

# The STONE curve: A ROC-derived model performance assessment tool

Michael W. Liemohn,<sup>1</sup> Abigail R. Azari,<sup>1,2</sup> Natalia Yu. Ganushkina,<sup>1,3</sup> and Lutz Rastätter<sup>4</sup>

<sup>1</sup>Department of Climate and Space Sciences and Engineering, University of Michigan, Ann Arbor, MI.

<sup>2</sup>Now at the Space Science Laboratory, University of California, Berkeley, CA

<sup>3</sup>Finnish Meteorological Institute, Helsinki, Finland

<sup>4</sup>Community Coordinated Modeling Center, NASA Goddard Space Flight Center, Greenbelt, MD

Corresponding author: Michael Liemohn ([liemohn@umich.edu](mailto:liemohn@umich.edu))

Submitted to *Earth and Space Science*

## Key Points:

- A new event-detection-based metric for model performance appraisal is given with sliding thresholds in both observational and model values
- The new metric is like the relative operating characteristic curve but uses continuous observational values, not just categorical status
- The new metric is used on real-time model predictions of common geomagnetic activity parameters, demonstrating its features and strengths

## AGU Index Terms:

- 1984 Statistical methods: Descriptive (4318)
- 4318 Statistical analysis (1984, 1986)
- 7924 Forecasting (1922, 2722, 4315)
- 0550 Model verification and validation
- 9820 Techniques applicable in three or more fields

## Keywords:

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: [10.1029/2020EA001106](https://doi.org/10.1029/2020EA001106)

**33 Abstract**

34 A new model validation and performance assessment tool is introduced, the sliding threshold of  
35 observation for numeric evaluation (STONE) curve. It is based on the relative operating  
36 characteristic (ROC) curve technique, but instead of sorting all observations in a categorical  
37 classification, the STONE tool uses the continuous nature of the observations. Rather than  
38 defining events in the observations and then sliding the threshold only in the classifier/model  
39 data set, the threshold is changed simultaneously for both the observational and model values,  
40 with the same threshold value for both data and model. This is only possible if the observations  
41 are continuous and the model output is in the same units and scale as the observations, i.e., the  
42 model is trying to exactly reproduce the data. The STONE curve has several similarities with the  
43 ROC curve – plotting probability of detection against probability of false detection, ranging from  
44 the (1,1) corner for low thresholds to the (0,0) corner for high thresholds, and values above the  
45 zero-intercept unity-slope line indicating better than random predictive ability. The main  
46 difference is that the STONE curve can be nonmonotonic, doubling back in both the x and y  
47 directions. These ripples reveal asymmetries in the data-model value pairs. This new technique is  
48 applied to modeling output of a common geomagnetic activity index as well as energetic electron  
49 fluxes in the Earth's inner magnetosphere. It is not limited to space physics applications but can  
50 be used for any scientific or engineering field where numerical models are used to reproduce  
51 observations.

52

**53 Plain Language Summary**

54 Scientists often try to reproduce observations with a model, helping them explain the  
55 observations by adjusting known and controllable features within the model. They then use a  
56 large variety of metrics for assessing the ability of a model to reproduce the observations. One  
57 such metric is called the relative operating characteristic (ROC) curve, a tool that assesses a  
58 model's ability to predict events within the data. The ROC curve is made by sliding the event-  
59 definition threshold in the model output, calculating certain metrics and making a graph of the  
60 results. Here, a new model assessment tool is introduced, called the sliding threshold of  
61 observation for numeric evaluation (STONE) curve. The STONE curve is created by sliding the  
62 event definition threshold not only for the model output but also simultaneously for the data  
63 values. This is applicable when the model output is trying to reproduce the exact values of a  
64 particular data set. While the ROC curve is still a highly valuable tool for optimizing the  
65 prediction of known and pre-classified events, it is argued here that the STONE curve is better  
66 for assessing model prediction of a continuous-valued data set.

67

68

69 **1. Introduction**

70 Numerical models are a fundamental feature of research in the natural sciences. Models  
71 are often used to explain strange and interesting features in an archival data set in order to assess  
72 the physical processes responsible for that observational signature. They are also used for  
73 prediction, using some estimate of future initial and boundary conditions to determine the state  
74 of the system, or even a particular observational quantity, ahead of time. These are typical uses  
75 of models in every discipline of Earth and space sciences.

76 There exists a large collection of metrics to assess the goodness of fit for these models to  
77 a particular data set. These metrics, for the most part, can be sorted into several major groupings,  
78 two of which are fit performance metrics and event detection metrics (e.g., Wilks, 2019; Joliffe  
79 and Stephenson et al., 2012; Liemohn, McCollough, et al., 2018). The former group, also called  
80 continuous metrics, is usually based on differencing each data-model value pair and includes  
81 many well-known assessment equations such as root mean square error, correlation coefficient,  
82 mean error, and prediction efficiency (e.g., Hogan and Mason, 2012; Morley et al., 2018). The  
83 second group, also called categorical metrics, is based on categorizing the observations into  
84 events and non-events and then assessing a model's ability to reproduce this classification. This  
85 is done through a contingency table (also commonly called a confusion matrix) in which each  
86 data-model pair gets two designations: determining if the observation is in the event state or not  
87 and similarly if the model value is in the event state or not. The similarity or difference of the  
88 data and model values is irrelevant, only the event/non-event designation matters. This second  
89 group includes other well-known assessment equations such as the probability of detection, false  
90 alarm rate, frequency bias, and Heidke skill score (see, e.g., Muller et al., 1944; Wilks, 2019).  
91 Continuous metrics are sometimes defined as comparisons that assess model parameters on a  
92 continuous scale. This is broader than the definition above and can include any event-based  
93 categorical metric by sweeping the model event identification threshold from low to high values.

94 A feature of the event detection metrics is that the model does not have to cover the same  
95 range or even have the same units as the observations. The model could be anything that might  
96 predict the event state of the observations. Furthermore, the observations do not have to be a  
97 continuous-valued real number set, but could be pre-categorized into events and non-events (or a  
98 multi-level classification). The model could be a continuous-valued real number set or a discrete-  
99 valued categorized set. When the data or model happens to be a continuous-valued real number  
100 set, then a threshold value for event identification is chosen, a threshold value that could be  
101 different between the observational events and the modeled events.

102 An event detection metric that is often used for weather prediction (e.g., Mason, 1982),  
103 psychology (e.g., Swets, 1973), medical clinical trials (e.g., Ekelund, 2011), and machine  
104 learning (e.g., Fawcett et al., 2006) is the relative (or, originally, receiver) operating  
105 characteristic (ROC) curve (see review by Carter et al., 2016). This is an assessment tool that can  
106 be applied when the model values are continuous-valued real numbers, using not just one event  
107 identification threshold but many. The method is to sweep the event definition threshold for the  
108 model values from low to high, calculating two specific metrics, the probability of detection  
109 (POD) and the probability of false detection (POFD), and plotting these two arrays against each  
110 other. The threshold that yields the location on the ROC curve closest to the upper left corner  
111 (high POD and low POFD) can be considered a possible "best setting" for event prediction by

112 this model. This is not the only location for an optimum pick of a final threshold along a ROC  
113 curve. Often the final choice will depend on the application and problem specific details. For  
114 example, recent developments have discussed the use of skill scores for different solar and space  
115 physics applications (e.g., Bobra & Couvidat 2015) and their location on ROC diagrams (e.g.,  
116 Manzato, 2007; Azari et al., 2018). A further detailed discussion on skill scores and their relation  
117 to ROC diagrams can be found within Manzato (2005). An integral quantity sometimes used  
118 from the ROC curve is the area under the curve (AUC), which is an overall measure of goodness  
119 of fit for the model to the observational events across all of the possible model value event  
120 identification thresholds.

121 The ROC curve has recently been used quite often in the Earth and space sciences to  
122 assess model performance at detecting events in an observational data set. It is used regularly in  
123 the atmospheric sciences, such as for regional ozone ensemble forecasting (e.g., Delle Monache  
124 et al., 2006), deciphering the microphysical properties of clouds (e.g., Gabriel et al., 2009), and  
125 forecasting summer monsoons over India (e.g., Borah et al., 2013). Earth scientists also employ  
126 the ROC curve for a diverse set of modeling activities, including the distribution of rock glaciers  
127 (e.g., Brenning et al., 2007), assessing triggering mechanisms of earthquake aftershocks (e.g.,  
128 Meade et al., 2017), and snow slab instability physics (e.g., Reuter & Schweizer, 2018). This  
129 also includes land-air interactions, such as mapping of expected ash cloud locations after  
130 eruptions (e.g., Stefanescu et al., 2014), modeling rainfall-induced landslides (e.g.,  
131 Anagnostopoulos et al., 2015), and statistically forecasting extreme corn losses in the eastern  
132 United States (Mathieu & Aires, 2018). The fields of space and planetary science have also  
133 started to employ this technique, such as for oblique ionogram retrieval algorithm assessment  
134 (Ippolito et al., 2016), identifying energetic particle flux injections at Saturn (e.g., Azari et al.,  
135 2018), magnetic activity prediction (e.g., Liemohn, McCollough, et al., 2018), and identifying  
136 solar flare precursors (e.g., Chen et al., 2019). In short, the ROC curve has become an essential  
137 tool, among many that can and should be applied, for model assessment across many natural  
138 science disciplines.

