

1 **Community-wide validation of geospace model local**
2 **K-index predictions to support model transition to**
3 **operations**

A. Glocer¹, L. Rastätter¹, M. Kuznetsova¹, A. Pulkkinen¹, H. J. Singer⁶, C.

Balch⁶, D. Weimer², D. Welling⁴, M. Wiltberger³, J. Raeder⁵ and R. S.

Weigel⁷, J. McCollough⁸, S. Wing⁹

¹NASA Goddard Space Flight Center,

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as doi:

10.1029/2016SW001387

June 4, 2016, 5:15am

D R A F T

4 **Abstract.**

5 We present the latest result of a community-wide space weather model val-
6 idation effort coordinated among the Community Coordinated Modeling Cen-
7 ter (CCMC), NOAA Space Weather Prediction Center (SWPC), model de-
8 velopers, and the broader science community. Validation of geospace mod-
9 els is a critical activity for both building confidence in the science results pro-
10 duced by the models and in assessing the suitability of the models for tran-
11 sition to operations. Indeed, a primary motivation of this work is support-
12 ing NOAA SWPC's effort to select a model or models to be transitioned into
13 operations. Our validation efforts focus on the ability of the models to re-
14 produce a regional index of geomagnetic disturbance, the local K-index. Our
15 analysis includes six events representing a range of geomagnetic activity con-
16 ditions and six geomagnetic observatories representing mid- and high-latitude
17 locations. Contingency tables, skill scores, and distribution metrics are used
18 for the quantitative analysis of model performance. We consider model per-
19 formance on an event-by-event basis, aggregated over events, at specific sta-
20 tion locations, and separated into high- and mid-latitude domains. A sum-
21 mary of results is presented in this report, and an online tool for detailed anal-
22 ysis is available at the CCMC.

Greenbelt, MD 20771, USA.

1. Introduction

23 Forecasting geomagnetic disturbance levels on the ground is a critical step in mitigating
24 the potentially severe impact of geomagnetically induced currents (GICs) [e.g., *Boteler et*
25 *al.*, 1998; *Pirjola*, 2005; *North American Electric Reliability Corporation*, 2012; *National*
26 *Research Council*, 2008]. The science community has responded with both first principles
27 and empirical models capable of forecasting these potentially hazardous disturbances.
28 Before such models can be transitioned into an operational setting, a comprehensive model
29 validation effort is required to determine the model quality and capabilities for improving
30 services. The Community Coordinated Modeling Center (CCMC), NOAA Space Weather
31 Prediction Center (SWPC), model developers, and the broader science community have
32 joined together to carry out this important validation effort. This report represents the
33 latest model validation findings in support of geospace model transition to operations.

34 This study builds on the prior studies of geospace model validation [*Pulkkinen et al.*,
35 2010, 2013; *Rastätter et al.*, 2011], and in particular is a direct follow on to [*Pulkkinen et*
36 *al.*, 2013]. That study focused on the ability of models to reproduce dB/dt (the variation
37 of ground magnetic field) at specific magnetometer locations. We encourage the reader
38 to refer to that work, as this study is a direct follow on to that effort. As the work of
39 [*Pulkkinen et al.*, 2013] was coming to completion, work was initiated on the present study,
40 to consider the ability of models to reproduce a local index of geomagnetic disturbance.
41 While the magnetic field fluctuations on short times, examined in the prior study, is more

²Center for Space Science and

42 directly tied to GIC prediction, a local index of variability is also useful as a convenient
43 measure of the local risk of GIC. Moreover, it is possible that a model would have more
44 skill in predicting the scaled range of magnetic field variability over a wider window than
45 over a relatively short-term variation.

46 The Kp index is a commonly used global measure of geomagnetic disturbances. It is a
47 measure on a scale of 0-9 of the average level of disturbance as measured by a scaled range
48 of delta-B at selected geomagnetic observatories For a detailed description of how Kp is
49 calculated see *Rostoker* [1972]. Local predictions of K , however, may differ significantly
50 from the global Kp -index. The interest in predicting potential GICs and geomagnetic
51 disturbances on a regional or local level, and the convenience of an activity index instead
52 of a raw prediction, provides part of the motivation for this study. Additionally, we will
53 be able to determine if the local value of the model derived K better represents the level
54 of activity at a particular location than the global Kp index.

55 The layout of the paper is as follows. Section 2 describes the organization of the
56 validation effort, Section 3 presents the metrics used to measure the model performance
57 and Section 4 details the models. Validation results are described in Section 5, and Section
58 6 discusses the findings.

2. Validation setting

59 As noted in the previous section, the present work builds on the validation study pre-
60 sented by *Pulkkinen et al.* [2013]. To avoid repeating the very complete description of the

Engineering Research, Virginia Polytechnic

61 validation setting provided previously, we will only provide an overview here as well as
62 new features particular to the current study.

63 Six events were chosen for the study consisting of the four events from the earlier GEM
64 Challenges [*Pulkkinen et al.*, 2010, 2013; *Rastätter et al.*, 2011] as well as two “surprise
65 events” chosen after the modelers delivered their models to CCMC for evaluation. CCMC
66 and NOAA SWPC scientists together choose the these two surprise events. The event list
67 is given in Table 1.

68 Three high-latitude (PBQ/SNK, ABK and YKC) and mid-latitude (WNG, NEW,
69 OTT) locations were selected. Table 2 and Fig. 1 show the locations of these stations. In
70 the case of the global MHD models, the magnetic field variations at each magnetometer
71 location were computed by a Biot-Savart integral over the entire domain. The integration
72 includes all currents in the magnetosphere, as well as the field-aligned currents in the
73 gap region between the MHD model’s inner boundary and the ionosphere, and the high-
74 latitude ionospheric currents. The CCMC tool used for the integration is described in
75 detail by *Rastätter et al.* [2013] and is applied to each of the Global MHD models used in
76 the study. The two empirical models (see Table 4) directly give the magnetic field at the
77 coordinates of the station. All model runs and ground magnetic field calculations (with
78 the exception of WingKp) were carried out at CCMC.

79 For every event under consideration (see Table 1), we evaluate the performance of the
80 model by comparing the observed vs predicted local K -values at the specific magnetometer
81 locations listed above. Throughout the paper K is calculated in the following way. First
82 we find the maximum “Range” of ΔB in the two horizontal directions.

Institute and State University, Blacksburg,

$$\text{Range} = \max [(\Delta B_{x,max} - \Delta B_{x,min}), (\Delta B_{y,max} - \Delta B_{y,min})] \quad (1)$$

Virginia, USA.

³High Altitude Observatory, National Center for Atmospheric Research, Boulder, Colorado, USA.

⁴Department of Atmospheric, Oceanic, and Space Sciences, University of Michigan, USA.

⁵Space Science Center & Physics Department, University of New Hampshire, USA.

⁶Space Weather Prediction Center, NOAA, Boulder, Colorado, USA.

⁷Department of Computational and Data Sciences, George Mason University, Fairfax, Virginia, USA.

⁸Air Force Research Laboratory, Kirtland AFB, NM, USA.

