

1 **Tropical Cyclones in Global Storm-Resolving**  
2 **Models**

3 **Falko Judt**

4 *National Center for Atmospheric Research, Boulder,*  
5 *Colorado, USA*

6 **and**

7 **Daniel Klocke<sup>8</sup>, Rosimar Rios-Berrios<sup>1</sup>, Benoit**  
8 **Vanniere<sup>13</sup>, Florian Ziemer<sup>3</sup>**

9 **and**

10 Ludovic Auger<sup>2</sup>, Joachim Biercamp<sup>3</sup>, Christopher  
11 Bretherton<sup>4</sup>, Xi Chen<sup>5</sup>, Peter Düben<sup>6</sup>, Cathy  
12 Hohenegger<sup>10</sup>, Marat Khairoutdinov<sup>7</sup>, Chihiro  
13 Kodama<sup>9</sup>, Luis Kornbluh<sup>10</sup>, Shian-Jiann Lin<sup>5</sup>, Masuo  
14 Nakano<sup>9</sup>, Philipp Neumann<sup>14</sup>, William Putman<sup>11</sup>,  
15 Niklas Röber<sup>3</sup>, Malcom Roberts<sup>15</sup>, Masaki Satoh<sup>12</sup>,  
16 Ryosuke Shibuya<sup>12</sup>, Bjorn Stevens<sup>10</sup>, Pier Luigi  
17 Vidale<sup>13</sup>, Nils Wedi<sup>6</sup>, Linjiong Zhou<sup>5</sup>,

18 <sup>1</sup>*National Center for Atmospheric Research, P.O. Box 3000,*  
19 *Boulder, CO 80307, USA*

20 <sup>2</sup>*CNRM Meteo-France, 42 av. G. Coriolis, 31057 Toulouse,*  
21 *France*

22 <sup>3</sup>*German Climate Computing Center, DKRZ, Bundesstraße*  
23 *45a, 20146, Hamburg, Germany*

24 <sup>4</sup>*Department of Atmospheric Sciences, University of*  
25 *Washington, Box 351640, Seattle, WA, 98195, USA*

26 <sup>5</sup>*Geophysical Fluid Dynamics Laboratory, Princeton*  
27 *University, Forrestal Campus/U.S. Route 1, P.O. Box 308,*  
28 *Princeton, NJ, 08542, USA*

29 <sup>6</sup>*ECMWF, Shinfield Road, Reading, RG2 9AX, UK*

30 <sup>7</sup>*School of Marine and Atmospheric Sciences, Stony Brook*  
31 *University, Stony Brook, NY, 11794, USA*

32 <sup>8</sup>*Hans-Ertel-Zentrum für Wetterforschung, Deutscher*  
33 *Wetterdienst, Frankfurter Straße 135, 63067 Offenbach,*  
34 *Germany*

35 <sup>9</sup>*Japan Agency for Marine-Earth Science and Technology,*  
36 *3173-15, Showa-machi, Kanazawa-ku, Yokohama, Kanagawa,*  
37 *236-0001, Japan*

38 <sup>10</sup>*Max Planck Institute for Meteorology, Bundesstraße 53,*  
39 *20146 Hamburg, Germany*

40 <sup>11</sup>*NASA Global Modeling and Assimilation Office, Goddard*  
41 *Space Flight Center, Greenbelt, Maryland, USA*  
42 <sup>12</sup>*Atmosphere and Ocean Research Institute, The University of*  
43 *Tokyo, 5-1-5 Kashiwanoha, Kashiwa, Chiba, 277-8564, Japan*  
44 <sup>13</sup>*National Centre for Atmospheric Science, University of*  
45 *Reading, Reading, UK*  
46 <sup>14</sup>*Helmut-Schmidt-Universität, Fakultät für Maschinenbau,*  
47 *High Performance Computing, Holstenhofweg 85, 22043*  
48 *Hamburg, Germany*  
49 <sup>15</sup>*UK Met Office, Fitzroy Rd, Exeter EX1 3PB, United*  
50 *Kingdom*

51 June 23, 2020

---

Corresponding author: Falko Judt, National Center for Atmospheric Research,  
P.O. Box 3000, Boulder, CO 80307, USA.  
E-mail: fjudt@ucar.edu

## Abstract

52

53 Recent progress in computing and model development has initiated the era  
54 of global storm-resolving modeling and with it the potential to transform  
55 weather and climate prediction. Within the general theme of vetting this  
56 new class of models, the present study evaluates nine global-storm resolving  
57 models in their ability to simulate tropical cyclones. Results show that,  
58 broadly speaking, the models produce realistic tropical cyclones and remove  
59 longstanding issues known from global models such as the deficiency to  
60 accurately simulate TC intensity. However, TCs are strongly affected by  
61 model formulation, and all models suffer from unique biases regarding the  
62 number of cyclones, intensity, size, and structure. Some models simulated  
63 TCs better than others, but no single model was superior in every way. The  
64 overall results indicate that global storm-resolving models are able to open  
65 a new chapter in tropical cyclone prediction, but they need to be improved  
66 to unleash their full potential.



67 **Keywords** tropical cyclone; typhoon; hurricane; global cloud-resolving  
68 model; tropical meteorology; numerical modeling; model verification; pre-  
69 dictability

## 70 **1. Introduction**

71 Tropical cyclones (TCs) are among the most destructive natural haz-  
72 ards, and predicting TCs is an important task of weather and climate  
73 models. Moreover, TCs are optimal testbeds for assessing the quality of  
74 numerical models, because their unique dynamics reveal deficiencies in the  
75 model formulation through artifacts such as unrealistic structure. The over-  
76 all purpose of the present study is to evaluate a new class of atmosphere  
77 models—global storm-resolving models (Satoh et al. 2019)—in their ability  
78 to simulate TCs. Specifically, we report on TC-related achievements, defi-  
79 ciencies, and biases in nine global storm-resolving models, and we hope that  
80 our findings will pave the way for improving the next generation of weather  
81 and climate models.

82 Global models have been a vital instrument in TC prediction although  
83 they have not been able to accurately predict TC intensity. A decade ago,  
84 Hamill et al. (2011) reported that global weather models, which at that  
85 time had mesh spacings between 50–150 km, were plagued by wind speed  
86 biases of down to  $-30 \text{ m s}^{-1}$ . Even though some progress has been made,

87 the most recent model with mesh spacings of 10–25 km still fail to capture  
88 the high winds of TCs (e.g., Magnusson et al. 2019; Hodges and Klingaman  
89 2019; Roberts et al. 2020). One of the main reasons for this shortcoming is  
90 insufficient horizontal resolution (Davis 2018). In fact, years of research with  
91 regional models have documented that storm-resolving resolution, here de-  
92 fined as  $<5$  km, is necessary to accurately simulate the inner-core structure  
93 of TCs (e.g., Chen et al. 2007; Gentry and Lackmann 2010), which in turn  
94 is necessary to predict TC intensity (e.g., Davis et al. 2008; Gopalakrishnan  
95 et al. 2012; Fox and Judt 2018).

96 The preceding arguments suggest that global storm-resolving models  
97 are ideal tools for TC prediction, because they combine the advantages  
98 of current-generation global and regional models, that is, they offer global  
99 coverage *and* storm-resolving horizontal resolution. Indeed, there has been  
100 some qualitative evidence that global storm-resolving models capture the  
101 inner-core structure of TCs quite realistically (e.g., Fudeyasu et al. 2008;  
102 Zhou et al. 2019). Other studies have demonstrated that models with 7–  
103 10 km mesh spacings reduce some of the biases found in coarser-resolution  
104 models (Manganello et al. 2012; Nakano et al. 2017). However, the immense  
105 computational resources needed to run global models with mesh spacings  
106 of  $\leq 5$  km have so far precluded a detailed, TC-focused evaluation of those  
107 models. The present study attempts to fill this gap by evaluating the models

108 that participated in the DYAMOND initiative (Stevens et al. 2019), and  
109 it expands on the brief overview of TCs already presented in Stevens et al.  
110 (2019).

111 Given computational limitations and the general purpose of DYAMOND,  
112 each participating model provided only one 40-day simulation. This means  
113 that it was not possible to evaluate the models as usually done in the weather  
114 prediction community, i.e., by computing errors of metrics such as maximum  
115 wind speed from a large number of short-range forecasts (e.g., DeMaria et al.  
116 2014; Nakano et al. 2017). It was also not possible to evaluate long-term  
117 TC climatologies as in climate studies (e.g., Camargo et al. 2005; Bengts-  
118 son et al. 2007; Manganello et al. 2012; Roberts et al. 2020). Instead, we  
119 focused on answering the following questions:

- 120 • What are the biases in TC number, tracks, intensity, and size over  
121 those 40 days?
- 122 • Do the models produce TCs with a realistic structure?
- 123 • Do the models have similar biases, or does each model have its own?

124 The validity of the study rests on three important assumptions, namely  
125 (i) the 40-day period of DYAMOND is sufficient to draw general conclusions  
126 about the TC characteristics in each model, (ii) objects identified as TCs  
127 by the tracking software (see section 2) would also be identified as TCs

128 by human forecasters (and vice versa), and (iii) the observations used to  
129 evaluate the models are sufficiently accurate.

130 We are confident that (i) holds true because the discrepancies between  
131 the models were substantial and almost certainly caused by different model  
132 formulations. Furthermore, even though 40 days is relatively short, we have  
133 global statistics, and the sampling is not as sparse as one might intuit. It is  
134 more difficult to judge the validity of (ii) and (iii), but given the amount of  
135 past studies that relied on those assumptions, we assumed they would hold  
136 for this work, too.

137 Lastly, we emphasize that high horizontal resolution is necessary but not  
138 sufficient for accurately simulating TC structure and intensity. Advances in  
139 ocean coupling and model physics are critical as well (e.g., Lee and Chen  
140 2014; Mogensen et al. 2017; Magnusson et al. 2019). One area that seems  
141 to be particularly important is the parameterization of the boundary layer  
142 (Kanada et al. 2012; Kepert 2012; Zhang et al. 2015) and the surface layer,  
143 especially the drag (Zeng et al. 2010; Green and Zhang 2013; Magnusson  
144 et al. 2019).

145 The remainder of the paper is structured as follows: in section 2, we  
146 present the data and methods. Section 3 contains the results, organized  
147 into subsections on (i) TC number and tracks, (ii) intensity, (iii) size, (iv)  
148 structure, and (v) the sensitivity of TCs on resolution and parameterized

149 convection. The findings are discussed in section 4 and the paper ends with  
150 a summary and conclusions in section 5.

## 151 **2. Data and Methods**

152 This study leverages the vast data repository of DYAMOND, which  
153 contains the output from the following nine global models: ARPEGE, FV3,  
154 GEOS, ICON, IFS<sup>1</sup>, MPAS, NICAM, SAM, and UM. For details about  
155 the DYAMOND experiment and the participating models see Stevens et  
156 al. (2019) and references therein. The models were run on meshes with  
157 maximum spacings between 2.5 km (ARPEGE, ICON) and 7.8 km (UM).  
158 All models except GEOS were initialized with the 00 UTC 1 August 2016  
159 analysis from the *European Centre for Medium-Range Weather Forecasts*  
160 (*ECMWF*) and integrated for 40 days (1 August–10 September 2016). The  
161 sea surface temperature and sea ice fields were prescribed using 7-day run-  
162 ning mean analyses from *ECMWF*.

163 To identify TCs in the model output, we employed the *GFDL vortex*  
164 *tracker* (Marchok 2002; Biswas et al. 2018). This software searches for TCs  
165 based on spatial minima and maxima in the following fields: (i) relative  
166 vorticity at 10 m, (ii) sea-level pressure, (iii) wind speed at 10 m, 850

---

<sup>1</sup>The IFS model considered here is an experimental version of the operational IFS model with 4-km mesh spacing and explicitly simulated deep convection

167 hPa, and 700 hPa, and (iv) layer-mean temperature between 300-500 hPa  
168 (Marchok 2002). The tracker produces *track files* with 6-hourly *records*  
169 that contain TC location (latitude/longitude), maximum 10-m wind speed  
170 ( $v_{max}$ ), minimum sea-level pressure ( $p_{min}$ ), and wind radii  $r_{17}$ ,  $r_{25}$ , and  $r_{32}$ ,  
171 i.e., the maximum radial extent of 17 m s<sup>-1</sup>, 25 m s<sup>-1</sup>, and 32 m s<sup>-1</sup> winds  
172 in each compass quadrant (northeast, southeast, southwest, and northwest).

173 To evaluate the models, we used best track data from the International  
174 Best Track Archive for Climate Stewardship [IBTrACS version 4; Knapp  
175 et al. (2010, 2018)]. Specifically, we used the data from the WMO agency  
176 responsible for a given storm, and we accounted for wind speed reporting  
177 differences by converting all  $v_{max}$  values to 1-min sustained winds following  
178 Harper et al. (2008). Note that the IBTrACS data does not contain di-  
179 rect observations or objective analyses, but subjective analyses from human  
180 forecasters based on available but limited observations. For simplicity, we  
181 will nevertheless refer to the IBTrACS data as “observations”.

182 For a number of reasons, the workflow was not trivial. For example, some  
183 groups provided the output on their native model mesh, which rendered the  
184 data unreadable for the tracker. Furthermore, the high-resolution output  
185 caused the tracker to falsely identify hundreds of convective objects as TCs.  
186 To overcome those issues, we carried out the following three-step process:

187 1. Interpolate the output from each model to a common longitude/latitude

188 grid with  $0.5^\circ$  resolution.