139 The ROC curve, however, only assesses the model's ability to predict a single  
140 observational event identification threshold. While this is desirable if the data were pre-classified  
141 as events or non-events, this imposes a simplification of the data set when the observations are  
142 also continuous-valued real numbers. That is, the ROC curve does not test the model's ability to  
143 predict events across the full range of the data. A family of ROC curves can be produced using  
144 different data-value event identification thresholds (and sweeping the model-value event  
145 identification threshold to produce each ROC curve), which is acceptable if the model is only  
146 being used to maximize the prediction of events. If the model, however, is trying to reproduce  
147 the exact values of the observations, then it is useful to conduct an assessment for which the data  
148 and model have the same threshold setting. The ROC curve, unfortunately, cannot easily test the  
149 model's ability to reproduce the observed events at the same threshold setting, sweeping through  
150 all possible event identification thresholds.

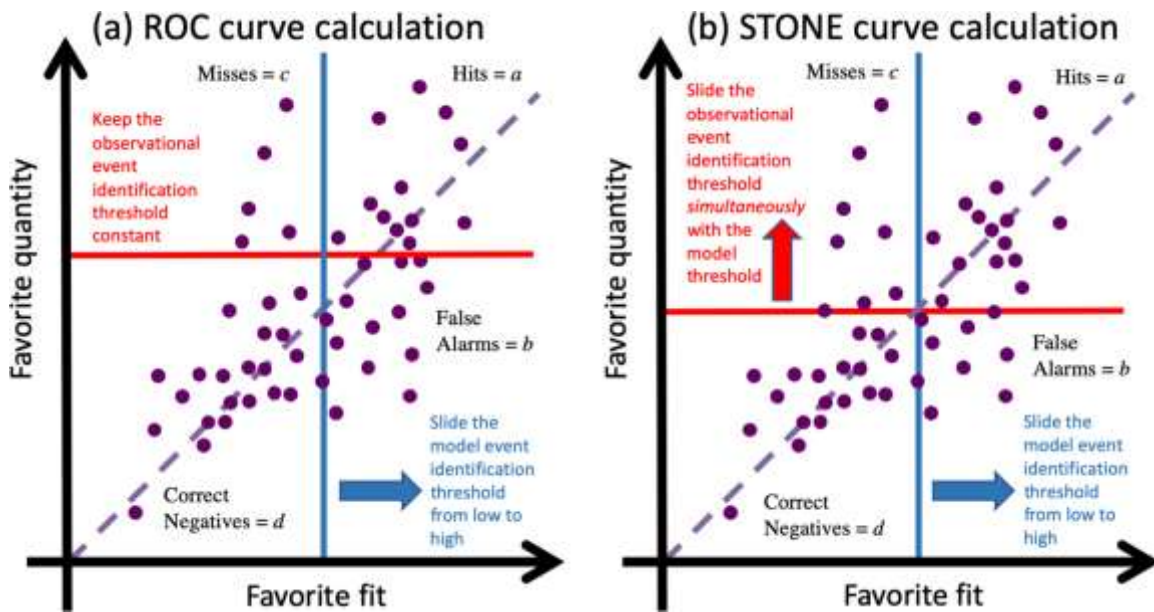
151 One could argue the need for a new metric that provides a more complete evaluation  
152 across multiple observation thresholds. Like the ROC curve, this new metric should test a  
153 model's ability to predict observed events across the full range of possible model-value event  
154 identification settings, but rather than using a single observational event categorization, it should  
155 sweep through the same range of event identification thresholds as used for the model. Such a  
156 metric is proposed below, called the sliding threshold of observation numeric evaluation, or

157 STONE, curve. This is based on the ROC curve but includes the desirable features described  
 158 above. The work then presents an application of the STONE curve to two space physics data  
 159 sets, the prediction of a geomagnetic activity index and energetic electron fluxes in near-Earth  
 160 space. Similarities and differences between the ROC and STONE curves are discussed, as well  
 161 as the interpretive meaning of features in the STONE curve.

## 162 2. Method of Calculation

163 The calculation of a STONE curve is rather similar to that of a ROC curve, with one  
 164 major exception – both thresholds slide together, incrementing the two event identification  
 165 thresholds simultaneously so that the same threshold value is used for both the data and the  
 166 model at each setting from low to high across the range. Because this tool is for continuous-  
 167 valued observations and model results, for which an “event” is an arbitrary designation, there  
 168 does not have to be a pre-defined event threshold in the observations. In fact, it is desired that the  
 169 model match the observations for all levels of “event” definition. Therefore, in the STONE tool,  
 170 the two thresholds move together. This is illustrated in Figure 1, showing an arbitrary data set  
 171 plotted against a model output that is trying to reproduce these values.

172



173

174 **Figure 1.** Idealized examples of how to calculate (a) the ROC curve and (b) the STONE curve.  
 175 In (a), only the blue curve shifts while the red curve is fixed at some level. In (b), both the red  
 176 and blue thresholds shift together. As these lines shift, data points are converted from one  
 177 quadrant to another. The purple dashed curve is the zero-intercept unity-slope line, for reference.

178

179 Figure 1a shows the calculation scenario for the ROC curve, with the event identification  
 180 threshold for the observations set to a fixed value and the threshold for the model results  
 181 sweeping from low to high values. Annotations label the four quadrants of the chart, as defined  
 182 by these two thresholds. As the model threshold changes, the points in the chart change quadrant.  
 183 Specifically, two shifts occur: points in the “hits” quadrant (variable  $a$ ) move to the “misses”

184 quadrant (*c*) and points in the “false alarms” quadrant (*b*) move to the “correct negatives”  
 185 quadrant (*d*).

186 The ROC curve is defined from two metrics in the “discrimination” category (Murphy &  
 187 Winkler, 1987) of data-model comparison techniques. Discrimination metrics are assessments  
 188 that only use a portion of the data values within a specified range (and the corresponding model  
 189 values). For event detection metrics, the usual practice is to use the event state of the  
 190 observations to define the subsets of the data. In particular, the ROC curve uses POD and POFD,  
 191 which have the following formulas:

$$192 \quad \text{POD} = \frac{a}{a+c} \quad (1)$$

$$193 \quad \text{POFD} = \frac{b}{b+d} \quad (2)$$

194 Where *a*, *b*, *c*, and *d* are point counts from the quadrants in the scatter plot. It is seen that these  
 195 two formulas are mutually exclusive, POD only uses the hits and misses quadrants while POFD  
 196 only uses the false alarms and correct negatives quadrants. Because the data threshold remains  
 197 fixed for the ROC curve, the points either contribute to POD or POFD, regardless of the model  
 198 threshold designation. For a very low model threshold setting, all of the points are in either the  
 199 hits or false alarms quadrants, which sets both POD and POFD to one. As the model threshold is  
 200 increased, points are converted from hits to misses and from false alarms to correct negatives,  
 201 which monotonically decreases POD and POFD. For a very high model threshold, all of the  
 202 points will then be misses or correct negatives, and both POD and POFD will be zero.

203 Figure 1b shows the calculation scenario for the STONE curve. In this situation, both  
 204 event identification thresholds move simultaneously. The four quadrants are still defined as with  
 205 the ROC curve, but with both thresholds changing, the shift of points from one quadrant to  
 206 another is not so simple. For a very low threshold setting, nearly all points will be hits and  
 207 perhaps a few will be false alarms. Thus, like the ROC curve, the STONE curve also begins in  
 208 the (1,1) corner of POFD-POD space (assuming a “low” starting threshold value). Also similarly,  
 209 for a very large threshold setting, nearly all points will be correct negatives and perhaps a few  
 210 will be misses, with the STONE curve ending in the (0,0) corner of POFD-POD space. Another  
 211 similarity is that false alarms are converted into correct negatives as the threshold setting  
 212 increases.

213 The choice of POD and POFD for the STONE curve calculation is purely for  
 214 convenience and direct comparison with the ROC curve. Because the STONE curve assumes that  
 215 the data are continuous rather than categorical, with event status being defined by a sweeping  
 216 threshold, other metrics could have been chosen. For example, rather than basing its calculation  
 217 on discrimination, the equations for reliability could have been used. These are defined using the  
 218 quadrants on either side of the model threshold rather than either side of the data event threshold.  
 219 In this case, metrics such as the correct alarm ratio and the miss ratio would provide similar  
 220 information to the proposed usage of POD and POFD.

221 The big difference between the ROC and STONE curve calculations, however, is that as  
 222 the event identification threshold increases, a hit event can shift to any of the other three  
 223 quadrants. If it is far above the data threshold but close to the model threshold, then the threshold  
 224 increase will cause the point to shift from being a hit to a miss. If it is close to the data threshold  
 225 but far away from the model threshold, then it will shift from being a hit to being a false alarm. If

226 it is close to both thresholds, then there is a chance it will cross both lines during the incremental  
227 shift and jump from the hits regions to the correct negatives zone. Only the first of these three  
228 moves (hits to misses) occurs with the ROC curve calculation. In addition, misses are shifting to  
229 become correct negatives as the observational threshold is incremented to higher values, another  
230 move that is not part of the ROC curve calculation. The behavior of the POD and POFD values  
231 as a function of threshold, therefore, are not intuitively known and the STONE curve does not  
232 have to be monotonic between its (1,1) and (0,0) endpoints.

233 There are case where the ROC and STONE curves require special handling. For the  
234 limiting case of no observed events, then POD is undefined (zero divided by zero). The converse  
235 – when there are no non-events in the observation set – leaves POFD undefined. When  
236 calculating a ROC curve it is necessary, therefore, to check the comparison interval to ensure  
237 that both events and non-events occur in the data. As long as there is at least one data value in  
238 each of these event status categories, then the ROC curve will monotonically vary from (0,0) to  
239 (1,1). For the STONE curve, the issue is slightly different. If the event identification threshold,  
240 the same for both data and model, is set to a value below the smallest value, then all points are in  
241 the hits quadrant (POD = 1) but POFD is undefined. In this extreme threshold choice, POFD  
242 should be set to 1. Similarly, if the threshold is above the highest value, then all points are correct  
243 negatives (POFD = 0) and POD is undefined. The remedy is to set POD to zero when the  
244 threshold is swept beyond the end of the values. These corrections yield a STONE curve that  
245 extends to the two corners of (0,0) and (1,1), like the ROC curve.

