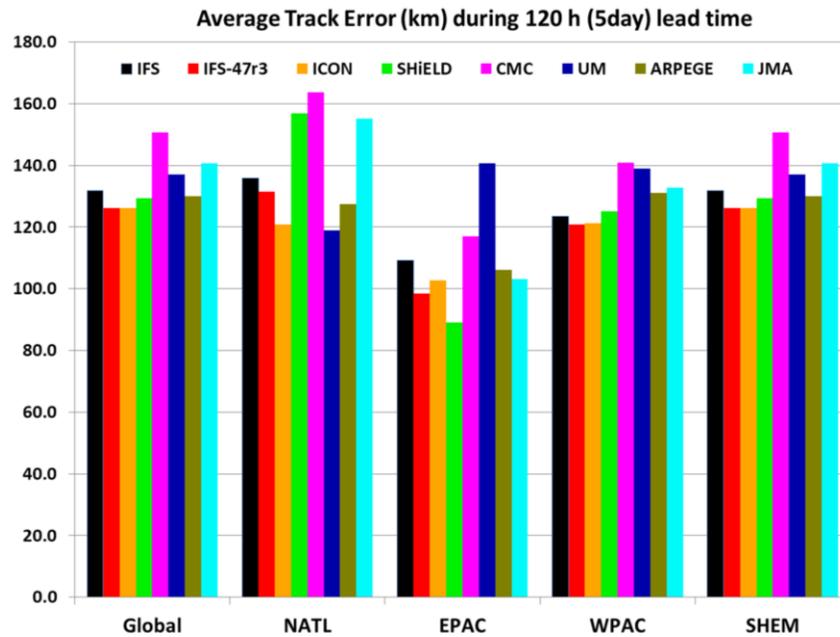


177

178 FIGURE 2. (a) Global mean TC track forecast errors (km) at every 12 hour forecast lead time for IFS (black), IFS-
 179 47r3 (red), SHiELD (green), ICON (yellow), UM (blue), CMC (magenta), ARPEGE (grass green), and JMA (light
 180 blue). The 95% confidence levels for IFS are indicated by the gray color shading. Numbers of homogeneous cases
 181 for individual lead times are listed in the brackets at the bottom of each abscissa. Vertical gray dotted lines indicate
 182 72 and 120 hour forecast lead times. (b) Global mean TC track forecast error differences of IFS-47r3 (red), SHiELD
 183 (green), ICON (yellow), UM (blue), and ARPEGE (grass green) comparing to IFS.

184 Figure 3 shows the 5-day average TC track errors for all models in the entire globe and in
 185 the four major sub-regions individually. For the two IFSs, IFS-47R3 shows lower track errors
 186 than IFS globally and in all major sub-regions. ICON shows competitive low track errors to IFS-
 187 47R3 in the WPAC and the SHEM, and a very low track error in the NATL. SHiELD performs
 188 the lowest track error in the EPAC, and low track errors besides the two IFSs and ICON in the
 189 SHEM and the WPAC. However, SHiELD has a much larger track error in the NATL, which
 190 results from the slow moving bias shown in the forecasts of Hurricane Florence and the bias of
 191 direction of motion shown in the forecasts of Hurricane Leslie. Detailed investigations can be
 192 found in Text S1 and Figs. S1-S3. The performance of the TC track forecast in UM is notable.

193 The model shows the lowest track error in the NATL but the largest track error in the EPAC,
 194 while its track errors in the WPAC and the SHEM are also at the high end compared to other
 195 models. JMA performs competitively to IFS-47R3 and ICON in EPAC, but not in other sub-
 196 regions. With regard to CMC, it shows larger track errors than other models in most sub-regions
 197 except for in the EPAC during this targeted year. From Fig. 3, we can also see that the globally-
 198 averaged track forecast errors among the models is dominated by the errors in the WPAC and the
 199 SHEM since the majority of the TCs were found in these two sub-regions.



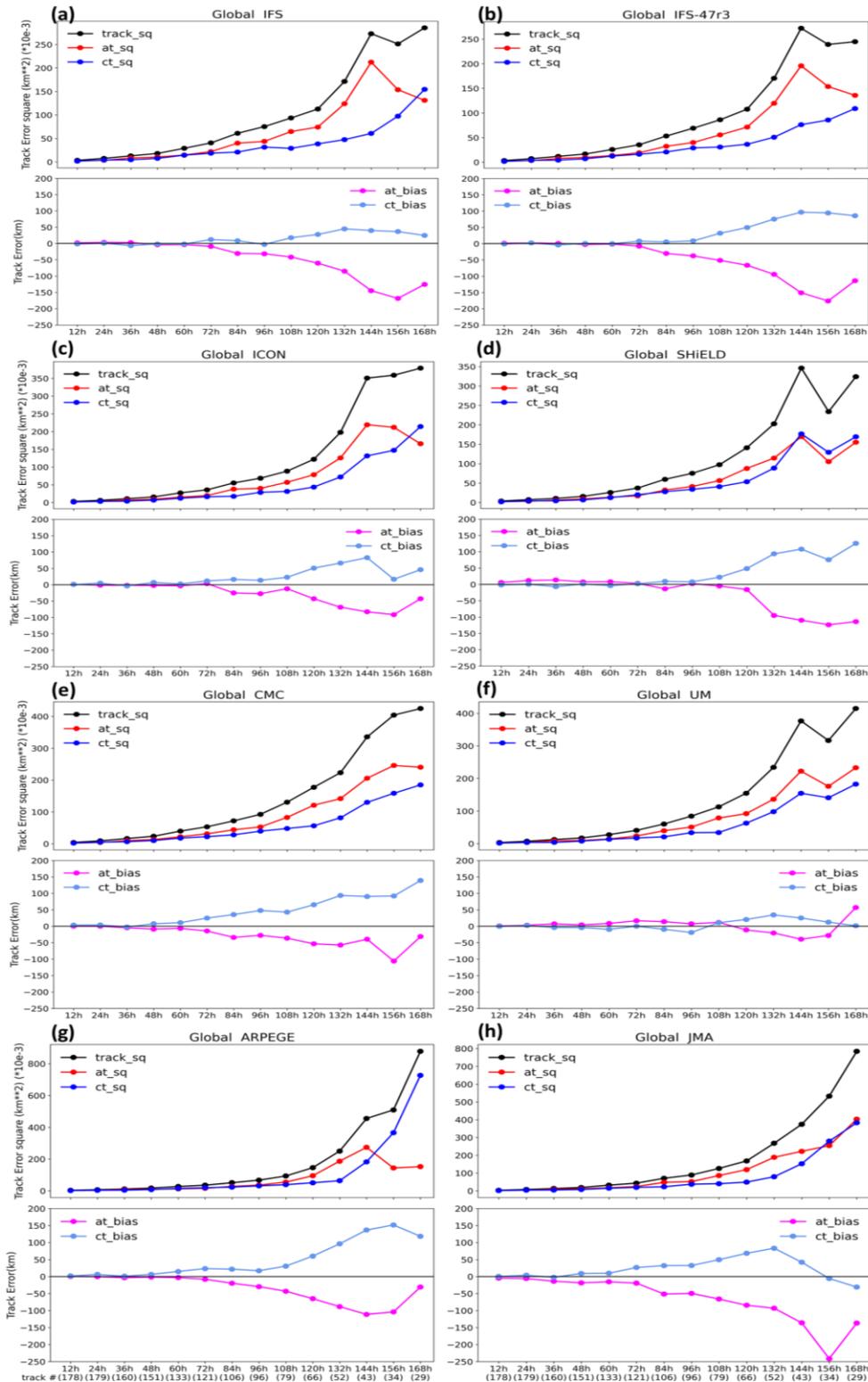
200

201 FIGURE 3 Averaged Track errors (km) in globe and 4 sub-regions during the 120-hour lead time for the 8 models.
 202 Abbreviations and colors used for the models are the same as in Fig. 2a. Abbreviations used for the sub-regions on
 203 the abscissa are the same as in Fig. 1.

204 The sources of track errors can be due to biases either in forecasts of the TC translational
 205 speed or the TC direction of motion. The lower sub-panels in Fig. 4 show the globally averaged
 206 along-track (AT) errors and cross-track (CT) errors (perpendicular to the track) for all models.
 207 Both AT and CT errors are calculated as great circle distances. Most of the models start to show
 208 negative AT biases and positive CT biases during 72 to 120-hour lead time, which indicates that
 209 the TC track errors during the later lead times are mostly due to the slow and northward (for an
 210 easterly moving TC) moving biases. In general, UM shows the smallest AT and CT biases
 211 among all models, and the biases of SHIELD take place at longer forecast lead times than other
 212 models.

213 By the Pythagorean Theorem the square of the total error equals the squares of the AT
 214 and CT errors (Chen et al. 2019a). The squares of total track errors, CT errors, and AT errors are
 215 plotted in upper sub-panels in Fig. 4 to illustrate the proportion of contributions from the AT and
 216 CT errors respectively to the total error. From Figs. 4a,b, we can see that the AT error
 217 contributes more than the CT error to the total track error in the two IFSs, which indicates that
 218 the track errors in the IFSs are dominated by the slow-moving bias (negative AT biases). The
 219 characteristic of the consistently larger AT error square than the CT error square can also be
 220 found in ICON, CMC, and UM, but the differences between their AT and CT error squares are

221 smaller than those in the IFSs. SHIELD, ARPEGE, and JMA show relatively closer AT and CT
222 error squares, especially SHIELD. This indicates that the contributions of the slow and
223 northward moving biases to the total track errors are similar in these models. For ARPEGE, the
224 CT error is the dominant track error in late lead times, while the AT error contributes more to the
225 JMA's total TC track errors during the 120 to 144-hour lead time.



226

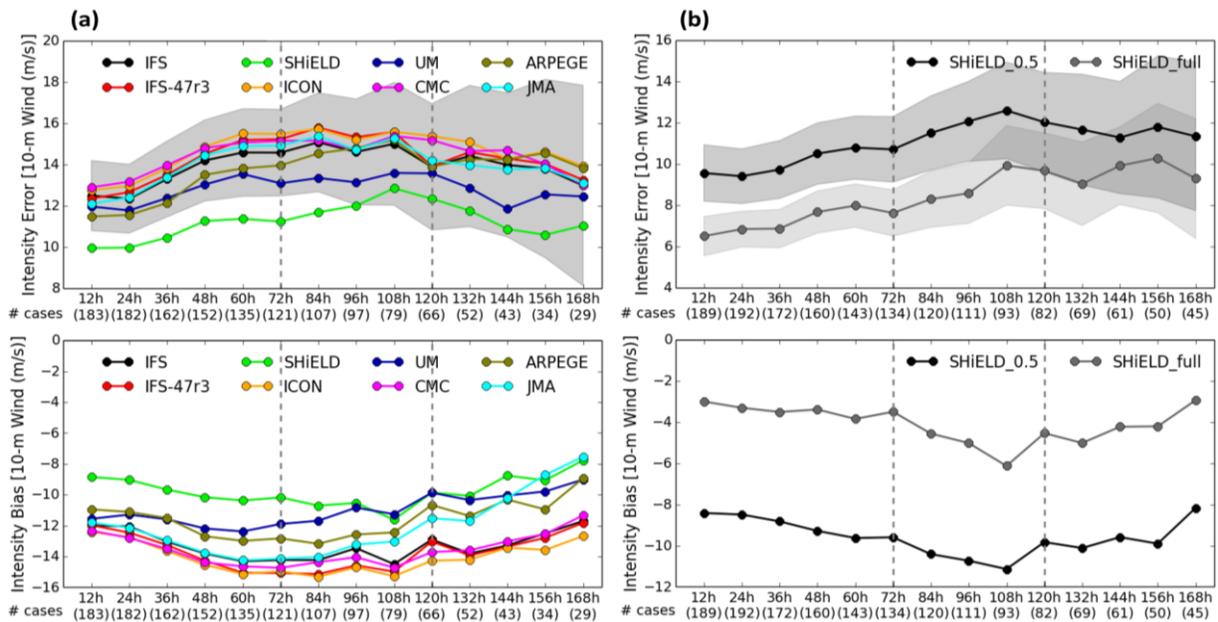
227 FIGURE 4 Global analyses of along-track (AT) error and cross-track (CT) error for (a) IFS, (b) IFS-47R3, (c)
 228 ICON, (d) SHiELD, (e) CMC, (f) UM, (g) ARPEGE, and (h) JMA. The squares of total track errors (black), along-
 229 track errors (red), and cross-track errors (blue) are in the upper panels for each model. The biases of along-track
 230 (magenta) and cross-track (light blue) errors are in the lower panels. Numbers of homogeneous cases for each lead
 231 time are listed at the bottom of lower panels.