189 2. Run the tracker on the interpolated grids. Keep in mind that the  
190 track files contain information from the smoothed data.

191 3. Use the storm center information from step 2 to search for the actual  
192  $v_{max}$ ,  $p_{min}$ , and  $r_{17}$ ,  $r_{25}$ ,  $r_{32}$  in the native model files, and overwrite  
193 the data in the track files with these new values.

194 Even after this process, the software tracked objects that human meteo-  
195 rologists would not identify as TCs, such as disorganized convective systems  
196 and heat lows over the deserts of Iran and central Asia. To reduce the num-  
197 ber of falsely-identified objects as much as possible, the track files were  
198 quality-controlled using the following criteria:

- 199 • drop all storms that form inland over Arabia and Iran,
- 200 • drop all storms with lifetimes under 48 h,
- 201 • drop all storms that never achieved a  $v_{max}$  of  $7.5 \text{ m s}^{-1}$ ,
- 202 • drop all records poleward of  $\pm 40^\circ$  latitude (i.e., remove storms that  
203 become extratropical).

204 The IBTrACS data were quality-controlled using the same criteria to ho-  
205 mogenize model data and observations.

## 206 3. Results

### 207 3.1 *Number of Tropical Cyclones and Tracks*

Fig. 1

Fig. 2

208 Meteorological services observed a global total of 24 TCs during the 40-  
209 day DYAMOND period, while the models simulated between 12–31 TCs,  
210 i.e., 50–140% of the observed value (Fig. 1). Most of the models simulated  
211 fewer TCs than observed; specifically, six of the nine models simulated less  
212 than 24 TCs (ARPEGE, FV3, ICON, IFS, MPAS, SAM; Figs. 1b,c,e,f,g,i),  
213 and only NICAM and UM simulated more TCs than observed (Figs. 1h,j).  
214 GEOS simulated exactly 24 TCs (Fig. 1d), however, given the limited  
215 sample size and the likelihood that a different tracker may have yielded  
216 slightly different numbers, we do not wish to emphasize the exact number  
217 of TCs each model produced.

218 According to the observations, the Western Pacific was the most active  
219 basin during the DYAMOND time period, followed by the Eastern Pacific,  
220 Atlantic, and Indian Ocean. All models agreed that the Western Pacific was  
221 going to be the most active basin, and the simulated tracks were generally  
222 oriented from south to north like in the observations (Fig. 1). A plausible  
223 reason for the track agreement is that all models were able to capture the  
224 large-scale steering flow over the Western Pacific. However, the models were  
225 not as successful in the other basins. For example, in the Eastern Pacific,



226 all models except MPAS (Fig. 1g) simulated fewer TCs than observed, and  
227 there was less agreement between observed and simulated tracks. FV3 seems  
228 to have done best in terms of tracks in this basin (Fig. 1c). TC activity  
229 in the Atlantic proved to be particularly difficult to capture, and some  
230 models simulated a very active basin while others simulated a very quiet  
231 one. Specifically, NICAM produced 11 Atlantic TCs (Fig. 1h), whereas  
232 FV3 and IFS only produced one (Figs. 1c,f).

233 TC formation events during the DYAMOND period were not spread out  
234 uniformly over time but occurred in more or less well-defined periods (Fig.  
235 2). The models simulated the temporal modulation of activity in rough  
236 agreement with the observations. For example, in the Western Pacific, most  
237 models correctly simulated a greater number of formation events before  
238 22 August than after that date (Fig. 2a). In the Eastern Pacific, the  
239 models missed some of the formation events in early August, but they agreed  
240 with the observations on a second round of activity in late August/early  
241 September (Fig. 2b). In the Atlantic, about half of the models suggested  
242 a relatively active period in mid/late August, around the same time four  
243 formation events were observed (Fig. 2c). On the other hand, the models  
244 struggled with capturing the timing of TC formation in the Indian Ocean  
245 (Fig. 2d); however, with only two observed events, this basin is likely not  
246 representative.

247 At this point we can only speculate why the models were able to capture  
248 the temporal modulation of activity beyond the typical predictability limit  
249 of weather prediction, which is around two weeks. One possible reason is  
250 that the models were able to capture the modulating effect of intraseasonal  
251 variability as previously shown by Nakano et al. (2015). Another possible  
252 reason is that the pre-scribed sea-surface temperatures artificially impart  
253 longer predictability on the atmosphere.

254 Perhaps most importantly, Figure 2 demonstrates that no model suffered  
255 from a climate drift, that is, no model showed the number of TC formation  
256 events to unrealistically increase or decrease over the 40-day period. This  
257 highlights the quality of the DYAMOND models, which were not tuned for  
258 the experiment.

259 As a final remark, we note that UM produced three ensemble members  
260 in addition to the official 40-day DYAMOND run. The differences in TC  
261 numbers and tracks within that ensemble were as large as (or at times larger  
262 than) than inter-model differences (not shown). This indicates that more  
263 simulations and ensemble runs are needed to properly assess the predictive  
264 skill of each model beyond the broad statements made above.

### 3.2 Tropical Cyclone Intensity

Fig. 3

Timeseries of  $v_{max}$  in Fig. 3 provide a broad overview of the intensity of the TCs and allow for a cursory model evaluation. Some biases are clearly evident; for example, ICON and SAM produced storms that were generally too weak (Figs. 3e,i), whereas ARPEGE produced a few storms that were much too strong. In fact, ARPEGE produced storms with unrealistically high  $v_{max}$  of  $>100$  m s<sup>-1</sup> (Fig. 3b), most likely because the evaporation coefficient was set to a wrong value (Stevens et al. 2019).

Fig. 4

Fig. 5

Fig. 6

According to the observations, the TCs during the first two weeks of August remained relatively weak with only two storms reaching hurricane intensity ( $v_{max} \geq 33$  m s<sup>-1</sup>; Fig. 3a). On the other hand, some of the TCs that formed in the second half of August became quite intense with four storms reaching major hurricane intensity ( $v_{max} \geq 50$  m s<sup>-1</sup>). Most models had issues with capturing this pattern. Specifically, a number of models simulated storms in the first half of August that were too intense (ARPEGE, GEOS, NICAM, UM; Figs. 3b,d,h,j). From all models, MPAS seems to have best captured the overall pattern (Fig. 3g).

To evaluate the models regarding intensity in more depth, we compared the observed and modeled frequency distributions of  $v_{max}$  (Fig. 4) and  $p_{min}$  (Fig. 5). We chose to compare frequency distributions instead of  $v_{max}$  and  $p_{min}$  errors, because the models did not simulate all observed TCs and not

286 all simulated TCs were observed. We present the frequency distributions  
287 by way of *kernel density estimates* (Silverman 2018), because this method  
288 yields smooth curves that make a comparison easier. The kernel density  
289 estimates were implemented using the python seaborn library.

290 The observed  $v_{max}$  distribution has a broad primary peak centered near  
291  $20 \text{ m s}^{-1}$ , a secondary peak near  $50 \text{ m s}^{-1}$ , and a fat tail towards higher  
292 values (Fig. 4). All models were able to produce this bi-modal distribution  
293 to some degree, but certain models deviated more from the observations  
294 than others. ICON and SAM deviated most dramatically: both models  
295 produced a narrow primary peak, mainly because they were not able to  
296 simulate high intensities (Figs. 4d,h). FV3 and GEOS shifted the secondary  
297 peak to higher values (Figs. 4b,c), whereas IFS and MPAS shifted it to lower  
298 values (Figs. 4e,f). ARPEGE produced a very broad distribution, partly  
299 related to its over-intensification issue (Fig. 4a). NICAM reproduced the  
300 observed distribution for  $v_{max} > 25 \text{ m s}^{-1}$  better than the other models,  
301 but missed some of the weaker intensities with  $v_{max} < 20 \text{ m s}^{-1}$  (Fig. 4g).

302 The observed  $p_{min}$  distribution has a well-defined primary peak around  
303  $1000 \text{ hPa}$ , and a fat tail extending towards lower pressures with hint of  
304 a secondary maximum near  $950 \text{ hPa}$  (Fig. 5). All models were able to  
305 capture the general shape of the observed distribution, with MPAS and  
306 UM matching the observations best (Figs. 5f,i). Most of the other models

307 produced storms that were too deep, although in different ways. In FV3,  
308 the distribution showed the same shape as the observation but shifted to  
309 deeper values (Figs. 5b); in IFS, the secondary maximum was much more  
310 pronounced than in the observations (Figs. 5f); and GEOS was somewhere  
311 in between FV3 and IFS (Figs. 5c). In ARPEGE and NICAM, some  
312 storms were much deeper than the observations, causing the tail to stretch  
313 too far to the left (Figs. 5a,g). SAM is unique in that the main peak was  
314 shifted to much higher values. We shall note here that SAM's  $p_{min}$  values  
315 are ambiguous, because SAM uses the anelastic equations and the quantity  
316 *pressure* can only be determined to within a function proportional to the  
317 base-state density field with arbitrary amplitude (Bannon et al. 2006).

318 Lastly, we evaluated the overall TC activity by means of *accumulated*  
319 *cyclone energy* (ACE), a quantity that estimates the wind energy produced  
320 by one or multiple TCs over their lifetime. It is computed according to  
321  $ACE = 10^{-4} \sum v_{max}^2$ , where  $v_{max}$  is in units of knots (1 knot = 0.51 m  
322  $s^{-1}$ ). According to the observations, the ACE during the DYAMOND pe-  
323 riod was 169 (Fig. 6). Since the wind speed enters the ACE calculation  
324 as a squared value, ACE is quite sensitive to uncertainty in the analyzed  
325  $v_{max}$  values. We therefore estimated a lower and upper bound by assum-  
326 ing that all observed  $v_{max}$  records have an error of  $\pm 5$  m  $s^{-1}$ , an estimate  
327 based on Torn and Snyder (2012) and Landsea and Franklin (2013). This

328 assumption yielded a lower bound of 118 ACE units and an upper bound  
329 of 230 ACE units. Most models were within these uncertainty bounds or  
330 slightly above, indicating that the DYAMOND models produced realistic  
331 amounts of ACE, even without tuning. Only three models were outside the  
332 uncertainty bounds: GEOS overestimated ACE, whereas ICON and SAM  
333 produced less ACE than observed.

### 334 3.3 Tropical Cyclone Size

Fig. 7

335 Size is an important TC parameter because it correlates with the risk  
336 for storm surge, but it is often neglected and infrequently used for model  
337 validations. We examined the radius of gale-force winds ( $r_{17}$ ) and present  
338 the median of all  $r_{17}$  records as our metric of choice (Fig. 7). Results for  $r_{25}$   
339 and  $r_{32}$  were qualitatively similar (not shown), indicating that the results  
340 are not sensitive to a particular wind speed threshold. The observational  
341 error bars were computed by increasing/decreasing each  $r_{17}$  record by 50%  
342 (Landsea and Franklin 2013).

343 In general, the models overestimated TC size. TCs in ARPEGE, FV3,  
344 ICON, and NICAM were substantially larger than observed (Figs. 7a,b,d,g).  
345 In fact, ARPEGE and ICON produced very expansive wind fields, and  
346 their median  $r_{17}$  reached radially outward to 300 km (more than double  
347 the observations). In contrast, the median size of TCs in GEOS matched

348 the observations remarkably well (Fig. 7c), and UM came in as a clear  
349 second (Fig. 7i). Storms in IFS and SAM were somewhat smaller than  
350 observed, but still within the uncertainty estimates (Figs. 7e,h). A common  
351 bias in the models was associated with the asymmetry of the wind field.  
352 Concretely, the observed  $r_{17}$  was largest in the northeast quadrant, but in  
353 FV3, ICON, MPAS, and NICAM, it was largest in the southeast quadrant  
354 (Figs. 7b,d,f,g). This result suggests that the models are deficient in their  
355 representation of TC structure; the prospect of which will be examined in  
356 the next section.

### 357 3.4 Tropical Cyclone Structure

358 The TC wind-pressure relationship, i.e., the function that relates  $p_{min}$   
359 to  $v_{max}$ , is often used to inform whether models simulate TC structure re-  
360 alistically. The DYAMOND models produced a variety of wind pressure  
361 relationships, with some models being closer to the observation than others  
362 (Fig. 8). FV3 and GEOS stand out for reproducing the observed rela-  
363 tionship remarkably well (Fig. 8b, c). Most other models have a tendency  
364 to produce a relationship that drops off too fast, or in other words, for a  
365 given  $p_{min}$ , the  $v_{max}$  is too low. This behavior was most pronounced in  
366 ICON (Fig. 8d), and less noticeable in ARPEGE and MPAS (Fig. 8a,f). A  
367 possible explanation for this behavior will be discussed in section 4. SAM

Fig. 8

Fig. 9

Fig. 10

Fig. 11

Fig. 12

368 was unique and had an unrealistic wind-pressure relationship that bended  
369 upward (Fig. 8h). This phenomenon was not due to a single outlier but  
370 likely related to the the surface pressure field being an ambiguous quantity  
371 in this model (see also section 3.2).