⁹Johns Hopkins University Applied Physics Laboratory, Laurel, Maryland, USA

83 over a three-hour window sliding by 15 minutes, where $B_{x,max}$, $B_{x,min}$, $B_{y,max}$, and $B_{y,min}$
84 indicate the max and min values in the window of the two horizontal components of the
85 magnetic field (North and East in geomagnetic dipole coordinates). Strictly speaking, the
86 quiet day variation should be subtracted before the range is calculated. However, neglect-
87 ing this only introduces a relatively small error when geomagnetic activity is disturbed
88 The Range is then divided by a station specific scaling factor. Scaling factors for stations
89 used in this validation study are specified by IAGA through International Service of Ge-
90 omagnetic Indices (ISGI) and is, generally speaking, a function of geomagnetic latitude.
91 Those values are given in Table 2. K is then found from the scaled range using a lookup
92 table given in Table 3. The same approach was used for both models and observations.
93 As stated before, we follow the earlier GEM Challenges and the earlier validation study
94 using the magnetometer stations listed in Table 2 and shown in Fig. 1. Three high- as
95 well as three mid-latitude stations (the same as for [Pulkkinen et al., 2013]) were included
96 in the present study (Table 2). Station PBQ was no longer available in late 2007 and
97 was therefore SNK was used. We therefore use station SNK for the 5th and 6th events.
98 We use the results from the model and observations from [Pulkkinen et al., 2013] for the
99 time series used to calculate K in this study. No new models runs or data processing was
100 carried out to get the time series from which we calculate the local K value. An exception
101 to this is a rerun of the 5_WEIMER empirical model to account for errors in how that
102 model was run in the previous study. The new results from that model (referred to as
103 6_WEIMER here and in the online plotting tool) are used in this analysis. 6_WEIMER
104 has the outputs correctly rotated to geomagnetic dipole coordinates whereas 5_WEIMER
105 does not. In addition, the CCMC had run the 5_WEIMER model with the Y component

106 of the IMF always set to zero, due to a program error in the CCMC run scripts. The
107 model developer found the problem which was subsequently fixed by CCMC for the rerun
108 named 6_WEIMER. The previous dB/dt study has not yet been corrected.

3. Metrics

109 The model validation is largely built on event-based analyses, as described in *Pulkkinen*
110 *et al.* [2013], and a distribution metric that provides new insight into model performance.
111 The event-based analysis determines where K exceeds a threshold of k_{thres} in a three-hour
112 sliding window. We then generate a contingency table that presents the number of correct
113 hits, false alarms, missed events and correct no events [e.g., *Lopez et al.*, 2007]. In this work
114 the thresholds for K were chosen to roughly correspond to the moderate ($K = 6$), and
115 severe ($K = 8$) geomagnetic storm levels as defined by the NOAA Space Weather Scales
116 (see e.g. <http://www.swpc.noaa.gov/noaa-scales-explanation>). The selected thresholds
117 are chosen with the idea that higher K values representing stronger events are of more
118 interest for space weather applications.

119 The contingency tables presented in the results section contain four entries per model
120 evaluated: The number of times the threshold crossing was accurately predicted H (hits),
121 the number of false predictions where a threshold crossing was predicted but not observed
122 F , the number of observed threshold crossings missed by a model M and the number of
123 times the model correctly predicted that no crossing occurred N . These entries are used
124 to compute the metrics used to quantify model performance. NOAA SWPC proposed
125 three metrics for use in the final analyses: Probability of Detection (POD), Probability of
126 False Detection (POFD) and Heidke Skill Score (HSS). For interest, we also include the
127 Critical Success Index (CSI) as an additional skill score; however, it is not used for model

128 ranking. For HSS, a 1 indicates a perfect score, a 0 demonstrates no skill as compared
129 to random chance, and negative values mean that random chance has more skill than
130 the model prediction. For POD, a 1 indicates a perfect score, while a 0 indicates that a
131 model never makes a correct detection. For POFD, a 0 indicates a perfect score, while
132 a 1 indicates that a model always makes false detections. For detailed descriptions of
133 these metrics, we refer the interested reader back to the previous study by *Pulkkinen et*
134 *al.* [2013].

135 In addition to the event tables and skill scores, we also consider a newly defined distri-
136 bution metric. In this metric, we consider the distribution of model predictions when the
137 observations are a particular value of $k = k_0$. A model that performs well in this metric
138 would show a distribution peaked around k_0 with very little spread in the distribution. A
139 model with significant random error would exhibit broadening of the distribution around
140 k_0 . A model with systematic error would have the distribution shifted so the peak is above
141 or below k_0 . A model with both systematic and random errors would exhibit both a shift
142 and broadening of the distribution around k_0 . In this study, we consider the distribution
143 metric for three values of $k = 4, 6, 8$, and qualitatively compare the results to examine for
144 the relative presence of random and systematic error in model predictions. This compar-
145 ison could potentially be made more rigorous in future studies by using autocorrelation
146 peaks.

4. Models

147 We include the same five models used in [*Pulkkinen et al.*, 2013]. These included em-
148 pirical models by *Weimer* [2013] and *Weigel et al.* [2003] and major US global magneto-
149 hydrodynamic (MHD) models from University of Michigan [*Tóth et al.*, 2012], the Center

150 for Integrated Space Weather Modeling (CISM) [Willberger *et al.*, 2004], and University
151 of New Hampshire [Raeder *et al.*, 2008]. In addition to these models, we also include the
152 WingKp model of Global Kp prediction [Wing *et al.*, 2005]. This last model was added
153 in order to determine the “value added” of models that can predict regional K values,
154 compared with a model currently used to predict a single global magnetic disturbance
155 level that is assumed to apply everywhere.

156 As with the prior evaluation study, each model that participated in the current study was
157 provided to CCMC. Communications with the model developers was essential to assure
158 each model was installed correctly with correct settings and used appropriately. The
159 WingKp model was treated differently because it is already operational at NOAA/SWPC
160 and hence, the model was evaluated by the NOAA/SWPC staff with minimal involvement
161 of its developer. We used the same model settings as in the previous study with final
162 settings determined in August 2011. No model could participate if it could not run
163 at least twice real-time on a 64 processor super computer. In otherwords, one hour
164 of simulated time could be completed in a half hour of wall time. This is critical to
165 ensuring models evaluated could operate in a realistic operational environment. Detailed
166 model descriptions and milestones of model deliveries and run executions are presented in
167 [Pulkkinen *et al.*, 2013]. All simulations, except for WingKp, were performed at CCMC
168 using identical computational resources and were driven by ACE level 2 data for Events
169 2-6. As reported by [Skoug *et al.*, 2004], only low resolution data could be constructed
170 for event 1. Additionally, the plasma density data for the event were derived from the
171 Plasma Wave Instrument on board the Geotail Satellite.

172 The WingKp model was run at AFRL since it was not one of the models in the CCMC
173 inventory. Details of this output can be found in the report by [McCollough *et al.*, 2014].
174 Additionally, AFRL was not able to provide results for event 3 which was outside their
175 run window. While the other models were all driven by identical ACE level 2 data, the
176 WingKp model was run with the real-time ACE data, and occasionally was not able to
177 supply a prediction due to missing data. Such predictions show up as a no data flag
178 ($K=-1$) in the online plotting and are excluded from our metrics analysis. The different
179 input data should be kept in mind when comparing model performance. WingKp was
180 handled differently than the other models because, when available, its purpose was to
181 compare the local prediction of K by the models under evaluation with a Kp prediction
182 that is currently available to SWPC forecasters.

183 Table 4 presents some of the features of each model. Some of these models, such as the
184 Weimer model and each of the global MHD models can be accessed through the CCMC
185 for runs-on-request.

5. Results

186 All of the time series of local K values are posted online and visualizations can be made
187 through the CCMC (http://ccmc.gsfc.nasa.gov/challenges/dBdt/metrics_results.php).
188 Figure 2 shows an example time series of the observed vs modeled K for the event 2
189 (Table 1). Each model is shown in a separate panel (red line) together with the observa-
190 tions (black line). We chose a random mid-latitude station for this demonstration.

191 Event-based metrics are broken out in several different ways. First, all the events and
192 stations are combined, as presented in Figure 3 and tables 5 and 6, to obtain an overall
193 view of model performance. The models are ordered from left to right by the HSS,

194 although all the event-based skill scores, previously discussed, are presented. It is also
195 of interest to examine the performance for different latitudes. Therefore we report the
196 results summed over all events and high-latitude (PBQ/SNK, ABK, YKC) stations and
197 mid-latitude (WNG, NEW, OTT) stations. Figures 4 and 5 show the performance for
198 high-latitude stations and mid-latitude stations respectively. Other configurations were
199 also considered such as grouping the results by the first four events that were known to
200 the modelers ahead of the study, and the two events added later. However, in the interest
201 of brevity the associated tables are not included here. We note that caution must be
202 taken when determining groupings or setting thresholds to ensure that there are enough
203 threshold crossing events. To that end we do not focus on individual magnetometer but
204 rather the groupings specified above. The smallest number of threshold crossings in any
205 grouping considered is 171 out of 1422 total events for midlatitude magnetometers with
206 a threshold of 8.

