1 **Quantifying the storm-time thermospheric neutral density variations using model**

2 **and observations**

3 **E. Ceren Kalafatoglu Eyiguler[1], J. S. Shim[2], M. Kuznetsova[3], Z. Kaymaz[1], B. R. Bowman[4],**

4 **M. V. Codrescu[5], S. C. Solomon[6], T. J. Fuller-Rowell[5], A. J. Ridley[7], P. M. Mehta[8] and E.**

5 **K. Sutton[9]**

6 [1] İstanbul Technical University, Faculty of Aeronautics and Astronautics, Department of

7 Meteorology, İstanbul, TR

8 [2] The Catholic University of America, NASA GSFC, Greenbelt, MD, USA

9 [3] NASA, Goddard Space Flight Center, Greenbelt, MD, USA.

10 [4] Air Force Space Command Space Analysis Division, Peterson AFB, CO, USA

11 [5] Space Weather Prediction Center, NOAA, Boulder, CO, USA

12 [6] High Altitude Observatory, National Center for Atmospheric Research, Boulder, CO, USA

13 [7] School of Engineering, University of Michigan, Ann Arbor, MI, USA

14 [8] Department of Mechanical and Aerospace Engineering, West Virginia University, WV, USA.

15 [9] Air Force Research Laboratory, Kirtland AFB, NM, United States

16

17 Corresponding author: E. Ceren Kalafatoglu Eyiguler (ceren.kalafatoglu@itu.edu.tr)

18    **Key Points:**

19    • Using the average and maximum values of neutral densities to determine the model

20      performances can be misleading

21    • Removing the quiet-time trend from the neutral density reveals the actual performance of

22      the model in simulating the storm-time variations

23    • Mean absolute error, prediction efficiency and normalized root mean square error should

24      be considered together for the evaluations

25

26

27

28

29

30

31

32

33

34

35

36 **Abstract**

37 Accurate determination of thermospheric neutral density holds crucial importance for satellite

38 drag calculations. The problem is two-fold and involves the correct estimation of the quiet-time

39 climatology and storm-time variations. In this work, neutral density estimations from two

40 empirical and three physics-based models of the ionosphere-thermosphere are compared with the

41 neutral densities along the CHAMP satellite track for six geomagnetic storms. Storm time

42 variations are extracted from neutral density by 1) subtracting the mean difference between

43 model and observation (bias), 2) setting climatological variations to zero, and 3) multiplying

44 model data with the quiet time ratio between the model and observation. Several metrics are

45 employed to evaluate the model performances. We find that the removal of bias or climatology

46 reveals actual performance of the model in simulating the storm-time variations. When bias is

47 removed, depending on event and model, storm-time errors in neutral density can decrease by an

48 amount of 113% or can increase by an amount of 12% with respect to error in models with quiet

49 time bias. It is shown that using only average and maximum values of neutral density to

50 determine the model performances can be misleading since a model can estimate the averages

51 fairly-well but may not capture the maximum value or vice versa. Since each of the metrics used

52 for determining model performances provide different aspects of the error, among these, we

53 suggest employing mean absolute error, prediction efficiency and normalized root mean square

54 error together as standard set of metrics for the neutral density.

55 **Plain Language Summary**

56 Thermospheric neutral density is the largest source of uncertainty in atmospheric drag

57 calculations. Consequently, mission and maneuver planning, satellite lifetime predictions,

58 collision avoidance and orbit determination depend on the accurate estimation of the

59 thermospheric neutral density. Thermospheric neutral density varies in different time scales. In

60 short time scales, the largest variations occur due to the geomagnetic storms. Several empirical

61 and physics-based models of the ionosphere-thermosphere system are used for estimating the

62 variations in the neutral density. However, the storm-time responses from the models are clouded

63 by the climatology (background variations), upon which the effect of geomagnetic storms are

64 superimposed. In this work, we show that it is critical to use reference levels for the neutral

65 density to extract the true performance of the models for the evaluation of the storm-time

66 performances. We demonstrate that mean absolute error, prediction efficiency and normalized

67 root mean square error should be considered together for the performance evaluations, since each

68 of them provides different aspects of the error.

69 **1 Introduction**

70 It is known that the atmospheric drag acting on satellites is significant between the altitudes 160

71 and 800 km (Zesta and Huang, 2016). Consequently, in atmospheric drag calculations, in orbit

72 determination, the largest uncertainty comes from the thermospheric neutral density (Hejduk and

73 Snow, 2018; Bussy-Virat et al., 2017). The effects of the uncertainty in neutral density are not

74 only limited to orbit prediction, accurate density estimates are also needed for mission and

75 maneuver planning and collision avoidance (Storz et al., 2005). Low Earth orbit (LEO) satellites

76 are under the influence of the thermospheric environment and their lifetimes depend on the

77 variation of the neutral density (Prölss, 2011). Consequently, real-time estimation of the

78 atmospheric drag, which is important for satellite operations, heavily relies on the correct

79 estimation of the thermospheric neutral density.

80 Variations in thermospheric density can be decomposed into three main components: 1) the

81 variations, which are governed by the solar irradiance (solar-cycle dependent, seasonal, diurnal)

82    (Qian and Solomon, 2012), 2) the variations due to upward propagating tides and waves from the

83    mesosphere (Sutton et al., 2005), and 3) the storm-time variations, which are largely influenced

84    by the heat sources that come into play during geomagnetic activity, such as Joule heating

85    (Fedrizzi et al, 2011; Kim et al., 2006), auroral particle precipitation (Deng et al., 2013) and

86    heating due to small scale field-aligned currents (FACs) (Lühr et al., 2004). The former two

87    components control the quiet-time variation in neutral density, which is referred to as

88    climatological (background) variations in this study. In addition, the thermospheric composition

89    modulates the changes in thermospheric neutral density (Qian et al., 2009). In some geomagnetic

90    storm cases, the damping of the thermospheric density by NO cooling is significantly stronger

91    than expected. Those cases are classified as problem storms by Knipp et al. (2013) and it is

92    shown that the thermosphere's response is strongly associated with the pre-storm properties of

93    the solar wind. Different drivers of geomagnetic storms, such as the Coronal Mass Ejections

94    (CME) and Corotating Interaction Regions (CIRs) cause different environmental responses in the

95    thermosphere (McGranaghan et al., 2014). CIR and CME effects on thermospheric densities

96    were investigated in several studies (Chen et al., 2014; Chen et al., 2012; McGranaghan et al.,

97    2014; Lei et al., 2011; Thayer et al., 2008). Even though less geoeffective in terms of Dst

98    magnitude, the total effect of CIR storms was found to be comparable to CME induced

99    enhancements in thermospheric neutral density (Chen et al., 2014).