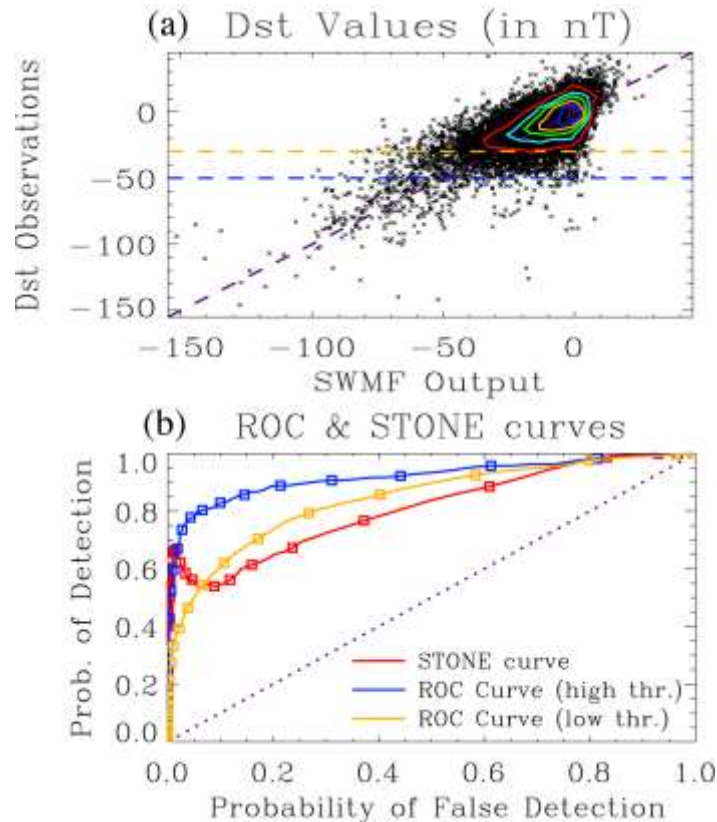
### 246 **3. Application of the STONE tool**

247 With this definition for the STONE curve, it can be used on a few example data-model  
248 comparisons to illustrate the similarities and differences with the ROC curve. Here, two  
249 comparisons will be shown. The first is for a model prediction of a geomagnetic activity index,  
250 originally presented by Liemohn, Ganushkina, et al. (2018), and the second is for energetic  
251 electrons in near-Earth space, originally presented by Ganushkina et al. (2019).

#### 252 **3.1. Predicting a geomagnetic activity index**

253 Liemohn, Ganushkina, et al. (2018) compared the output from experimental real-time  
254 simulations of the Space Weather Modeling Framework (SWMF) against the disturbance storm-  
255 time index, Dst (Rostoker et al., 1972). The SWMF is a collection of space physics numerical  
256 models simulating the Sun-Earth space environment (Toth et al., 2012), and in many other  
257 planetary environments (e.g., Jia et al., 2012; Ma et al., 2013; Dong et al., 2014; Liemohn et al.,  
258 2017). This geospace environment simulation has a very similar setup to that of Pulkkinen et al.  
259 (2013), using the Block Adaptive Tree Roe-type Upwind Scheme (BATS-R-US)  
260 magnetohydrodynamic model coupled to the Rice Convection Model (RCM) and the Ridley  
261 Ionosphere Model (RIM). Real-time solar wind and interplanetary magnetic field input was  
262 taken from the Advanced Composition Explorer (ACE) satellite. The simulated Dst time series  
263 from the SWMF was calculated with the method from Yu et al. (2010) and compared against the  
264 real-time version of the Dst index as produced by the Kyoto World Data Center for  
265 Geomagnetism. The interval of comparison spans from 19 April 2015 until 17 July 2017, which  
266 is 27 months of 1-hour resolution measurements and corresponding model output values (just  
267 under 300,000 data-model pairs).

268 Figure 2a shows a scatter plot of the SWMF Dst values against the observed Dst values.  
 269 While the individual points are analyzed as unique contributions, they are binned to produce the  
 270 colored curves on the plot, demarking contours of 50 points within a 5-by-5 nT grid. Note that,  
 271 because Dst is near zero for quiet times and shifts to negative values during storm times, events  
 272 are defined as values below (i.e., more negative) a chosen threshold. As defined by Gonzalez et  
 273 al. (1994), a typical designation for the Dst index measuring a storm situation is -30 nT or below  
 274 for a weak storm and -50 nT or below for a moderate storm, so these two settings are used for the  
 275 ROC curve observational threshold setting. These two thresholds are indicated in Figure 2a as  
 276 horizontal dashed lines.  
 277



278

279 **Figure 2.** (a) Scatter plot of the observed real-time Dst time series (y-axis values) against a  
 280 prediction Dst time series from the SWMF (x-axis values). The contours are drawn every 50  
 281 points per 5x5 nT bin. Also drawn are horizontal dashed lines at the ROC event thresholds of -30  
 282 and -50 nT, with events defined as the points below these lines. A purple dashed zero intercept  
 283 unity slope line is also drawn, for reference. (b) STONE (red) and ROC curves (blue for -50 nT,  
 284 orange for -30 nT observed event threshold) calculated from the scatter plot. Symbols are shown  
 285 along all three curves at every 5 nT threshold increment. The diagonal dotted line with zero  
 286 intercept and unity slope is shown for reference.

287

288 The ROC and STONE curves are calculated as follows and shown in Figure 2b. To create  
 289 a ROC curve, the model threshold setting is initially set to +10 nT and then swept in 1 nT  
 290 increments to -120 nT. The data threshold for events is held fixed, at -50 nT for the blue curve



291 and -30 for the orange curve. To create STONE curve (red line), this same model threshold  
292 variation is followed, but the data threshold is also swept from +10 to -120 nT. Symbols are  
293 shown along each of the plots every 5 nT of threshold increment.

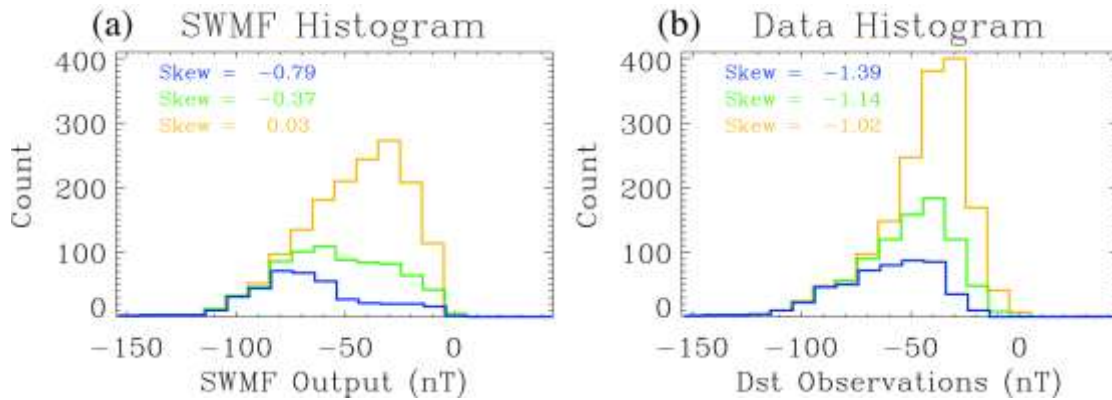
294 Some features of Figure 2b should be noted. It is seen that the ROC curves monotonically  
295 increase from (0,0) to (1,1). The ROC curve with a -50 nT event threshold is well above the  
296 zero-intercept, unity-slope line (the diagonal purple dotted line on Figure 2b), indicating that the  
297 model is reasonably good at reproducing moderate and stronger storm events recorded by the  
298 real-time Dst index. The closest approach to the upper left corner occurs at a threshold of -37 nT  
299 for the -50 nT threshold ROC curve and -17 nT for the -30 nT ROC curve, which indicates that  
300 the model somewhat underpredicts the strength of such storms.

301 The STONE curve lies both above and below these two ROC curves, depending on the  
302 threshold. The STONE curve is coincident with each ROC curve at the locations where the ROC  
303 curve model threshold setting is equal to the observational threshold setting (-30 nT for the  
304 orange curve, -50 nT for the blue curve). They cross elsewhere, too, such as in the low-threshold  
305 (i.e., a threshold of near and above zero) region in the upper right region of the plot. It is seen  
306 that the STONE curve is not monotonic but includes a local maximum and local minimum at the  
307 “high threshold” settings (minimum at -28 nT threshold and maximum at -52 nT threshold). The  
308 nonmonotonicity is because POD increases at these threshold values. An increase in POD is  
309 achieved by more points leaving the misses quadrant than leaving from the hits quadrant.

310 This is better understood by considering the distribution of points beyond a few threshold  
311 choices. Figure 3 shows histograms of the points above a particular data or model threshold  
312 setting. In particular, three threshold settings are displayed – -30 nT, -40 nT, and -50 nT –  
313 showing the points at “higher” (more negative) Dst values in both the data and model (left and  
314 right columns, respectively). For Figure 3a, the counts are for all points below some horizontal  
315 line of an event identification threshold setting of the observations. For Figure 3b, the counts are  
316 for all points to the left of some event identification threshold setting for the model values. The  
317 calculated skew for these distributions is listed in each panel.