232 It is also found that models have different AT and CT errors in different sub-regions. All
 233 models showed a slow-moving bias and a poleward bias in the NATL and the WPAC, but in the
 234 EPAC, except for JMA, most models show a fast-moving bias. In contrast, there are no
 235 consistently slow or fast moving biases among models in the SHEM. Detailed analyses of AT
 236 and CT errors for all eight models in the four major sub-regions can be found in Text S2 and
 237 Figs. S4-S7.

238

239 **3.2 TC intensity forecasts**

240 It has been more challenging to predict TC intensity than track, especially for global
 241 models which usually cannot resolve fine scale interactions between thermal dynamics and
 242 dynamics due to insufficient resolutions. As outlined in Section 2, the 10-day forecast outputs
 243 were interpolated to a common 0.5 degree for each model. Therefore, the model-predicted TC
 244 intensities found by the tracker are underestimated due to the low data resolution. However, it is
 245 still of interest to compare the relative differences of TC intensities among the models for the
 246 interpolated output data. The global mean TC intensity errors and biases based on the maximum
 247 10-m wind speed are presented in Fig. 5a. SHIELD predicts a much stronger TC intensity than
 248 other models, followed by UM and ARPEGE. Figure 5b uses SHIELD as an example to
 249 demonstrate the differences between the TC intensities obtained from the native resolution (13
 250 km grid) outputs and in the interpolated 0.5 degree resolution data. The average differences of
 251 total error and bias are between 3 to 5 ms^{-1} .

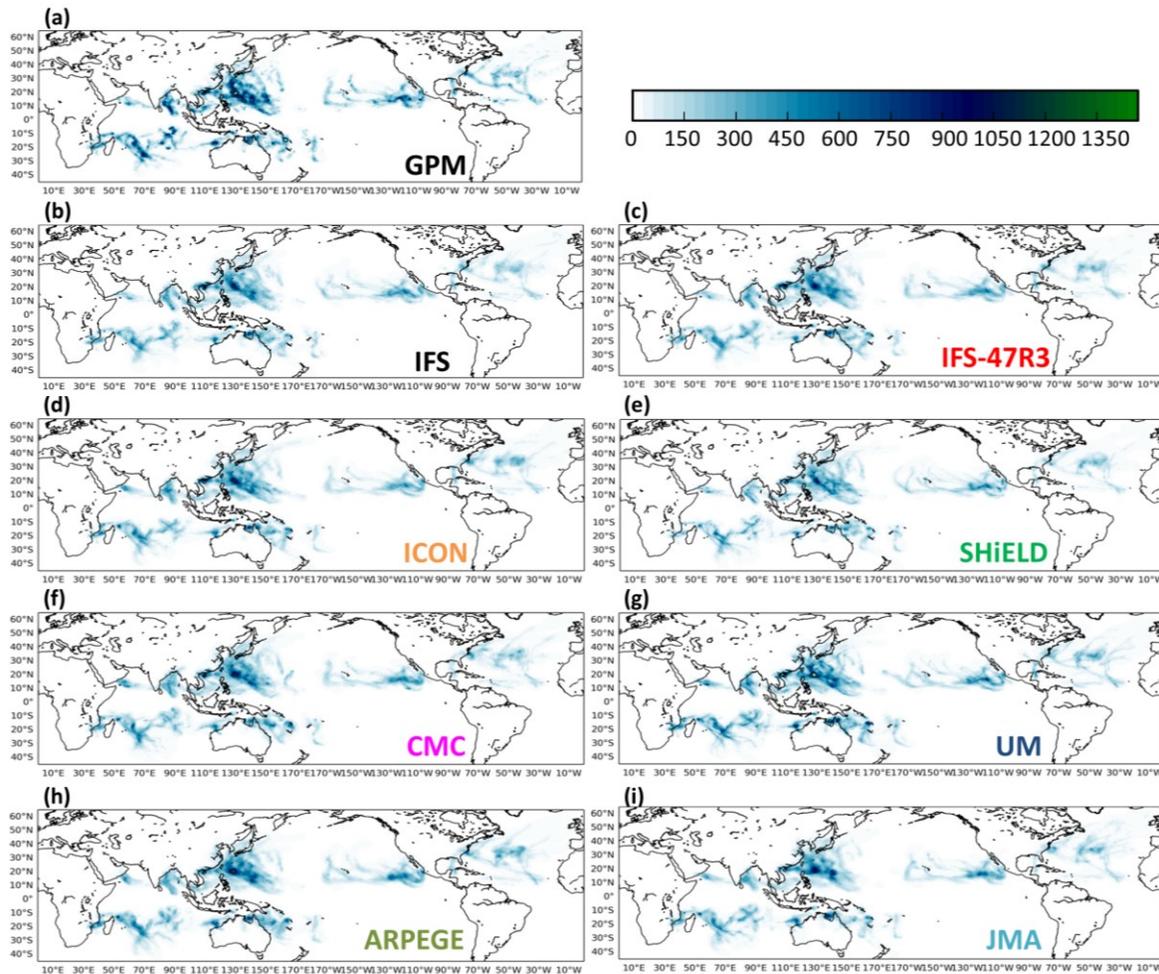


252

253 **FIGURE 5.** (a) Global mean TC intensity errors and biases. Upper panel: Absolute error of the maximum 10-m wind
 254 speed (m s^{-1}) along with the model forecast lead time for 8 models. Abbreviations and colors used for the models are
 255 the same as in Fig. 2a. The 95% confidence levels for IFS are indicated by the gray color shading. Numbers of
 256 homogeneous cases for individual lead times are listed in the brackets at the bottom. Vertical grey dotted lines
 257 indicate 72 hour and 120 hour lead times. Lower panel: As in the upper panel, but for the bias of the maximum 10-m
 258 wind speed (m s^{-1}). (b) As in (a), but for SHIELD native resolution data (black; SHIELD_full) and SHIELD 0.5
 259 degree interpolated data (gray; SHIELD_0.5). The 95% confidence levels for each resolution data are indicated by
 260 the same medium and light transparent grey shading areas, with their overlapping region denoted by dark grey
 261 shading.

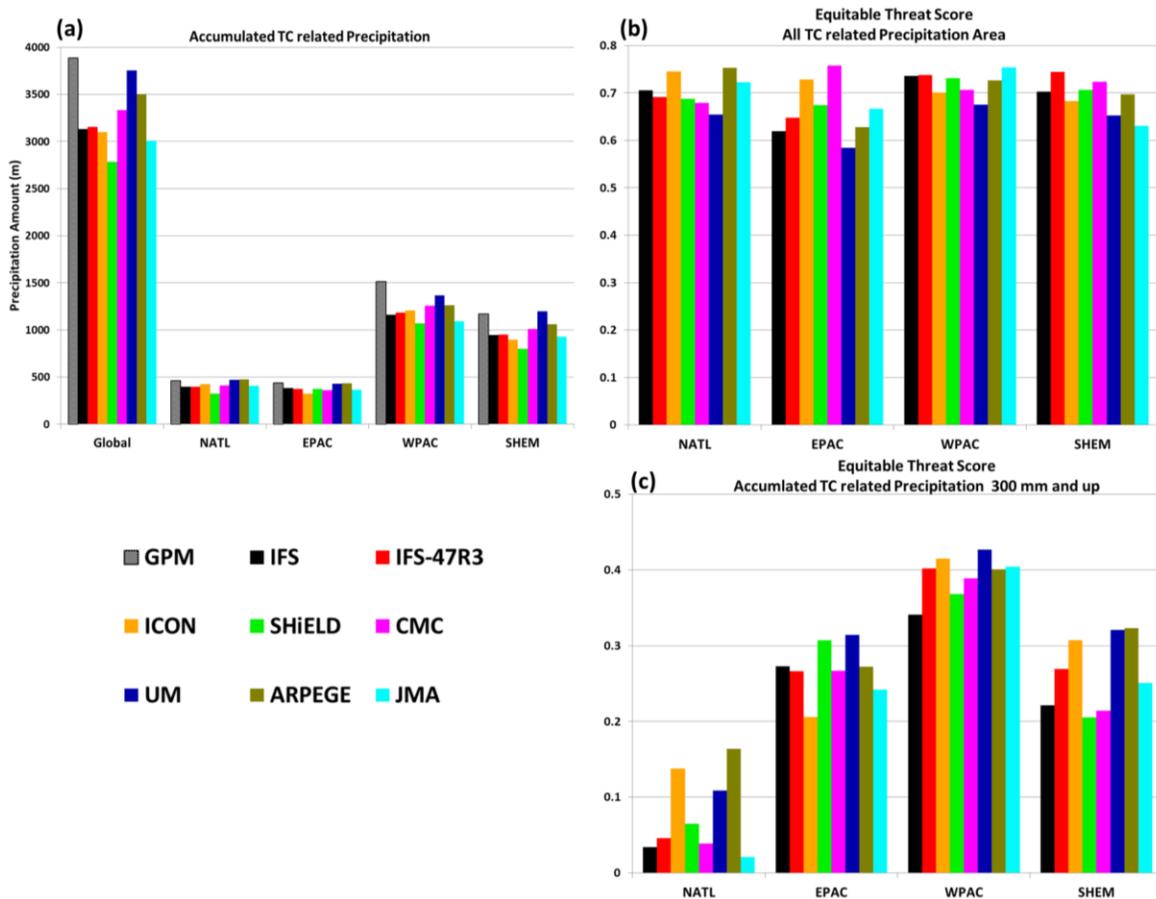
262 3.3 Forecasts of TC-related precipitation

263 Since the performance of TC intensity forecasts cannot be fully represented by the
 264 interpolated data, here, the TC-related precipitation is evaluated to provide another perspective
 265 on the forecasted TC characteristics in the models. Using the TC track information, the
 266 precipitation within 350 km of each TC center is used to investigate the TC-related precipitation
 267 for each model. Figure 6 shows the accumulated total precipitation for all TCs during the
 268 DIMOSIC period in each model compared to the Global Precipitation Measurement (GPM)
 269 observational data (Fig. 6a). The comparison shows that all models under-predict the amount of
 270 precipitation, especially in the most active areas of the WPAC and the EPAC. From a broad
 271 visual comparison, UM and SHIELD appear to have produced the highest and lowest amounts of
 272 precipitation among all the models, respectively. This can be confirmed by comparing the
 273 accumulated precipitation of models to the GPM data presented in Fig. 7a. The UM shows a
 274 much larger amount of precipitation than other models, followed by ARPEGE and CMC. In
 275 contrast, SHIELD shows the least TC-related precipitation among all models, except in the
 276 EPAC. The ranks of models are similar in different sub-regions, while the global precipitation
 277 amounts are dominated by those in the WPAC and the SHEM, as expected.