372 Since the 10-m winds in a TC and therefore  $v_{max}$  are strongly affected  
373 by the surface layer parameterization, we also investigated the relationship  
374 between  $p_{min}$  and 850-hPa  $v_{max}$ . The graphs were qualitatively similar to  
375 Fig. 8 (not shown), indicating that the wind-pressure relationships in Fig. 8  
376 are not merely a product of each model’s boundary layer and surface layer  
377 parameterizations, but stem from differences in the overall model imple-  
378 mentation including the dynamical cores.

379 Snapshots of 10-m wind speed demonstrate the diversity of the models  
380 in simulating the surface wind field (Fig. 9). There were striking differences  
381 in eyewall shape, size, and symmetry, as well as in the radial extent of the  
382 wind field. Some models produced unrealistic wind fields, either too large  
383 and too strong (ARPEGE; Fig. 9a), or too faint and with peculiar waviness  
384 (SAM; Fig. 9h). The wind fields of FV3, GEOS, and MPAS were arguably  
385 most similar to that of a canonical intense TC, with a distinct eyewall that  
386 contained multiple convective- and mesoscale asymmetries (Figs. 9b,c,f).

387 The ICON example was unique in that it did not reveal a distinct eyewall  
388 with sharp gardients; its wind field was rather diffuse and spread out over



389 a large area (Fig. 9d). In contrast, the IFS example was a very small TC  
390 with a radially constrained wind field (Fig. 9e). The NICAM example, Fig.  
391 9g, had an even larger hurricane-force (wind speed  $\geq 33 \text{ m s}^{-1}$ ) wind field  
392 than ICON, but it also had a distinct eyeall like most other models—albeit  
393 somewhat smoother than the eyewalls in FV3, GEOS, and MPAS. The wind  
394 field from the UM example exhibited the smoothest structure, the widest  
395 eyewall, and the clearest imprint of the model mesh—all consistent with  
396 UM being the model with the lowest resolution (Fig. 9i).

397 A closer look at the kinematic structure of the modeled TCs was achieved  
398 by creating composites of the azimuthally-averaged circulation (Fig. 10).  
399 Each model’s composite includes the individual cases where  $v_{max} \geq 33 \text{ m}$   
400  $\text{s}^{-1}$ . Broadly speaking, all models produced a typical kinematic structure,  
401 that is, a well-defined *primary circulation* with a tangential wind maximum  
402 in the lower troposphere near the storm center, and a well-defined *secondary*  
403 *circulation* manifested by strong radial inflow in the boundary layer, rising  
404 motion in the eyewall region, and radial outflow in the mid- to upper tro-  
405 posphere. Despite the overall agreement, there were noteworthy differences  
406 between the models, which will be discussed next. Note that we will assume  
407 that the inter-model structure differences are due to model formulation and  
408 not due the varying intensity of the composite storms.

409 The differences in the overall tangential wind structure can be eluci-

410 dated by comparing the size of the radius of maximum tangential wind, the  
411 compactness of the wind maximum (specifically, the radial extent of the  
412  $35 \text{ m s}^{-1}$  isotach), and the decay of the tangential wind in the radial and  
413 vertical direction. The composite storms had radii of maximum tangential  
414 wind roughly between 30–70 km, with ARPEGE and IFS on the lower end  
415 (Figs. 10a,e) and ICON on the upper end (Fig. 10d). In FV3 and MPAS,  
416 the wind maximum was comparatively narrow and confined, and the radial  
417 extent of the  $35 \text{ m s}^{-1}$  isotach was less than 20 km (Figs. 10b,f). On the  
418 other hand, in ICON and NICAM, the wind maximum was rather broad,  
419 and the radial extent of the  $35 \text{ m s}^{-1}$  isotach was greater than 50 km (Figs.  
420 10e,g). Differences in the radial and vertical decay rates mirror the previ-  
421 ous discussion of storm size, that is, models in which the tangential wind  
422 decayed more slowly, such as in ICON and NICAM, were the ones that  
423 produced comparatively larger storms.

424     Given the lack of an equivalent observational dataset, it is difficult to  
425 assess what model produced a particularly realistic tangential wind struc-  
426 ture. The observational composites of Gao et al. (2019, their Fig. 5c) and  
427 Komaromi and Doyle (2017, their Fig. 7a) at least suggests that no model  
428 produced a particularly unrealistic structure.

429     As for the vertical motion, ARPEGE and IFS had the steepest eyewall  
430 slopes (Figs. 10a,e). The other extreme was UM, which had the most

431 pronounced eyewall tilt (Fig. 10i). In ICON and NICAM, the eyewall  
432 updraft was spread out and diffuse (Figs. 10d,g), but in IFS and MPAS it  
433 was relatively narrow and confined (Figs. 10e,f). Besides these differences  
434 in the eyewall region, there were differences in the rainband region, too.  
435 Specifically, the vertical motion between  $r = 100\text{--}250$  km was noticeably  
436 stronger in ICON, MPAS, and NICAM than in GEOS, IFS, and SAM (Figs.  
437 10d,f,g versus Figs. 10c, e,h). This difference may be a reflection of more  
438 or stronger rainbands in the former models.

439 Again, it is difficult to say which models produced a particularly realistic  
440 structure because no equivalent observational dataset exists for the TCs  
441 observed during the DYAMOND period. Stern and Nolan (2009) showed  
442 that the slope of the eyewall depends on the size of the radius of maximum  
443 wind, which would explain why the eyewall updraft in IFS has a steeper  
444 slope than in IFS, but it cannot explain the differences between models  
445 with similarly sized radii of maximum wind, such as MPAS and UM.

446 The upper-tropospheric outflow also differed between the models, espe-  
447 cially with regard to the altitude of the outflow maximum and the depth of  
448 the outflow layer. For instance, the outflow was comparatively deep in FV3  
449 (Fig. 10b) and comparatively shallow in IFS (Fig. 10e). In ARPEGE and  
450 ICON, the outflow maximum occurred at a height of 15 km (Figs. 10a,d),  
451 but in most of the other models, it occurred mostly below 15 km.

452 One particularly noteworthy feature, produced somewhat more promi-  
453 nently by FV3, GEOS, and IFS, is the descending flow above the outflow  
454 layer that merges with the ascending outflow from below (Figs. 10b,c,e). We  
455 are not aware of either observational or modeling studies that show such a  
456 feature in TCs; to the contrary, there is reasonable evidence to suggest that  
457 at least in intense TCs, it may be common to have a shallow layer of weak  
458 inflow atop the upper-level outflow layer (e.g., Kieu et al. 2016; Komaromi  
459 and Doyle 2017; Heng et al. 2017; Duran and Molinari 2018).

460 Inter-model differences in the boundary layer inflow were mostly in the  
461 form of variations of inflow layer depth and strength (Fig. 11). Specifically,  
462 IFS and SAM produced comparatively shallow inflow layers that did not  
463 extend much above 1 km height (Figs. 11e,h). In GEOS and ICON, the  
464 inflow layer had a maximum depth of 1.5 km (Figs. 11c,d), and in the  
465 other models, its maximum depth extended slightly above 1.5 km. The  
466 observational composite of Zhang et al. (2011, their Fig. 5b) shows that  
467 the inflow layer depth increases from 900 m at the radius of maximum wind  
468 to 1.5 km roughly 200 km from the center, which is in broad agreement  
469 with most of the models.

470 From basic TC dynamics one would expect that the inflow strength  
471 correlate with the average intensity of the TCs simulated by the models.  
472 However, this was not the case. For example, ICON, which simulated mostly

473 weak TCs, produced stronger inflow than FV3, MPAS, and NICAM, which  
474 simulated much stronger TCs (Fig. 11d vs. Figs. 11b,f,g). In fact, with  
475 inflow magnitudes of  $9 \text{ m s}^{-1}$ , the inflow in FV3, MPAS, and NICAM was  
476 relatively weak not only compared to the other models, but also compared  
477 to observations, which show an inflow magnitude of  $20 \text{ m s}^{-1}$  (Zhang et al.  
478 2011).

479 Besides the kinematic structure, we also explored the thermodynamic  
480 TC structure in our set of global storm-resolving simulations. To this  
481 end, we examined the TC warm core, here represented by the tempera-  
482 ture anomaly with respect to the mean temperature between  $r=300\text{--}700$   
483 km (Fig. 12). All models produced a warm core, and all models agreed  
484 on its general structure (expansive in the upper levels, radially confined be-  
485 low). Differences emerged mostly in the vertical structure of the warming  
486 inside the TC eye, and in the upper and lower level temperature anomalies  
487 outside the eye.

488 Most models agreed that the warm anomaly peaks at a height of just  
489 under 10 km. More pronounced differences between the models appeared  
490 in the vertical structure of the warm core, which ranged from a single, ver-  
491 tically confined maximum in FV3 and GEOS (Figs. 12b,c), to an extended  
492 vertical column in NICAM (Fig. 12g), to a clear double maximum of anoma-  
493 lously warm air in UM (Fig. 12i). The other models fell somewhere in

494 between these three distinct cases. Most observational studies indicate that  
495 the warm core is maximized in the upper troposphere (Frank 1977; Bram-  
496 mer and Thorncroft 2017; Komaromi and Doyle 2017), in agreement with  
497 most of the DYAMOND models. However, Stern and Nolan (2012) claimed  
498 that the maximum warming should be between 4–8 km, with a potential  
499 secondary maximum at higher altitudes. Kieu et al. (2016) also claimed  
500 that a double-warm core structure is the norm rather than the exception.  
501 According to those studies, UM had a particularly realistic thermodynamic  
502 structure, even though it was an outlier among the DYAMOND models.

503 Compared to the model differences in terms of the warm core, the dif-  
504 ferences above the outflow layer were equally if not more striking. Above  
505 15-km height, the models did not even agree on the sign of the temperature  
506 anomaly. In particular, IFS and ARPEGE produced a strong cool anomaly  
507 ( $< -3$  K; incidentally, IFS and ARPEGE were the only spectral models),  
508 whereas NICAM, SAM, and UM produced a warm anomaly. FV3, GEOS,  
509 ICON, and MPAS were somewhere in between the extremes and produced  
510 a weak cool anomaly ( $> -1$  K). Observational composites generally show a  
511 weak cold anomaly above the outflow layer (Frank 1977; Komaromi and  
512 Doyle 2017), although instantaneous snapshots of intense TCs may also  
513 show strong cold anomalies (Komaromi and Doyle 2017).

514 Temperature differences were also found in the boundary layer, although

515 less dramatic: NICAM was anomalously cool (Fig. 12g), and IFS was  
516 anomalously warm (Fig. 12e). The other models had weak cool anomalies  
517 or no clear signal. Note that IFS and NICAM were polar opposites of each  
518 other (NICAM: warm in the upper levels, cool in the lower levels, IFS: vice  
519 versa).

### 520 *3.5 Sensitivity of Tropical Cyclone Formation and Intensity* 521 *to Model Resolution and Parameterized Deep Convection*

Fig. 13

522 In addition to the primary high-resolution simulation, some DYAMOND  
523 models produced sensitivity runs with lower resolution. For example, ICON  
524 produced six simulations with mesh spacings of 2.5, 5, 10, 20, 40, and 80  
525 km, all **without** parameterized convection (ICON no-conv), and an ad-  
526 ditional three simulations with mesh spacings of 20, 40, and 80 km **with**  
527 parameterized convection (ICON conv). These nine simulations provided  
528 an opportunity to investigate the sensitivity to model resolution and pa-  
529 rameterized convection in a controlled way (Fig. 13, Fig. 14) .

Fig. 14

530 As for sensitivity to resolution, there was a clear inverse relationship and  
531 the number of simulated TCs **increased** when resolution was **decreased**  
532 (Fig. 13, left column). Concretely, the highest resolution run produced  
533 the fewest TCs (15; Fig. 13a), and the lowest resolution run produced the  
534 most TCs (50; Fig. 13h). In the simulations with intermediate resolution,

535 the number of TCs was relatively constant (around 20). The sensitivity  
536 to resolution seemed to be basin dependent. In the Atlantic and Eastern  
537 Pacific, the 80-km ICON produced five to six times as many TCs as the  
538 2.5-km ICON (Figs. 13a,h), but in the Western Pacific, the 80-km ICON  
539 produced only two times as many TCs as the 2.5-km ICON. In the Indian  
540 Ocean, the number of events seemed to be insensitive to resolution, and  
541 each run produced either one or two TCs.

542 As for sensitivity to parameterized convection, the model produced dra-  
543 matically fewer TCs once the parameterization was turned on (Fig. 13, left  
544 vs. right column). This effect was most pronounced at lower resolution.  
545 Specifically, the number of TCs dropped from 23 to 17 in the 20-km runs  
546 (Figs. 13d,e), from 21 to 14 in the 40-km runs (Figs. 13f,g), and from 50  
547 to a mere 9 in the 80-km runs (Figs. 13h,i).

548 The runs with parameterized convection also featured substantially lower  
549 ACE (Fig. 14). Again, the effect was most dramatic at lower resolution,  
550 but even for an intermediate resolution of 20 km, the ACE was reduced  
551 by 65%. This result suggests that convection parameterization did not just  
552 reduce the number of TCs, but it also made them weaker and their lifetime  
553 shorter.

554 Interestingly, the ICON no-conv runs produced more or less the same  
555 amount of ACE at all resolutions (Fig. 14). Evidently, the lack of intense



556 storms in the lower-resolution runs was compensated by a larger number  
557 of weak storms. An interesting follow-up question would be to investigate  
558 whether this compensation was pure luck or whether the amount of back-  
559 ground available potential energy that is converted to kinetic energy by TCs  
560 is a resolution-independent quantity, such as mean precipitation (Hoheneg-  
561 ger et al. 2020).