318 In Figure 3a, it is evident, both qualitatively from the histograms and quantitatively from  
319 the skew values, that the distribution of model output values is significantly changing across  
320 these three observational threshold settings. For the more negative threshold, there are far fewer  
321 model values between zero and -50 nT. That is, across these threshold settings, many of the  
322 points in the misses quadrant were converted into correct negatives. In Figure 3b, the three  
323 distributions have essentially the same shape, with a large negative skew. These distributions do  
324 not undergo the same systematic alteration in their shape the way that the distributions in Figure  
325 3a did. Putting these two features together, it means that more misses were removed than hits,  
326 and so POD increased as the STONE threshold was swept to more negative values between -30  
327 nT and -50 nT. This resulted in a nonmonotonic wiggle in the STONE curve at these thresholds.

328

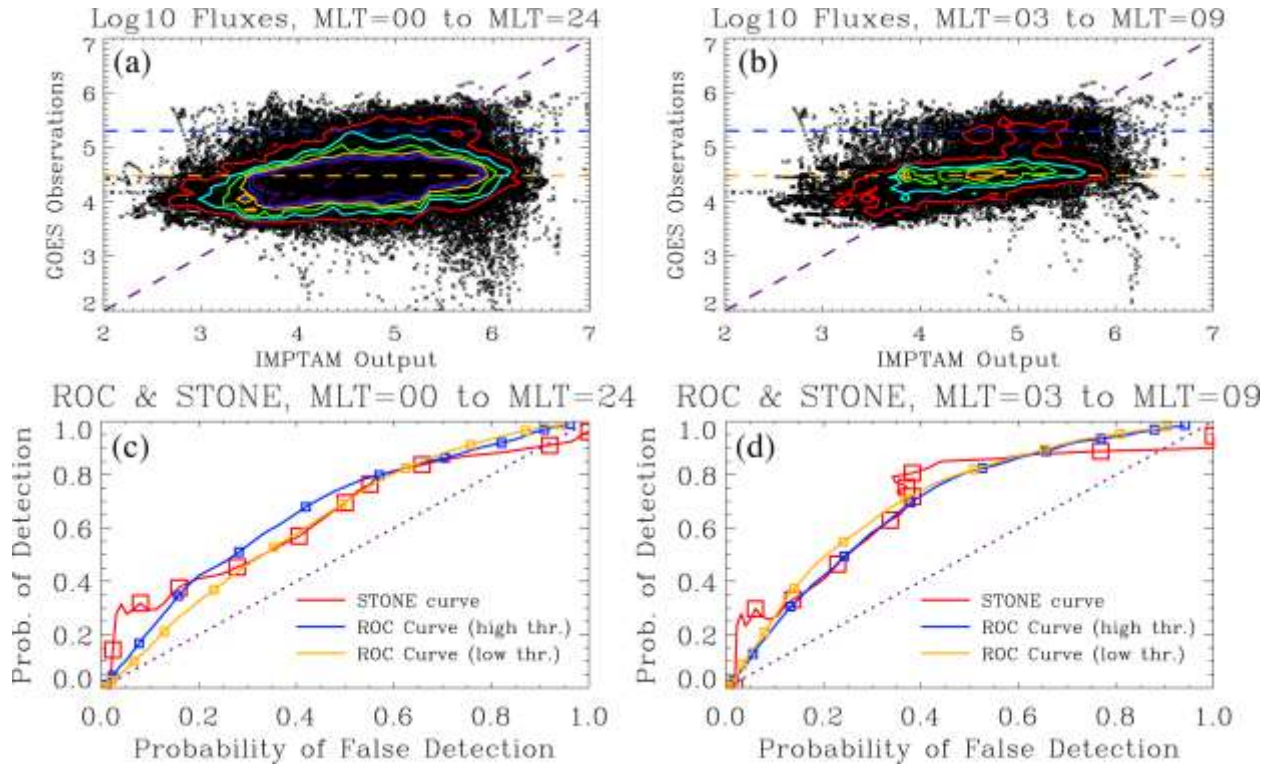


**Figure 3.** (a) Histogram of model values for all data values below three different thresholds: -30 nT (orange curve), -40 nT (green curve), and -50 nT (blue curve). (b) Histogram of data values for all model values below the same three thresholds. The bin sizes for each histogram is 10 nT. The calculated skew for each distribution is listed in each plot.

### 3.2. Predicting energetic electrons in near-Earth space

Ganushkina et al. (2019) compared real-time output from the inner magnetosphere particle transport and acceleration model (IMPTAM) with measurements from the magnetosphere electron detector (MAGED) on the geosynchronous orbiting environmental satellites (GOES) in geostationary orbit at 6.62 Earth radii geocentric distance over the American sector (Rowland & Weigel, 2012; Sillanpaa et al., 2017), specifically, with data from GOES-13, -14, and -15. IMPTAM, initially developed by Ganushkina et al. (2001) and used regularly for investigating the physics of plasma sheet electron transport (e.g., Ganushkina et al., 2013, 2014), has been running in a real-time operational mode since February 2013, first in Europe and then a mirror site at the University of Michigan. Ganushkina et al. (2015) made an initial comparison of these model output values against a few months of GOES data, while Ganushkina et al. (2019) provided a far more robust validation analysis of the model, covering over 18 months (September 20, 2013 through March 31, 2015). It is this second interval that will be used again for this study.

Figures 4a and 4b show two scatter plots comparing the IMPTAM and GOES electron differential number fluxes at 40 keV. The colored contours show the point density, with a new curve every 50 points within a bin (defined, for these contours, with 10 bins per decade in both the data and model values). Figure 4a presents the full data set while Figure 4b only shows the comparison for those values in the 03 to 09 magnetic local time (MLT) range, the region found by Ganushkina et al. (2019) to have a “good comparison” between the data and model values. On each of these plots, two observational event thresholds are shown as the horizontal dashed lines, drawn at  $5 \times 10^4$  and  $2 \times 10^5$  electrons  $\text{cm}^{-2} \text{s}^{-1} \text{sr}^{-1} \text{keV}^{-1}$  in green and blue, respectively.



358

359

360

361

362

363

364

365

366

367

368

369

370

371

372

373

374

375

376

377

378

379

380

381

382

**Figure 4.** Scatter plot comparing GOES and IMPTAM 40 keV electron differential number fluxes (log base 10 of electrons  $\text{cm}^{-2} \text{s}^{-1} \text{sr}^{-1} \text{keV}^{-1}$ ) for (a) all MLTs and (b) the 03-09 MLT range. Color contours are shown every 50 points per bin (10 bins per decade in both data and model). The horizontal dashed lines show the ROC thresholds of  $3 \times 10^4$  and  $2 \times 10^5$ . A purple dashed zero intercept and unity slope line is shown for reference. The lower panels show STONE curves (red) and ROC curves (blue for  $2 \times 10^5$  and orange for  $3 \times 10^4$ ) for (c) the full MLT comparison and (d) the 03-09 MLT range. Symbols are shown every factor of 2 increase in threshold value. The diagonal dotted line with zero intercept and unity slope is shown for reference.

Figures 4c and 4d show the ROC and STONE curves for these two data-model comparisons, the full set with values at all MLTs and the subset from 03 to 09 MLT, respectively. In both Figures 4c and 4d, the STONE curve again has a nonmonotonic shape at high threshold settings (above  $4 \times 10^5$ ). Like the similar case for the Dst STONE curve in Figure 2b, this shows that, for these thresholds, more points are being removed from the misses quadrant than being removed from the hits quadrant.

Figure 4d has another unusual feature in the STONE curve, seen as a nonmonotonicity in the x-axis values. This is from the POFD values increasing with increasing threshold (rather than decreasing, as they always do with a ROC curve). This is occurring for thresholds between  $1 \times 10^4$  and  $4 \times 10^4$ , just as the STONE curve crosses the orange ROC curve. Considering equation (2) above, the correct negatives in the denominator are always increasing with increasing threshold, as points convert to this quadrant from any of the other three quadrants. For POFD to increase, the false alarms had to increase faster than the correct negative point count. This is seen in Figure 4b as the points have a horizontal peak (highlighted by the flat, elongated color contours). Many

383 points are being converted from the hits quadrant into the false alarms quadrant and, for these  
384 threshold settings, this conversion to false alarms outpaces the conversion of points into the  
385 correct negatives quadrant. This results in a ripple in the STONE curve at these thresholds.

386 Figures 4c and 4d, the STONE curve is quite close to the two ROC curves, which are  
387 very similar to each other. This can be understood from the “flatness” of the cloud of points in  
388 the scatter plots in Figures 4a and 4b. The points are not well aligned with the zero intercept and  
389 unity slope line, revealing less than perfect agreement between the observations and model  
390 output. However, in terms of physics-based real-time modeling of near-Earth magnetospheric  
391 electron fluxes, this is actually quite good, arguably the best that is currently available. This  
392 means that all ROC curves will be close to each other, as any observational event identification  
393 threshold will have a relatively similar transfer of points between the quadrants. However,  
394 because the model is trying to exactly reproduce the observed flux values, the STONE curve can  
395 be calculated, and this new curve includes several nonmonotonicities. The wiggles and ripples in  
396 the STONE curve reveal thresholds where the distribution of points, in either the vertical or  
397 horizontal direction, are asymmetric, bi-modal, or otherwise non-Gaussian. The ROC curves  
398 cannot reveal this kind of information about the distribution of points in the scatter plot the way  
399 that the STONE curve can.

#### 400 **4. Discussion**

401 The STONE curve introduced above is a new tool for assessing the ability of a model  
402 with a continuous-valued output to exactly match a continuous-valued data set. As illustrative  
403 example usages, it was applied to two recently-published data-model comparisons, a prediction  
404 of the disturbance storm-time index Dst and a prediction of energetic electron fluxes in near-  
405 Earth space.