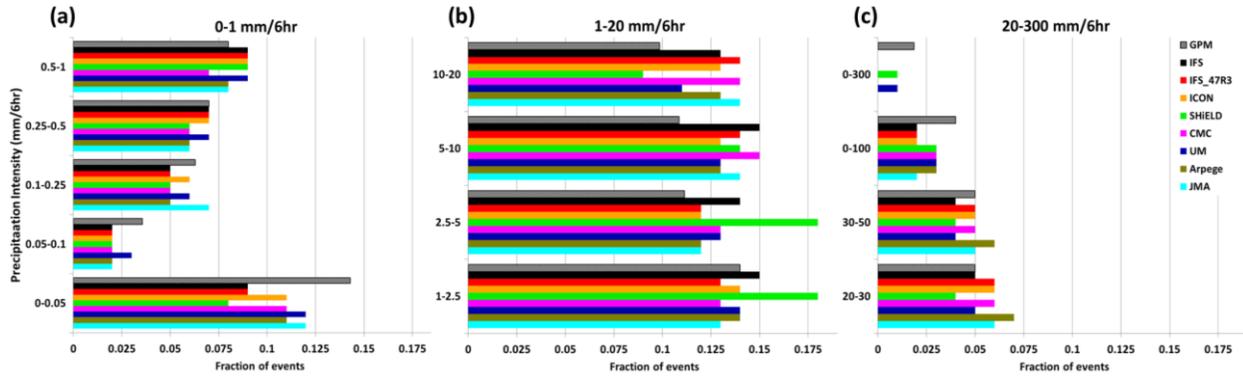


278
 279 FIGURE 6. Accumulated TC-related precipitation (unit: mm) for all TCs in the DIMOSIC period in (a) Global
 280 Precipitation Measurement (GPM) observations, (b) IFS, (c) IFS-47R3, (d) ICON, (e) SHIELD, (f) CMC, (g) UM,
 281 (h) ARPEGE, and (i) JMA.

282 To more objectively compare the forecasted locations of TC-related precipitation in each
 283 model to the GPM observations, the equitable threat scores (ETSs; Schaefer 1990) are computed.
 284 The ETSs for all TC-related precipitation areas in the four sub-regions for all eight models are
 285 compared in Fig. 7b. Note that although UM shows the closest precipitation amount to the GPM
 286 observation data (Fig. 7a), its ETS (skill) is lower than other models when considering all TC-
 287 related precipitation areas. This could be related to its relatively larger track errors (Fig.3) that
 288 cause the displacement of precipitation locations. However, for the areas with at least 300 mm of
 289 accumulated TC-related precipitation, the ETSs of UM are generally higher than those of other
 290 models (Fig. 7c). This is likely due to its relatively better prediction of precipitation amounts
 291 (Fig. 7a). In contrast, SHIELD under-predicts the precipitation amounts, but owing to its better
 292 track forecast in the EPAC (Fig. 3), it is able to achieve relatively higher ETSs in this sub-region
 293 (Figs. 7b,c). When comparing the two IFSs, Fig. 7 shows that their accumulated precipitation
 294 amounts are similar, but the newer IFS-47R3 generally had higher ETS scores (Figs. 7b,c).
 295 Finally, JMA shows relatively higher ETSs in the WPAC in both categories, while in the NATL,
 296 the highest ETSs in the two categories are achieved by ICON and ARPEGE (Figs. 7b,c).



297
 298 FIGURE 7. (a) Global accumulated TC-related precipitation (unit: m) and for the four sub-regions for all TCs in the
 299 DIMOSIC period. The equitable threat scores (ETSs) for (b) all TC-related precipitation areas and (c) the areas with
 300 300 mm and up accumulated TC-related precipitation. The GPM analysis data in (a) is shown in the bars with the
 301 grey checkerboard pattern. Abbreviations and colors used for the models and abbreviations used for the sub-regions
 302 on the abscissa are the same as in Fig. 3.



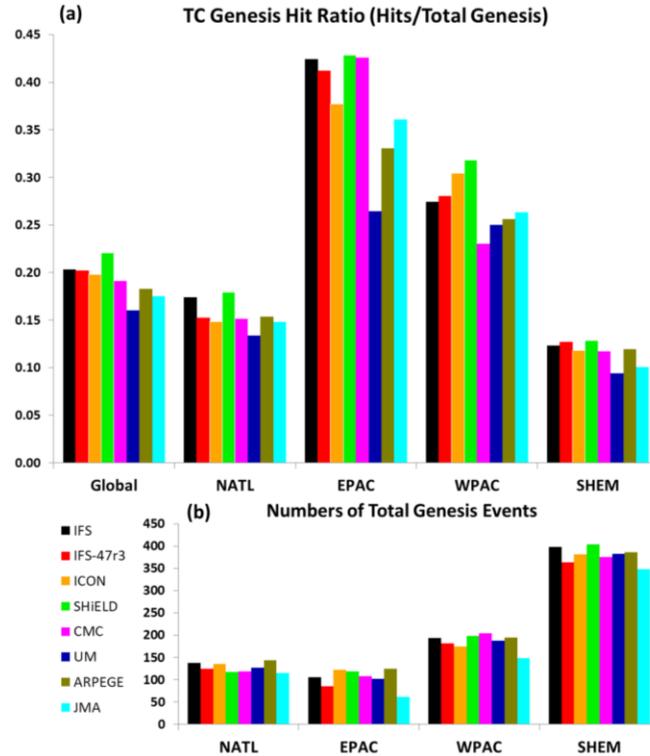
303
 304 FIGURE 8. Fractions of precipitation events (GPM observations and model forecasts) in each precipitation intensity
 305 bin. (a) 5 bins for precipitation intensities from 0 to 1 mm(6 h)⁻¹. (b) four bins for precipitation intensities from 1 to
 306 20 mm(6 h)⁻¹. (c) four bins for precipitation intensities from 20 to 300 mm(6 h)⁻¹. Abbreviations and colors used for
 307 the model are the same as in Fig. 7.

308 TC-related precipitation based on different precipitation intensities was also analyzed.
 309 Figure 8 shows the fractions of precipitation events in different precipitation intensity bins. Most
 310 of the models under-predicted light (weaker than 0.5 mm (6 h)⁻¹; Fig. 8a) and heavy (stronger
 311 than 50 mm (6 h)⁻¹; Fig. 8c) precipitations, but over-predicted medium precipitation events (Fig.
 312 8b). Although SHIELD and UM were able to predicts some heavy precipitation events in the bin
 313 of 100-300 mm (6 h)⁻¹, SHIELD significantly over-predicted the events between 1-5 mm (6 h)⁻¹.
 314 The SHIELD development team at GFDL will closely examine the precipitation forecasts in the
 315 model in the near future, particularly to better isolate the possible reasons for these excessive
 316 precipitation amounts.

317

318 3.4 Forecasts of TC genesis

319 When a timeline contains an observed TC genesis, the track and intensity forecasts of the
 320 TC are verified based on the model forecasts initialized at or after the observed TC genesis time.
 321 In contrast, to investigate the models' performance for TC genesis, the 10-day forecast runs
 322 initialized before the observed TC genesis time which is based on the first "TD (tropical
 323 depression)" recorded in the ATCF best track data are considered. All TCs found by the GFDL
 324 simple tracker in these forecasts but not existing as TCs in the initial conditions for the forecast
 325 are counted as genesis events in the models. If a TC genesis has a track that "matches" an
 326 observed TC track, the genesis case is categorized as a "hit event". Otherwise, it is a "false
 327 alarm". The same criteria is used as in Chen et al. (2019b) to judge a model storm was a
 328 successful prediction of an observed genesis event.

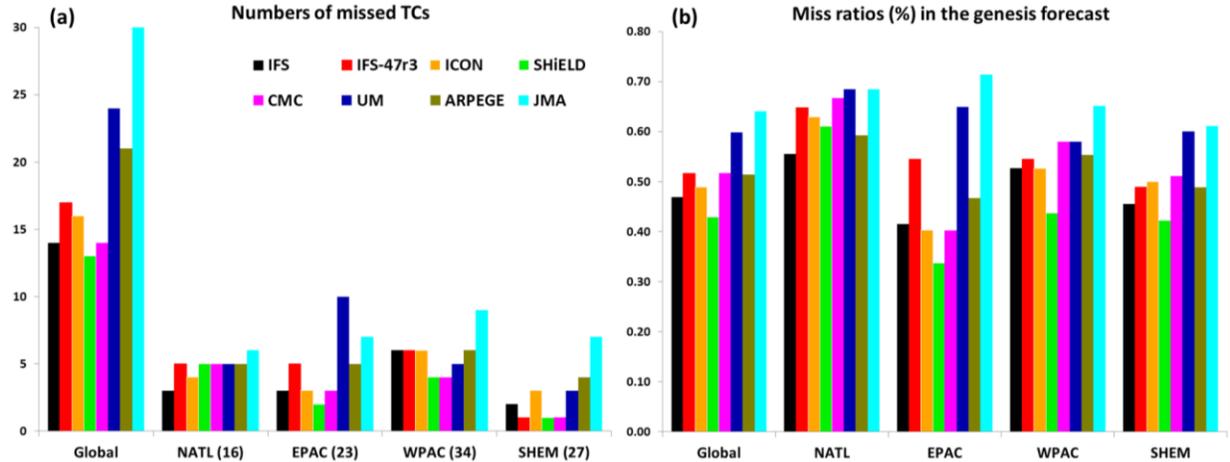


329

330 FIGURE 9. (a) Ratios of hit events to the total number of genesis events and (b) Numbers of total genesis events
 331 (sum of hit events and false alarms) for all models. Abbreviations and colors used for the models and abbreviations
 332 used for the sub-regions on the abscissa are the same as in Fig. 7.

333 Figure 9 shows the TC genesis ratios (hits to total predicted genesis events) and the
 334 number of total genesis events for each of the eight models in different regions. The sum of hit
 335 events and false alarms is equal to the number of total forecasted genesis events. We first find
 336 that all models show the highest hit ratios (Fig. 9a) with the fewest total genesis events in the
 337 EPAC. This indicates that models can predict TC genesis more skillfully in the EPAC than in
 338 other sub-regions. In contrast, models show the lowest hit ratios but the largest numbers of total
 339 genesis events in the SHEM, which indicates that models generate more false alarms in this sub-
 340 region than in the others. In general, SHIELD demonstrates the highest hit ratios both globally
 341 and in all sub-regions, followed by the two IFSs and ICON. CMC also shows high hit ratios in
 342 the EPAC. UM shows the lowest hit ratios in most regions with the exception of the WPAC.

343 During the DIMOSIC period, there were 16, 23, 34, and 27 TCs generated in the NATL,
 344 the EPAC, the WPAC, and the SHEM, respectively. However, not all observed TC geneses were
 345 predicted by the models. Figure 10a lists the number of TCs which were completely missed by
 346 models in each sub-region. It shows that JMA missed the genesis of total 30 TCs globally which
 347 is the most among all models. Most TC missed by JMA were in the WPAC and the SHEM. UM
 348 also missed many TC geneses in the EPAC. SHIELD had the least number of missed TCs
 349 globally followed by IFS and CMC. The newer IFS-47R3 missed more TCs than IFS in the
 350 NATL and the EPAC.