100   LEO satellite observations and empirical and physics-based models are employed in the

101   investigations of thermospheric neutral density (Lathuillère et al., 2008; Sutton et al., 2006; Liu

102   et al., 2005; Pardini et al., 2012; Codrescu et al., 2012; Deng et al., 2013; Solomon et al, 2011).

103   The Challenging Micro-Satellite Payload (CHAMP) and Gravity Recovery and Climate

104   Experiment (GRACE) satellites are the most used satellites for the investigations of the neutral

105   density and the associated atmospheric drag acting on satellites (Anderson et al., 2009; Picone et

106   al., 2002; Bruinsma and Forbes, 2010; Liu et al., 2010; Xu et al., 2011; Huang et al., 2014;

107   Bruinsma, 2015; Bruinsma et al., 2018). Recently, data from Swarm constellation has also been

108   employed to derive the thermospheric neutral densities (Huang and Zesta, 2016; Siemes et al.,

109   2016; Kodikara et al., 2018). In this kind of approach, the densities are calculated from the

110   accelerometers on the spacecraft (Sutton, 2005).

111   However, in-situ measurements from satellites only provide the current state of the

112   thermosphere. Hence, the empirical models involving semi-physical relations, which take

113   geomagnetic and solar indices as input and the physics-based models of the ionosphere-

114   thermosphere (IT) are employed to nowcast and forecast of the future state of the IT system in

115   global scales. The nowcast and forecast of neutral density are necessities for early-action and

116   response and orbit determination of the LEO spacecraft.

117   Comparisons between the model and observations are made in different time scales: daily global

118   mean (Solomon et al., 2011; Qian et al., 2008), orbit averaged (Bowman et al., 2008) and along

119   the satellite track (Connor et al., 2016; Shim et al., 2012). Comparisons for longer time scales

120   that are associated with the periodicities in neutral density such as the 27-day, 81-day and yearly

121   variations were also carried out in several studies (Rhoden et al., 2000; Qian and Solomon, 2012;

122   Bruinsma et al., 2018).

123   Several metrics are employed to assess the model performances. For the neutral density studies,

124   the most used metrics are the mean absolute error (MAE), bias (B), correlation (R), root mean

125   square error (RMSE), standard deviation (Std), prediction efficiency (PE), ratio of maximum and

126   ratio of average (Pardini et al., 2012, Shim et al., 2012; Elvidge et al., 2014; Bruinsma, 2015;

127   Elvidge et al., 2016; Emmert et al., 2017; Kodikara et al., 2018) and the version of the metrics in

128    log space (Picone et al., 2002; Sutton, 2018; Bruinsma et al., 2018). Each of these metrics has

129    advantages and disadvantages (Hyndman et al., 2012; Shcherbakov et al., 2013). For example,

130    the MAE provides the average difference between the model and observation and it is easy to

131    use. However, it does not offer any information on the amount of the error when compared to the

132    variations at large with respect to the event in percentage. Likewise, "ratio"s provide the

133    difference between the observation and estimate at an instant, but they do not deliver information

134    on the properties of the temporal evolution of the error. Std and RMSE are highly sensitive to

135    outliers and may lead to the overestimation of errors in some cases. Among the metrics, the PE is

136    becoming increasingly used by the space weather community. PE is a dimensionless quantity and

137    represents the measure of success in reproducing a time series. PE basically compares the order

138    of magnitude of model errors with the magnitude of variations of the measurements/reference

139    data. However, one handicap of PE is that, it does not provide the actual value of difference

140    between the observation and estimations. It is also worth to note that in the literature, same

141    equations are used in the calculations of all metrics given above, except the bias metric. Bias

142    may have different definitions based on the study. Bias is sometimes calculated as the difference

143    between the model and observation in percentage (Pardini et al., 2012) and sometimes as the

144    mean difference between the model and observation (Elvidge et al., 2016). In our work, we

145    define model bias as the quiet-time mean difference between the model and observation (mean of

146    model minus mean of observation). Additionally, we do not use it as a metric, but rather, use the

147    quiet-time model bias to extract the storm-time variations from the neutral density. The

148    definitions of the metrics that we use in our study are given in Section 2.3.

149    As a summary, all metrics provide different aspects of the error. Hence, Chai et al. (2014)

150    suggests using not only one, but several metrics together, especially in studies involving the

151   assessment of more than one model when the error distribution becomes important.

152   Consequently, this is the case for the neutral density studies and a variety of metrics are

153   employed together in comparisons. However, there are not any consensus on what to use as a

154   standard set of metrics. The community need at the current time is to be able to run the models

155   for real-time calculations of atmospheric drag in support of real-time satellite operations. For this

156   purpose, there is a need to assess the performances of the models and to specify the conditions

157   when they perform satisfactorily and when they do not (Shim et al., 2014; Shim et al., 2015).

158   This study is a continuation of the GEM-CEDAR challenge for the assessment and

159   benchmarking of the empirical and coupled models of the ionosphere-thermosphere and is a

160   deliverable of the International Forum on Space Weather Modeling Capabilities Assessment. In

161   the first study of the series, Shim et al. (2011) compared the model results with the local

162   measurements available from EISCAT radars for the ionospheric parameters $N_mF2$, $h_mF2$ and

163   vertical drift with limited latitudinal coverage. Shim et al. (2012) focused on the space-borne

164   measurements of the $N_mF2$, $h_mF2$, ionospheric electron density and thermospheric neutral density

165   along the satellite track at the measurement locations. $N_mF2$ and $h_mF2$ from the models were

166   compared with the observations from the Constellation Observing System for Meteorology,

167   Ionosphere and Climate (COSMIC) while ionospheric electron density and thermospheric neutral

168   density were compared using the measurements from CHAMP. In both studies, root mean square

169   error (RMSE), prediction efficiency (PE), ratio of (max-min) and ratio of maxima were

170   employed to assess the model performances. They reported that the model performances depend

171   on the metrics used and varied with latitude and geomagnetic levels. No models outperformed

172   others in estimating the thermospheric and ionospheric parameters in all cases.