## 562 4. Discussion

Fig. 15

563 One of the drawbacks of global storm-resolving models is their immense  
564 computational cost, which poses questions about cost versus benefit. One  
565 may, for example, postulate that regional high-resolution models suffice  
566 for TC prediction. Although a practical alternative, regional models have  
567 disadvantages such as determining the ideal domain size and placement  
568 for a regional domain. More importantly, regional domains require lateral  
569 boundary conditions, which have “serious negative effects” (Warner et al.  
570 1997). One of those effects is that errors creep in through the boundaries  
571 and render longer-range forecasts less skillful than those made by global  
572 models. Putting it slightly differently, regional models are very dependent  
573 on the global model forcing being “good enough”.

574 One may also follow Manganello et al. (2012) and argue that hydrostatic  
575 models with mesh spacings of 10 km and parameterized convection are

576 sufficient for producing realistic TCs. Nonetheless, mesh spacings of  $<5$  km  
577 are still required for realistically simulating  $v_{max}$  and the dynamic processes  
578 in the TC inner core (e.g., Chen et al. 2007; Gentry and Lackmann 2010;  
579 Judt and Chen 2010; Gopalakrishnan et al. 2012; Davis 2018). Observations  
580 and numerical models indicate that such processes are important for rapid  
581 intensification (e.g., Miyamoto and Takemi 2015; Guimond et al. 2016; Judt  
582 and Chen 2016). In fact, a case study by Fox and Judt (2018) suggested that  
583 simulating extreme cases of rapid intensification requires  $\leq 1$  km horizontal  
584 grid spacing. Since extreme storms are highly disruptive to society, being  
585 able to reliably predict or project intense TCs has great value.

586 As a potential easy target for bias reduction in the models, we examined  
587 whether models with similar biases used similar parameterization schemes.  
588 For example, we investigated whether the models with a TKE-like boundary  
589 layer parameterization produced similar intensity biases versus models that  
590 used a diagnostic eddy diffusivity. However, no such relationships were  
591 found. In the end, there are variety of reasons for the model diversity,  
592 including but not limited to: cloud microphysics, boundary layer processes,  
593 and the dynamical cores (with differences in effective resolutions).

594 In agreement with other studies, this paper also demonstrates that high  
595 resolution is necessary yet not sufficient to capture the  $v_{max}$  of TCs. For  
596 example, ICON was tied with ARPEGE for highest resolution (2.5 km),

597 yet ICON struggled to produce intense TCs while ARPEGE produced un-  
598 realistically strong TCs. These intensity biases are likely a consequence of  
599 the respective model’s surface flux formulation, as demonstrated by Fig.  
600 15, which shows the surface fluxes of momentum and latent heat over an  
601 area  $300\times 300$  km centered on the strongest TC in each model. The drag  
602 in ICON increased much faster with wind speed than in ARPEGE (Fig.  
603 15a), which means that there was a comparatively stronger “break” on the  
604 surface wind in ICON. ICON also had significantly weaker latent heat fluxes  
605 for a given wind speed, providing less amount of “fuel” (Fig. 15b).

606 The monotonically increasing momentum flux in Fig. 15a also indicates  
607 that the models did not account for the saturation of the drag at wind speeds  
608 above  $25\text{ m s}^{-1}$  (e.g., Powell et al. 2003; Donelan 2004; Chen et al. 2013;  
609 Curcic and Haus 2020). This shortcoming was found in other models as well  
610 (not shown), and it may be the reason why the wind-pressure relationship  
611 in several models deviated from observations at higher winds (Fig. 8). In  
612 fact, the wind-pressure relationship in IFS seems to improve when drag  
613 is computed in a more realistic three-way coupled atmosphere-wave-ocean  
614 model (Magnusson et al. 2019).

615 Lastly, there is much evidence that the storm count (and storm-count-  
616 related model biases) are sensitive to the tracker and to the model formula-  
617 tion/resolution (Roberts et al. 2020; Vanniere et al. 2020). This can be an

618 issue when comparing models as weak TCs might be over- or under-detected  
619 depending on the threshold used. In the end, only further studies can speak  
620 to the robustness of the results presented in this paper.

## 621 **5. Summary and Conclusions**

622 We evaluated nine global storm-resolving models that participated in the  
623 DYAMOND initiative (Stevens et al. 2019) in their ability to simulate TCs.  
624 Specifically, we validated and compared the number of TCs each model pro-  
625 duced, their tracks, intensity, size, and structure. With mesh spacings be-  
626 tween 2.5–7.8 km, the DYAMOND models are the highest-resolution global  
627 models that have so far been analyzed for this purpose.

628 The results suggest that global storm-resolving models are able to sim-  
629 ulate the structure and intensity of TCs more realistically than previous  
630 generations of global models. However, we found that TCs are strongly  
631 affected by model formulation, and essentially all models had biases. We  
632 found that no model did best in all regards, although some models did,  
633 generally speaking, better than others. For instance, GEOS produced the  
634 observed number of TCs, captured TC size better than any other model,  
635 and produced a realistic wind-pressure relationship. But GEOS also pro-  
636 duced too many strong storms and had the largest ACE bias of all models  
637 (it is unclear if ocean coupling would reduce this bias). Other models that

638 did generally well were FV3, MPAS, and UM.

639 On the other hand, ICON, IFS, and SAM had some issues with size,  
640 structure, and intensity. For example, ICON and SAM produced storms  
641 that were too weak. ICON, IFS, and SAM were also not able to capture  
642 the wind-pressure relationship as realistically as GEOS, FV3, and MPAS,  
643 pointing to deficiencies in their numerical formulations. We also found that  
644 parameterized convection strongly reduces the number and intensity of TCs  
645 in comparison to simulations without convection parameterization (at least  
646 for simulations with a mesh spacing of  $>20$  km). This sensitivity high-  
647 lights the problems and ambiguities that come with parameterizing deep  
648 convection.

649 In a nutshell, we believe that the ability to realistically simulate TCs  
650 in global models is critical for weather and climate prediction. This study  
651 demonstrates that global-storm resolving models are an optimal tool to  
652 advance TC prediction; however, they need to be improved to unleash their  
653 full potential. Surface layer, pbl are targets for improvements.

## 654 **Acknowledgements**

655 This material is based upon work supported by the National Center for  
656 Atmospheric Research, which is a major facility sponsored by the National  
657 Science Foundation under Cooperative Agreement No. 1852977. Support

658 by the Centre of Excellence ESiWACE is acknowledged. ESiWACE has  
659 received funding from the European Unions Horizon 2020 research and in-  
660 novation programme under grant agreement numbers 675191 and 823988.  
661 The authors thank the German Climate Computing Center for hosting,  
662 providing access to and supporting at evaluating the DYAMOND data  
663 sets. We acknowledge high-performance computing support from Cheyenne  
664 (doi:10.5065/D6RX99HX) provided by NCAR's Computational and Infor-  
665 mation Systems Laboratory, sponsored by the National Science Foundation.  
666 RS, MS, MN and CK were supported by the FLAGSHIP2020 within the  
667 priority study4 (Advancement of meteorological and global environmental  
668 predictions utilizing observational "Big Data") and the Integrated Research  
669 Program for Advancing Climate Models (TOUGOU) Grant Number JP-  
670 MXD0717935457 from the Ministry of Education, Culture, Sports, Science,  
671 and Technology of Japan. The NICAM simulation was performed on the  
672 Earth Simulator of the Japan Agency for Marine-Earth Science and Tech-  
673 nology. MR, PLV and BV acknowledge funding from the EU H2020 PRI-  
674 MAVERA project under Grant Agreement no. 641727. The authors would  
675 also like to thank David Ahijveych for his assistance with the GFDL tracker,  
676 and Dan Stern, Will Komaromi, and Kevin Hodges for valuable comments.  
677 The comments from XXX reviewers helped improving the manuscript.

## References

678

679 Bannon, P. R., J. M. Chagnon, and R. P. James, 2006: Mass conservation  
680 and the anelastic approximation. *Monthly Weather Review*, **134(10)**,  
681 2989–3005.

682 Bengtsson, L., K. I. Hodges, M. Esch, N. Keenlyside, L. Kornblueh, J.-J.  
683 Luo, and T. Yamagata, 2007: How may tropical cyclones change in a  
684 warmer climate? *Tellus A: Dynamic Meteorology and Oceanography*,  
685 **59(4)**, 539–561.

686 Biswas, M. K., D. Stark, and L. Carson, 2018: *GFDL Vortex Tracker Users’*  
687 *Guide V3.9a*. National Center for Atmospheric Research, Develop-  
688 mental Testbed Center.

689 Brammer, A., and C. D. Thorncroft, 2017: Evaluation of reanalysis tropical  
690 cyclone structure with global hawk 1 dropsonde observations.

691 Camargo, S. J., A. G. Barnston, and S. E. Zebiak, 2005: A statistical assess-  
692 ment of tropical cyclone activity in atmospheric general circulation  
693 models. *Tellus A: Dynamic Meteorology and Oceanography*, **57(4)**,  
694 589–604.

695 Chen, S. S., J. F. Price, W. Zhao, M. A. Donelan, and E. J. Walsh, 2007:  
696 The CBLAST-hurricane program and the next-generation fully cou-

697       pled atmosphere–wave–ocean models for hurricane research and pre-  
698       diction. *Bulletin of the American Meteorological Society*, **88(3)**, 311–  
699       318.

700   Chen, S. S., W. Zhao, M. A. Donelan, and H. L. Tolman, 2013: Directional  
701       wind–wave coupling in fully coupled atmosphere–wave–ocean mod-  
702       els: Results from CBLAST-hurricane. *Journal of the Atmospheric*  
703       *Sciences*, **70(10)**, 3198–3215.

704   Curcic, M., and B. K. Haus, 2020: Revised estimates of ocean surface drag  
705       in strong winds. *Geophysical Research Letters*, **47(10)**.

706   Davis, C., W. Wang, S. S. Chen, Y. Chen, K. Corbosiero, M. DeMaria,  
707       J. Dudhia, G. Holland, J. Klemp, J. Michalakes, H. Reeves, R. Ro-  
708       tunno, C. Snyder, and Q. Xiao, 2008: Prediction of landfalling hur-  
709       ricanes with the advanced hurricane WRF model. *Monthly Weather*  
710       *Review*, **136(6)**, 1990–2005.

711   Davis, C. A., 2018: Resolving tropical cyclone intensity in models. *Geo-*  
712       *physical Research Letters*, **45(4)**, 2082–2087.

713   DeMaria, M., C. R. Sampson, J. A. Knaff, and K. D. Musgrave, 2014:  
714       Is tropical cyclone intensity guidance improving? *Bulletin of the*  
715       *American Meteorological Society*, **95(3)**, 387–398.



- 716 Donelan, M. A., 2004: On the limiting aerodynamic roughness of the ocean  
717 in very strong winds. *Geophysical Research Letters*, **31(18)**.
- 718 Duran, P., and J. Molinari, 2018: Dramatic inner-core tropopause vari-  
719 ability during the rapid intensification of hurricane patricia (2015).  
720 *Monthly Weather Review*, **146(1)**, 119–134.
- 721 Fox, K. R., and F. Judt, 2018: A numerical study on the extreme intensifi-  
722 cation of hurricane patricia (2015). *Weather and Forecasting*, **33(4)**,  
723 989–999.
- 724 Frank, W. M., 1977: The structure and energetics of the tropical cyclone i.  
725 storm structure. *Monthly Weather Review*, **105(9)**, 1119–1135.
- 726 Fudeyasu, H., Y. Wang, M. Satoh, T. Nasuno, H. Miura, and W. Yanase,  
727 2008: Global cloud-system-resolving model NICAM successfully sim-  
728 ulated the lifecycles of two real tropical cyclones. *Geophysical Re-*  
729 *search Letters*, **35(22)**.
- 730 Gao, K., L. Harris, J.-H. Chen, S.-J. Lin, and A. Hazelton, 2019: Improv-  
731 ing AGCM hurricane structure with two-way nesting. *Journal of*  
732 *Advances in Modeling Earth Systems*, **11(1)**, 278–292.
- 733 Gentry, M. S., and G. M. Lackmann, 2010: Sensitivity of simulated tropical

734 cyclone structure and intensity to horizontal resolution. *Monthly*  
735 *Weather Review*, **138(3)**, 688–704.

736 Gopalakrishnan, S. G., S. Goldenberg, T. Quirino, X. Zhang, xf Marks,  
737 K.-S. Yeh, R. Atlas, and V. Tallapragada, 2012: Toward improving  
738 high-resolution numerical hurricane forecasting: Influence of model  
739 horizontal grid resolution, initialization, and physics. *Weather and*  
740 *Forecasting*, **27(3)**, 647–666.

741 Green, B. W., and F. Zhang, 2013: Impacts of air–sea flux parameterizations  
742 on the intensity and structure of tropical cyclones. *Monthly Weather*  
743 *Review*, **141(7)**, 2308–2324.

744 Guimond, S. R., G. M. Heymsfield, P. D. Reasor, and A. C. Didlake, 2016:  
745 The rapid intensification of hurricane karl (2010): New remote sens-  
746 ing observations of convective bursts from the global hawk platform.  
747 *Journal of the Atmospheric Sciences*, **73(9)**, 3617–3639.