207 As described in Section 3, we also incorporate a “distribution” metric. The concept
208 behind this metric is as follows: We examine the distribution of model predictions at a
209 particular station for an observed K at that same station. Although we do not employ
210 a mathematically rigorous analysis of the model performance in the distribution metric,
211 a great deal can still be learned by visual inspection of the distributions. For instance,
212 a peak shifted to the left represents a systematic under-prediction while a peak shifted
213 to the right represents a systematic over prediction. When taken in conjunction with the
214 contingency tables and skill scores the results can be quite illuminating. A model that
215 has a high-probability of false detection, for instance, could have those false detections
216 as a result of a systematic error causing the model to consistently predict higher values,

217 random errors causing the model to result in more false detections, or a combination of
218 both. The contingency tables alone cannot pinpoint the type of error, but including the
219 distribution metric can provide insight into the cause for, in this case, the false detection.

220 When evaluating results from using the distribution metric, we consider the results
221 station-by-station to gain a more granular picture of model performance. One important
222 factor to keep in mind is that the number of events decreases for $K = 8$ and may be
223 very small when considering the distribution on a station-by-station basis (on the order
224 of 50 events). To be concise, here we only present a single example of the distribution
225 metric; however, all the figures are made available in the online supplementary material.
226 Figure 6 shows an example of the distribution metric for the 6_WEIMER Model. The
227 figure presents results for $K=4$ (left column), $K=6$ (middle column), or $K=8$ (right
228 column). Additionally, each row presents results for a different magnetometer station.
229 In the following paragraphs we will summarize the results of this distribution metric for
230 each model, starting with the 6_WEIMER and 9_SWMF models which were the top
231 performers in the event-based metrics.

232 For both mid-latitude stations (OTT and NEW), for observed $K=4$ and $K=8$, the dis-
233 tribution of model predictions for the 6_WEIMER Model is peaked below the observations.
234 For $K=6$ the distribution of model predictions is peaked right at 6 for the mid-latitude
235 stations. For high-latitude stations for all observed values of K the distribution is seen
236 to be shifted to the left representing a systematic under prediction. This pattern seems
237 consistent with the event based studies when the model showed low POFD (apparently
238 due to the systematic under-prediction) and the strongest performance among models for

239 mid-latitude stations when the K threshold is set to 6, but worse performance for higher
240 K threshold and high-latitude.

241 The 9_SWMF Model distribution results for mid-latitude stations are typically peaked
242 at or near the correct values of K . Some moderate spread in the distributions are present
243 indicating the presence of some random error. The same largely holds true for high-
244 latitude results with the spreading a bit more pronounced. Also a slight systematic shift
245 towards under-prediction is seen when the observed $K=8$. This is consistent with the
246 trend seen in the event studies that performance for 9_SWMF was stronger for mid-
247 latitude compared to high-latitude. It is also consistent with the finding from the event
248 table that 9_SWMF has higher skill for threshold of $K=8$ (compared to $K=6$) for mid-
249 latitude, but the reverse is true for high-latitude. Note that virtually identical results are
250 found for 9a_SWMF, which is expected, as it is the same model run, but the magnetometer
251 timeseries from which K is calculated is provided by the model's internal tools rather than
252 the CCMC tool. This provides an independent check of the CCMC tool for calculating
253 the magnetometer timeseries.

254 For the 2_LFM-MIX Model the distribution of model predictions for an observed K tend
255 to peak below the observed value of K for both mid- and high-latitude stations. This shift
256 in the peak of the distribution relative to the observed K is indicative of a systematic
257 under-prediction by the model. The 2_LFM-MIX model was found to have extraordinarily
258 low POFD in the event based analysis which is likely a result of this systematic shift. Some
259 modest evidence of random error is visible in the spreading of the distribution, but it is
260 not enough to result in significant false detections for the K thresholds considered.

261 The 4_OPENGGCM Model demonstrates a large number of occurrences in the model
262 predictions of K values greater than then observed K . Sometimes this is a systematic
263 shift in the distribution (e.g., WNG and NEW , $K=4$), and sometimes it appears to be
264 more random error (e.g., OTT $K=4$ and NEW $K=6$). Regardless of whether the shift
265 is systematic or random, the high-occurrence of predictions significantly exceeding the
266 observations, particularly for mid-latitude stations and lower K values, results in a large
267 rate of false detection (even if true detections are plentiful). This finding is consistent
268 with the high-POFD and high POD exhibited by 4_OPENGGCM in the event studies.

269 For the 2_WEIGEL Model, for both mid- and high-latitude stations, and for all choices
270 of observed K , the distribution of model predictions is peaked below the observations.
271 Such a shift represents a systematic under-prediction of the model. As a result, the model
272 is likely to have a low POFD. These findings are consistent with the event-based analysis
273 which demonstrates that the 2_WEIGEL model has low POFD.

274 Finally, the WingKp Model demonstrates a very large spread indicating significant
275 random error when trying to predict K using the global Kp prediction. For $K=8$, the
276 results are more peaked at the correct value of K although some random error is still
277 visible. The results are similar for high-latitude which is consistent with the event based
278 analysis. However, not including the strongest storm for this model may introduce some
279 bias in the analysis for larger K values. The results for station PBQ are particularly good
280 with peaks at the correct values of K , albeit with some spread. However, the results for
281 stations YKC and ABK exhibit significant random error for all values of K . As WingKp
282 produces a single global prediction of Kp , and we are using that prediction for local K

283 predictions, some error is to be expected. From this type of analysis we can see that the
284 error is mostly random in nature.

285 In summary, the distribution metric, is quite useful in understanding and interpreting
286 the results of the event based metrics. The distribution metric reveals the presence of sys-
287 tematic and random errors and how that can affect the POD and POFD (either positively
288 or negatively).

6. Discussion

289 This work describes another phase of the geospace model validation effort building
290 on the earlier GEM modeling challenges and the dB/dt validation study summarized in
291 *Pulkkinen et al.* [2013]. The work was carried out in coordination among the CCMC,
292 NOAA SWPC, modelers and the science community. The focus of the effort was to
293 evaluate the ability of geospace models to predict the local K index and moreover to
294 evaluate the potential value added of a local prediction over the global prediction.

295 We considered two types of metrics in evaluating the model K prediction: skills scores
296 calculated from event-based contingency tables and a distribution metric. The skills scores
297 (POD, POFD and HSS) from event-based contingency tables for different K thresholds
298 were the primary metric used to rank the models. In particular, the HSS, reflects how
299 much better a model skill is compared to random chance. The derived contingency tables
300 were compiled by grouping all the stations and events together, by separating high-latitude
301 stations and mid-latitude stations for all events, and by separating events into those
302 known to the model developer ahead of time (first four events) and the surprise events
303 selected after models were delivered to CCMC for evaluation (last two events). These
304 different groupings allow us to draw more detailed conclusions about model performance

305 and suitability for forecasting K values at mid-latitude vs high-latitude and for strong
306 events vs very strong events. The distribution metric was an additional tool used to gain
307 insight into aspects of model performance such as revealing random error and systematic
308 errors.

309 In terms of actual model performance, the 9_SWMF and 9a_SWMF models were con-
310 sistently strong performers in all the metrics almost always ranking near the top in all
311 categories. The model had relatively high-POD and low POFD resulting in a HSS that
312 was always among the best. The distribution metric revealed the presence of a moderate
313 amount of random error and limited systematic error. We reiterate that similar perfor-
314 mance is expected for 9_SWMF and 9a_SWMF since they are actually the same model
315 except for how the ground-magnetic field perturbation is calculated.