173    In model comparison and validation, the absence of a standard set of metrics complicates the

174    evaluation and synthesis of the results of different studies. As a part of the systematic evaluation

175    of the models in this study, our aims are to present ways to facilitate the comparison of the

176    storm-time performances of the models and to provide a useful set of metrics for the neutral

177    density studies. We present methods to remove the quiet time variations from the neutral density,

178    so that the storm-time changes are revealed. Accordingly, direct comparisons can be made

179    between the model estimations and observations from the CHAMP satellite for the disturbed

180    periods. The climatology removal methods are called as baseline shifts, since they match the

181    level of quiet-time neutral density estimated from the models with the quiet-time level of neutral

182    density variations observed by CHAMP. Orbital averages of thermospheric neutral density along

183    the CHAMP satellite track are used to evaluate the model performance. We show that baseline

184    shifts are a necessity in order to correctly assess the storm-time performances of the models and

185    the climatology and storm-time variations should be evaluated separately as the dominant

186    mechanisms and their time-scales are different in each. In Section 2, the events selected for the

187    case studies are introduced and baseline shifting methods are described. Section 3 presents the

188    results and involves the comparison of baseline shifting methods and the neutral density

189    estimations from the empirical and physics-based models of the IT. Lastly, we conclude the

190    study and discuss the future needs of the community in Section 4.

191    **2 Data and Methodology**

192    Two empirical and three physics-based models are employed in this study. The empirical models

193    are Naval Research Laboratory Mass Spectrometer and Incoherent Scatter Extended

194    (NRLMSISE-00, will be referred to as MSIS, hereafter) (Picone et al., 2002) and Jacchia-

195    Bowman-2008 (JB2008) (Bowman et al., 2008), whereas the physics-based models are

196    Thermosphere-Ionosphere-Electrodynamics General Circulation Model (TIEGCM1.95)

197    (Richmond et al., 1992), Coupled Thermosphere Ionosphere Plasmasphere electrodynamics

198    (CTIPe) (Millward et al., 2001; Codrescu et al., 2008) and GITM (Ridley et al., 2006). The

199    models were run using the NASA Community Coordinated Modeling Center Runs-on-Request

200    system. The results can be found by searching the simulation IDs that are given in Table S1.

201    Additionally, Table S1 provides information on the version and the resolution of the models for

202    each run. For each run and model, the initial parameters and model input are the same. Table S2

203    shows the input parameters to the models. For physics-based models, ionospheric electric

204    potentials have to be specified to describe the interaction of the solar wind and magnetosphere

205    with the ionosphere. This is handled by selecting a high-latitude driver, which describes the

206    electrodynamic input from the magnetosphere and solar wind into the high-latitude ionosphere

207    under different solar wind conditions. In this study, Weimer-2005 (Weimer, 2005) ionospheric

208    potentials are employed as the high-latitude driver for each physics-based model for consistency.

209    Details on the models and their standard configurations for the runs can be found in (Shim et al.,

210    2011; Shim et al., 2012).

211    The model results are compared against the newly updated thermospheric neutral density data set

212    from CHAMP by Mehta et al. (2017), which is referred to as M2017, hereafter. Previous studies

213    of systematic assessment (Shim et al., 2012; Shim et al., 2014) used older versions of neutral

214    density data that were also derived from CHAMP accelerometer measurements (Sutton et al.,

215    2005). Besides, prior to the M2017, the most recent version of neutral density data which had

216    been widely used in comparisons was the Version 2.3 of Sutton (2009). This version is also

217    detailed on a report by (Sutton, 2011). The differences between the previous versions of neutral

218    density data sets and the M2017 are associated with the modeling of the drag coefficient ($C_D$),

219 which is a coefficient in the equation of satellite drag. The drag coefficient is a number that

220 depends on the geometry of the spacecraft and the properties of the impinging particles. Precise

221 calculations of the drag coefficient are necessary for accurate neutral density estimations, since

222 the neutral density is calculated using accelerometer data, hence the $C_D$. The M2017 considers a

223 more complicated geometry and uses the most recent advances in the modeling of gas-surface

224 interactions and the modeling of physical $C_D$. In their work, Mehta et al. (2017) reported

225 differences up to 20% for some cases with respect to the neutral density estimates of Sutton

226 (2008). In this study, to give the difference between the newly derived and old data sets, the

227 Version 2.3 data set (Sutton, 2009) is also included in the comparisons. The (Sutton, 2009)

228 Version 2.3 is represented as SV2.3 throughout the paper.

229 In this work, we investigate the storm-time performances of the IT models for six geomagnetic

230 storms, which were particularly chosen by the GEM-CEDAR community for the systematic

231 evaluation of the models. According to the NOAA classification based on the Kp index, the

232 intensity of selected events ranges from weak to severe. Table 1 presents the extreme values of

233 geomagnetic and solar indices along with the solar wind drivers for the events Hemispheric

234 Power (HP) index is also given in Table 1 since it is an input to the physics-based models. In the

235 Table, HSS denotes the high speed streams.

236 Figure 1 shows the storm-time maximum neutral density on the left, storm-time average neutral

237 density from the models and M2017 in the middle, and the timing difference between the neutral

238 density maximum in M2017 and the maximum in models in the right panel, for each

239 geomagnetic storm case. As evident from the plot, the storm-time maximum and average neutral

240 densities from M2017 display a decreasing trend with weaker geomagnetic storms. Even though

241 SV2.3 always shows higher values than M2017, it follows the same trend in neutral densities.

242 For the neutral density maximum, all models show the same tendency as in CHAMP

243 observations, except the 2005-243 event, which is due to an HSS. TIEGCM and JB2008

244 overestimate the neutral density peak in each event, whereas GITM slightly underestimates in

245 four of the six events (2005-135, 2005-243, 2007-142 and 2007-91). MSIS neutral density

246 maxima are higher than M2017 for events with Kp<6, but lower than M2017 for events with

247 Kp>6, except the 2006-348 event. CTIPe estimates are slightly higher than but very close to

248 M2017 in most of the events. Overall, CTIPe and GITM are the two models that generally show

249 the closest neutral density maxima to M2017.

250 These patterns in the modeled neutral density maxima change in the average neutral densities. A

251 model overestimating the neutral density maxima in M2017 can give a lower average than the

252 M2017 or vice versa for the same events. For example, JB2008 and GITM for 2005-135,

253 TIEGCM for the 2005-243 and MSIS for the 2006-348 and 2007-142 show the opposite

254 behavior in terms of storm-time neutral density average and maximum. In the figure, it is seen

255 that MSIS underestimates the neutral density average in all selected events except the 2007-91.