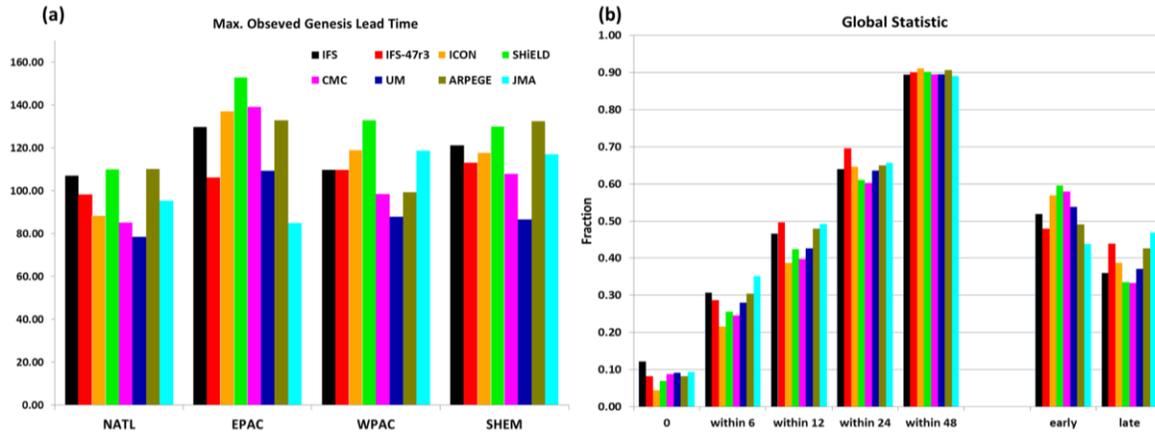
351
 352 FIGURE 10. (a) Numbers of missed TCs and (b) miss ratios (%) in the genesis forecasts from the eight models.
 353 Abbreviations and colors used for the models and abbreviations used for the sub-regions on the abscissa are the
 354 same as in Fig. 7. The numbers on the abscissa in (a) indicate the observed TC numbers in each sub-regions.

355 In the DIMOSIC period, models were initialized every three days. Hence, during the 10
 356 days before an observed TC genesis event, a model could have three or four 10-day forecasts
 357 initialized and these runs are expected to predict this genesis event. Therefore, besides counting
 358 the number of completely missed TCs, the “miss ratio” can be computed as the number of
 359 missing cases compared to the number of expected genesis hit events (Chen et al. 2019b). The
 360 miss ratios for each model in the different sub-regions are shown in Fig. 10b to better reveal the
 361 differences among the models. It shows that SHIELD shows the lowest miss ratios generally,
 362 except for in the NATL, while JMA and UM still struggle with relatively high miss ratios
 363 globally. We note that IFS-47R3 shows higher miss ratios than IFS in all sub-regions (Fig. 10b),
 364 including in the WPAC and the SHEM where IFS-47R3 shows the same or fewer numbers of
 365 completely missed TCs than IFS (Fig. 10a).

366 Following Chen et al. (2019b), beyond the scores of hit events, false alarms, and missing
 367 cases, we also investigate how precisely a model could predict the timing of TC genesis by
 368 comparing the “length of lead time” (see Fig. 8 in Chen et al. 2019b). The observed genesis lead
 369 time (OLT) is defined as the difference in time between the model initial time and the time at
 370 which observed TC genesis occurred. On the other hand, the time span from the model initial
 371 time to the model-predicted TC genesis lead time is referred as the model genesis lead time
 372 (MLT). The differences between the MLT and OLT (DMO) can indicate how accurate a model
 373 is in generating storms at the observed genesis time. If a model-predicted TC genesis occurred
 374 exactly at the observed TC genesis time, the DMO of this hit event is “zero”. A positive DMO
 375 means that the model hit event occurs later than the observed TC genesis time, while negative
 376 DMO values are associated with early initiation of the TC in the model

377 For each observed TC, it is expected that more than one hit event will happen in the set of
 378 10-day forecasts that cover the observed genesis time. To assess the predictive skill of each
 379 model, we only consider the maximum OLT, corresponding to the integration that identified the
 380 observed TC at the longest lead time. Figure 11a shows the mean values of the maximum OLT
 381 of all observed TCs in the four major sub-regions. In the NATL, both SHIELD and ARPEGE
 382 show a 110-hour OLT which is longer than the OLTs of other models, e.g. 78-hour OLT of UM.
 383 This indicates that SHIELD and ARPEGE could, on average, predict a hit TC genesis event 32

384 hours earlier than UM in the NATL. SHIELD also shows the earliest hit events in the EPAC and
 385 the WPAC, while ARPEGE shows the earliest hit events in the SHEM. The models in this study
 386 generally predict hit events earlier in the EPAC than in other sub-regions, except for JMA, which
 387 performs better in the WPAC and the SHEM than in other sub-regions. It is also interesting to
 388 see that IFS shows earlier hit events than the newer IFS-47R3 in most sub-regions.



389
 390 FIGURE 11. (a) Mean values of maximum observed genesis lead time (in hours) of all storms in each sub-region for
 391 the eight models. Abbreviations used for the sub-regions on the abscissa are the same as in Fig. 7. (b) Fractions of
 392 global total hit events in each model that occurred within a certain DMO length. On the abscissa, “0” is for hit
 393 events which happened at the observed genesis time. “Within 6 (12, 24, or 48)” is for hit events with DMO lengths
 394 in 6 (12, 24, or 48) hours. “Early” is for all hit events with negative DMOs and “late” is for all hit events with
 395 positive DMOs. Abbreviations and colors used for the models in both (a) and (b) are the same as in Fig. 7.

396 Figure 11b shows the fraction of global total hit events in each model which occurred
 397 within a certain length of DMO (indicated on the abscissa). From the definition of DMO, a
 398 model showing more genesis cases with short DMOs indicates that the model could predict more
 399 accurate genesis timings of its hit events. It can be found that IFS shows the highest fraction
 400 among all of the eight models in the “zero” DMO length category. The results indicate that IFS
 401 shows the highest ratio of its hit events forecasted at the observed TC genesis time among the
 402 models. Besides the two IFSSs, JMA and ARPEGE also accurately predict TC genesis timings
 403 within the first three categories (“zero”, “within 6”, and “within 12”) which is the 24-hour
 404 window (12 hours before or after) centered on the observed TC genesis time. In contrast, ICON
 405 shows the smallest fractions in the first three DMO length categories, which implies that the
 406 accuracy of TC genesis timing of ICON is relatively low compared to the other models.

407 From the result of the “within 48” DMO in Fig. 11b, we can see that more than 89% of
 408 hit events in each of the models occur within the 48 hours before or after the observed TC
 409 genesis time. When comparing the results in the “early” and “late” categories, it is seen that most
 410 models forecast their hit events before the observed TC genesis time, except for JMA. The ratios
 411 of “early” to “late” cases are larger in SHIELD and CMC compared to other models, while the
 412 IFS-47R3 shows relatively even number of cases of hit events generated before or after the
 413 observed TC genesis time.

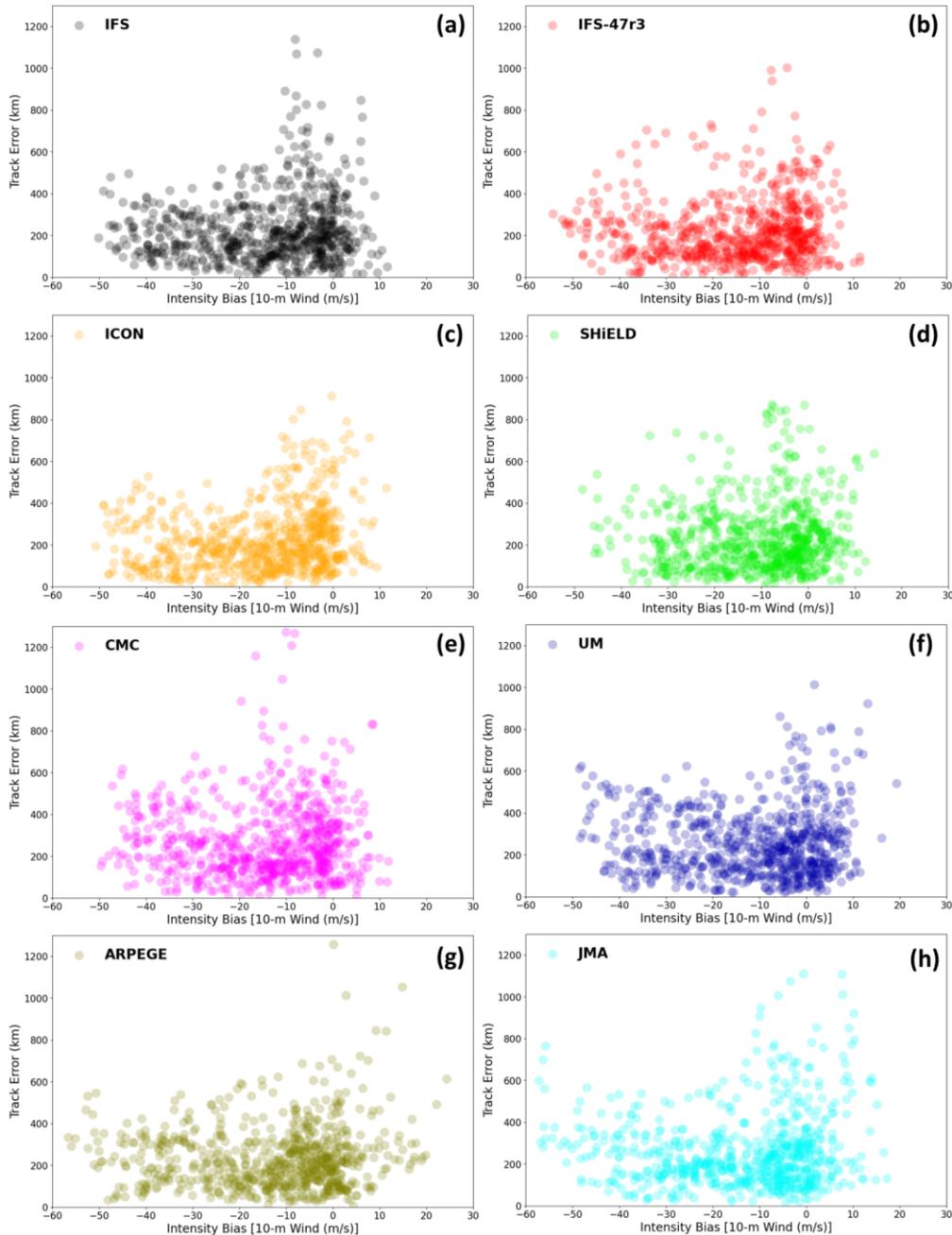
414 **4 Summary and Discussion**

415 The DIMOSIC project provides a great opportunity to engage the worldwide community
416 of medium-range modeling centers on cooperative model research and development. This study
417 investigated TC forecast skills in the eight participating global medium-range forecast models
418 during the year-long DIMOSIC period (June 2018 to June 2019). All models conducted 10-day
419 forecasts from the same initial conditions based on the ECMWF IFS model cycle 45R1. The
420 horizontal resolutions of the eight models ranged from 5 to 25km, and there were different
421 choices of dynamical cores and physics parameterizations across the models (Magnusson et al.
422 2022). The forecast skills of TC track and intensity have been presented for the eight models.
423 The TC-related precipitation and the performance of TC genesis forecasts have also been
424 evaluated.

425 Comparing the model forecasts to the observations for the 109 TCs in the DIMOSIC
426 period, IFS (45R1) and the updated version IFS-47R3 shows the best global averaged TC track
427 forecasts, followed by ICON and SHIELD. CMC shows a relatively higher error than others
428 before the 72-hour lead time. Based on our preliminary investigations, it could be related to the
429 initializing moisture shock given that the CMC has much moister analyses than IFS which may
430 induce convection collapses after the initialization with IFS ICs. For the TC track forecasts in
431 different sub-regions, UM and ICON show the lowest track errors in the NATL, while SHIELD
432 had the best track forecasts in the EPAC. In the WPAC and the SHEM, both IFS-47R3 and
433 ICON show the lowest TC track errors, followed by SHIELD. From the analyses of along-track
434 (AT) and cross-track (CT) errors, the models behave differently in different sub-regions. All
435 models showed a slow-moving bias and a poleward bias in the NATL and the WPAC, but in the
436 EPAC, except for JMA, most models show a fast-moving bias. In contrast, there are no
437 consistently slow- or fast-moving biases among models in the SHEM.