316 The 2_LFM-MIX model typically had lower performance compared to other models
317 as measured by the HSS. The exception was the last two events for mid-latitude where
318 the model performance was in the middle of the pack. The model typically exhibited
319 lower POD and POFD. The distribution metric shows a clear tendency of this model to
320 under-predict K and that likely results in the lower POD, POFD, and HSS. We note that
321 these results are consistent with the earlier dB/dt study in which the 2_LFM-MIX model
322 performed worse for larger thresholds of magnetic perturbation. It is possible that the
323 model would perform better for lower K thresholds for calculating the contingency tables,
324 just as the model did better in the dB/dt study for lower thresholds. However, the present
325 study is focused primarily on model ability to detect strong and very strong disturbances,
326 not small or moderate disturbances. A cursory examination of a lower threshold of $K=4$
327 did not result in a significant change in the ordering of models by performance (although

328 the HSS increased). Another factor contributing to the poor model performance during
329 storm-time is the lack of of ring current model. More recent version of the LFM include
330 coupling with the Rice Convection Model (RCM) [Pembroke *et al.*, 2012] and are likely
331 to improve performance on these metrics.

332 The 6_WEIMER statistical model performed exceptionally well for mid-latitudes for a
333 threshold of $K=6$, the top performer in this category. The model performance decreased
334 significantly for mid-latitudes with a threshold of $K=8$, but the performance was still
335 strong. In contrast to mid-latitudes the model performance dropped significantly at high-
336 latitude for both K thresholds.

337 The 4_OPENGGCM model had mixed performance. It generally had very good POD,
338 but it also had a consistently elevated POFD. As seen from the distribution metric results,
339 the model had a tendency to over predict, leading to a high POD and high POFD. As
340 a result, sometimes the model has a good HSS and sometimes worse depending on how
341 strongly the POD outweighed the POFD. Significant random and systematic error was
342 likely the cause of the the higher POFD. Regardless of the cause, and overall result on
343 the HSS, an elevated POFD is a concern that needs to be considered in an operational
344 setting. The model did perform better in the last two events compared to the first four.

345 The 2_WEIGEL model was never the top-performing model, but it was also never the
346 worst performing model as measured by HSS. The distribution metric results showed that
347 the model typically underpredicted the observations, and as a result, have an exceedingly
348 low POFD with a reasonable POD.

349 One of the key questions this study addresses is: “How well do geospace models predict
350 local geomagnetic activity (K) compared to representing that activity by the global Kp

351 index?" To answer that question we included in our analysis the WingKp model, which
352 is currently used by SWPC as one method for predicting short-term Kp . The WingKp
353 model never ranked at the bottom or the top of the model rankings based on its HSS.
354 Interestingly, the model used in this way was also often not the lowest performing model,
355 indicating that using the WingKp prediction of global Kp (as a local K prediction) would
356 actually exhibit higher skill than using the local K predicted by some models. However,
357 the POFD was typically elevated compared to other models. An elevated POFD raises
358 concerns for using the global Kp prediction from WingKp for local forecasts of K , but it
359 also demonstrates the potential value of a local K forecast. All local K forecasts (except
360 for 4_OPENGGCM) consistently had much lower POFD than WingKp. However, the
361 POD score is near the top in some cases. One caution when interpreting these results
362 is that the WingKp model used different solar wind inputs than the other models. It is
363 possible that the results could have been somewhat different had the same input solar
364 wind parameters been used.

365 One consideration for transition to operations is lead times for model prediction. The
366 main constraint in this regard is the input data from ACE which arrives at most one
367 hour ahead of the event. The empirical models in this study can provide a practically
368 instantaneous prediction with very modest computing resources while the MHD models
369 are more resource intensive. As noted earlier, one requirement for the MHD models was
370 they could run in twice real-time on a moderately sized supercomputing cluster. If larger
371 computational resources are available these models could run faster. Nevertheless, the
372 empirical models will always be more computationally efficient than the MHD models.

373 All the models had positive HSS demonstrating better prediction skill than random
374 chance. Moreover, we found most results consistent with the dB/dt study of *Pulkkinen*
375 *et al.* [2013]. When considering all events, a POD of around 70% is found for the top
376 performing models for mid-latitude stations, even with a K threshold of 8. For high-
377 latitude stations, the POD possible for top performing models drops to around 50%. In
378 either case, the POFD for most models is exceedingly low for the thresholds considered.
379 Whether this performance is sufficient for current space weather prediction needs, or if
380 further improvement is required is not a question addressed in this study. We also note
381 that this study only evaluates model prediction of K and therefore cannot be used to
382 draw conclusions about how those models would perform when predicting other quantities,
383 even closely related ones. Indeed, it is entirely possible to that a model can produce a
384 value of K that is very close to that determined from the measurements, while having
385 ΔB predictions with signs that are mostly opposite of the measured value. As a result
386 of the model evaluation conducted by CCMC in coordination with modelers and NOAA-
387 SWPC, NOAA-SWPC has decided to transition the SWMF model to space weather
388 operations and to give further consideration to the Weimer model. As the models continue
389 to improve and evolve, it is likely that more geospace models will transition to operations
390 for purposes of addressing specific user needs, for incorporating improved models, and for
391 ensemble modeling. Indeed, this validation is just one step on the path of operationalizing
392 state-of-the-art codes for space weather forecasting.

393 **Acknowledgments.**

394 The data from the ground based magnetic observatories was critical to this study. As
395 such, we thank the institutions that support those observatories as well asl INTERMAG-

396 NET for promoting high-standards of practice (www.intermagnet.org). The National
397 Center for Atmospheric Research is supported by the National Science Foundation. All
398 model output used in the analysis is available through the CCMC as described in the
399 manuscript.

References

- 400 Boteler, D.H., R.J. Pirjola, and H. Nevanlinna (1998), The Effects of Geomagnetic Dis-
401 turbances on Electrical Systems at the Earth's Surface, *Adv. Space Res.*, 22, 17-27.
- 402 Lopez, R. E., S. Hernandez, M. Wiltberger, C.-L. Huang, E. L. Kepko, H. Spence,
403 C. C. Goodrich, and J. G. Lyon (2007), Predicting magnetopause crossings at
404 geosynchronous orbit during the Halloween storms, *Space Weather*, 5, S01005,
405 doi:10.1029/2006SW000222.
- 406 McCollough, J. P., S. L. Young, and W. R. Frey (2014), Real-Time Validation of the *Kp*
407 Predictor Model, *AFRL Tech. Rep.*, AFRL-RV-PS-TR-2015-0073.
- 408 National Research Council (2008), Severe Space Weather Events-Understanding Societal
409 and Economic Impacts: A Workshop Report, *The National Academies Press*, Washing-
410 ton, DC.
- 411 North American Electric Reliability Corporation GMD Task Force (2012), 2012 Special
412 Reliability Assessment Interim Report: Effects of Geomagnetic Disturbances on the
413 Bulk Power System, *NERC*, February 2012.
- 414 Pembroke, A., F. Toffoletto, S. Sazykin, M. Wiltberger, J. Lyon, V. Merkin, and P.
415 Schmitt (2012), Initial results from a dynamic coupled magnetosphere-ionosphere-ring
416 current model, *J. Geophys. Res.*, 117(A2), doi:10.1029/2011JA016979.