256 JB2008 overestimates the storm-time neutral density in four of the six events and underestimates

257 in two events. Neither the MSIS, nor the JB2008 display the decreasing trend with weakening

258 geomagnetic activity in average neutral density average that is illustrated in M2017 for the

259 selected event set. Despite, TIEGCM, and GITM display the decreasing trend also for the neutral

260 density averages, except the 2005-243 event as in neutral density maxima case. None of the

261 models are found to be consistently closer to M2017 in terms of neutral density average.

262 Timing differences between the models and M2017 also change with respect to event.

263 Interestingly, most of the models performed the best in capturing the timing of maximum in

264 2005-190 event, which is due to a CME during an HSS. The variations in timing differences

265      seem to be random. The timing difference between the maxima of M2017 and the models are

266      found to be between $\pm 7.5$ hours.

267      In Figure 1, the storm-time neutral density maxima and averages include not only the storm-time

268      neutral density variations but also the climatological variations. That is, the model biases are also

269      included in evaluations. In the following sections, we show that removing the climatology or

270      quiet-time model bias reveals the actual performance of the models in simulating the

271      thermospheric neutral density variations during geomagnetic activity. Our approach for assessing

272      the storm-time model performances consists of three steps, such as orbit averaging,

273      climatology/bias removal and assessment of the results. In the following sections, we describe

274      the tools designed for each step. The codes were written in MATLAB and are in transition to

275      Python language.

276      **2.1 Orbit Averaging Tool (OAT)**

277      The orbit averaging tool (OAT) is used for taking orbital averages of thermospheric neutral

278      density from CHAMP and models. Comparisons along the track involve local time effects,

279      small-scale structures, and diurnal and seasonal variations (Qian and Solomon, 2012; Liu et al.,

280      2005; Lühr et al., 2004; Kwak et al., 2009), which make it hard to specify the reason behind the

281      difference in model estimations and observations. On the other hand, taking orbital averages

282      smooths out the temporal and spatial variations due to the spacecraft position on a single orbit

283      and provides the globally averaged response to the geomagnetic storm. It was also shown

284      previously by Burke et al. (2007) that the change in orbit-averaged densities occurs

285      systematically whereas the local density exhibits large variations.

286    The OAT works with CHAMP ephemeris data. First, the beginning and end times of each orbit

287    are determined: an orbit starts at the highest northern latitude, crosses the highest southern

288    latitude and ends at the highest northern latitude. One orbit lasts approximately 92 minutes.

289    There are typically ~15 orbits in a day. Neutral density observations from CHAMP and

290    estimations from each model are averaged over every single orbit of the spacecraft.

291    **2.2 Baseline Shifting Tool (BAST)**

292    In this study, we are concerned with the storm-time performances of the models. Thus, to

293    compare only the storm-time responses, the baseline shifting tool (BAST) is used. BAST adjusts

294    the quiet-time neutral density level of the models to match the quiet time level of M2017. The

295    adjustment is handled by assuming that unless there is a geomagnetic storm, the neutral density

296    variations will continue to fluctuate around the quiet time level of neutral density. Consequently,

297    three types of adjustment are employed:1) subtracting the average quiet-time difference between

298    the models and observation (Shift1-SH1), 2) setting off the climatology to zero by subtracting

299    the quiet-time neutral density average from the models and the observation (Shift2: SH2), and 3)

300    multiplying the model results with the quiet time average ratio between the model and

301    observation (Shift3:SH3). All adjustments are applied separately to the model results. Hereafter,

302    we call the adjustments as baseline shifts, since they shift the quiet time reference level of the

303    model results to the observation or to the zero level. In the shifting procedure, the "quiet time"

304    refers to the neutral density variations, which are only due to the changes in the solar irradiance

305    and tides. Subsequently, any additional changes in the neutral density that are due to the

306    geomagnetic disturbances are referred to as storm-time variations. The storm-time variations are

307    considered to be superimposed on the quiet time neutral density variations (Lühr et al., 2011).

308 All three shifts work with the quiet time average of thermospheric neutral density from the

309 model and observations. Hence, the correct identification of the quiet time intervals is important.

310 To determine the quiet time intervals, we select a threshold for the Kp index and the neutral

311 density fluctuations as observed by the CHAMP satellite. An interval is defined as quiet when

312 Kp < 3- and the orbit-averaged neutral density difference between two consecutive orbits of

313 CHAMP is less than or equal to $1.25 \times 10$-13 kg/m$^3$. The threshold, $1.25 \times 10^{-13}$ kg/m$^3$, was

314 selected by inspecting the orbit-averaged neutral density variations on quiet day cases (2007-79,

315 2007-190, 2007-341) used in (Shim et al., 2012) (see Figure S1). We define it as the start of the

316 storm when the increase in CHAMP neutral density is more than $1.25 \times 10^{-13}$ kg/m$^3$ and there is

317 an increasing trend in orbit-averaged neutral density in two consecutive orbits. The end of the

318 storm is marked as the time when CHAMP neutral densities return to quiet-time average neutral

319 density level. Table 2 details the shifts that are applied to the thermospheric neutral density.

320 As a result of the shifting processes, we estimate the errors to be as high as the selected

321 threshold: $\pm 1.25 \times 10^{-13}$ kg/m$^3$, which is about 5% to 7% of the quiet-time neutral density of the

322 selected events.

323 Figure 2 shows the 2006-348 event, which is classified as 'severe' according to the NOAA

324 geomagnetic storm scale based on Kp, as an example event for baseline shifts. The selected quiet

325 time interval for the event, which was determined according to thresholds for Kp and neutral

326 density level is between 13/12/2006 15:00 UT and 14/12/2006 14:00. The original time series

327 from the model and observations are displayed on the left and the shifts 1, 2 and 3 are found on

328 the right panels. It is seen that most of the models overestimate the neutral densities during the

329 quiet-time interval. Appropriately, the shifts remove the bias from the models, so that we can

330 compare the storm-time variations directly between the models and M2017.