748 Hamill, T. M., J. S. Whitaker, M. Fiorino, and S. G. Benjamin, 2011: Global  
749 ensemble predictions of 2009’s tropical cyclones initialized with an  
750 ensemble kalman filter. *Monthly Weather Review*, **139(2)**, 668–688.

751 Harper, B. A., J. D. Kepert, and J. D. Ginger, 2008: , Guidelines for convert-  
752 ing between various wind averaging periods in tropical cyclone con-

753           ditions, world meteorological organization. Technical report, World  
754           Meteorological Organization.

755 Heng, J., Y. Wang, and W. Zhou, 2017: Revisiting the balanced and un-  
756           balanced aspects of tropical cyclone intensification. *Journal of the*  
757           *Atmospheric Sciences*, **74(8)**, 2575–2591.

758 Hodges, K. I., and N. P. Klingaman, 2019: Prediction errors of tropical  
759           cyclones in the western north pacific in the met office global forecast  
760           model. *Weather and Forecasting*, **34(5)**, 1189–1209.

761 Hohenegger, C., L. Kornblueh, D. Klocke, T. Becker, G. Cioni, J. F. Engels,  
762           U. Schulzweida, and B. Stevens, 2020: Climate statistics in global  
763           simulations of the atmosphere, from 80 to 2.5 km grid spacing. *Jour-*  
764           *nal of the Meteorological Society of Japan. Ser. II*, **98(1)**, 73–91.

765 Judt, F., and S. S. Chen, 2010: Convectively generated potential vorticity  
766           in rainbands and formation of the secondary eyewall in hurricane rita  
767           of 2005. *Journal of the Atmospheric Sciences*, **67(11)**, 3581–3599.

768 Judt, F., and S. S. Chen, 2016: Predictability and dynamics of tropical  
769           cyclone rapid intensification deduced from high-resolution stochastic  
770           ensembles. *Monthly Weather Review*, **144(11)**, 4395–4420.

771 Kanada, S., A. Wada, M. Nakano, and T. Kato, 2012: Effect of plane-

772 tary boundary layer schemes on the development of intense tropical  
773 cyclones using a cloud-resolving model. *Journal of Geophysical Re-*  
774 *search: Atmospheres*, **117(D3)**, n/a–n/a.

775 Kepert, J. D., 2012: Choosing a boundary layer parameterization for tropi-  
776 cal cyclone modeling. *Monthly Weather Review*, **140(5)**, 1427–1445.

777 Kieu, C., V. Tallapragada, D.-L. Zhang, and Z. Moon, 2016: On the devel-  
778 opment of double warm-core structures in intense tropical cyclones.  
779 *Journal of the Atmospheric Sciences*, **73(11)**, 4487–4506.

780 Knapp, K. R., H. J. Diamond, J. P. Kossin, M. C. Kruk, and C. J. Schreck,  
781 2018: International best track archive for climate stewardship (ib-  
782 tracs) project, version 4.

783 Knapp, K. R., M. C. Kruk, D. H. Levinson, H. J. Diamond, and C. J.  
784 Neumann, 2010: The international best track archive for climate  
785 stewardship (IBTrACS). *Bulletin of the American Meteorological*  
786 *Society*, **91(3)**, 363–376.

787 Komaromi, W. A., and J. D. Doyle, 2017: Tropical cyclone outflow and  
788 warm core structure as revealed by HS3 dropsonde data. *Monthly*  
789 *Weather Review*, **145(4)**, 1339–1359.

790 Landsea, C. W., and J. L. Franklin, 2013: Atlantic hurricane database

791 uncertainty and presentation of a new database format. *Monthly*  
792 *Weather Review*, **141(10)**, 3576–3592.

793 Lee, C.-Y., and S. S. Chen, 2014: Stable boundary layer and its impact  
794 on tropical cyclone structure in a coupled atmosphere–ocean model.  
795 *Monthly Weather Review*, **142(5)**, 1927–1944.

796 Magnusson, L., J.-R. Bidlot, M. Bonavita, A. R. Brown, P. A. Browne, G. D.  
797 Chiara, M. Dahoui, S. T. K. Lang, T. McNally, K. S. Mogensen,  
798 F. Pappenberger, F. Prates, F. Rabier, D. S. Richardson, F. Vitart,  
799 and S. Malardel, 2019: ECMWF activities for improved hurricane  
800 forecasts. *Bulletin of the American Meteorological Society*, **100(3)**,  
801 445–458.

802 Manganello, J. V., K. I. Hodges, J. L. Kinter, B. A. Cash, L. Marx, T. Jung,  
803 D. Achuthavarier, J. M. Adams, E. L. Altshuler, B. Huang, E. K. Jin,  
804 C. Stan, P. Towers, and N. Wedi, 2012: Tropical cyclone climatol-  
805 ogy in a 10-km global atmospheric GCM: Toward weather-resolving  
806 climate modeling. *Journal of Climate*, **25(11)**, 3867–3893.

807 Marchok, T. P., 2002: How the ncep tropical cyclone tracker works. In *AMS*  
808 *Conference on Hurricanes and Tropical Meteorology, San Diego, CA,*  
809 *P1.13*, American Meteorological Society.

- 810 Miyamoto, Y., and T. Takemi, 2015: A triggering mechanism for rapid  
811 intensification of tropical cyclones. *Journal of the Atmospheric Sci-*  
812 *ences*, **72(7)**, 2666–2681.
- 813 Mogensen, K. S., L. Magnusson, and J.-R. Bidlot, 2017: Tropical cyclone  
814 sensitivity to ocean coupling in the ECMWF coupled model. *Journal*  
815 *of Geophysical Research: Oceans*, **122(5)**, 4392–4412.
- 816 Nakano, M., M. Sawada, T. Nasuno, and M. Satoh, 2015: Intraseasonal  
817 variability and tropical cyclogenesis in the western north pacific sim-  
818 ulated by a global nonhydrostatic atmospheric model. *Geophysical*  
819 *Research Letters*, **42(2)**, 565–571.
- 820 Nakano, M., A. Wada, M. Sawada, H. Yoshimura, R. Onishi, S. Kawahara,  
821 W. Sasaki, T. Nasuno, M. Yamaguchi, T. Iriguchi, M. Sugi, and  
822 Y. Takeuchi, 2017: Global 7-km mesh nonhydrostatic model inter-  
823 comparison project for improving TYphoon forecast (TYMIP-g7):  
824 experimental design and preliminary results. *Geoscientific Model*  
825 *Development*, **10(3)**, 1363–1381.
- 826 Powell, M. D., P. J. Vickery, and T. A. Reinhold, 2003: Reduced drag coef-  
827 ficient for high wind speeds in tropical cyclones. *Nature*, **422(6929)**,  
828 279–283.

829 Roberts, M. J., J. Camp, J. Seddon, P. L. Vidale, K. Hodges, B. Van-  
830 niere, J. Mecking, R. Haarsma, A. Bellucci, E. Scoccimarro, L.-P.  
831 Caron, F. Chauvin, L. Terray, S. Valcke, M.-P. Moine, D. Putrasa-  
832 han, C. Roberts, R. Senan, C. Zarzycki, and P. Ullrich, 2020: Impact  
833 of model resolution on tropical cyclone simulation using the High-  
834 ResMIP–PRIMAVERA multimodel ensemble. *Journal of Climate*,  
835 **33(7)**, 2557–2583.

836 Satoh, M., B. Stevens, F. Judt, M. Khairoutdinov, S.-J. Lin, W. M. Putman,  
837 and P. Düben, 2019: Global cloud-resolving models. *Current Climate*  
838 *Change Reports*, **5(3)**, 172–184.

839 Silverman, B., 2018: *Density Estimation for Statistics and Data Analysis*.  
840 Routledge.

841 Stern, D. P., and D. S. Nolan, 2009: Reexamining the vertical structure  
842 of tangential winds in tropical cyclones: Observations and theory.  
843 *Journal of the Atmospheric Sciences*, **66(12)**, 3579–3600.

844 Stern, D. P., and D. S. Nolan, 2012: On the height of the warm core in  
845 tropical cyclones. *Journal of the Atmospheric Sciences*, **69(5)**, 1657–  
846 1680.

847 Stevens, B., M. Satoh, L. Auger, J. Biercamp, C. S. Bretherton, X. Chen,  
848 P. Düben, F. Judt, M. Khairoutdinov, D. Klocke, C. Kodama,

849 L. Kornbluh, S.-J. Lin, P. Neumann, W. M. Putman, N. Röber,  
850 R. Shibuya, B. Vanniere, P. L. Vidale, N. Wedi, and L. Zhou, 2019:  
851 DYAMOND: the DYnamics of the atmospheric general circulation  
852 modeled on non-hydrostatic domains. *Progress in Earth and Plane-*  
853 *tary Science*, **6(1)**.

854 Torn, R. D., and C. Snyder, 2012: Uncertainty of tropical cyclone best-track  
855 information. *Weather and Forecasting*, **27(3)**, 715–729.

856 Vanniere, B., M. J. Roberts, P. L. Vidale, K. I. Hodges, M.-E. Demory,  
857 L.-P. Caron, E. Scoccimarro, L. Taurent, and R. Senan, 2020: The  
858 moisture budget of tropical cyclones in highresmp models : large-  
859 scale environmental balance and sensitivity to horizontal resolution.  
860 *submitted to Journal of Climate*.

861 Warner, T. T., R. A. Peterson, and R. E. Treadon, 1997: A tutorial on lat-  
862 eral boundary conditions as a basic and potentially serious limitation  
863 to regional numerical weather prediction. *Bulletin of the American*  
864 *Meteorological Society*, **78(11)**, 2599–2617.

865 Zeng, Z., Y. Wang, Y. Duan, L. Chen, and Z. Gao, 2010: On sea sur-  
866 face roughness parameterization and its effect on tropical cyclone  
867 structure and intensity. *Advances in Atmospheric Sciences*, **27(2)**,  
868 337–355.



- 869 Zhang, J. A., D. S. Nolan, R. F. Rogers, and V. Tallapragada, 2015: Eval-  
870 uating the impact of improvements in the boundary layer parame-  
871 terization on hurricane intensity and structure forecasts in HWRF.  
872 *Monthly Weather Review*, **143(8)**, 3136–3155.
- 873 Zhang, J. A., R. F. Rogers, D. S. Nolan, and F. D. Marks, 2011: On the  
874 characteristic height scales of the hurricane boundary layer. *Monthly*  
875 *Weather Review*, **139(8)**, 2523–2535.
- 876 Zhou, L., S.-J. Lin, J.-H. Chen, L. M. Harris, X. Chen, and S. L. Rees, 2019:  
877 Toward convective-scale prediction within the next generation global  
878 prediction system. *Bulletin of the American Meteorological Society*,  
879 **100(7)**, 1225–1243.

## List of Figures

881	1	TC tracks and numbers from observations (black/grey) and models (orange) for the DYAMOND period (1 Aug–10 Sep 2016). Numbers are given for each basin (Indian Ocean, Western Pacific, Eastern Pacific, Atlantic); the global total number of TCs is shown in the lower right. . . . .	45
882	2	Timeseries of TC formation events in the Western Pacific (a), Eastern Pacific (b), Atlantic (c), and Indian Ocean (d) from observations (black) and models (orange). . . . .	46
883	3	Timeseries of maximum surface wind speed ( $v_{max}$ ) for each TC from observations (black, grey) and models (orange). . . . .	47
884	4	Kernel density estimates of maximum wind speed from observations (black) and models (orange). . . . .	48
885	5	Kernel density estimates of minimum sea-level pressure from observations (black) and models (orange). . . . .	49
886	6	Accumulated cyclone energy (ACE) from observations (grey) and models (orange). The lower bound of the uncertainty range in observed ACE assumes that all $v_{vmax}$ observations have an error of $-5 \text{ m s}^{-1}$ or $+5 \text{ m s}^{-1}$ (upper bound). . . . .	50
887	7	Average storm size as measured by the median $17 \text{ m s}^{-1}$ wind radius for each storm quadrant from observations (black) and models (orange). Dashed grey circles indicate radius intervals of 100 km. The error bars in the observations are based on an error estimate of 50%. . . . .	51
888	8	TC wind-pressure relationships from observations (black) and models (orange). The curves are least-squares fitted quadratic functions. Note: the peculiar shape of the fit line in SAM (h) is not caused by the obvious outlier at $65 \text{ m s}^{-1}$ and 950 hPa. Excluding this outlier will not change the fit substantially. . . . .	52
889	9	Snapshots of 10-m wind speed of the strongest storm from each model at the time of peak intensity. . . . .	53
890	10	Radius-height composites of azimuthally-averaged tangential wind speed (grey shading) and radial/vertical flow (colored streamlines) from each model. The $20 \text{ m s}^{-1}$ -contour is annotated. The composites include all snapshots where a storm's $v_{max} \geq 33 \text{ m s}^{-1}$ . . . . .	54
891			
892			
893			
894			
895			
896			
897			
898			
899			
900			
901			
902			
903			
904			
905			
906			
907			
908			
909			
910			
911			
912			
913			
914			
915			