438 For TC intensity forecasts, based on the TC tracker results using the interpolated 0.5
439 degree resolution data, SHIELD performs relatively better than other models, followed by UM
440 and ARPEGE. From Table 2, we can see that the resolutions of the models range between 5 and
441 25 km, and the resolution of SHIELD is in the middle of that range. Therefore, the outperforming
442 TC intensity by SHIELD may imply that the resolution is not the only major factor limiting TC
443 intensity in global models. The use of dynamics and physics in the model also plays important
444 role. The performance of TC track and intensity forecasts could reveal some of the
445 characteristics of a model especially related to its dynamics and physics interactions. In Chen et
446 al. (2019b), it has been demonstrated that updating the GFS dynamical core to the nonhydrostatic
447 FV3 (Lin 2004; Putman and Lin 2007; Harris et al. 2020) can largely improve TC intensity
448 forecasts, and additional improvements in TC intensity and genesis forecasts were seen when
449 replacing the Zhao-Carr cloud microphysics scheme with the advanced GFDL cloud
450 microphysics scheme (Zhou et al. 2019).

451 Here, we attempt to probe into the characteristics of the models based on their biases of
452 TC track and intensity. Figure 12 shows the scatter plots of TC track and intensity errors for all
453 of the forecasts during the 72-120-hour lead time in each model. Some similarities can be found
454 in the scattered distributions of the two IFSs and ICON, including both ranges of intensity bias
455 and track error. This is consistent with the findings in Magnusson et al. (2022) that IFS and
456 ICON behave relatively similarly due to the sharing of partial physical parameterizations sharing
457 between ECWMF and DWD. In contrast, SHIELD, ARPEGE, and JMA show rather unique
458 patterns their own.



459

460 FIGURE 12. Scatter plot distribution of track errors (unit: km; ordinate) and intensity biases (the maximum 10-m
 461 wind speed; unit: m s^{-1} ; abscissa) of all forecasts during the lead time of 72-120 hour for (a) IFS, (b) IFS-47r3, (c)
 462 ICON, (d) SHIELD, (e) CMC, (f) UM, (g) ARPEGE, and (h) JMA.

463 Figure 12 also shows that in most models, forecasts with larger track error (>700 km) are
 464 usually accompanied by smaller intensity biases ($<10 \text{ m s}^{-1}$). In contrast, for those forecasts with
 465 larger intensity biases, their track errors are not consistently larger. At GFDL, it has been noticed
 466 that when the performance of TC track forecasts was improved by using an advection scheme
 467 with a stronger damping in the dynamics, a degradation of TC intensity was observed. The two-
 468 delta filter in the non-monotonic advection scheme and the monotonicity constraint in the tracer
 469 advection affect the model diffusivity which can also impact the diabatic heating and the location

470 of the TC deep convection relative to the eye (Gao et al. 2021). This was attributed to the impact
471 of stronger damping, which suppresses finer-scale features and activities, e.g. grid-scale
472 convection in the TCs, which further suppresses the TC intensities. An in-depth study of the
473 impact of grid-scale convection activity on TC track forecasts in SHIELD is in preparation.

474 Since the interpolated data cannot fully represent the performance of TC intensity
475 forecasts in the models, the TC-related precipitation was also evaluated to provide another
476 perspective on forecasted TC characteristics. Compared to the GPM observation data, all models
477 under-predict the amount of TC-related precipitation, especially in the Pacific Ocean. UM better
478 captures the regions with annual accumulated TC-related precipitation of more than 300mm
479 compared to the other models. However, when considering all TC-related precipitation areas, the
480 ETS of UM is generally lower than that of other models. This could be related to the relatively
481 large track errors of UM, especially in the EPAC. In contrast, SHIELD shows the largest dry bias
482 in TC-related precipitation among all models, but it still achieves relatively high ETSs in the
483 EPAC due to its better track forecasts in this sub-region. As to the intensity of TC-related
484 precipitation, most models over-predict medium precipitation events but under-predict light and
485 heavy precipitation events. Among all models, SHIELD noticeably over-predicts precipitation
486 events at the intensity of 1-5 mm (6 h)⁻¹. The SHIELD development team at GFDL will take a
487 close look at its low precipitation amount and over-predicted medium intensity precipitation
488 events in the future.

489 The assessment of TC genesis forecast skill was based here on hits and misses, measures
490 that showed significant inter-model variability across the different sub-regions of interest. All
491 models show the highest hit ratios with the fewest total genesis events in the EPAC, which
492 indicates that models can better predict TC genesis in the EPAC. In contrast, models generate
493 more false alarms in the SHEM. SHIELD shows the highest hit ratios globally, followed by the
494 two IFSs and ICON. CMC also achieves high hit ratios in the EPAC. In contrast, UM shows
495 lower hit ratios than other models. As for the missed TC genesis cases, JMA missed 30 of the
496 100 observed TC geneses during the target period, which is the most among all models. Note that
497 JMA uses the coarsest model resolution among all participating models, which could impair its
498 TC genesis performance. In contrast, SHIELD shows the least number of missed TCs globally,
499 which may be benefiting from its better TC intensity forecast.

500 We also investigate how well the participating models predict the timing of TC genesis
501 by comparing the “length of lead time” proposed by Chen et al. (2019b). The results show that
502 models can generally accurately predict TC formation earlier in the EPAC than in other sub-
503 regions, except for JMA which predicts the WPAC and the SHEM hit events earlier than in other
504 sub-regions. SHIELD generally predicts the earliest hit events globally, while ARPEGE also
505 predicts TC genesis events earlier than other models in the NATL and the SHEM. Based on the
506 differences between the model genesis lead time and the observed genesis lead time, IFS shows
507 the most accurate timing of the TC genesis forecast, followed by JMA and ARPEGE. In contrast,
508 the accuracy of genesis timing in ICON is relatively lower than that of other models. We also
509 found that most models develop TCs earlier than observed in the best track, except for JMA. In
510 addition, more than 89% of hit events in each of the models occur within 48 hours of the
511 observed genesis time.

512 The comparison between IFS (version 45R1) and IFS-47R3 provides an opportunity to
513 examine the incremental change obtained for an upgrade of one model. The upgrade from
514 version 45R1 to 47R3 includes many changes in data assimilation and model physics. The

515 changes and meteorological impacts have been documented by ECMWF on the website:
516 <https://www.ecmwf.int/en/forecasts/documentation-and-support/changes-ecmwf-model>. One of
517 the listed impacts from the upgrade is the improvement of TC position errors. From our analysis,
518 the average track errors during the first 120 hours of IFS-47R3 are 2.5 to 10.8 km less than IFS
519 (Fig. 3) in the four major sub-regions, which is consistent with the ECMWF implementation
520 report. However, from the analyses of the along-track and cross-track errors, the biases of slow
521 and poleward movement are similar in these two model versions. Possibly associated with the
522 major upgrade to moist physics (Bechtold et al. 2020), IFS-47R3 shows a slightly larger negative
523 TC intensity biases than the older IFS version which was also found in Magnusson et al. (2021).
524 Although total precipitation predictions remain similar across the two IFS versions, the newer
525 IFS-47R3 achieves higher ETS scores for large accumulations, especially in the WPAC and the
526 SHEM. We also found that IFS-47R3 has more precipitation events with stronger precipitation
527 intensity (10-50 mm (6 h)⁻¹) than IFS. However, as to the TC genesis forecast, IFS-47R3 shows
528 some degradation from IFS, including more missed TC genesis event, higher miss ratios, shorter
529 genesis lead times, and less accuracy on TC genesis timing. These degradations in TC genesis
530 performance may be related to the weaker TC intensities in the newer version.

531 To summarize, in this study, extensive evaluation was made of the performance of TC
532 forecasts in the DIMOSIC models based on one year of model predictions initialized with the
533 same initial conditions. Although it is hard to precisely isolate the influence of individual
534 components in different model formulations on TC forecast skill in such an overview, the
535 comparisons based on different evaluation metrics highlight important similarities and
536 differences between the models. The results will be valuable for model developers in
537 participating centers as a benchmark of TC forecast skill with the impact of the initial condition
538 quality removed. Also, common forecast biases of the TC movement and TC-related
539 precipitations indicate general deficiencies in DIMOSIC models and point out a direction for
540 model developers for further model improvement.

541

542

543 **Acknowledgments**

544 The authors thank Kun Gao, Jie Chen, Morris Bender, and Tom Knutson for GFDL internal
545 review, and Lucas Harris, James Doyle, and Simon Lang for their comments helped to improve
546 this article. Authors also would like to thank other DIMOSIC participants, Duncan Ackerley,
547 Yves Bouteloup, K. C. Kwon, Yoonjin Lim, Mio Mastueda, Takumi Matsunobu, and Yamaguchi
548 Munehiko for their contribution to the DIMOSIC project.

549

550 **Open Research**

551 All DIMOSIC model interpolated data can be requested from Linus Magnusson. All TC
552 analyses are archived in the GFDL Tape Archive System at /archive/jhc/DIMOSIC/Analysis/TC
553 and can be requested from Jan-Huey Chen.

554

555 **References**

- 556 Bechtold, P., R. Forbes, I. Sandu, S. Lang, & M. Ahlgrimm (2020), A major moist physics
557 upgrade for the IFS. *ECMWF Newsletter*, 164, URL:
558 <https://www.ecmwf.int/en/newsletter/164/meteorology/major-moist-physics-upgrade-ifs>
- 559 Chen, J.-H., S.-J. Lin, L. Magnusson, M. A. Bender, X. Chen, L. Zhou, B. Xiang, S. L. Rees, M.
560 J. Morin, & L. Harris (2019a), Advancements in hurricane prediction with NOAA's next
561 generation forecast system. *Geophysical Research Letters*, 46 (8), 44954501,
562 doi:10.1029/2019GL082410.
- 563 Chen, J.-H., Lin, S., Zhou, L., Chen, X., Rees, S., Bender, M., & Morin, M. (2019b), Evaluation
564 of Tropical Cyclone Forecasts in the Next Generation Global Prediction System, *Monthly*
565 *Weather Review*, 147(9), 3409-3428, doi: 10.1175/MWR-D-18-0227.1
- 566 DWD (2022), ICON : Icosahedral Nonhydrostatic Weather and Climate Model. *Technical*
567 *Report, DWD*. URL: <https://code.mpimet.mpg.de/projects/iconpublic/wiki/Documentation>
- 568 ECMWF (2018), IFS Documentation CY45R1. *Technical Report, ECMWF*. URL:
569 <https://www.ecmwf.int/en/publications/ifs-documentation>
- 570 ECMWF (2021), IFS Documentation CY47R3. *Technical Report, ECMWF*. URL:
571 <https://www.ecmwf.int/en/publications/ifs-documentation>