- 417 Pirjola, R., Effects of space weather on high-latitude ground systems, *Advances in Space*
418 *Research*, 36, 2231-2240, 2005.
- 419 Pulkkinen, A., L. Rastätter, M. Kuznetsova, M. Hesse, A. Ridley, J. Raeder, H.J. Singer,
420 and A. Chulaki (2010), Systematic evaluation of ground and geostationary magnetic
421 field predictions generated by global magnetohydrodynamic models, *Journal of Geo-*
422 *physical Research*, 115, A03206, doi:10.1029/2009JA014537.
- 423 Pulkkinen, A., M. Kuznetsova, A. Ridley, J. Raeder, A. Vapirev, D. Weimer, R. S. Weigel,
424 M. Wiltberger, G. Millward, L. Rastätter, M. Hesse, H. J. Singer and A. Chulaki
425 (2013), Geospace Environment Modeling 2008-2009 Challenge: ground magnetic field
426 perturbations, *Space Weather*, Vol. 9, S02004, doi:10.1029/2010SW000600.
- 427 Pulkkinen, A., L. Rastätter, M. Kuznetsova, H. Singer, C. Balch, D. Weimer, G. Toth,
428 A. Ridley, T. Gombosi, M. Wiltberger, J. Raeder, R. Weigel(2011), Community-wide
429 validation of geospace model ground magnetic field perturbation predictions to support
430 model transition to operations, *Space Weather*, Vol. 11, doi:10.1002/swe.20056.
- 431 Raeder, J., D. Larson, W. Li, E. L. Kepko, and T. Fuller-Rowell (2008),
432 OpenGGCM simulations for the THEMIS mission, *Space Sci. Rev.*, 141, 535,
433 doi:10.1007/s11,21400894215.
- 434 Rastätter, L., M. Kuznetsova, A. Vapirev, A. Ridley, M. Wiltberger, A. Pulkki-
435 nen, M. Hesse and H.J. Singer (2011), Geospace Environment Modeling 2008-
436 2009 Challenge: geosynchronous magnetic field, *Space Weather*, Vol. 9, S04005,
437 doi:10.1029/2010SW000617.
- 438 Rastätter, L., G. Toth, M. M. Kuznetsova, and A. A. Pulkkinen (2014), CalcDeltaB:
439 An efficient postprocessing tool to calculate ground-level magnetic perturbations from

- 440 global magnetosphere simulations, *Space Weather*, 11, doi:10.1002/2014SW001083
- 441 Rostoker, G. (1972), Geomagnetic Indices, *Reviews of Geophysics*, 10, 935–950,
442 doi:10.1029/RG010i004p00935.
- 443 Skoug, R. M., Gosling, J. T., Steinberg, J. T., McComas, D. J., Smith, C. W., Ness, N. F.,
444 Hu, Q., Burlaga, L. F., (2004), Extremely high-speed solar wind: 2930 October 2003,
445 *J. Geophys. Res.*, 109, A09102, doi:10.1029/2004JA010494.
- 446 Tóth, G. et al (2012), Adaptive numerical algorithms in space weather modeling, *J.*
447 *Comput. Phys.*, 231.
- 448 Weigel, R. S., A. J. Klimas, D. Vassiliadis (2003), Solar wind coupling to and predictability
449 of ground magnetic fields and their time derivatives, *J. Geophys. Res.*, 108 (A7), 1298,
450 doi:10.1029/2002JA009627.
- 451 Weimer, D. R. (2013), An empirical model of ground-level geomagnetic perturbations,
452 *Space Weather*, 11, 107-120, doi: :10.1002/swe.20030.
- 453 Wiltberger, M., W. Wang, A. G. Burns, S. C. Solomon, J. G. Lyon, and C. C.
454 Goodrich (2004), Initial results from the coupled magnetosphere ionosphere thermo-
455 sphere model: magnetospheric and ionospheric responses, *Journal of Atmospheric and*
456 *Solar-Terrestrial Physics*, 66(1), 1411-1423, doi:10.1016/j.jastp.2004.03.026.
- 457 Wing, S., J. R. Johnson, J. Jen, C.-I. Meng, D. G. Sibeck, K. Bechtold, J. Freeman,
458 K. Costello, M. Balikhin, and K. Takahashi (2005), *Kp* forecast models, *Journal of*
459 *Geophysical Research*, 110, A04203, doi:10.1029/2004JA010500.

Table 1. Geospace events studied in the validation activity. The last two columns give the minimum Dst index and the maximum Kp index of the event, respectively.

Event #	Date and time	min(Dst)	max(Kp)
1	October 29, 2003 06:00 UT - October 30, 06:00 UT	-353 nT	9
2	December 14, 2006 12:00 UT - December 16, 00:00 UT	-139 nT	8
3	August 31, 2001 00:00 UT - September 1, 00:00 UT	-40 nT	4
4	August 31, 2005 10:00 UT - September 1, 12:00 UT	-131 nT	7
5	April 5, 2010 00:00 UT - April 6, 00:00 UT	-73 nT	8-
6	August 5, 2011 09:00 UT - Aug 6, 09:00 UT	-113 nT	8-

Table 2. The locations of the geomagnetic observatories used in the study.

Station name	Station code	Geomagnetic lat	Geomagnetic lon	Scaling Factor
Yellowknife	YKC	68.9	299.4	3.0
Newport	NEW	54.9	304.7	1.4
Poste-de-la-Baleine	PBQ	65.5	351.8	3.0
Sanikiluaq	SNK	66.4	356.1	3.0
Ottawa	OTT	55.6	355.3	1.5
Abisko	ABK	66.1	114.7	3.0
Wingst	WNG	54.1	95.0	1.0



Figure 1. The locations and the station codes of the geomagnetic observatories used in the study. Geomagnetic dipole coordinates are used. Red and blue circles indicate high-latitude and mid-latitude stations, respectively, used in the final analyses in Section 5.

Table 3. Look up table to determine K from scaled range of ΔB .

K -index	nT range
0	$0 \leq \text{Range of } \Delta B < 5$
1	$5 \leq \text{Range of } \Delta B < 10$
2	$10 \leq \text{Range of } \Delta B < 20$
3	$20 \leq \text{Range of } \Delta B < 40$
4	$40 \leq \text{Range of } \Delta B < 70$
5	$70 \leq \text{Range of } \Delta B < 120$
6	$120 \leq \text{Range of } \Delta B < 200$
7	$200 \leq \text{Range of } \Delta B < 330$
8	$330 \leq \text{Range of } \Delta B < 500$
9	$500 \leq \text{Range of } \Delta B$

Author Manuscript

Table 4. Models analyzed in the validation effort. Each model is assigned a unique model identifier given by the leftmost column of the table. The table indicates the model description, and if applicable, the number of cells and the minimum spatial resolution used in the global MHD part of the model. See text in Section 4 for details.

Identifier	(model version) Model	Grid (# of cells, min. res.)
2_LFM-MIX	(LTR-2.1.1) LFM coupled with ionospheric electrodynamics	163,000, $0.4 R_E$
3_WEIGEL	empirical model	N/A
4_OPENGGCM	(OpenGGCM 4.0) global MHD coupled with CTIM	3.9 million, $0.25 R_E$
6_WEIMER	empirical model	N/A
9_SWMF	(SWMF 2011-01-31) BATS-R-US coupled with RIM and RCM	1 million, $0.25 R_E$
9a_SWMF	Same as 9_SWMF but using internal SWMF calculation for magnetometer timeseries	

Acronyms:

RIM	Ridley Ionosphere Model
RCM	Rice Convection Model
CTIM	Coupled Thermosphere Ionosphere Model