406 The STONE curve is quite similar to the ROC curve. It is based on the same contingency  
407 table calculations of POD and POFD, plotting these two values against each other for a range of  
408 event threshold settings. Like the ROC curve, it starts at (1,1) for low threshold settings and  
409 moves to (0,0) for high threshold settings. Also like the ROC curve, being above the zero-  
410 intercept, unity-slope line indicates a prediction that is better than random chance. Curves are  
411 better when they are closer to the upper left corner in POFD-POD space, and a common choice  
412 for the best optimization point along a ROC or STONE curve is that closest to this corner as this  
413 point reveals the best model threshold setting for optimizing discrimination performance. That is,  
414 both curves reveal a possible best model threshold setting for event prediction, the ROC curve  
415 revealing the best settings for a specified observational event identification threshold and the  
416 STONE curve revealing the best setting against the an identically defined observational event. Of  
417 course, this is “best” only if discrimination is what should be optimized for the particular  
418 application. A different threshold settings might be most favorable if other considerations  
419 outweigh discrimination, such as minimizing false alarms or maximizing a particular skill score.  
420 Focusing on user needs during the validation of a model is a foundational element of Application  
421 Usability Levels (Halford et al., 2019) and should always be considered when assessing a  
422 model’s performance. If the user is most concerned about optimizing one of these other features  
423 of the data-model comparison for their particular decision-making needs, then a model event  
424 identification threshold should be chosen that best addresses that need.

425 Another similarity is that the integral area of the ROC curve, AUC, is equally applicable  
426 to the STONE curve. AUC, a synthesis of the entire threshold-setting range into a single number,

427 indicates the quality of the chosen model to predict the events identified in the observational data  
428 (see the detailed explanation of AUC in Fawcett (2006) or Ekelund (2011)). Being an integrated  
429 quantity, AUC is a complementary metric to the “best” model threshold setting for event  
430 prediction mentioned in the preceding paragraph because AUC uses information from all model  
431 threshold settings, even those with POFD-POD coordinates far from the “best setting” upper-left  
432 corner of the graph. Comparing AUCs for several STONE curves (i.e., using different models  
433 against the same data set) will provide a quantitative assessment of which of the models has the  
434 best system-level predictive capability against that data set. It could be that the model with the  
435 highest AUC is not the model with a point along its STONE curve closest to (0,1) in POFD-POD  
436 space. Such a case reveals that the first model, with the higher AUC, has the best model physics  
437 for reproducing the data set as a whole, but that the second model is actually best at predicting  
438 events with a particular threshold setting. Because it is calculated the same way, AUC can be  
439 used to compare STONE curves just like it is for ROC curves.

440 A key difference between the STONE and ROC curves is that the STONE curve can have  
441 nonmonotonicities. These features, which can be wiggles with respect to either POD or POFD,  
442 reveal features of the model prediction of events that are not easily extracted from a ROC curve.  
443 This makes the STONE curve somewhat like a fit performance metric, even though it is an  
444 event-detection metric that disregards the difference between the data-model pairs.

445 The nonmonotonicities in the STONE curve reveal information about the distribution of  
446 points in the data-model comparison. Specifically, they show the existence of an asymmetry,  
447 perhaps a non-Gaussian point spread like a skewed or bimodal distribution, for the pairs above  
448 that threshold setting. The nonmonotonicities might also prove useful in assessing changes in  
449 model bias; if the model is unbiased in part of the value range but biased in another, this shift  
450 relative to the observations could be revealed in the nonmonotonic features of the STONE curve.  
451 Combined with a histogram or even fit-performance data-model comparison formulas for this  
452 subset of either the data or model values, the nature of this distribution can be explored.

453 Why not just start out by calculating fit performance metrics on these subsets? The  
454 answer is because the subset of interest would not have been known; the STONE curve revealed  
455 the thresholds where the distribution had a changing or non-Gaussian distribution. That is, it  
456 could be used to optimize the fit performance analysis by identifying the subset of the data or  
457 model that should be considered in more detail. Also, the STONE curve includes information not  
458 just within a subset of the data (discrimination) or a subset of the model (reliability), but includes  
459 information about the entire data-model comparison set, because POD and POFD use all data-  
460 model pairs in the point counting in the quadrants. For one of the specific examples in the  
461 manuscript: continuous metrics will tell the user very little about the SWMF’s ability to predict  
462 magnetic storms of -50 nT or less. A ROC curve is far more suited to this, and a STONE curve  
463 one step farther, revealing the ability of the model to predict Dst levels below any threshold  
464 (which could be accomplished by a large family of ROC curves). No continuous metric that does  
465 this type of assessment. If the detection of events is desired, then the STONE curve is an  
466 advantageous assessment tool in addition to standard continuous fit performance metrics.

467 A useful follow-on study to this would be a detailed analysis of the features of the  
468 STONE curve to the underlying distribution of points in the data-model scatter plot. That is, by  
469 assuming known two-dimensional distributions of points of several different shapes and  
470 parameter settings, the connection between the distribution and the resulting features in the  
471 STONE curve can be isolated. The STONE curve features could also be put into perspective

472 relative to other event detection and fit performance metrics. Such an in-depth assessment of the  
473 STONE curve is beyond this initial description and illustrative usage of this metrics tool but is  
474 planned as a future project.

475 A key feature of the STONE curve is that it reveals the threshold (or range of thresholds)  
476 for which the model does well at reproducing similarly-defined events in the data. A single ROC  
477 curve cannot do this because it uses a fixed threshold for identifying events in the observations.  
478 When the data are continuously-varying values and the model is seeking to reproduce these exact  
479 values, then it is useful to examine the event detection capability of the model at the same  
480 threshold settings between data and model. A single ROC curve doesn't do this, except at one  
481 threshold setting. The STONE curve, therefore, is a better assessment tool for models that are  
482 trying to predict the exact value of a data set.

483 The ROC is still a highly useful tool for event prediction and this study does not seek to  
484 replace it with the STONE curve. Indeed, the ROC curve is optimal for categorical data sets  
485 where the observations have been pre-classified as events and non-events. In this case, the  
486 STONE curve cannot be used because the data and model are on different scales, the former  
487 being a binary yes-no designation and the latter being either a real number range or its own  
488 categorical designation. The ROC curve can handle this difference in units while the STONE  
489 curve cannot.

490 The two example data-model comparisons to which the STONE curve was applied are  
491 both from space physics. The first was an evaluation between a physics-based model of  
492 geospace, running in real time, with the real-time version of the Dst index, a measure of  
493 geospace activity (see its comparison with other similar indices in Katus & Liemohn (2013)).  
494 Many models exist for the prediction of Dst (see the review by Liemohn, McCollough, et al.,  
495 2018), with some models doing exceptionally well at reproducing the observed time series.  
496 While this chosen model for this comparison is arguably the best physics-based model for  
497 reproducing Dst (see, for comparison, the solar cycle storm-interval Dst comparison of Liemohn  
498 & Jakowski (2008)), it is not the best model available at predicting this index. In fact, many  
499 empirical models are substantially better at capturing the storm intervals of Dst. The second  
500 example was a comparison of a physics-based model of energetic electron fluxes in the near-  
501 Earth magnetosphere, running in real time, with real-time observations from a geosynchronous  
502 spacecraft. Magnetospheric charged particle fluxes are notoriously difficult to reproduce with  
503 physics-based modeling approaches (see, e.g., Morley et al., 2018), and even empirical models  
504 reduce the problem to remove the fast temporal dynamics, averaging over a day (e.g., Li, 2004)  
505 or an hour (e.g., Boynton et al., 2019). That is, these two examples represent state-of-the-art  
506 physics-based approaches to space weather nowcasting, but are not the best predictions of these  
507 two quantities across the field.

508 It is worth stating here that there are many other metrics in existence for evaluating a  
509 scatter plot of data-model values like that shown in Figure 1. No one metric equation or  
510 technique does everything; each was designed to assess only a specific aspect of the relationship.  
511 That is, neither the ROC curve nor the STONE curve should be used as the sole assessment tool  
512 for a model against a particular data set. In practice, many metrics, from both the continuous fit-  
513 performance grouping and from the categorical event-detection grouping, should be applied to  
514 examine the quality of the model from a number of perspectives.

515 It should be mentioned that this is not the first application of sliding both the  
516 observational and model event identification threshold. As one example of this, in their  
517 presentation and initial usage of the extreme dependency score (EDS), Stephenson et al. (2008)  
518 simultaneously moved both thresholds. Events become rarer with increasing threshold and that  
519 study examined the relationship of EDS as a function of this rarity – moving both thresholds  
520 together, as is done here for the STONE curve.

521 A final note to make here is that this is not the first usage of the STONE curve. Both  
522 Liemohn, Ganushkina, et al. (2018) and Liemohn, McCollough, et al. (2018) used STONE  
523 curves in the plots labeled as ROC curves. It is clear that these panels are mislabeled because  
524 nonmonotonicities are seen in these lines.

## 525 5. Conclusions

526 A new data-model comparison assessment tool has been introduced, described, used, and  
527 interpreted – the sliding threshold of observations for numeric evaluation curve. Based on the  
528 relative operating characteristic curve, the STONE curve is created by plotting POD against  
529 POFD for a wide range of threshold settings. The main difference with the ROC curve is that the  
530 STONE curve requires the data to be continuous-valued real numbers and the model to be  
531 attempting to reproduce these exact values. The threshold is moved not only for the model, as is  
532 done for the ROC curve, but also for the observational event identification threshold setting,  
533 which is moved simultaneously with the model threshold setting.