916	11	Radius-height composites of azimuthally-averaged radial wind	
917		speed in the lowest 2 km from each model. The dashed black	
918		line depicts the inflow layer height, here defined as the layer	
919		with radial wind $< -1 \text{ m s}^{-1}$ . The composites include all	
920		snapshots where a storm's $v_{max} \geq 33 \text{ m s}^{-1}$ . . . . .	55
921	12	Radius-height composites of the TC <i>warm core</i> from each	
922		model, computed as the azimuthally-averaged temperature	
923		anomaly with respect to the mean temperature between $r =$	
924		$300\text{--}700 \text{ km}$ . The composites include all snapshots where a	
925		storm's $v_{max} > 32 \text{ m s}^{-1}$ . . . . .	56
926	13	TC tracks and numbers from various ICON runs (orange)	
927		and observations (grey). Left: ICON runs <b>without</b> deep	
928		convective parameterization, right: ICON runs <b>with</b> deep	
929		convective parameterization. The model resolution, given in	
930		each panel, increases from top to bottom. . . . .	57
931	14	ACE from observations (grey) and various ICON runs <b>with</b>	
932		(light orange) or <b>without</b> (dark orange) deep convective pa-	
933		rameterization. Model resolution increases from top to bot-	
934		tom. The error bars are the lower an upper bounds assuming	
935		that all $v_{vmax}$ observations have an error of $\pm 5 \text{ m s}^{-1}$ . . . .	58
936	15	Momentum flux (top) and latent heat flux (bottom) from	
937		ARPEGE and ICON as a function of wind speed. The data	
938		are from the same time and domain as the snapshots in Fig.	
939		9. Instead of a raw scatter plot, the data are binned and the	
940		color saturation is a measure of points per bin. . . . .	59

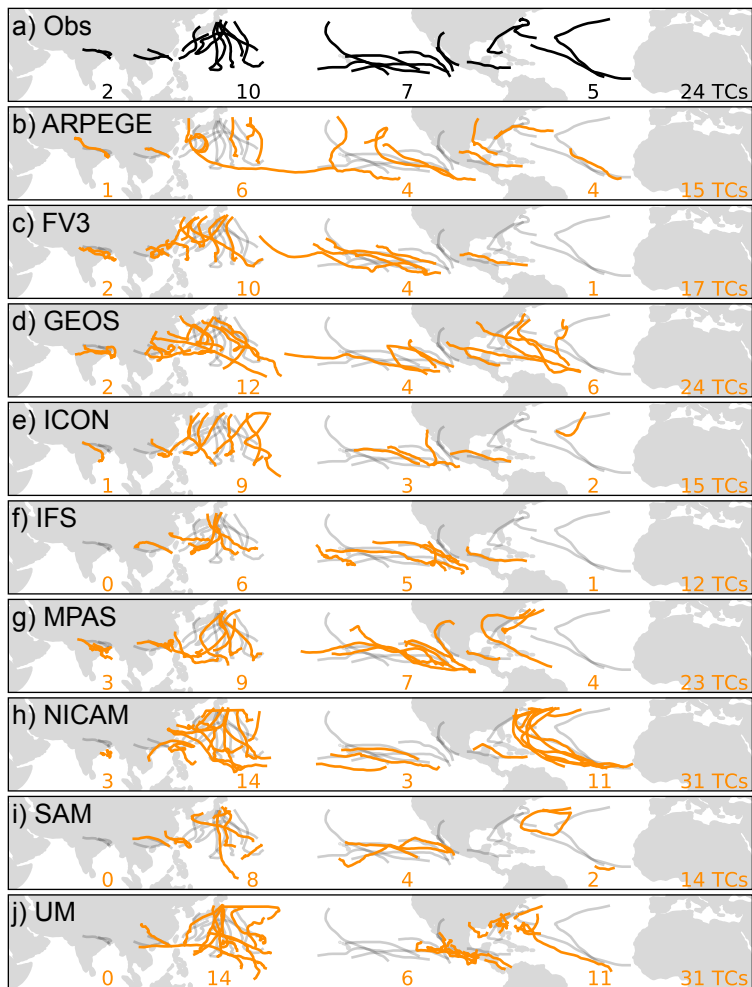


Fig. 1. TC tracks and numbers from observations (black/grey) and models (orange) for the DYAMOND period (1 Aug–10 Sep 2016). Numbers are given for each basin (Indian Ocean, Western Pacific, Eastern Pacific, Atlantic); the global total number of TCs is shown in the lower right.

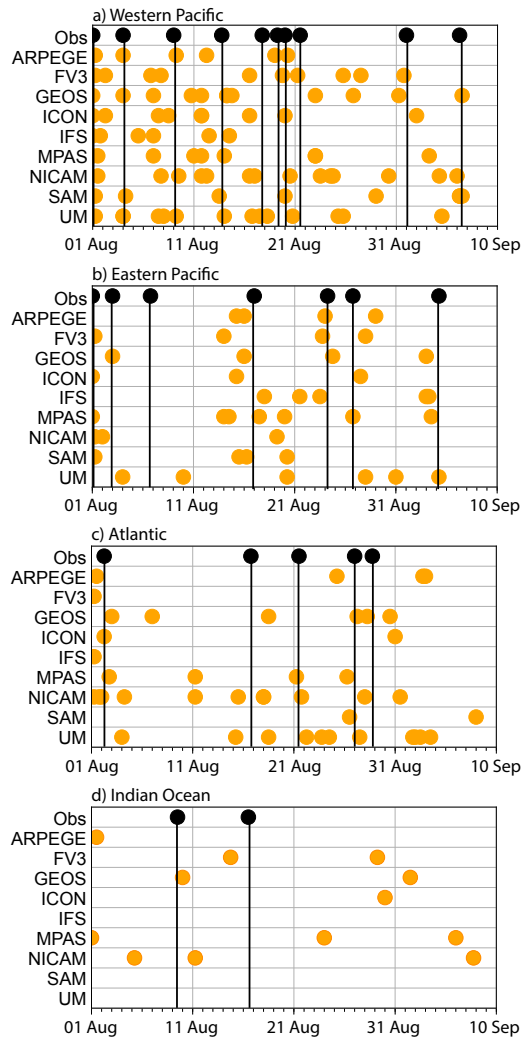


Fig. 2. Timeseries of TC formation events in the Western Pacific (a), Eastern Pacific (b), Atlantic (c), and Indian Ocean (d) from observations (black) and models (orange).

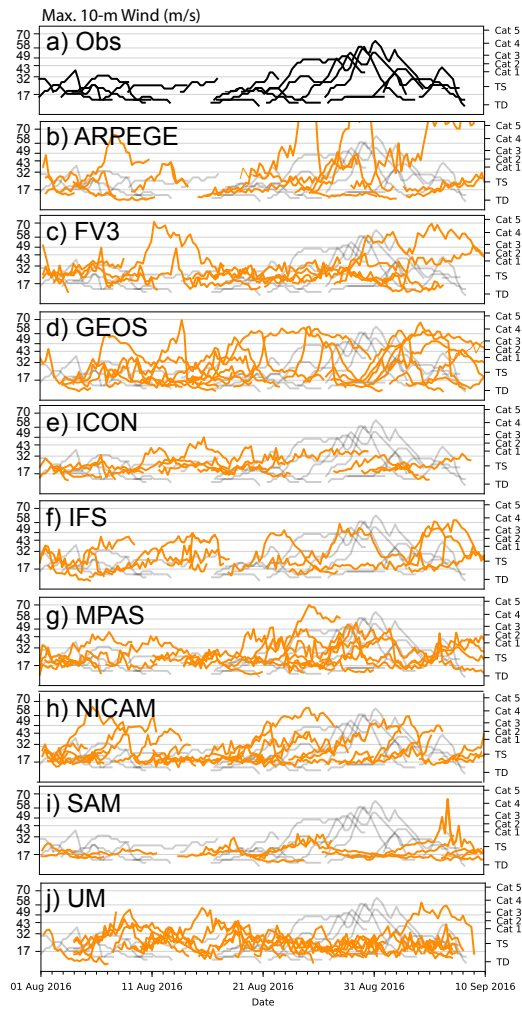


Fig. 3. Timeseries of maximum surface wind speed ( $v_{max}$ ) for each TC from observations (black, grey) and models (orange).

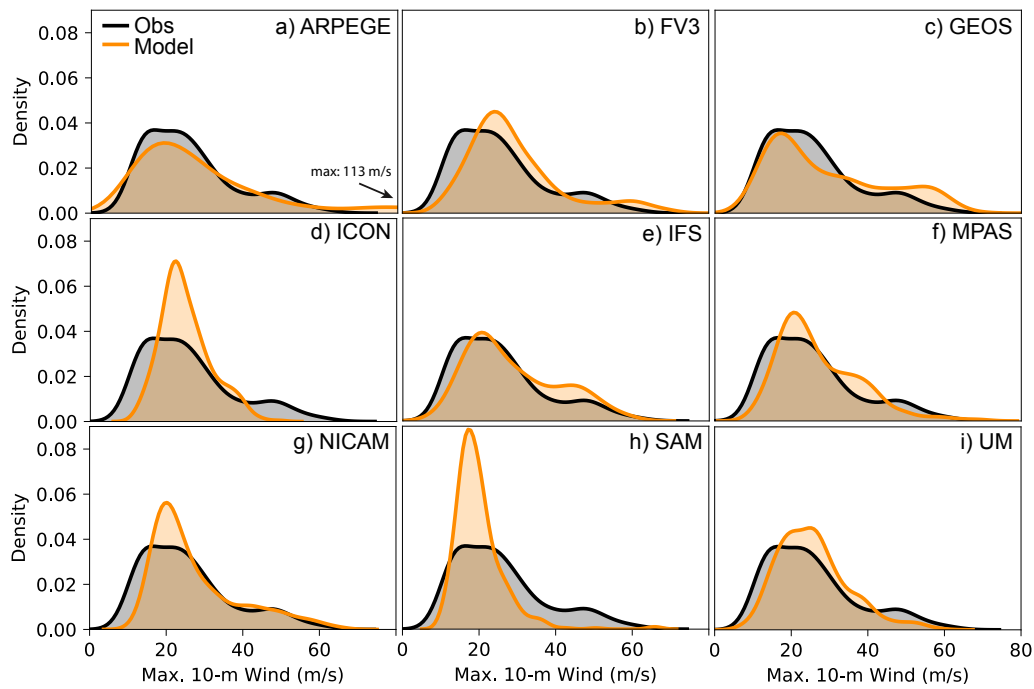


Fig. 4. Kernel density estimates of maximum wind speed from observations (black) and models (orange).

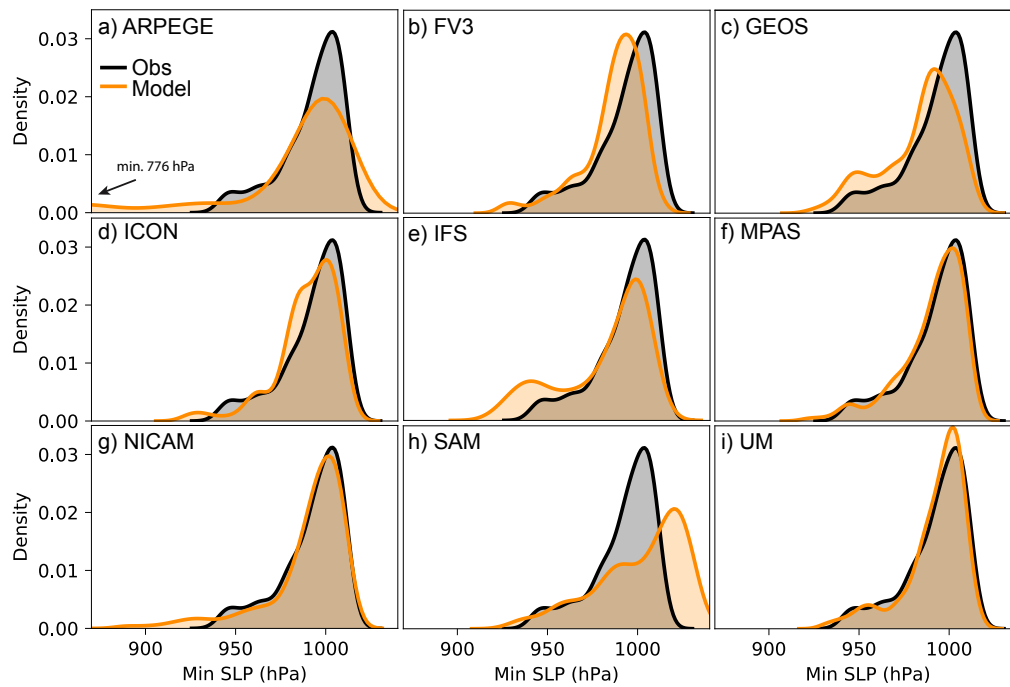


Fig. 5. Kernel density estimates of minimum sea-level pressure from observations (black) and models (orange).



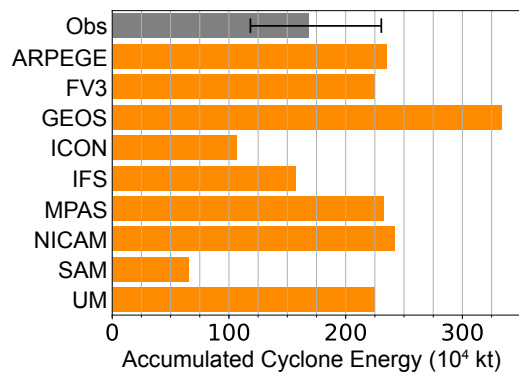


Fig. 6. Accumulated cyclone energy (ACE) from observations (grey) and models (orange). The lower bound of the uncertainty range in observed ACE assumes that all  $v_{vmax}$  observations have an error of  $-5 \text{ m s}^{-1}$  or  $+5 \text{ m s}^{-1}$  (upper bound).