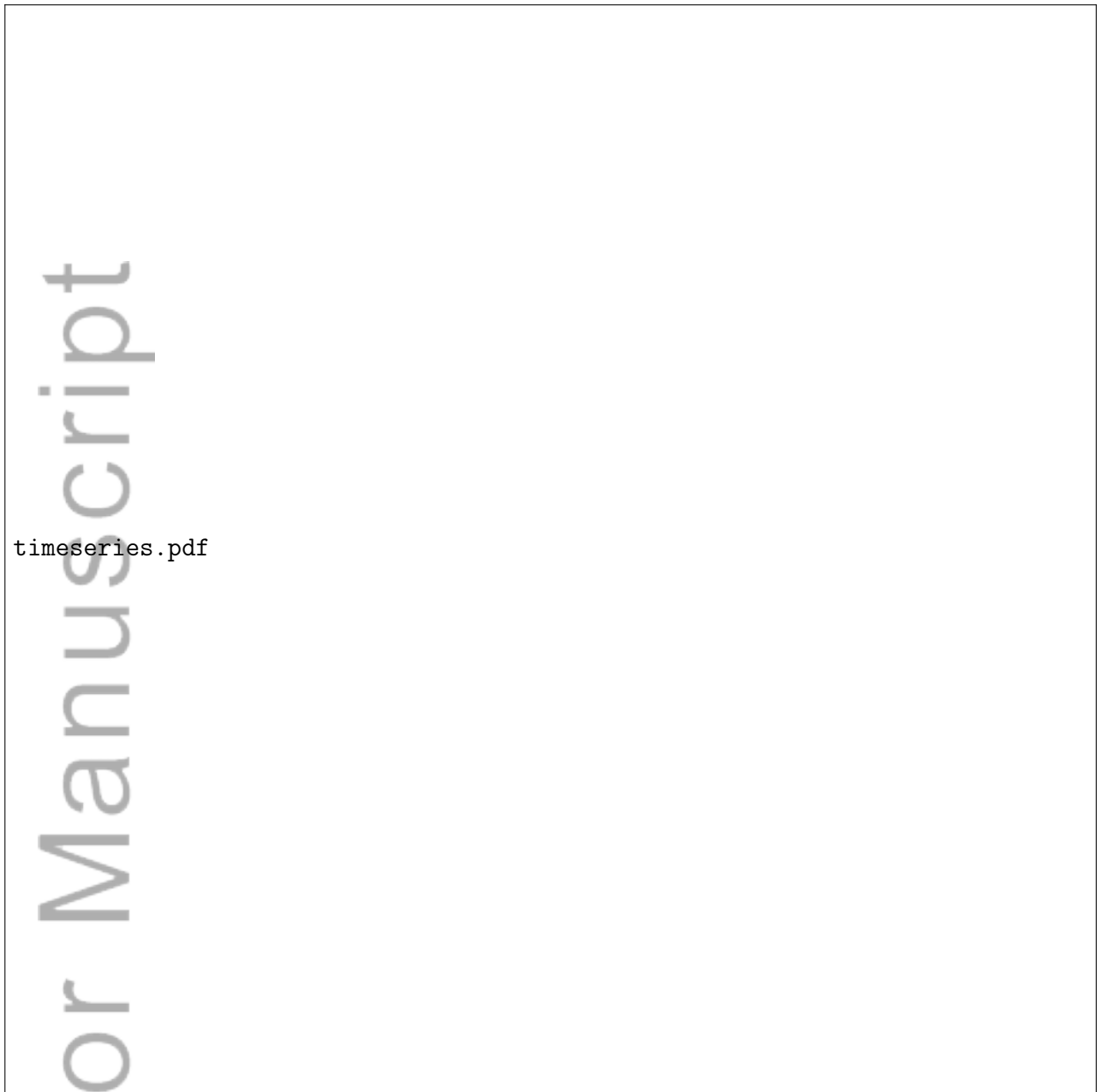


Figure 2. Time series of the observed (Black) and modeled (Red) K predictions for a particular mid-latitude station (OTT). Each panel shows a different model's prediction.

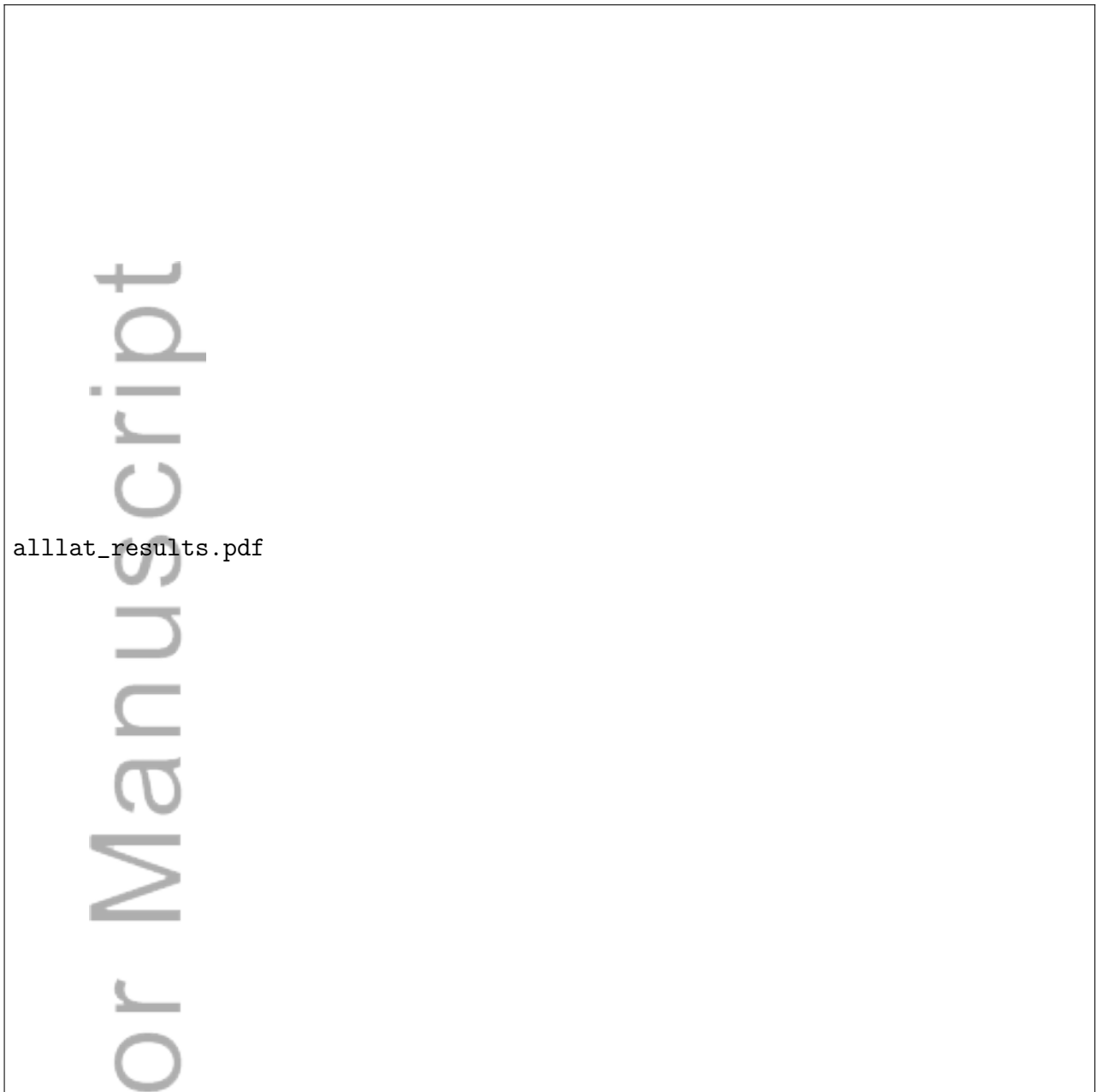


Figure 3. Heidke Skill Score (HSS), Critical Success Index (CSI), Probability of Detection (POD) (blue curve) and Probability of False Detection (POFD) (yellow curve) defined in Section 3 for the K thresholds 6 (left panel) and 8 (right panel). POD and POFD obtained by integrating over the three mid-latitude stations and the three high-latitude stations. The models (see Table 4) are ordered according to their HSS. The model with the largest HSS is the leftmost in all panels.

Run	n_event	n_noevent	H	F	M	N	HSS	CSI	POD	POFD
9_SWMF	1240	1532	801	74	439	1458	0.61	0.61	0.65	0.05
9a_SWMF	1240	1532	752	38	488	1494	0.60	0.59	0.61	0.02
6_WEIMER	1240	1532	605	20	635	1512	0.50	0.48	0.49	0.01
2_WEIGEL	1240	1532	537	25	703	1507	0.44	0.42	0.43	0.02
WingKp	1151	1117	722	279	429	838	0.38	0.50	0.63	0.25
4_OPENGGC	1240	1532	803	425	437	1107	0.37	0.48	0.65	0.28
2_LFM-MIX	1240	1532	353	26	887	1506	0.29	0.28	0.28	0.02

Table 5. Table for all stations, threshold 6

Run	n_event	n_noevent	H	F	M	N	HSS	CSI	POD	POFD
9a_SWMF	395	2377	201	55	194	2322	0.57	0.45	0.51	0.02
9_SWMF	395	2377	210	80	185	2297	0.56	0.44	0.53	0.03
2_WEIGEL	395	2377	116	41	279	2336	0.37	0.27	0.29	0.02
4_OPENGGC	395	2377	139	145	256	2232	0.33	0.26	0.35	0.06
WingKp	370	1898	121	137	249	1761	0.29	0.24	0.33	0.07
6_WEIMER	395	2377	79	18	316	2359	0.28	0.19	0.20	0.01
2_LFM-MIX	395	2377	42	11	353	2366	0.16	0.10	0.11	0.00

Table 6. Table for all stations, threshold 8



Figure 4. Heidke Skill Score (HSS) (red curve), Critical Success Index (CSI) (blue curve), Probability of Detection (POD) (green curve) and Probability of False Detection (POFD) (yellow curve) defined in Section 3 for the K thresholds 6 (left panel) and 8 (right panel). POD and POFD are obtained by integrating over the three high-latitude stations. The models (see Table 4) are ordered according to their HSS. The model with the largest HSS is the leftmost in all panels.