534 The STONE curve has many features in common with the ROC curve with one large  
535 exception – it can have nonmonotonicities in both the POD and POFD values. For the ROC  
536 curve, the points shift within the quadrants defining POD or within the quadrants used to define  
537 POFD, but not between these two mutually exclusive regions. The ROC curve is, therefore,  
538 always monotonic, sweeping from (1,1) to (0,0) in POFD-POD space. For the STONE curve, the  
539 motion of the observational threshold moves points from the POD regions to the POFD regions,  
540 allowing for these nonmonotonic features in the STONE curve.

541 These wiggles and ripples, however, reveal information about the underlying distribution  
542 of points in the data-model scatter plot. Specifically, if the distribution is shifted, asymmetric, or  
543 bi-modal, the STONE curve will have a nonmonotonicity. Further investigation of the  
544 distribution, through a histogram, skew calculation, or other metric assessment, can reveal the  
545 true nature of the data-model comparison for this threshold setting.

546 It is hoped that the STONE curve becomes a useful data-model comparison tool. It has  
547 been used with two space weather applications in this study but these are purely illustrative  
548 examples. A dozen studies using ROC curves across the Earth and space sciences were given in  
549 the Introduction above. Some of these studies were based on observations that were pre-  
550 classified yes/no as events or not, and so the ROC curve is the proper tool for assessing the  
551 model's ability to predict those events. Some of these studies, however, and others like them, are  
552 based on models trying to exactly predict the observed data values, in which case the STONE  
553 curve might be a useful assessment tool. For any continuous-valued model trying to reproduce  
554 the exact numbers of a continuous-valued data set, the STONE curve can be calculated, perhaps,  
555 as shown for the two examples here, revealing additional information about the data-model  
556 comparison than can be obtained from the ROC curve alone. The STONE curve is a general  
557 purpose metric for use whenever a model is trying to exactly reproduce a continuous-valued data

558 set. It can be used with both archival observations as well as for assessment of real-time  
559 nowcasting across the full breadth of science and engineering disciplines.

560

## 561 **Acknowledgments and Data**

562 The authors would like to thank the US government for sponsoring this research, in  
563 particular research grants from NASA (NNX14AC02G, NNX16AG66G, NNX17AI48G, and  
564 NNX17AB87G) and NSF (AGS-1414517). The part of the research done by M. Liemohn and N.  
565 Ganushkina received funding from the European Union Horizon 2020 Research and Innovation  
566 programme under grant agreement 637302 (PROGRESS) and 870452 (PAGER). A. Azari's  
567 contributions are based on work supported by the NSF Graduate Research Fellowship Program  
568 (DGE 1256260). The SWMF simulations were conducted on the computing facilities at NASA  
569 GSFC and the run output is freely available on their website ([https://ccmc.gsfc.nasa.gov/cgi-  
570 bin/SWMFpred.cgi](https://ccmc.gsfc.nasa.gov/cgi-bin/SWMFpred.cgi)) and the CCMC iSWA interactive tool (<https://ccmc.gsfc.nasa.gov/iswa/>).  
571 The real-time solar wind data was provided by NOAA SWPC  
572 (<http://www.swpc.noaa.gov/products/real-time-solar-wind>). The authors thank the World Data  
573 Center in Kyoto, Japan for the real-time Dst values ([http://wdc.kugi.kyoto-  
574 u.ac.jp/dst\\_realtime/presentmonth/index.html](http://wdc.kugi.kyoto-u.ac.jp/dst_realtime/presentmonth/index.html)). The IMPTAM simulations are available through  
575 the Finnish Meteorological Institute (<http://imptam.fmi.fi/>) and at the University of Michigan  
576 (<http://citrine.engin.umich.edu/imptam/>).

577 In addition to the archival repositories listed above, the specific observational data sets  
578 and the model output files used in this study are available at the University of Michigan Deep  
579 Blue Data repository,  
580 [https://deepblue.lib.umich.edu/data/concern/data\\_sets/02870v99r?locale=en](https://deepblue.lib.umich.edu/data/concern/data_sets/02870v99r?locale=en) .

581

## 582 **References**

- 583 Anagnostopoulos, G. G., Fatichi, S., & Burlando, P. (2015). An advanced process-based  
584 distributed model for the investigation of rainfall-induced landslides: The effect of  
585 process representation and boundary conditions. *Water Resources Research*, *51*, 7501–  
586 7523. <https://doi.org/10.1002/2015WR016909>
- 587 Azari, A. R., Liemohn, M. W., Jia, X., Thomsen, M. F., Mitchell, D. G., Sergis, N., et al. (2018).  
588 Interchange injections at Saturn: Statistical survey of energetic H<sup>+</sup> sudden flux  
589 intensifications. *Journal of Geophysical Research: Space Physics*, *123*, 4692– 4711.  
590 <https://doi.org/10.1029/2018JA025391>
- 591 Bobra, M. G., & Couvidat, S. (2015). Solar flare prediction using SDO/HMI vector magnetic  
592 field data with a machine learning algorithm. *The Astrophysical Journal*, *789*, 135.  
593 <http://dx.doi.org/10.1088/0004-637X/798/2/135>
- 594 Borah, N., Sahai, A. K., Chattopadhyay, R., Joseph, S., Abhilash, S., & Goswami, B. N. (2013).  
595 A self-organizing map–based ensemble forecast system for extended range prediction of  
596 active/break cycles of Indian summer monsoon. *Journal of Geophysical Research -  
597 Atmospheres*, *118*, 9022– 9034. <https://doi.org/10.1002/jgrd.50688>