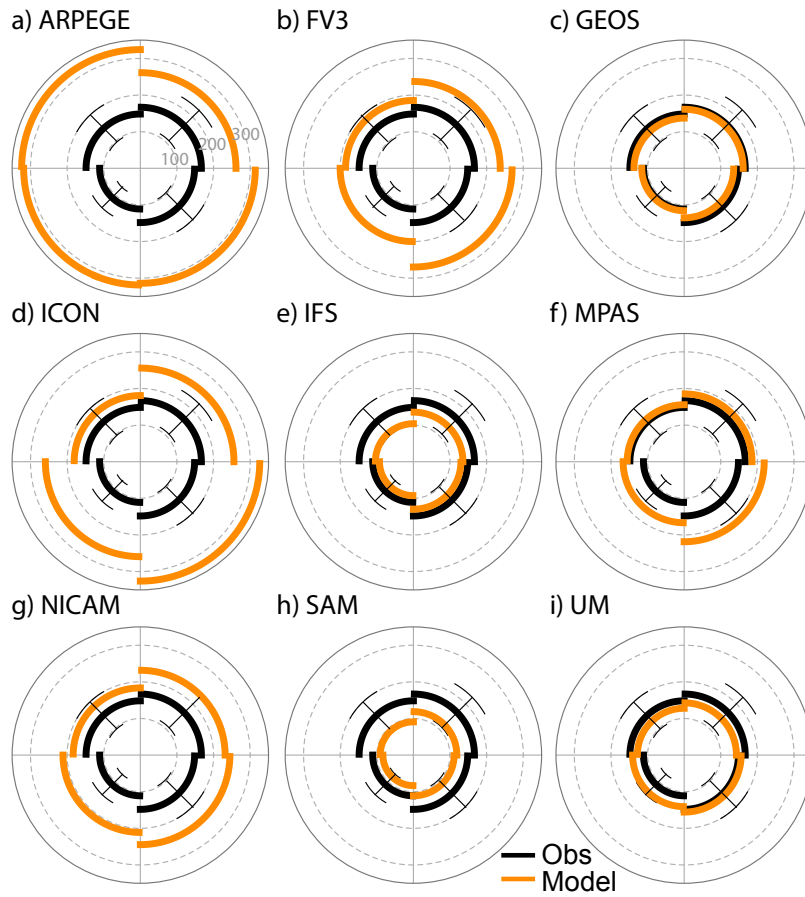


Fig. 7. Average storm size as measured by the median  $17 \text{ m s}^{-1}$  wind radius for each storm quadrant from observations (black) and models (orange). Dashed grey circles indicate radius intervals of 100 km. The error bars in the observations are based on an error estimate of 50%.

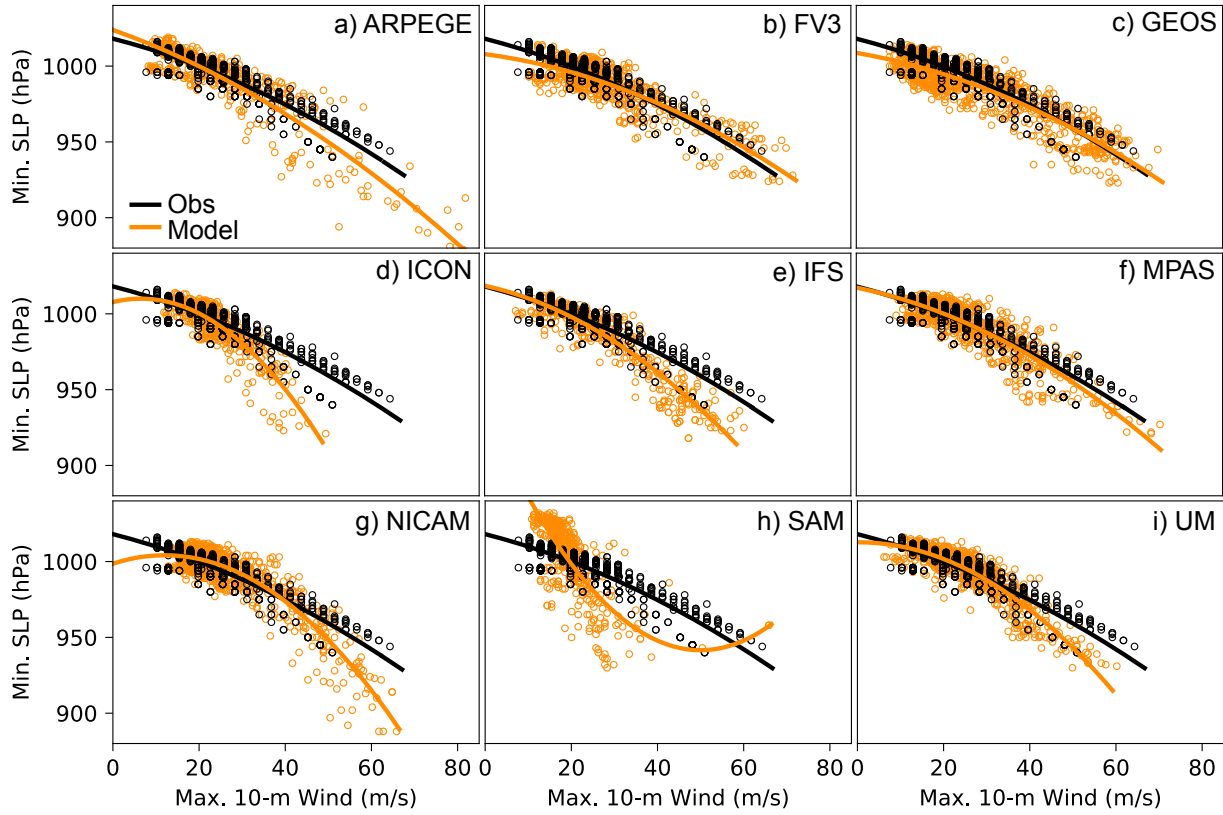
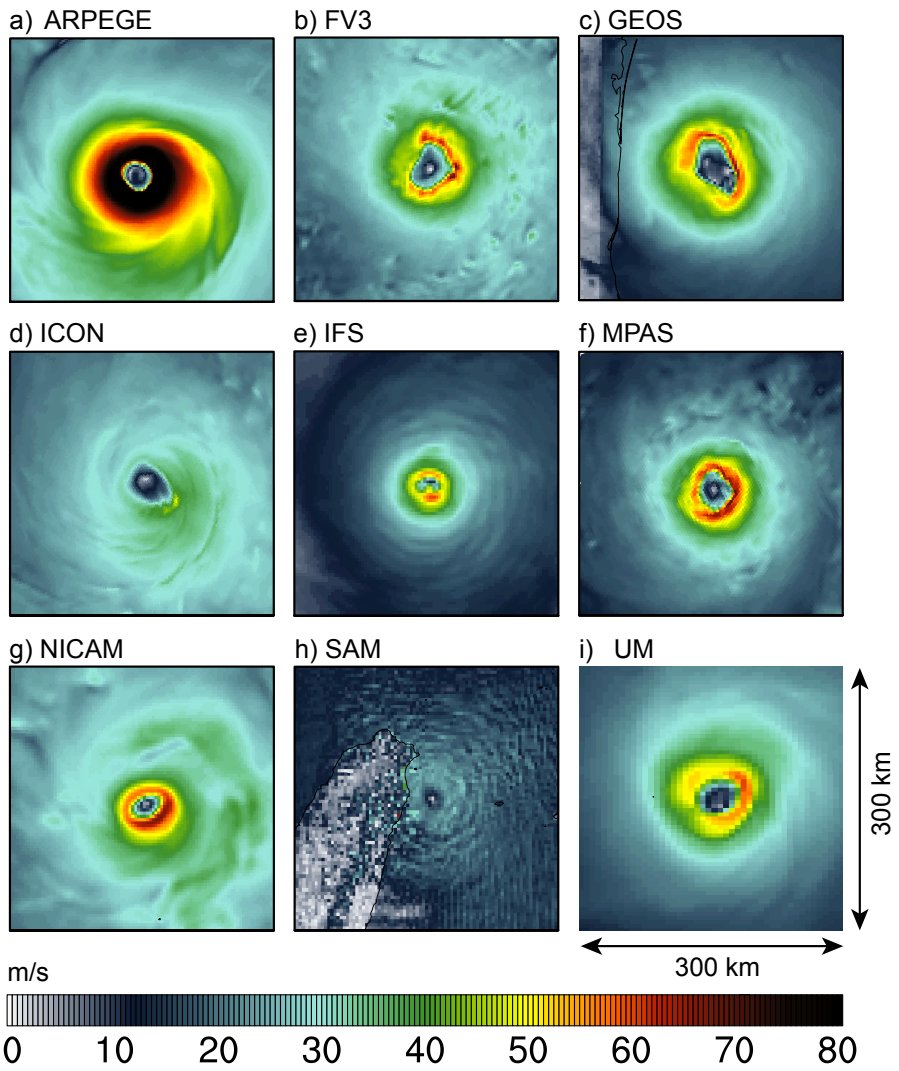


Fig. 8. TC wind-pressure relationships from observations (black) and models (orange). The curves are least-squares fitted quadratic functions. Note: the peculiar shape of the fit line in SAM (h) is not caused by the obvious outlier at  $65 \text{ m s}^{-1}$  and  $950 \text{ hPa}$ . Excluding this outlier will not change the fit substantially.



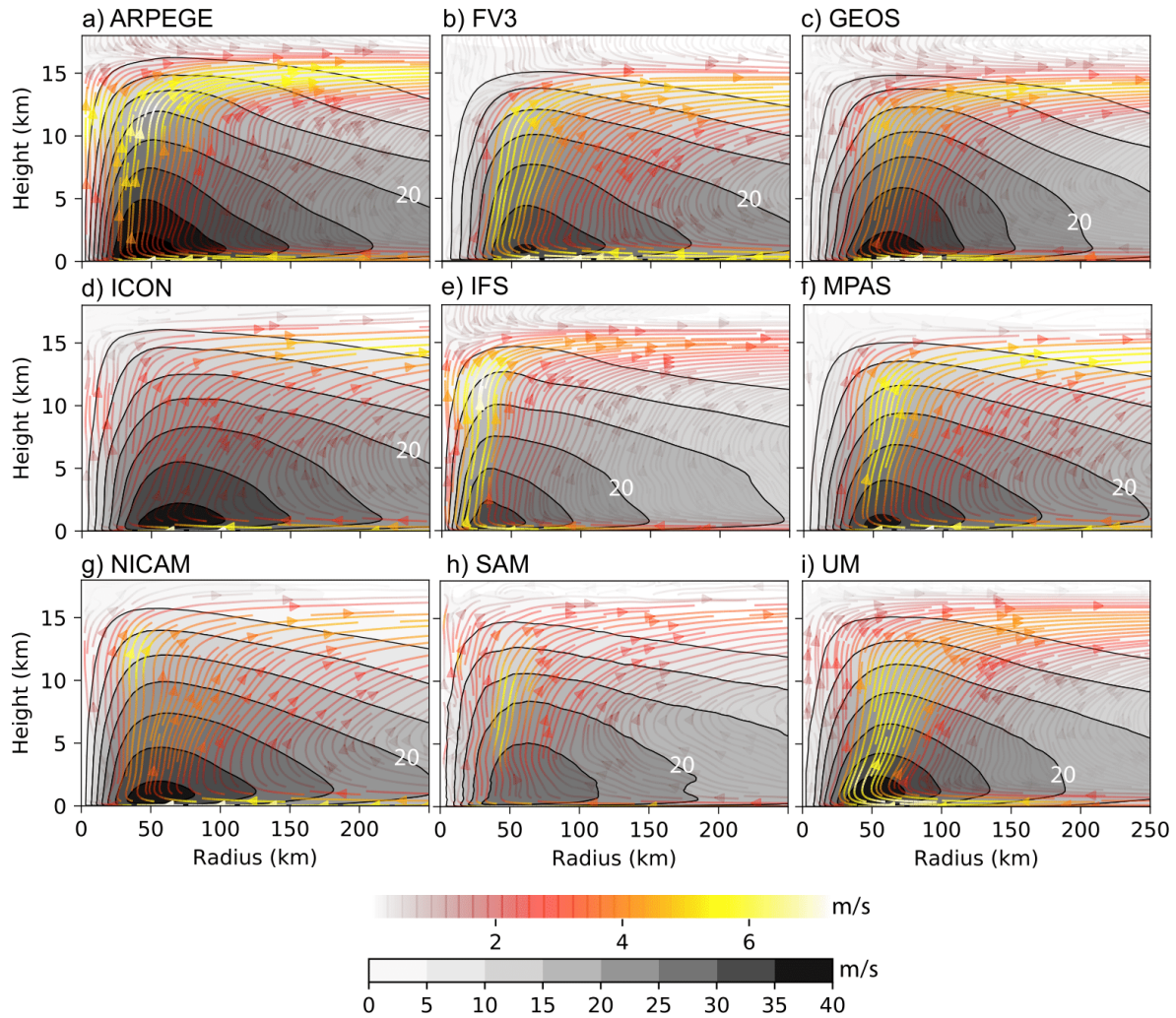


Fig. 10. Radius-height composites of azimuthally-averaged tangential wind speed (grey shading) and radial/vertical flow (colored streamlines) from each model. The  $20 \text{ m s}^{-1}$ -contour is annotated. The composites include all snapshots where a storm's  $v_{max} \geq 33 \text{ m s}^{-1}$ .