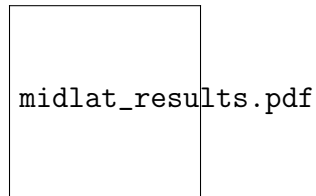
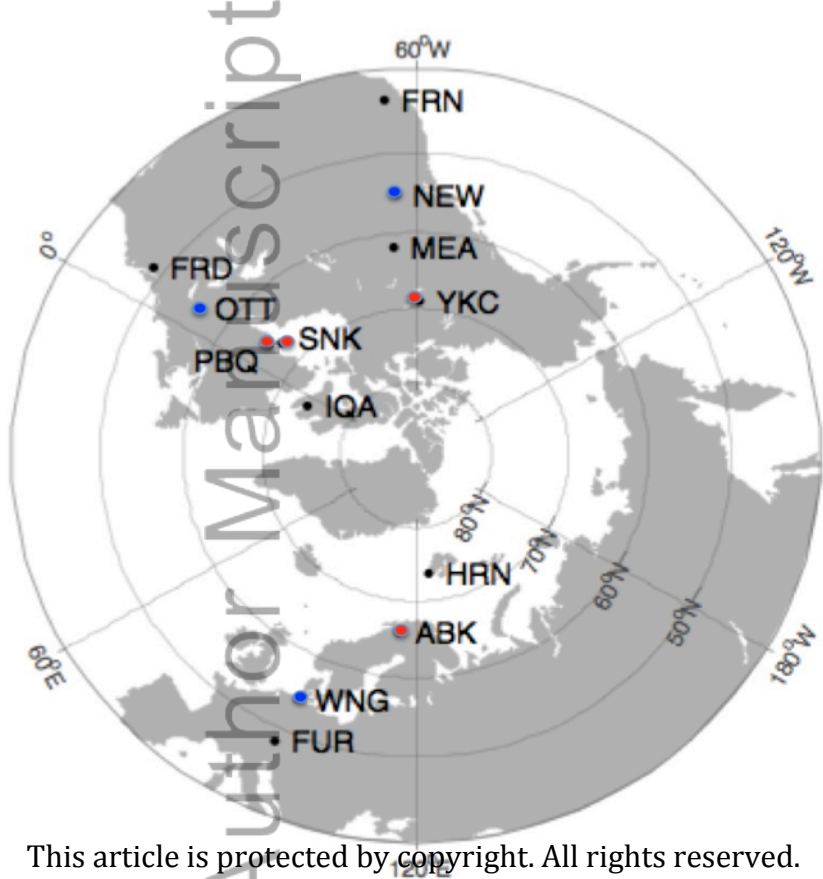


Figure 5. Heidke Skill Score (HSS) (red curve), Critical Success Index (CSI) (blue curve), Probability of Detection (POD) (green curve) and Probability of False Detection (POFD) (yellow curve) defined in Section 3 for the K thresholds 6 (left panel) and 8 (right panel). POD and POFD are obtained by integrating over the three mid-latitude stations. The models (see Table 4) are ordered according to their HSS. The model with the largest HSS is the leftmost in all panels.



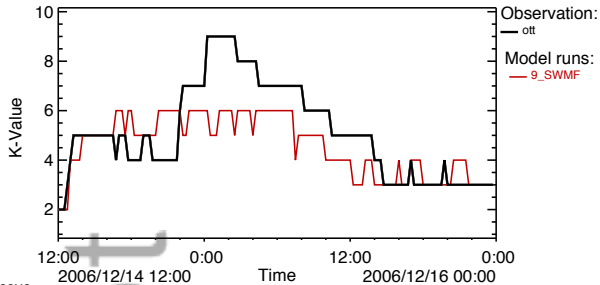
Figure 6. Distribution of 6_WEIMER Model predictions when $K=4$ (left column), $K=6$ (middle column), and $K=8$ (right column). Each row presents results for a different mid-latitude station.



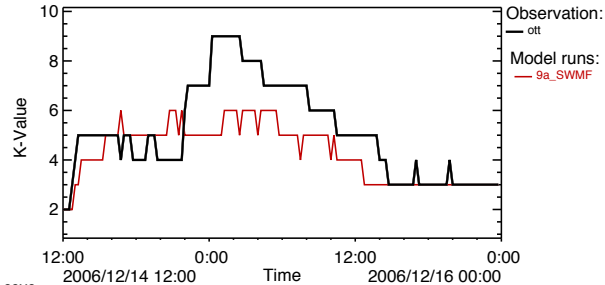
This article is protected by copyright. All rights reserved.

Sample Time Series

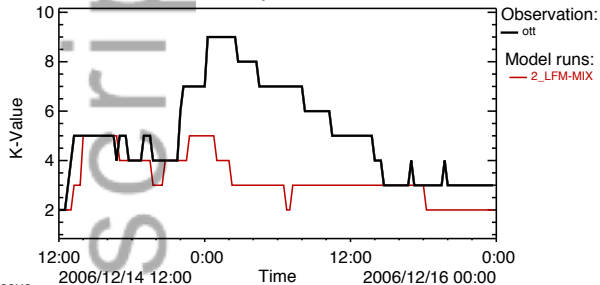
K-Value from observatory file: ott_OBS_20061214.txt



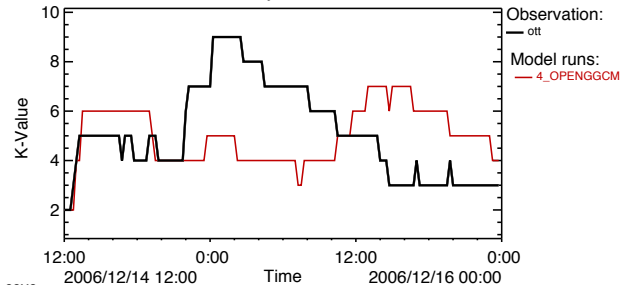
K-Value from observatory file: ott_OBS_20061214.txt



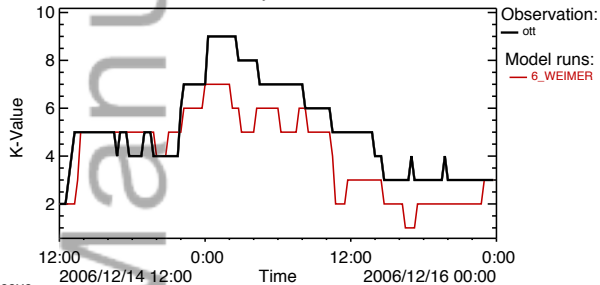
K-Value from observatory file: ott_OBS_20061214.txt



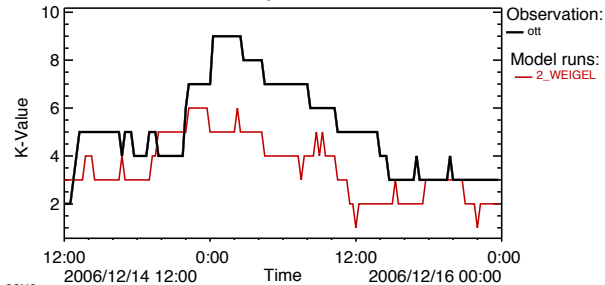
K-Value from observatory file: ott_OBS_20061214.txt