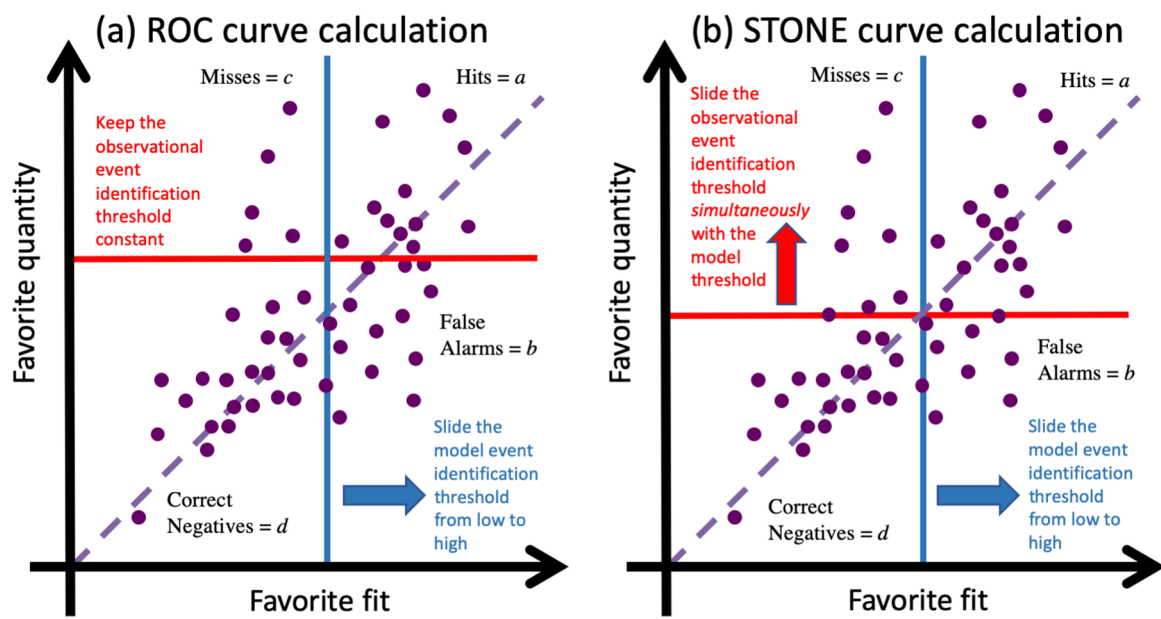


- 598 Boynton, R. J., Amariutei, O. A., Shprits, Y. Y., & Balikhin, M. A. (2019). The system science  
599 development of local time-dependent 40-keV electron flux models for geostationary  
600 orbit. *Space Weather*, 17, 894–906. <https://doi.org/10.1029/2018SW002128>
- 601 Brenning, A., Grasser, M., & Friend, D. A. (2007). Statistical estimation and generalized  
602 additive modeling of rock glacier distribution in the San Juan Mountains, Colorado,  
603 United States. *Journal of Geophysical Research – Solid Earth*, 112, F02S15.  
604 <https://doi.org/10.1029/2006JF000528>
- 605 Carter, J. V., J. Pan, S. N. Rai, & S. Galandiuk (2016). ROC-ing along: Evaluation and  
606 interpretation of receiver operating characteristic curves. *Surgery*, 159(6), 1638-1645,  
607 <https://doi.org/10.1016/j.surg.2015.12.029>
- 608 Chen, Y., Manchester, W. B., Hero, A. O., Toth, G., Dufumier, B., Zhou, T., et al. (2019).  
609 Identifying solar flare precursors using time series of SDO/HMI Images and SHARP  
610 Parameters. *Space Weather*, 17, 1404–1426. <https://doi.org/10.1029/2019SW002214>
- 611 Delle Monache, L., Hacker, J. P., Zhou, Y., Deng, X., & Stull, R. B. (2006). Probabilistic  
612 aspects of meteorological and ozone regional ensemble forecasts. *Journal of Geophysical  
613 Research - Atmospheres*, 111, D24307. <https://doi.org/10.1029/2005JD006917>
- 614 Dong, C., S. W. Bougher, Y. Ma, G. Toth, A. F. Nagy, & Najib, D. (2014). Solar wind  
615 interaction with Mars upper atmosphere: Results from the one-way coupling between the  
616 multifluid MHD model and the MTGCM model. *Geophysical Research Letters*, 41,  
617 2708–2715. <https://doi.org/10.1002/2014GL059515>
- 618 Ekelund, S. (2011). ROC Curves – What are They and How are They Used?, *Point of Care*,  
619 11(1), 16-21. <https://doi.org/10.1097/POC.0b013e318246a642>
- 620 Fawcett, T. (2006). An introduction to ROC analysis. *Pattern Recognition Letters*, 27, 861-874.  
621 <https://doi.org/10.1016/j.patrec.2005.10.010>
- 622 Gabriel, P., Barker, H. W., O'Brien, D., Ferlay, N., & Stephens, G. L. (2009). Statistical  
623 approaches to error identification for plane-parallel retrievals of optical and  
624 microphysical properties of three-dimensional clouds: Bayesian inference. *Journal of  
625 Geophysical Research - Atmospheres*, 114, D06207.  
626 <https://doi.org/10.1029/2008JD011005>
- 627 Ganushkina N. Yu., T. I. Pulkkinen, V. F. Bashkurov, D. N. Baker, & X. Li (2001). Formation of  
628 intense nose structures. *Geophysical Research Letters*, 28, 491-494.
- 629 Ganushkina, N. Y., Amariutei, O. A., Shprits, Y. Y., & Liemohn, M. W. (2013). Transport of the  
630 plasma sheet electrons to the geostationary distances. *Journal of Geophysical Research:  
631 Space Physics*, 118 (1), 82-98. <https://doi.org/10.1029/2012JA017923>
- 632 Ganushkina, N. Y., Liemohn, M. W., Amariutei, O. A., & Pitchford, D. (2014). Low-energy  
633 electrons (550 keV) in the inner magnetosphere. *Journal of Geophysical Research: Space  
634 Physics*, 119 (1), 246-259. <https://doi.org/10.1002/2013JA019304>
- 635 Ganushkina, N. Y., Amariutei, O. A., Welling, D., & Heynderickx, D. (2015). Nowcast model  
636 for low-energy electrons in the inner magnetosphere. *Space Weather*, 13 (1), 16-34.  
637 <https://doi.org/10.1002/2014SW001098>

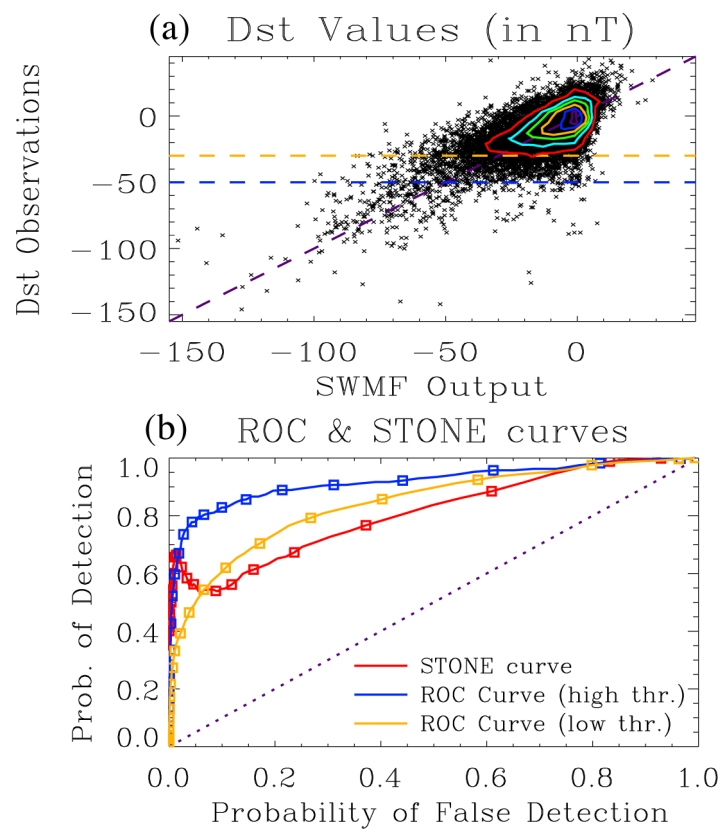
- 638 Ganushkina, N. Yu., Liemohn, M. W., & Dubyagin, S. (2018), Current systems in the Earth's  
639 magnetosphere, *Reviews of Geophysics*, 56(2), 309-332. [https://doi.org/  
640 10.1002/2017RG000590](https://doi.org/10.1002/2017RG000590)
- 641 Ganushkina, N. Y., Sillanpaa, I., Welling, D. T., Haiducek, J., Liemohn, M. W., Dubyagin, S., &  
642 Rodriguez, J., (2019). Validation of Inner Magnetosphere Particle Transport and  
643 Acceleration Model (IMPTAM) on the long-term GOES MAGED measurements of keV  
644 electron fluxes at geostationary orbit. *Space Weather*, 17, 687-708.  
645 <https://doi.org/10.1029/2018SW002028>
- 646 Halford, A., Kellerman, A., Garcia-Sage, K., Klenzing, J., Carter, B., McGranaghan, R., Guild,  
647 T., Cid, C., Henney, C., Ganushkina, N., Burrell, A., Terkildsen, M., Thompson, B. J.,  
648 Pulkkinen, A., McCollough, J., Murray, S., Leka, K. D., Fung, S., Bingham, S., Walsh,  
649 B., Liemohn, M., Bisi, M., Morley, S., & Welling, D. (2019), Application Usability  
650 Levels: A framework for tracking project product progress, *Journal of Space Weather  
651 and Space Climate*, 9, A34. <https://doi.org/10.1051/swsc/2019030>
- 652 Hogan, R. J. & Mason, I. B. (2012). Deterministic Forecasts of Binary Events. In *Forecast  
653 Verification* (eds I. T. Jolliffe and D. B. Stephenson). doi:[10.1002/9781119960003.ch3](https://doi.org/10.1002/9781119960003.ch3)
- 654 Ippolito, A., Scotto, C., Sabbagh, D., Sgrigna, V., & Maher, P. (2016). A procedure for the  
655 reliability improvement of the oblique ionograms automatic scaling algorithm. *Radio  
656 Science*, 51, 454– 460. <https://doi.org/10.1002/2015RS005919>
- 657 Jia, X., Hansen, K. C., Gombosi, T. I., Kivelson, M. G., Toth, G., DeZeeuw, D. L., & Ridley, A.  
658 J., (2012). Magnetospheric configuration and dynamics of Saturn's magnetosphere: A  
659 global MHD simulation. *Journal of Geophysical Research Space Physics*, 117, A05225.  
660 <https://doi.org/10.1029/2012JA017575>
- 661 Jolliffe, I.T., & Stephenson, D.B. (2012). *Forecast verification: A practitioner's guide  
662 in atmospheric science*. Wiley-Blackwell, Hoboken, NJ.
- 663 Katus, R. M., & Liemohn, M. W. (2013). Similarities and differences in low-to-mid-latitude  
664 geomagnetic indices. *Journal of Geophysical Research*, 118, 5149-5156.  
665 <https://doi.org/10.1002/jgra.50501>.
- 666 Li, X. (2004). Variations of 0.7–6.0 MeV electrons at geosynchronous orbit as a function of solar  
667 wind. *Space Weather*, 2, S03006. <https://doi.org/10.1029/2003SW000017>
- 668 Liemohn, M. W., & Jazowski, M. (2008). Ring current simulations of the 90 intense storms  
669 during solar cycle 23. *Journal of Geophysical Research*, 113, A00A17.  
670 <https://doi.org/10.1029/2008JA013466>
- 671 Liemohn, M. W., Xu, S., Dong, C., Bougher, S. W., Johnson, B. C., Ilie, R., & De Zeeuw, D. L.,  
672 (2017). Ionospheric control of the dawn-dusk asymmetry of the Mars magnetotail current  
673 sheet. *Journal of Geophysical Research Space Physics*, 122, 6397-6414.  
674 <https://doi.org/10.1002/2016JA023707>
- 675 Liemohn, M. W., Ganushkina, N. Y., De Zeeuw, D. L., Rastaetter, L., Kuznetsova, M., Welling,  
676 D. T., Toth, G., Ilie, R., Gombosi, T. I., & van der Holst, B. (2018). Real-time SWMF  
677 and CCMC: assessing the Dst output from continuous operational simulations. *Space  
678 Weather*, 16, 1583-1603. <https://doi.org/10.1029/2018SW001953>.