331     Before the baseline shifting procedure, MSIS is one of the best performing models with a

332     maximum close to the M2017 for the 2006-348 event. However, with the removal of its bias, it is

333     found that it actually underestimates the neutral density enhancement due to the geomagnetic

334     storm. In the case of TIEGCM, the model overestimates the quiet-time neutral density so much

335     that, the neutral density maximum and average during the storm are the highest among the

336     models. Consequently, the resulting differences between the model and observation are the

337     highest when the quiet-time bias is included. On the other hand, shifting the baseline to M2017

338     levels as seen in panels b and c indicate that the storm-time response as modeled by the

339     TIEGCM is closer to M2017 than they are before the shift. These cases demonstrate the

340     usefulness of the shifts in determining the actual storm-time response from the models.

341     Following the same assumptions as in case of SH1, SH2 and SH3, several other types of shifts

342     can also be applied to the data to remove the influence of the quiet time bias on the storm-time

343     performances. For example, an artificial time series can be produced using the quiet-time data by

344     assuming that the neutral density levels will remain the same on the following day. The easiest

345     way to produce an artificial time series is to sequentially iterate the neutral density during the

346     quiet time period to cover the entire event interval. Afterwards, this newly generated time series

347     can be used for point-to-point subtraction of 1) bias (Shift4, SH4) and 2) quiet time neutral

348     density at the same instant (Shift5-SH5) or for 3) point to point multiplication using the quiet

349     time ratios (Shift6-SH6). These procedures were also investigated in this work. However, since

350     the results of point-to-point shifts are similar to shifts based on quiet time averages, which are

351     described above, we chose to present only the results from SH1, SH2 and SH3. However, the

352     results of all shifts for the selected events are provided in the supplement from Figure S2 to

353     Figure S7. The figures demonstrate that, point-to-point shifting processes may lead to unphysical

354     variations in neutral density as in the case of GITM for weak events in this study.

355     **2.3 Performance Assessment Tool (PAT)**

356     After adjusting the baseline of the model and observations, storm-time model performances are

357     evaluated according to the M2017 data set. Performance Assessment Tool (PAT) measures the

358     model performances during individual events according to seven metrics. Those are: ratio

359     between the model maximum and CHAMP maximum (Ratio$_{max}$), ratio between the model mean

360     and CHAMP mean (Ratio$_{avg}$), time delay between the peak of the model and peak of the

361     CHAMP observation (TD), mean absolute error (MAE), normalized root mean square error

362     (NRMSE), prediction efficiency (PE) and integrated density change (IDC). Equations from 1 to

363     7 show the definitions of the metrics. The subscripts "i" and "j" represent the orbit number

364     during the quiet-time and entire event and "t", the time of the orbit, respectively. All calculations

365     are based on the storm-time variations after performing the baseline shifts.

366     $\text{Ratio}_{max} = \dfrac{\rho_{model,max}}{\rho_{M2017,max}}$     (1)

367     $\text{Ratio}_{avg} = \dfrac{\rho_{model,avg}}{\rho_{M2017,avg}}$     (2)

368     $\text{TD} = t_{model,max} - t_{M2017,max}$     (3)

369     $\text{MAE} = \sum |\rho_{M2017,i} - \rho_{model,i}| / N$     (4)

370     $\text{NRMSE} = RMSE / (\rho_{M2017,max} - \rho_{M2017,min}) = \sqrt{\sum \dfrac{(\rho_{M2017,i} - \rho_{model,i})^2}{N}} / (\rho_{M2017,max} - \rho_{M2017,min})$

371                                  (5)

372     $PE = 1 - RMS_{model} / RMS_{M2017} = 1 - \sqrt{\dfrac{\sum(\rho_{M2017,i} - \rho_{model,i})^2}{\sum(\rho_{M2017,i} - \overline{\rho_{M2017,i}})^2}}$     (6)

373 $$\text{IDC} = \sum_{j=1}^{n_{orbit}}\left(\sum_{t_{start}}^{t_{end}} \rho_{data,t} - \rho_{baseline}\right)_{j}; \rho_{baseline} = \sum_{i=1}^{q_{orbit}}\left(\sum_{t_{start}}^{t_{end}} \rho_{data,t}\right)_{i}/q_{orbit} \qquad (7)$$

374 Among the metrics, the IDC works with the orbit and storm-time integrated neutral densities.

375 The subscript "data" in Equation 7 denotes model or M2017 data. $q_{orbit}$ is the total number of

376 orbits during the quiet time interval. $n_{orbit}$ is the total number of orbits during the entire event

377 and "$t_{end}$" and "$t_{start}$" denote the start and end times of each orbit. Accordingly, $\rho_{baseline}$

378 represents the average of the orbit-integrated neutral density during the quiet time.

379 In contrast, other metrics use the orbit-averaged neutral densities. The perfect score for the ratios

380 (Ratio$_{max}$, Ratio$_{avg}$) is 1, whereas TD should be zero, meaning there is no lag between the peak of

381 the model and the time of the maximum from CHAMP. Determining the TD for less intense

382 events is different from determining the TD for intense events. In intense events, the maximum

383 of the neutral density is distinguishable, whereas in less intense events, there may be numerous

384 local maxima. Consequently, we first mark the timing of the maximum neutral density from

385 M2017, then detect the timing of the closest local maxima from the models. MAE gives the

386 average distance between the observation and model estimations. Values approaching to zero

387 indicate better agreement between the model and observations. Furthermore, MAE gives a

388 dimensioned skill score, that is, it has the same units with the neutral density (kg/m$^3$). On the

389 other hand, ratios, NRMSE and PE are dimensionless. PE varies between 1 and negative infinity.

390 PE equals to 1 indicates perfect agreement between the model and observations whereas PE=0

391 means the model errors are in the same order with the variations of the observations. Negative

392 PE values show that the observed mean is a better estimate for forecasts than the model (Shim et

393 al., 2012). The NRMSE, is the normalized version of RMSE. The NRMSE gives errors in

394 percentage. RMSE, consequently, NRMSE, vary with the variability of error magnitudes and the

395    mean absolute error (Wilmott and Matsuura, 2005). When interpreted together with the MAE,

396    NRMSE provides information on the variability of error magnitudes.

397    **3 Results and Discussion**

398    In this section, we present the storm-time performances of the models after the baseline shifting

399    methods are applied to the observation and model neutral density estimations from the models.