K-Value from observatory file: ott_OBS_20061214.txt

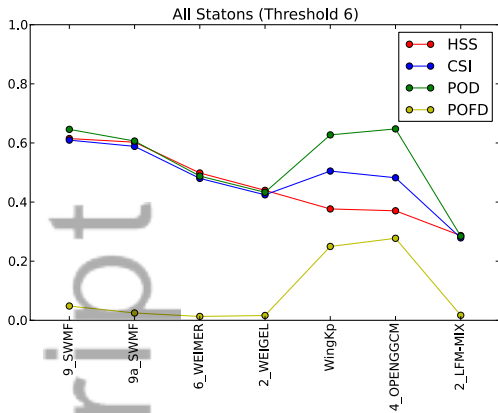


K-Value from observatory file: ott_OBS_20061214.txt

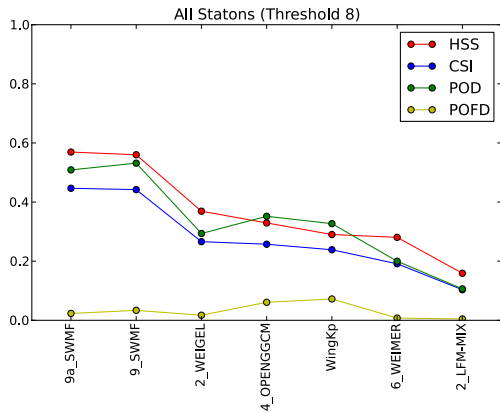


Skill Scores for All Stations

Threshold 6

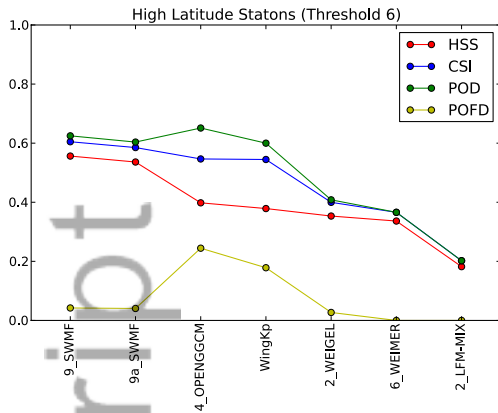


Threshold 8

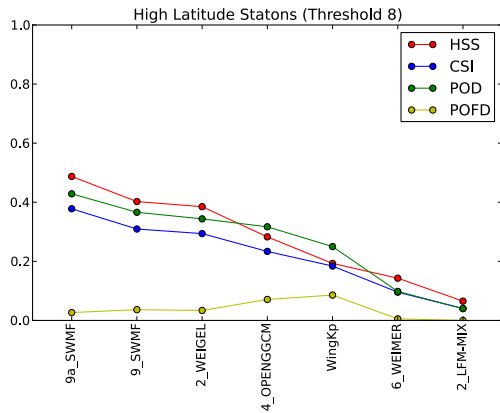


Skill Scores for High-Lat Stations

Threshold 6

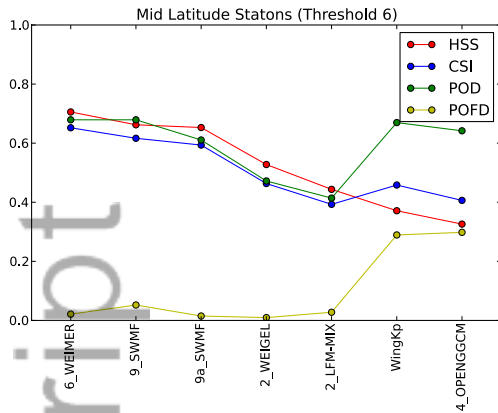


Threshold 8

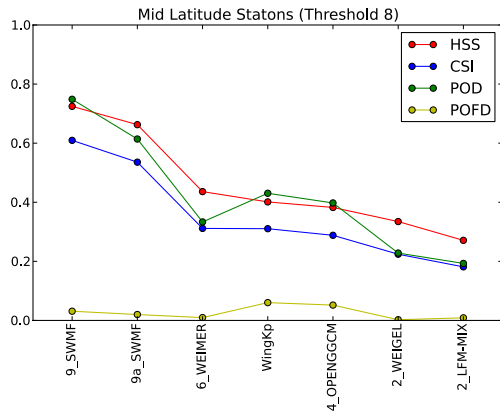


Skill Scores for Mid-Lat Stations

Threshold 6



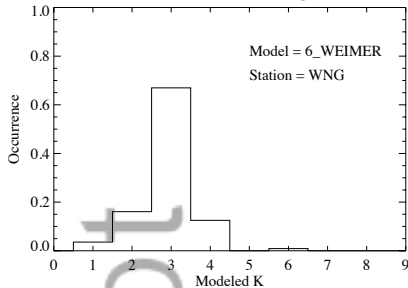
Threshold 8



Distribution Metric for 6-WEIMER (Mid-Lat)

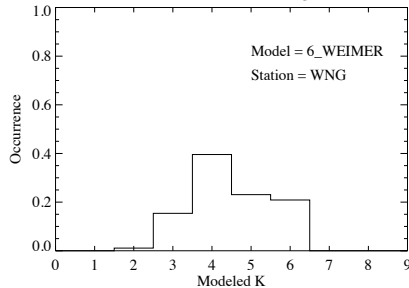
k=4

Distribution of Model K for Obs. K = 4



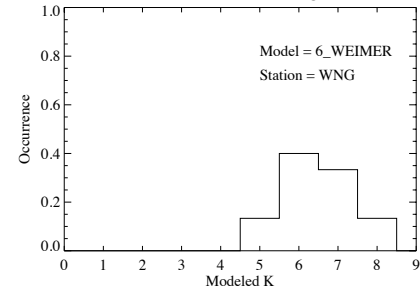
k=6

Distribution of Model K for Obs. K = 6

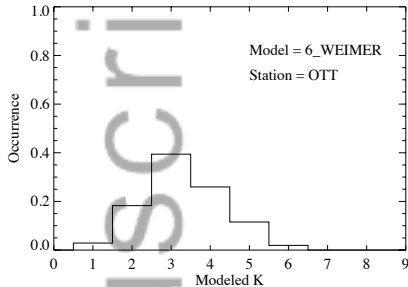


k=8

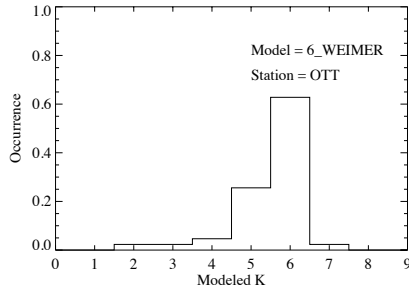
Distribution of Model K for Obs. K = 8



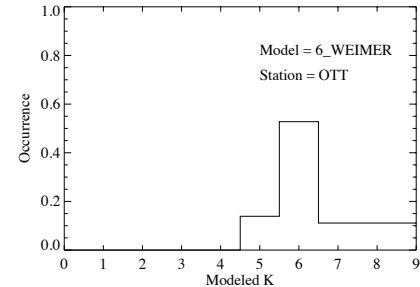
Distribution of Model K for Obs. K = 4



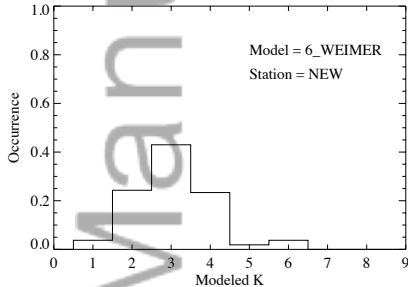
Distribution of Model K for Obs. K = 6



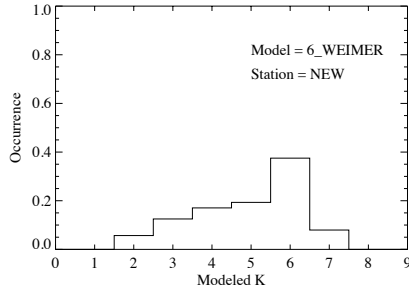
Distribution of Model K for Obs. K = 8



Distribution of Model K for Obs. K = 4



Distribution of Model K for Obs. K = 6



Distribution of Model K for Obs. K = 8