- 572 Gao, K., L. Harris, L. Zhou, M. A. Bender, and M. Morin (2021), On the sensitivity of hurricane
573 intensity and structure to horizontal tracer advection schemes in FV3. *Journal of the*
574 *Atmospheric Sciences*, 78(9), doi:10.1175/JAS-D-20-0331.13007-3021.
- 575 Girard, C., & Coauthors (2014), Staggered vertical discretization of the Canadian Environmental
576 Multiscale (GEM) model using a coordinate of the log-hydrostatic-pressure type. *Monthly*
577 *Weather Review*, 142 (3), 1183.
- 578 Harris, L. M., S.-J. Lin, & C. Tu (2016), High-resolution climate simulations using GFDL
579 HiRAM with a stretched global grid. *Journal of Climate*, 29, 4293–4314, doi:10.1175/JCLI-
580 D-15-0389.11196, doi:10.1175/MWR-D-13-00255.1.
- 581 Harris, L., & Coauthors (2020), GFDL SHIELD: A unified system for weather to seasonal
582 prediction. *Journal of Advances in Modeling Earth Systems*, 12 (10),
583 doi:10.1029/2020MS002223.
- 584 Hong, Y., K.-L. Hsu, S. Sorooshian, & X. Gao (2004), Precipitation estimation from remotely
585 sensed imagery using an artificial neural network cloud classification system. *Journal of*
586 *Applied Meteorology*, 43 (12), 1834-1853, doi:10.1175/JAM2173.1.
- 587 JMA (2019), Outline of the operational numerical weather prediction at the Japan
588 Meteorological Agency, *Technical Report, Japan Meteorological Agency*. URL
589 <http://www.jma.go.jp/jma/jma-eng/jma-center/nwp/outline2019-nwp/index.htm>.
- 590 Lin, S.-J. (2004), A “vertically Lagrangian” finite-volume dynamical core for global models,
591 *Monthly Weather Review*, 132, 2293-2307.
- 592 Magnusson, L., J.-H. Chen, S.-J. Lin, L. Zhou, & X. Chen (2019), Dependence on initial
593 conditions vs. model formulations for medium-range forecast error variations. *Quarterly*
594 *Journal of the Royal Meteorological Society*, 145, doi:10.1002/qj.3545

- 595 Magnusson, L., & co-authors (2021), Tropical cyclone activities at ECMWF. *ECMWF Technical*
596 *Memorandum 888*
- 597 Magnusson, L., D. Ackerley, Y. Bouteloup, J.-H. Chen, J. Doyle, P. Earnshaw, Y. C. Kwon, M.
598 Köhler, S. T. K Lang, Y.-J. Lim, M. Matsueda, T. Matsunobu, R. McTaggart-Cowan, A.
599 Reinecke, M. Yamaguchi1, & L. Zhou (2022), Skill of medium-range forecast models using
600 the same initial conditions. *Bulletin of the American Meteorological Society*, 103(9), E2050-
601 E2068, doi:10.1175/BAMS-D-21-0234.1
- 602 McTaggartCowan, R., & Coauthors (2019), Modernization of atmospheric physics
603 parameterization in Canadian NWP. *Journal of Advances in Modeling Earth Systems*, 11
604 (11), 35933635, doi:10.1029/2019MS001781.
- 605 Miller, R. J., A. J. Scrader, C. R. Sampson, & T. L. Tsui (1990), The Automated Tropical
606 Cyclone Forecast System (ATCF). *Weather and Forecasting*, 5, 653–660.
- 607 Mogensen, K. S. , L. Magnusson , & J.-R. Bidlot (2017), Tropical cyclone sensitivity to ocean
608 coupling in the ECMWF coupled model: Tropical cyclone sensitivity. *J. Geophys. Res.*
609 *Oceans*, 122, 4392–4412, doi:10.1002/2017JC012753.
- 610 Pollard, R. T. , P. B. Rhines , & R. O. R. Y. Thompson (1973), The deepening of the wind-mixed
611 layer. *Geophys. Fluid Dyn.* , 4, 381–404, doi:10.1080/03091927208236105.
- 612 Putman, W. M., & S.-J. Lin (2007) Finite-volume transport on various cubed-sphere grids.
613 *Journal of computational physics*, 227, 55-78, doi:10.1016/j.jcp.2007.07.022.
- 614 Rabier, F., H. Järvinen, E. Klinker, J.-F. Mahfouf, & A. Simmons (2000), The ECMWF
615 operational implementation of four-dimensional variational assimilation. I: Experimental
616 results with simplified physics. *Quarterly Journal of the Royal Meteorological Society*, 126
617 (564), 24 1143–1170, doi:10.1002/qj.49712656415

- 618 Roehrig, R., & Coauthors (2020), The cnrm global atmosphere model arpegeclimat 6.3:
619 Description and evaluation. *Journal of Advances in Modeling Earth Systems*, 12 (7),
620 doi:10.1029/2020MS002075.
- 621 Schaefer, J. T. (1990) The critical success index as an indicator of warning skill. *Weather and*
622 *Forecasting*, 5 , 570–575.
- 623 Sampson, C. R., & A. J. Schrader (2000), The Automated Tropical Cyclone Forecasting System
624 (version 3.2). *Bulletin of the American Meteorological Society*, 81, 1231–1240.
- 625 Walters, D., & Coauthors (2019), The met office unified model global atmosphere 7.0/7.1 and
626 jules global land 7.0 configurations. *Geoscientific Model Development*, 12(5), 1909-1963,
627 doi:10.5194/gmd-12-1909-2019.
- 628 Zeng, X. , & A. Beljaars (2005), A prognostic scheme of sea surface skin temperature for
629 modeling and data assimilation: Sea surface skin temperature scheme. *Geophysical*
630 *Research Letters*, 32, L14605, doi:10.1029/2005GL023030.
- 631 Zhou, L., S.-J. Lin, J.-H. Chen, L. M. Harris, X. Chen, & S. Rees (2019), Toward convective
632 scale prediction within the Next Generation Global Prediction System. *Bulletin of the*
633 *American Meteorological Society*. 100 (7): 1225-43, doi:10.1175/BAMS-D-17-0246.1



Journal of Advances in Modeling Earth Systems Supporting Information for

Tropical Cyclone Forecasts in the DIMOSIC Project

Jan-Huey Chen^{1,2}, Linjiong Zhou³, Linus Magnusson⁴, Ron McTaggart-Cowan⁵, and Martin Köhler⁶

¹National Oceanic and Atmospheric Administration/Geophysical Fluid Dynamics Laboratory, Princeton, NJ, USA

²University Corporation for Atmospheric Research, Boulder, CO, USA

³Atmospheric and Oceanic Sciences Program, Princeton University, Princeton, NJ, USA

⁴European Centre for Medium-Range Weather Forecasts, Reading, UK

⁵Environment and Climate Change Canada, Montreal, Canada

⁶Deutsche Wetterdienst, Offenbach, Germany

Contents of this file

Text S1
Figures S1 to S3
Text S2
Figures S4 to S7

Introduction

This supporting information file includes:

- 1) Analyses of TC track errors of GFDL SHIELD in the North Atlantic basin (NATL) (Text S1 and Figures S1-S3)
- 2) Analyses of along-track (AT) and cross-track (CT) errors for all 8 models in the 4 major sub-regions (Text S2 and Figures S4-S7)

Text S1.

During the DIMOSIC period, SHIELD is one of the top performance models which show low tropical cyclone (TC) track forecast errors in most sub-regions but not in the North Atlantic basin (NATL). We therefore further investigate the large track error of SHIELD in the NATL to supplement the results of track error analysis, and to provide references for the SHIELD development team at GFDL.

Figure S1a shows the mean TC track forecast errors for all of the 8 DIMOSIC models in the NATL, with the differences of TC track errors of the 7 models comparing to IFS (Fig. S1b). It can be found that SHIELD shows a much larger track error than most other models during the 72-120 lead times.

We investigated the track forecasts of the 16 TCs in the NATL during the target year individually and found that the basin-wide mean track errors are actually dominated by 2 TCs, Hurricanes Florence (2018) and Leslie (2018). Comparing to other leading models, SHIELD shows difficulty to forecast the storm movements of these two hurricanes. Figure S2 shows the forecasted tracks of Florence from IFS-47R3, SHIELD, ICON, and UM initialized at 00Z Sep. 1st (Fig. S2a) and 00Z Sep. 13th (Fig. S2b) comparing to the best track. It can be found that SHIELD shows slow moving biases on Florence's translation in both early and late stage of its lifetime, while IFS-47R3, ICON and UM do not show similar forecast biases. Different from Florence which had a steady westward movement crossing the Atlantic Ocean, Leslie had an irregular track with many loops in the middle of Atlantic Ocean before heading northeasterly to Portugal. From Fig. S3, we can see that SHIELD didn't well capture some of the sharp turns which are critical in the track forecast of this storm. In contrast, IFS-47R3, ICON, and UM can perform much better "v" shape turns in their forecasts initialized from 00Z 1st Oct. and from 00Z 7th Oct.

To summarize, the large TC track error of SHIELD in the NATL results from the slow moving bias shown in the forecasts of Hurricane Florence and the bias of direction of motion shown in the forecasts of Hurricane Leslie.

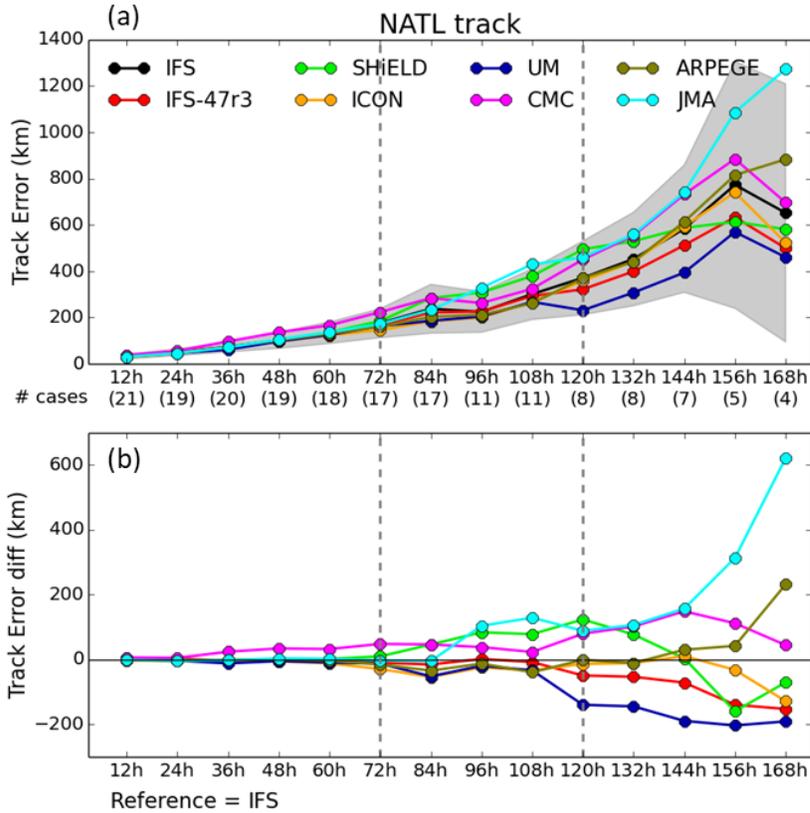


Figure S1. (a) Mean TC track forecast errors (km) along with the model forecast lead time for IFS (black), IFS-47r3 (red), SHIELD (green), ICON (yellow), UM (blue), CMC (magenta), APREGE (grass green), and JMA (light blue) in the North Atlantic basic (NATL). The 95% confidence levels for IFS are indicated by the gray color shading. Numbers of homogeneous cases for individual lead times are listed in the brackets at the bottom of each abscissa. Vertical gray dotted lines are indicated 72 and 120 h. (b) Mean TC track forecast error differences of IFS-47r3 (red), SHIELD (green), ICON (yellow), UM (blue), CMC (magenta), APREGE (grass green), and JMA (light blue) in the NATL comparing to IFS.