400    Figure 3 presents the ratio of maximum neutral density (top row) and ratio of average neutral

401    density (middle row) from each model to M2017. The best agreements are displayed between the

402    SV2.3 and M2017 for all events before and after the baseline shifting. The SH3 yields the best

403    results among the shifts for the SV2.3 and lead to one-to-one match between the M2017 and

404    SV2.3 for all events. This is because M2017 and SV2.3 are only different by a constant factor in

405    each event and SH3 finds and removes this factor by using the ratio between the SV2.3 and

406    M2017 during the selected quiet time interval.

407    For MSIS, CTIPe and GITM, baseline shifting causes the ratio of maximum to diverge from 1

408    for some events, whereas for TIEGCM and JB2008 the shifts cause performance enhancement in

409    capturing the maximum in M2017. MSIS and GITM are found to underestimate the maximum in

410    M2017 generally, after the shifts. For all models, SH1 produces the closest ratios to 1 among the

411    shifts for both the ratio of maximum and ratio of average neutral densities. SH2 causes the ratios

412    to be more spread for all events and models. Using SH3 leads to the underestimation of neutral

413    density average and maximum for all models except the JB2008. For JB2008, after the SH3, the

414    ratios approach closer to 1 with respect to other shifts for most of the events. However, there is

415    still overestimation in two of the events. Additionally, in TIEGCM, the 2005-135 event shows a

416    distinct behavior and captures the maximum in M2017 better after SH3. CTIPe overestimates for

417    events with Kp<7 and underestimates in with Kp≥7 before and after the shifts.

418    Qualitatively, the same conclusions mostly hold true for the ratio of neutral density averages and

419    maxima from the models; only the amount of underestimation or overestimation changes.

420    However, a model overestimating the neutral density maximum may underestimate the average

421    density as in JB2008 case for the 2006-348 event. Moreover, a model underestimating the neutral

422    density maximum may overestimate the average density as in CTIPe for the event 2005-243 and

423    GITM as in events with Kp≥7.

424    Timing differences between the maximum in M2017 and the models are shown on the bottom

425    row in Figure 3. SH1, SH2 and SH3 do not change the lags between the model maximum and

426    M2017. This is natural as only a constant value is used for the baseline shifts.

427    Figure 4 depicts the changes of the neutral density maximum (left panel) and average (middle

428    panel) from the quiet time values in percentage. Right panel shows the time and orbit-integrated

429    density change (IDC). The percentage change from the background variations and the IDC are

430    calculated around the zero-baseline level when all climatology is removed. Accordingly, SH2 is

431    used in the calculations of percentage change and the IDC. The percentages are calculated as

432    $\%Change = 100 \times (storm - quiet)/quiet..$

433    In M2017, the change in neutral density maximum due to the geomagnetic storm is found to be

434    nearly as twice as the change in neutral density average for the observations and models for all

435    events. The change in neutral density maximum ranges from 200% to 90% and the change in

436    neutral density average ranges from 100% to 45%. Both the change in maximum and average of

437    the observations (M2017 and SV2.3) show a decreasing trend with lower geomagnetic storm

438  intensity in terms of Kp. TIEGCM and CTIPe estimate the closest percentages to M2017 for

439  events with Kp≤7. CTIPe also performs reasonably well for events with Kp≥7.

440  In the right panel, geomagnetic storms with less Kp, which are due to HSSs (2007-142, 2005-

441  190) display IDCs as large as the events due to CMEs (2005-135, 2006-348). There is not any

442  model, which is consistently closer to the IDC from M2017. However, MSIS is closer to M2017

443  more times than the other models (4 of the 6 selected cases: 2006-348, 2005-243, 2007-142,

444  2007-91). TIEGCM overestimates in all events. Similar to TIEGCM, JB2008 and CTIPe are

445  higher than the M2017, except the 2005-135 and 2005-243 events, respectively. GITM shows a

446  distinction between Kp≥6+ and Kp<6+ events: it overpredicts the IDC in events with Kp≥6+ and

447  under predicts for events with Kp<6+ for the selected events.

448  Figure 5 presents mean absolute error (MAE), normalized root mean square error (NRMSE) and

449  prediction efficiency (PE) of the models for the selected events. MAE and NRMSE are

450  negatively-oriented skill scores, meanwhile PE is positively-oriented. This means that, lower

451  values of MAE and NRMSE are more desirable whereas PE closer to one shows the perfect

452  agreement between the models and M2017, in our case.

453  From Figure 5, the effect of baseline shifts on the storm-time performance of the models can be

454  distinguished. It is found that generally, the calculated errors after the baseline shifts are on the

455  same order for all models and range between 1% and 20%. In the figure, baseline shifts are

456  found to reduce the errors (MAE, NRMSE and PE) for the TIEGCM and SV2.3 for all cases.

457  Additionally, as in the case of the ratios, SV2.3 errors are more efficiently reduced using the SH3

458  compared to the other shifts.

459     The MAE provides information on the amount of mean error in dimensioned units (kg/m$^3$, in the

460     case of thermospheric neutral density). For the selected event set, MAE is found to be high for

461     strong events and low for weak events after the baseline shifts (except for GITM in 2005-243

462     and CTIPe in 2006-348), which is consistent with the findings of (Shim et al., 2012). Moreover,

463     Figure S8 shows that the behavior of RMSE is the same with MAE in all cases and models and

464     the amount of error grows with respect to event intensity. The amount of error increases with

465     stronger events because the temporal variability of the thermospheric neutral density is higher in

466     stronger geomagnetic storms. Normalization shows the errors are actually around the same order

467     of magnitude in terms of percentage for the events. A high MAE may account for a low NRMSE

468     based on the variation of the thermospheric neutral density during the event. On the other hand,

469     an increase in MAE after the shift, with respect to the original time series without shift, mirrors

470     itself as an increase in NRMSE with respect to the original time series, as well. Thus, basically,

471     the MAE and NRMSE provide the same information on the change in errors. However, NRMSE

472     gives the additional information that how much this error accounts for from the perspective of

473     the variability of the thermospheric neutral density based on the event.

474     The NRMSE from the models are confined between 60% and 10% after the shifts. Before the

475     shift, TIEGCM has the maximum NRMSE with ~125% for the event 2007-91. The shifts

476     revealed that its actual storm-time performance to be on the order of ~12% (SH1, SH2) to ~33%

477     (SH3) for the same event. In contrast, MSIS has a minimum error around ~25%, which increases

478     to ~37% (SH1, SH2) to ~41% (SH3) for the event 2005-135.