- 679 Liemohn, M. W., McCollough, J. P., Jordanova, V. K., Ngwira, C. M., Morley, S. K., Cid, C.,  
 680 Tobiska, W. K., Wintoft, P., Ganushkina, N. Y., Welling, D. T., Bingham, S., Balikhin,  
 681 M. A., Opgenoorth, H. J., Engel, M. A., Weigel, R. S., Singer, H. J., Buresova, D.,  
 682 Bruinsma, S., Zhelavskaya, I., Shprits, Y. Y., & Vasile, R. (2018). Model evaluation  
 683 guidelines for geomagnetic index predictions. *Space Weather*, 16, 2079–2102.  
 684 <https://doi.org/10.1029/2018SW002067>
- 685 Ma, Y. J., Nagy, A. F., Russell, C. T., Strangeway, R. J., Wei, H. Y., & Toth, G. (2013). A  
 686 global multispecies single-fluid MHD study of the plasma interaction around Venus.  
 687 *Journal of Geophysical Research Space Physics*, 118, 321–330.  
 688 <https://doi.org/10.1029/2012JA018265>
- 689 Manzato, A. (2005). An odds ratio parameterization for ROC diagram and skill score indices.  
 690 *Weather and Forecasting*, 20, 918–930. <https://doi.org/10.1175/WAF899.1>
- 691 Manzato, A. (2007). A note on the maximum Peirce skill score. *Weather and Forecasting*, 22,  
 692 1148–1154. <https://doi.org/10.1175/WAF1041.1>
- 693 Mathieu, J. A., & Aires, F. (2018). Using neural network classifier approach for statistically  
 694 forecasting extreme corn yield losses in Eastern United States. *Earth and Space Science*,  
 695 5, 622–639. <https://doi.org/10.1029/2017EA000343>
- 696 Meade, B. J., DeVries, P. M. R., Faller, J., Viegas, F., & Wattenberg, M. (2017). What is better  
 697 than Coulomb failure stress? A ranking of scalar static stress triggering mechanisms from  
 698  $10^5$  mainshock-aftershock pairs. *Geophysical Research Letters*, 44, 11,409–11,416.  
 699 <https://doi.org/10.1002/2017GL075875>
- 700 Morley, S. K., Brito, T. V., & Welling, D. T. (2018). Measures of model performance based on  
 701 the log accuracy ratio. *Space Weather*, 16, 69–88.  
 702 <https://doi.org/10.1002/2017SW001669>
- 703 Muller, R. H. (1944). Verification of short-range weather forecasts (a survey of the literature).  
 704 *Bulletin of the American Meteorological Society*, 25, 18–27.
- 705 Murphy, A. H. (1996). The Finley Affair: a signal event in the history of forecast verification.  
 706 *Weather and Forecasting*, 11, 3–20.
- 707 Murphy, A.H. & Winkler, R.L. (1987). A general framework for forecast verification. *Monthly*  
 708 *Weather Review*, 115, 1330–1338. [https://doi.org/10.1175/1520-0493\(1987\)115%3C1330:AGFFV%3E2.0.CO;2](https://doi.org/10.1175/1520-0493(1987)115%3C1330:AGFFV%3E2.0.CO;2)
- 710 Reiff, P. H. (1990). The use and misuse of statistics in space physics. *Journal of Geomagnetism*  
 711 *and Geoelectricity*, 42, 1145–1174. <https://doi.org/10.5636/jgg.42.1145>.
- 712 Reuter, B., & Schweizer, J. (2018). Describing snow instability by failure initiation, crack  
 713 propagation, and slab tensile support. *Geophysical Research Letters*, 45, 7019–7027.  
 714 <https://doi.org/10.1029/2018GL078069>
- 715 Rostoker, G. (1972). Geomagnetic indices. *Reviews of Geophysics and Space Physics*, 10, 935–  
 716 950.
- 717 Rowland, W., & Weigel, R. S. (2012). Intracalibration of particle detectors on a three-axis  
 718 stabilized geostationary platform. *Space Weather*, 10 (11).  
 719 <https://doi.org/10.1029/2012SW000816>

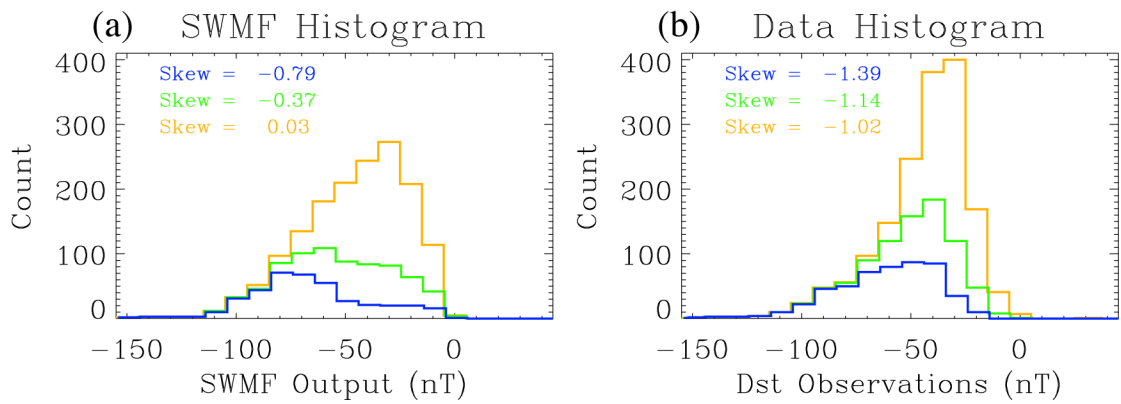
- 720 Pulkkinen, A., L. Rastätter, M. Kuznetsova, H. Singer, C. Balch, D. Weimer, et al. (2013).  
721 Community-wide validation of geospace model ground magnetic field perturbation  
722 predictions to support model transition to operations. *Space Weather*, *11*, 369–385.  
723 <https://doi.org/10.1002/swe.20056>
- 724 Sillanpaa, I., Ganushkina, N. Y., Dubyagin, S., & Rodriguez, J. V. (2017). Electron Fluxes at  
725 Geostationary Orbit From GOES MAGED Data. *Space Weather*, *15* (12), 1602-1614.  
726 doi: 10.1002/2017SW001698
- 727 Stefanescu, E. R., et al. (2014). Temporal, probabilistic mapping of ash clouds using wind field  
728 stochastic variability and uncertain eruption source parameters: Example of the 14 April  
729 2010 Eyjafjallajökull eruption. *Journal of Advances in Modeling Earth Systems*, *6*, 1173–  
730 1184. <https://doi.org/10.1002/2014MS000332>
- 731 Stephenson, D.B., Casati, B., Ferro, C.A.T. & Wilson, C.A. (2008), The extreme dependency  
732 score: a non-vanishing measure for forecasts of rare events. *Meteorological Applications*,  
733 *15*, 41-50. <https://doi.org/10.1002/met.53>
- 734 Swets, J. A. (1973). The relative operating characteristic in psychology. *Science*, *182*, 990-1000.
- 735 Toth, G., B. van der Holst, I. V. Sokolov, D. L. De Zeeuw, T. I. Gombosi, F. Fang, et al. (2012).  
736 Adaptive numerical algorithms in space weather modeling. *Journal of Computational*  
737 *Physics*, *231*, 870-903. <https://doi.org/10.1016/j.jcp.2011.02.006>
- 738 Wilks, D. S. (2019). *Statistical methods in the atmospheric sciences* (4th ed.). Oxford: Academic  
739 Press.



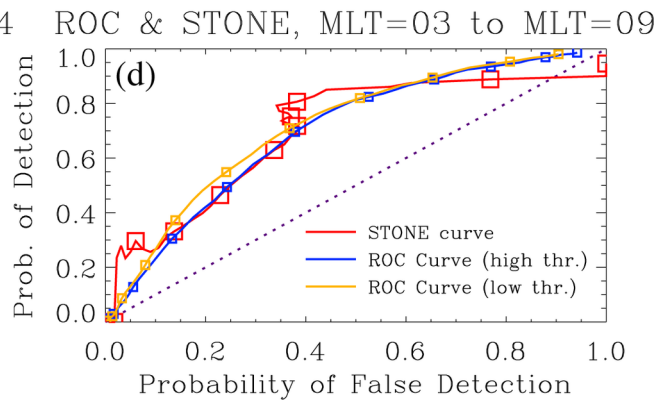
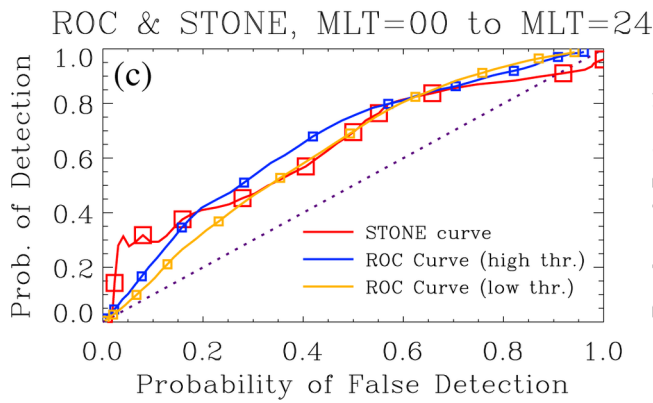
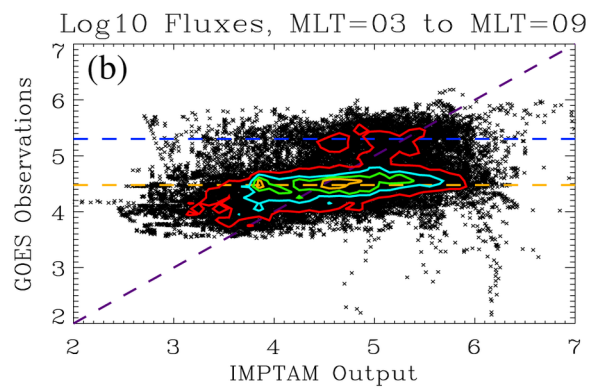
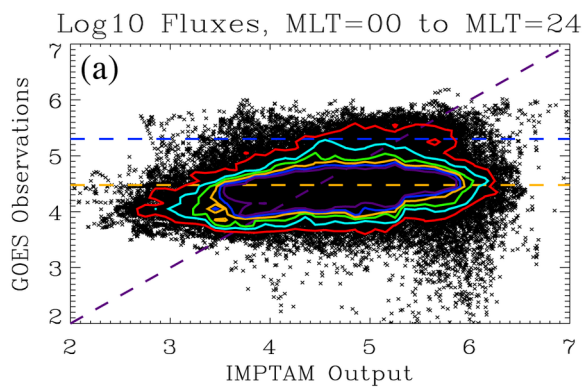
2020EA001106-f01-z-.tif



2020EA001106-f02-z.tif



2020EA001106-f03-z-.tif



2020EA001106-f04-z.tif