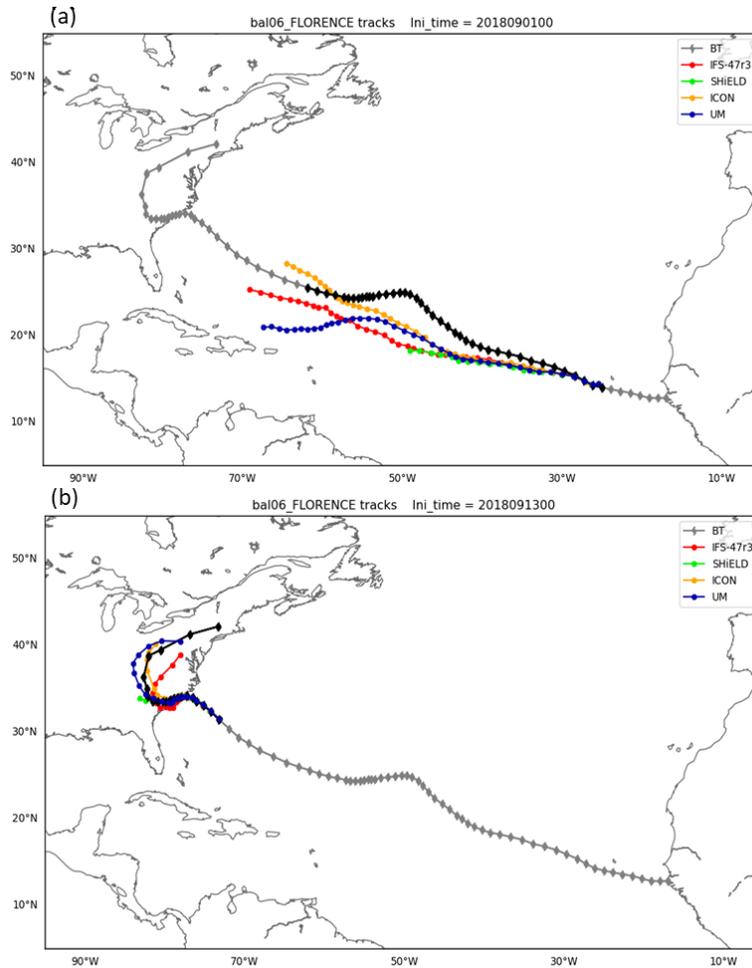


Figure S2. (a) Model forecasted tracks from IFS-47r3 (red), SHIELD (green), ICON (yellow), and UM (blue) initialized at 00Z 20180901 comparing to the Automated Tropical Cyclone Forecast (ATCF) best track for Hurricane Florence. Dots for model forecasts and typhoon symbols for the best track are at 6-hour interval. The best track of Hurricane Florence during the 10-day forecast period is plotted in black, and the rest part of the best track is plotted in grey. (b) As in (a), but for model forecasts initialized at 00Z 20180913.

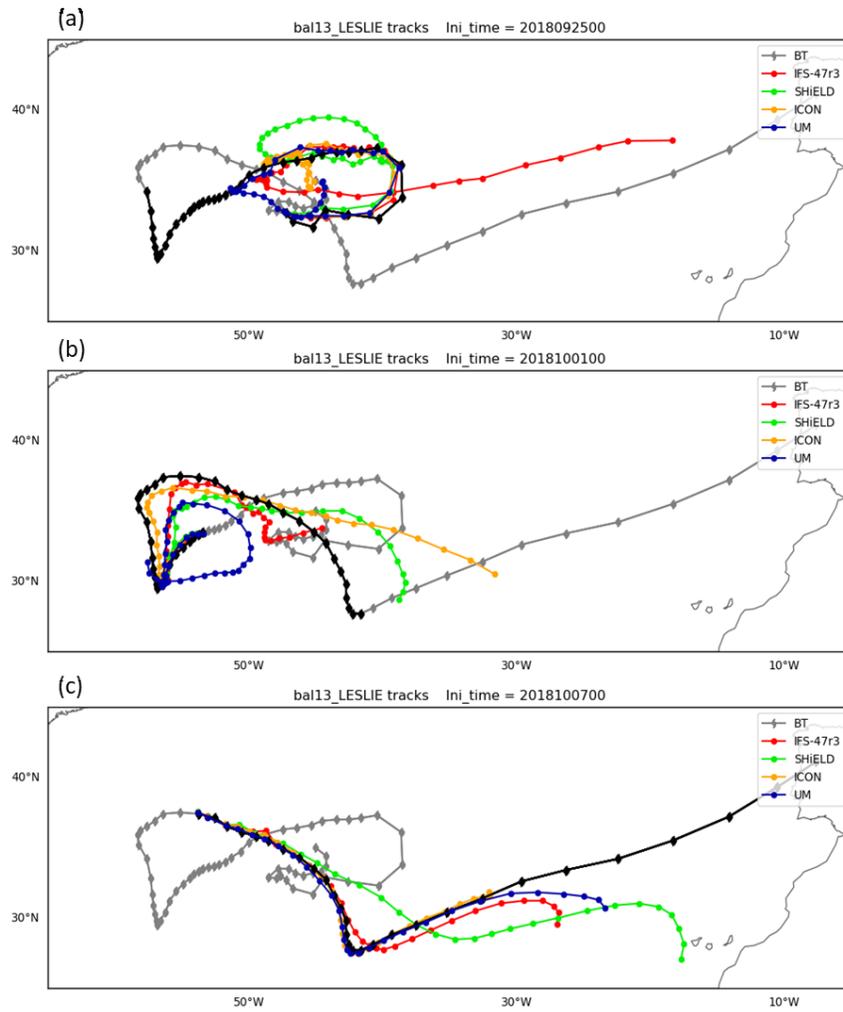


Figure S3. As in Figure S2, but for the model forecasts initialized at (a) 00Z 20180925, (b) 00Z 20181001, and (c) 00Z 20181007 for Hurricane Leslie comparing to the best track.

Text S2.

Analyses of along-track (AT) and cross-track (CT) errors for all eight models in the four major sub-regions are shown in Figs. S4-S7. We note that the models behave differently in different sub-regions. In the NATL (Fig. S4), the AT and CT errors during the first 5 days evenly contribute to the total track errors in the two IFSs. UM AT/CT skill is similar to that of the IFSs but with a larger AT error component on Day 3 to 4 (Fig. S4f). In contrast, ICON and CMC have a much larger AT error component than the CT error (Figs. S4c,e), which indicates that their total track errors in the NATL are mostly due to slow moving biases. The source of track error of SHIELD in the NATL is a combination of AT error during 96 to 132-hour lead time and CT error at Days 6 to 7 (Fig. S4d), which is consistent with the single case analyses (Text S1 and Figs. S1-S3).

In distinct behavior from the NATL, seven of the eight models show positive AT error biases in the EPAC during the first five days (Fig. S5). This indicates that the model-predicted TCs are usually moving too fast in this basin. In contrast, JAM suffers a slow-moving bias as shown by a dominant negative AT error bias (Fig. S5h). Also, the large track error of UM during the first five days in this region (Fig. 3) is mostly from a fast-moving bias, but a southward-moving bias also contributes. In the WPAC, all models consistently show negative AT error biases and positive CT error biases after Day 5 (Fig. S6). It is clear that a slow-moving bias is the main source of the total track errors in the two IFSs in this region (Figs. S6a,b). UM and JMA show large contributions of AT errors during Days 4-5, but northward-moving biases also contribute to the total error substantially (Figs. S6f,h). ICON, SHIELD, and CMC show even greater contributions of AT and CT errors to their total track errors than other models (Figs. S6c-e). In the SHEM (Fig. S7), except for SHIELD showing a relatively even contribution of AT and CT errors with small biases, other models generally exhibit much larger components of AT errors than CT errors. However, there are no consistently slow- or fast-moving biases among models in this basin. IFSs and CMC show consistently slow-moving biases after three days, while UM shows fast-moving biases in the later lead times. ICON, APREGGE, and JMA show mixed slow- or fast-moving biases in the seven days of lead time.

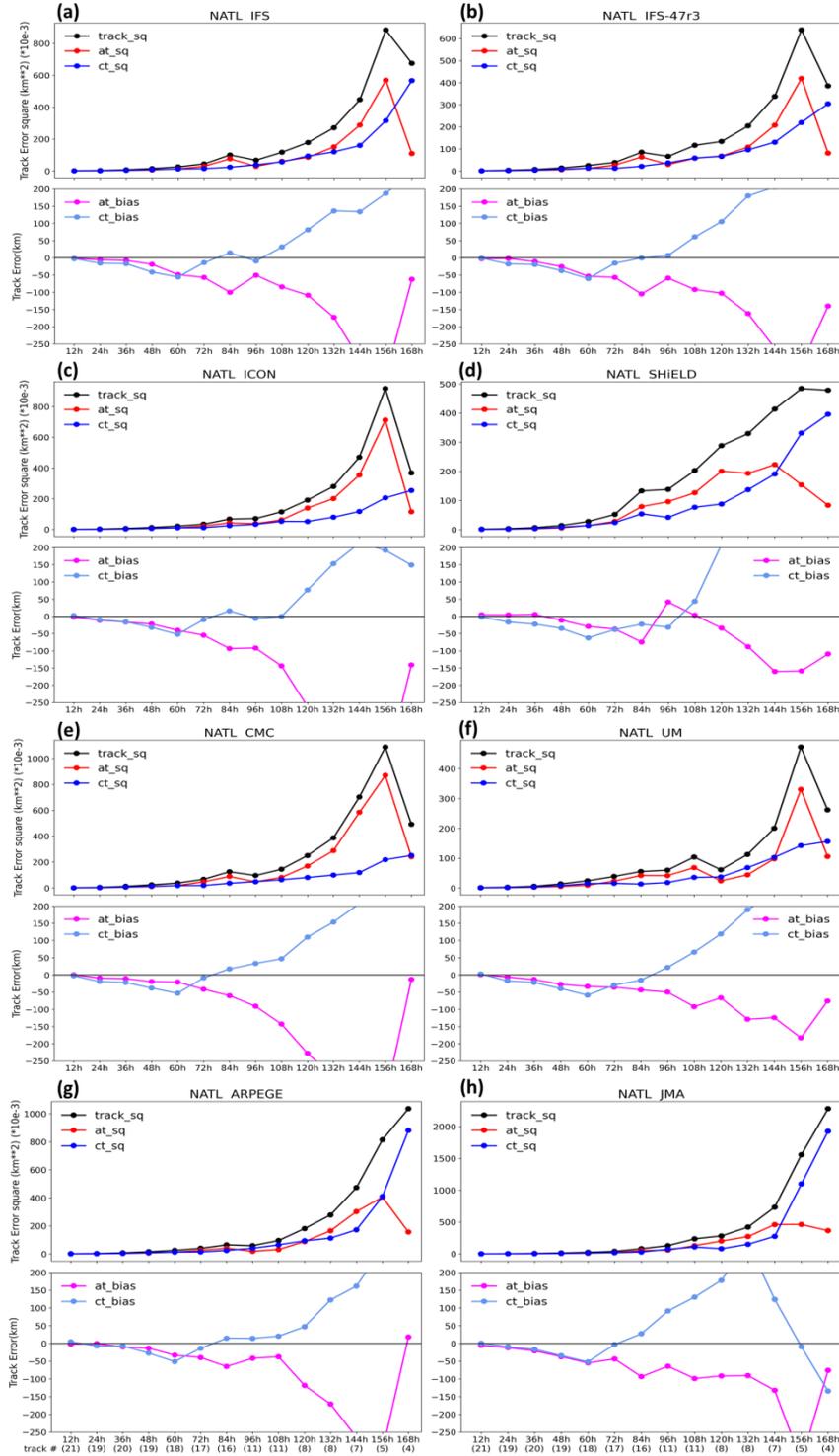


Figure S4. Analyses of long-track (AT) error and cross-track (CT) error for (a) IFS, (b) IFS-47R3, (c) ICON, (d) SHiELD, (e) CMC, (f) UM, (g) ARPEGE, and (h) JMA in the North Atlantic basin (NATL). The squares of total track errors (black), along-track errors (red), and cross-track errors (blue) are in the upper panels for each model. The biases of along-track (magenta) and cross-track (light blue) errors are in the lower panels. Numbers of homogeneous cases for each lead time are listed at the bottom of lower panels.