479     The 2005-135 is an exceptional case as can be seen from MAE, RMSE, NRMSE and PE of the

480     TIEGCM. Interestingly, only for TIEGCM among the other empirical and physics-based models

481     of the IT and only in this event, baseline removal via ratios (SH3) reduces the error more than

482  the shifts based on subtraction. In this case, the storm-time variation is so high and strong that it

483  is compensated by taking ratios. However, we argue that this is not the actual performance of the

484  model. Since storm effects are generally additions to the background neutral density (Lühr et al.,

485  2011), in the case of this event, the model, in fact, overestimates the storm-time variations so

486  much that the error is reduced via the SH3, which uses quiet-time ratios. On the other hand, for

487  other events, SH3 gives rise to the underestimation of the average and maximum values of the

488  neutral density from TIEGCM (Figures 2 and 3, Figures S2 to S6).

489  The PE on the right column shows the same variations with the NRMSE according to the event.

490  The PE increases when the NRMSE increases and vice versa. The PEs of TIEGCM for the

491  original, unshifted model neutral density are so low that the scales are compressed in the figure.

492  However, after the shifts, there is a clear improvement in model performances, which can be

493  seen from the frame interior to the figure.

494  The errors in TIEGCM seem to increase with the intensity of the geomagnetic storm. After

495  removing the climatology via the baseline shifts, the errors in CTIPe and JB2008 are also found

496  to decrease except the 2006-348 and 2005-135 events, respectively. The events with the most

497  errors in CTIPe model are found to be the problem cases, which Knipp et al., 2013 listed (2005-

498  190, 2005-243 and 2005-135). In the problem events, the damping of the thermospheric density

499  by NO cooling is more than expected, so that the density may not enhance as high as, that

500  estimated by the IT models. However, we should note that the version of CTIPe that is used in

501  this work does not include the correction to NO cooling at high Kp levels. From the selected

502  events, GITM appears to show a reduction in error for events with Kp≥6+ and growth in errors

503  for the events with Kp<6 after the baseline shifts. On the contrary, after the removal of the

504　climatology, for MSIS, the errors in the selected cases give the impression that they increase for

505　events with Kp≥6+, except the 2005-135 event, and decrease for events with Kp<6.

506　In our selected cases, after SH1 or SH2, TIEGCM performed the best for events with Kp<7

507　according to all metrics. Moreover, TIEGCM demonstrated the highest PE for most of the cases.

508　Lastly, it is found that the SH3 reduces the errors more than the other shifts for the SV2.3, since

509　neutral density is derived from the accelerometer on spacecraft and the error can be multiplied

510　during this process. All shifts and all events in terms of MAE, NRMSE, and PE show that the

511　SH3 works perfectly for the SV2.3 and the errors are on the order of ~1%, with a maximum of

512　~2.5%.

513　In addition to the errors from the models using the shifts SH1, SH2, and SH3 provided above, the

514　errors for the shifts with point-to-point subtraction and multiplication (SH4, SH5, SH6) are given

515　in the supplement (Figure S10, Figure S11). It can be seen from Figure S10 and Figure S11 that

516　the choice of the baseline shifting method does not affect the performance outcome of the

517　models. The errors obtained by using SH1 and SH4, SH2 and SH5 and SH3 and SH6 are very

518　close to each other.

519　Furthermore, additional metrics may be utilized serving to the special purposes of the studies.

520　For example, since their technique for data assimilation aims to reduce the errors in logarithmic

521　densities, Sutton (2018) used mean, normalized standard deviation and root mean square (rms)

522　errors of the log density ratio ($ln(\rho_{model}/\rho_{observation})$) in their work. For the sake of

523　comparison, we also tested these metrics for our events. Figure S9 presents the results. The

524　logarithmic mean gives similar results to the Ratio$_{avg}$ for all shifts, whereas for SH2, the errors

525　from the models are amplified in standard deviation and rms relying on the

526     *ln*(model/observation) ratio. The rms of log density ratio in SH1 and SH3 are found to be very

527     close to NRMSE.

528     **4 Summary and Conclusion**

529     In this study, we had two aims: 1) to find methods to facilitate the evaluation of the storm-time

530     performance of models and 2) to suggest a standard set of metrics to determine the model

531     performances.

532     For the first part, we presented methods to remove the quiet-time bias/climatology from the

533     models and referred to these methods as "baseline shifts". Shifts are based on subtraction of bias

534     from the models (SH1), subtraction of climatology from model and observation (SH2) and

535     multiplication of the quiet-time ratio between the model and observation with the model to match

536     the quiet-time neutral density level of observation (SH3). It was shown that defining the quiet-

537     time reference level is very critical in determining the actual storm-time performances. In some

538     events and models, the shifts were found to reduce the errors due to climatology in evaluating the

539     storm-time performances up to 113% (TIEGCM-2007-91: 125% to 12%) whereas in some

540     events, they increased the errors by 13% (MSIS-2005-135: 12% to 25%).

541     For the storm-time performance assessment of the models, SH1 and SH2 are found to work

542     equally well. The choice of different baseline levels (shifting the models to the level of CHAMP

543     observations or shifting observations and models to zero level by removing all the climatology)

544     does not change the amount of error associated with a model. Besides, SH3 increases the

545     variability of the errors from the models when compared to the other shifts. This is due to the

546     fact that the storm-time effects are generally superimposed upon the background (climatological)

547     variations and their nature is not multiplicative. Hence, modifying the original time series using

548     ratios does not work as efficiently as the subtraction process for the empirical and physics-based

549     models.

550     On the other hand, SH3 is efficient when comparing M2017 and SV2.3 as it depends on the quiet

551     time ratios. The difference between these two data sets is only a constant number, which depends

552     on the modeling of the $C_d$ and the geometry of the spacecraft. Therefore, the SH3 works the best

553     for SV2.3 when compared to other shifts. Hence, when neutral density is derived from

554     accelerometer data, systematic error and bias can be multiplied, so it is reasonable to divide to

555     remove them. It follows that the findings of the past model validation studies which used SV2.3

556     can be re-evaluated and calibrated using the SH3.

557     From the selected cases, it appears that, TIEGCM is more successful in low Kp events, and its

558     success rate decreases with the intensity of the storm. GITM shows a reduction in error for

559     events with Kp≥6+ and increase in errors for the events with Kp<6. On the contrary, the model

560     errors increase for MSIS for events with Kp≥6+, except the 2005-135 event, and decrease for

561     events with Kp<6 in this event set. JB2008 does not show any systematic errors for the selected

562     events. After the removal of the quiet time bias/climatology between M2017 and the models,

563     TIEGCM seems to perform the best in terms of all metrics for most of the selected events,

564     followed by CTIPe and GITM. For the selected cases, JB2008 was closer to M2017 than MSIS

565     for more of the events.