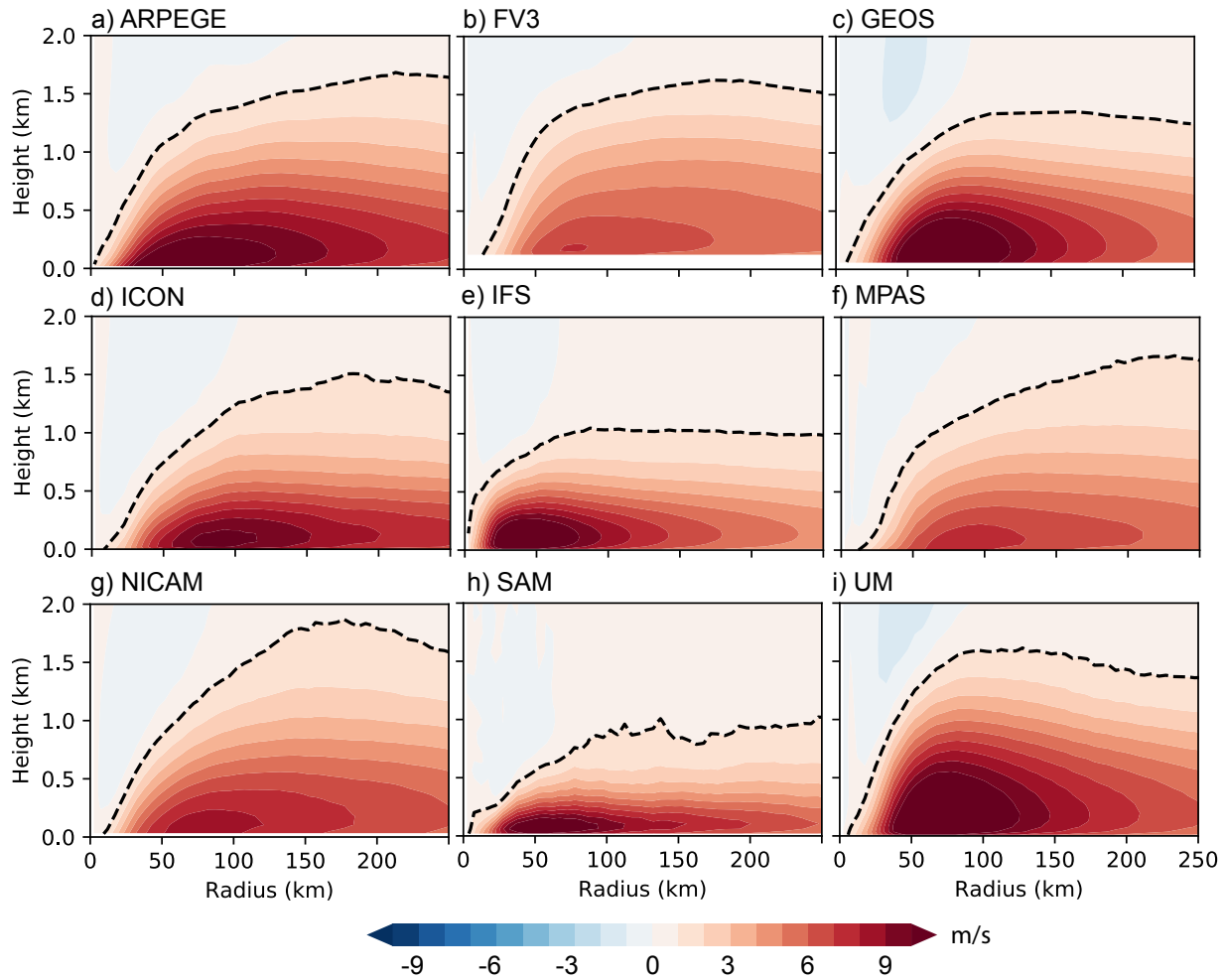


Fig. 11. Radius-height composites of azimuthally-averaged radial wind speed in the lowest 2 km from each model. The dashed black line depicts the inflow layer height, here defined as the layer with radial wind  $< -1 \text{ m s}^{-1}$ . The composites include all snapshots where a storm's  $v_{max} \geq 33 \text{ m s}^{-1}$ .

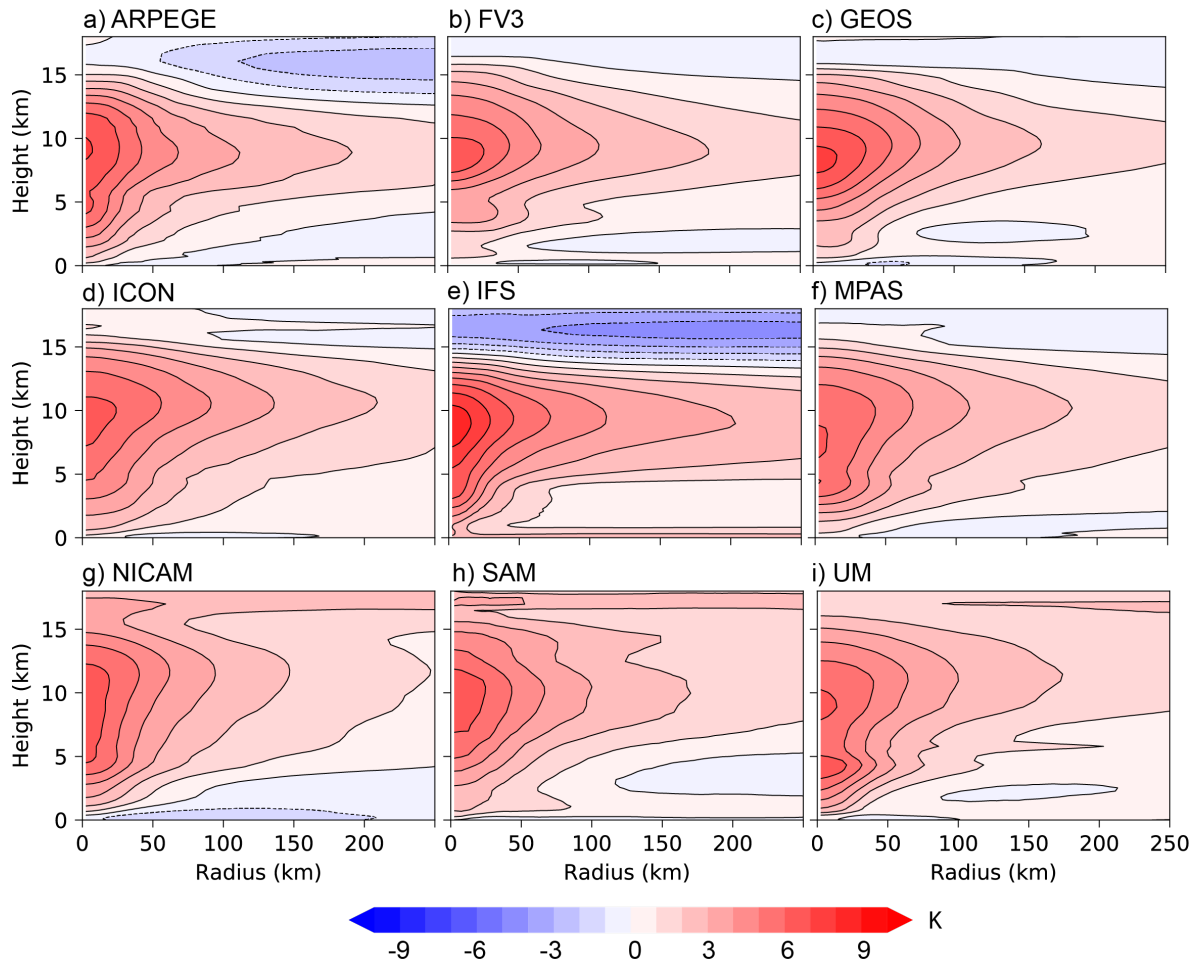


Fig. 12. Radius-height composites of the TC *warm core* from each model, computed as the azimuthally-averaged temperature anomaly with respect to the mean temperature between  $r = 300\text{--}700$  km. The composites include all snapshots where a storm's  $v_{max} > 32 \text{ m s}^{-1}$ .

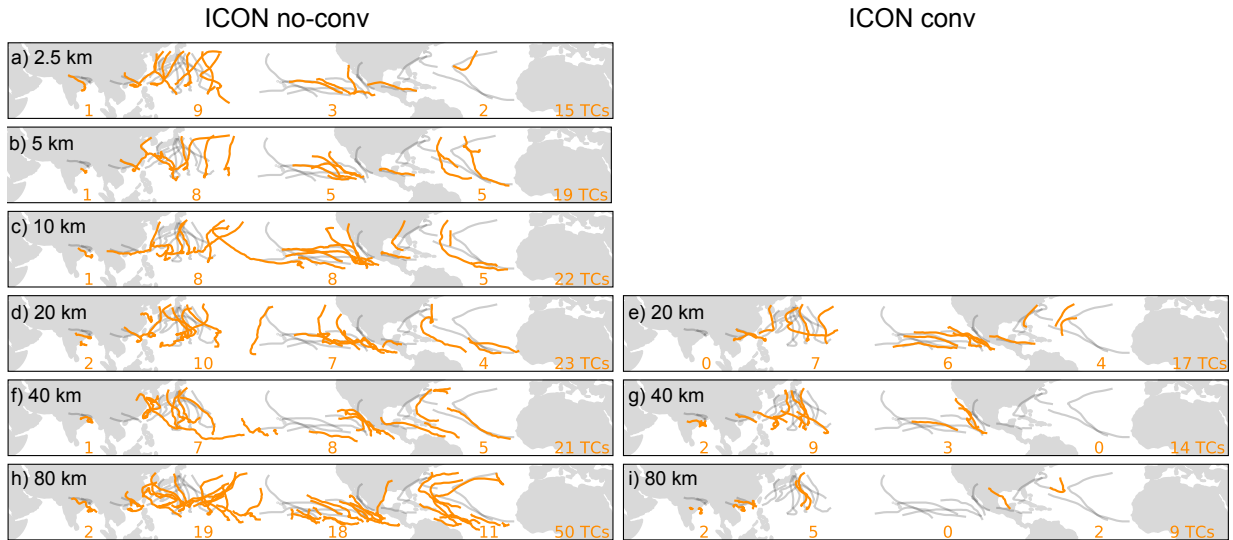


Fig. 13. TC tracks and numbers from various ICON runs (orange) and observations (grey). Left: ICON runs **without** deep convective parameterization, right: ICON runs **with** deep convective parameterization. The model resolution, given in each panel, increases from top to bottom.



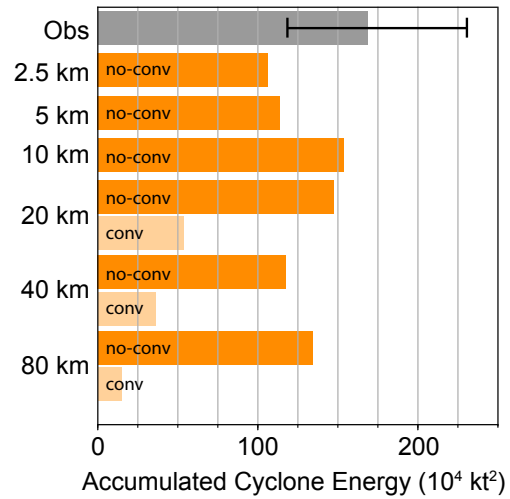


Fig. 14. ACE from observations (grey) and various ICON runs **with** (light orange) or **without** (dark orange) deep convective parameterization. Model resolution increases from top to bottom. The error bars are the lower and upper bounds assuming that all  $v_{max}$  observations have an error of  $\pm 5 \text{ m s}^{-1}$ .

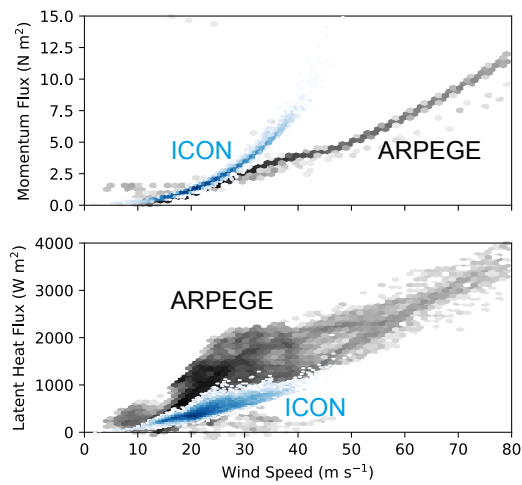


Fig. 15. Momentum flux (top) and latent heat flux (bottom) from ARPEGE and ICON as a function of wind speed. The data are from the same time and domain as the snapshots in Fig. 9. Instead of a raw scatter plot, the data are binned and the color saturation is a measure of points per bin.