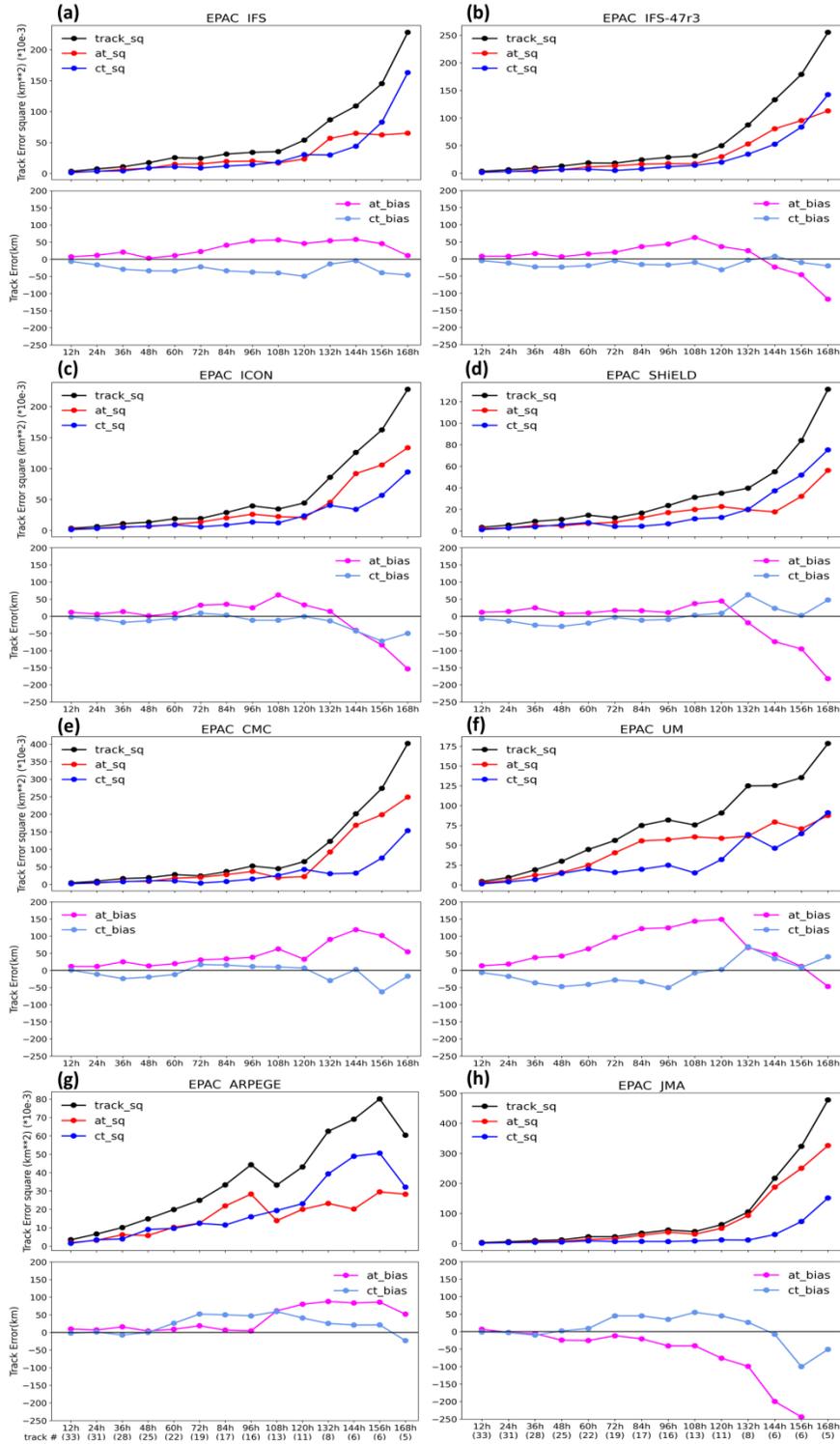


Figure S5. As in Figure S4, but for the analyses in the northeast Pacific basin (EPAC).

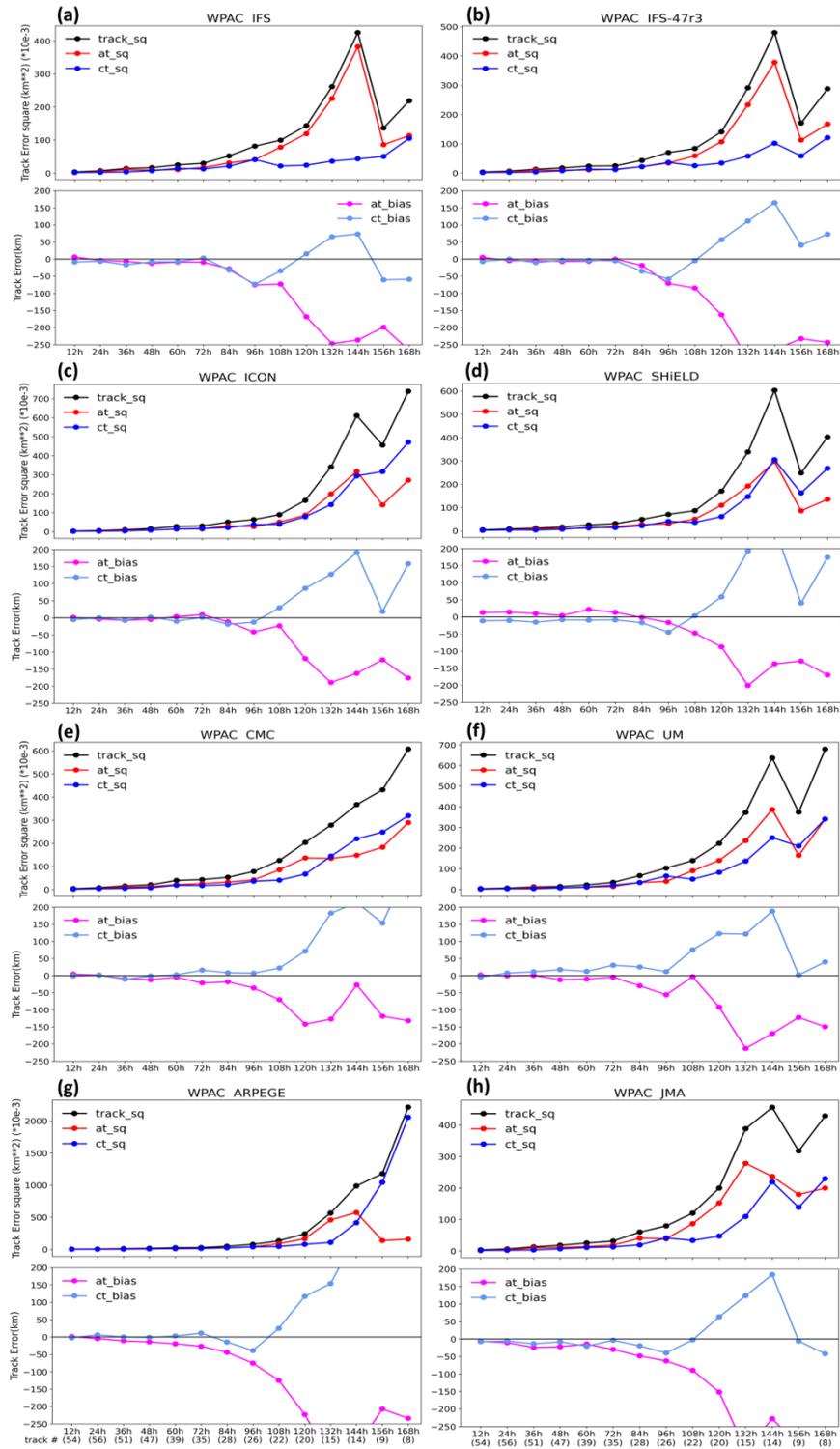


Figure S6. As in Figure S4, but for the analyses in the northwest Pacific basin (WPAC).

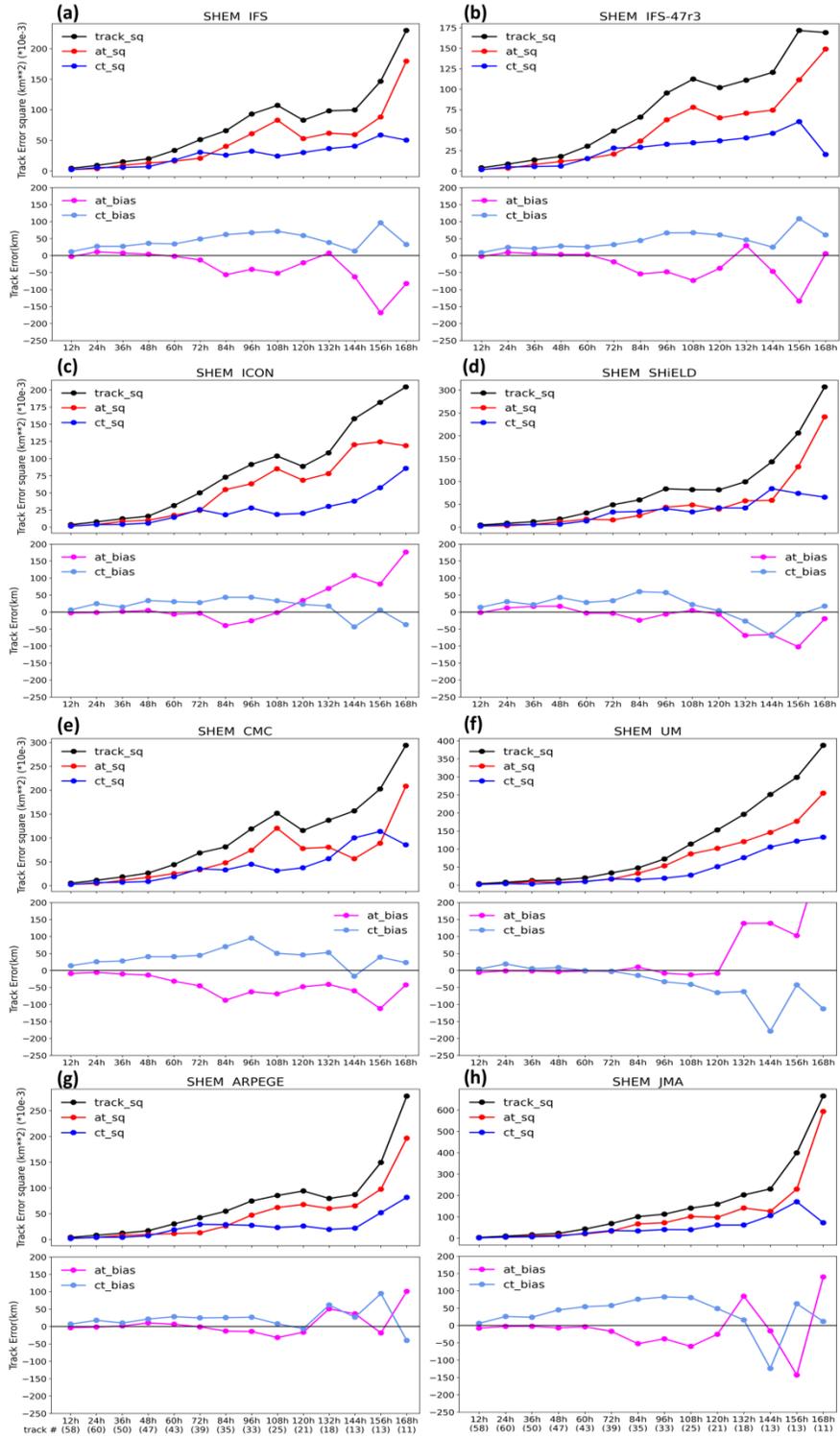


Figure S7. As in Figure S4, but for the analyses in the Southern Hemisphere (SHEM).