566     Three of the six events selected in this study were listed as problem storms by (Knipp et al.,

567     2013). They reported that the modeling of these storms is more difficult with respect to several

568     other events with less NO production. The NO cooling during these events restrict the neutral

569     density enhancement and neutral density does not increase as high as expected from the models.

570     In our event set, for these storms, the range of errors from the models are between 13% and 40%

571　and do not greatly differ from the other cases. Thus, we do not see any distinction among model

572　performances with respect to the solar wind drivers of the events. On the other hand,

573　performances of the MSIS, TIEGCM and GITM suggest differences based on Kp.

574　Furthermore, it is possible to estimate integrated neutral density change (IDC) during the storm

575　via SH2, which shifts the baseline to zero level. IDC is important as drag has a cumulative effect

576　on orbit determination and prediction (Emmert et al., 2017). For the evaluations in drag

577　calculations, we suggest using the upper limits for IDC that are calculated after SH2 to stay on

578　the safe side. In terms of the IDC metric, MSIS was found to be the closest to the M2017 in more

579　events than the other models studied here. This may be due to the fact that MSIS is trained with

580　the integrated neutral density (Picone et al., 2002).

581　The second part of this study involves selecting a standard set of metrics to quantify the errors in

582　neutral density. Seven metrics were investigated for this purpose: the ratio between the model

583　maximum and CHAMP maximum ($Ratio_{max}$), ratio between the model mean and CHAMP mean

584　($Ratio_{avg}$), time delay between the peak of the model and peak of the CHAMP observation (TD),

585　mean absolute error (MAE), normalized root mean square error (NRMSE), prediction efficiency

586　(PE) and integrated density change (IDC). In this study, we show that $Ratio_{max}$ and $Ratio_{avg}$ may

587　not be consistent with each other even after the baseline shifting procedure. A model

588　overestimating the ratio of maximum may predict the $Ratio_{avg}$ well. This is due to the shape of

589　the response curve and is controlled by how fast the growth and decay rates of the neutral density

590　are within the model. Thus, neither the neutral density maximum nor the neutral density average

591　is definitive in model performance assessment when used alone. In this study, consistency is

592　achieved between the skill scores MAE, RMSE, NRMSE and PE after the baseline shifts.

593　Consequently, we suggest using MAE, NRMSE and PE together for the neutral density

594  evaluations. MAE will provide the mean amount of error, NRMSE, the error percentage with

595  respect to the event and PE will provide how efficient the model is in capturing the variability

596  and mean of the neutral density observations.

597  To conclude, we have shown that baseline shifting is useful in assessing the storm-time model

598  performance when models have bias against the data during the quiet-time. Removal of the

599  baseline allows for the detection of actual storm-time response and performances from the

600  models. We emphasize that quiet-time climatology and storm-time performances of the models

601  should be evaluated separately, and after baseline shifts, especially for the models with quiet-

602  time bias. Even though, we focused on the storm-time performances of the models in this work,

603  we emphasize that for the long-term estimations of satellite drag, it is important to provide the

604  background neutral density precisely.

605  For satellite drag calculations, the accuracy of neutral density estimations is important. This

606  study shows, the IT models present variable errors depending on the event. None of the models

607  perform perfectly for all cases. In such cases, the uncertainty in thermospheric neutral density in

608  an event can be represented well by using an ensemble of models and iterating the results

609  (Elvidge et al., 2016). In an operational scenario, the ensemble method and baseline shifts using

610  the previous, quiet-day estimations can be used together to tune the models and their output, so

611  that the storm-time variations can be better estimated. Murray (2018) demonstrated the

612  usefulness of ensembles in space weather forecasting to determine the uncertainty and (Knipp,

613  2016) reported the studies, which use the ensemble method for space weather forecasting. We

614  also point out that multi-model ensemble forecasts can be of great use and are candidates for

615  future work, especially in respect of the integrated density change, maximum and average neutral

616 density which are found to be highly variable among the models and are important in satellite

617 drag calculations and for real-time operations.

638    **References**

639    Anderson, R. L., Born, G. H., & Forbes, J. M. (2009). Sensitivity of Orbit Predictions to Density

640    Variability. Journal of Spacecraft and Rockets, 46(6), 1214–1230.

641    https://doi.org/10.2514/1.42138

642    Bowman, B. R., Tobiska, W. K., Marcos, F. A., Huang, C. Y., Lin, C. S., & Burke, W. J. (2008).

643    37 TH COSPAR Scientific Assembly 2008 – Proposal for CIRA 2008, A New Empirical

644    Thermospheric Density Model JB2008 Using New Solar and Geomagnetic Indices, (August).

645    Bruinsma, S. L., & Forbes, J. M. (2010). Anomalous behavior of the thermosphere during solar

646    minimum observed by CHAMP and GRACE. Journal of Geophysical Research: Space Physics,

647    115(11), 1–8. https://doi.org/10.1029/2010JA015605.

648    Bruinsma, S., Sutton, E., Solomon, S. C., Fuller                                 -Rowell, T., & Fedr

649    weather modeling capabilities assessment: Neutral density for orbit determination at low Earth

650    orbit. Space Weather, 16. https://doi.org/10.1029/2018SW002027.

651    Burke, W. J., Huang, C. Y., Marcos, F. A., and Wise, J. O. (2007). Interplanetary control of

652    thermospheric densities during large magnetic storms, J. Atmos. Solar-Terr. Phys., 69, 279–287.

653    Bussy-Virat, C. D., Ridley, A. J., & Getchius, J. W. (2018). Effects of Uncertainties in the

654    Atmospheric Density on the Probability of Collision Between Space Objects. Space Weather,

655    519–537. https://doi.org/10.1029/2017SW001705

656    Chai, T., & Draxler, R. R. (2014). Root mean square error (RMSE) or mean absolute error

657    (MAE)? – Arguments against avoiding RMSE in the literature. Geoscientific Model

658    Development, 7(3), 1247–1250. https://doi.org/10.5194/gmd-7-1247-2014

659    Chen, G., Xu, J., Wang, W., Lei, J., & Burns, A. G. (2012). A comparison of the effects of CIR-

660    and CME-induced geomagnetic activity on thermospheric densities and spacecraft orbits: Case

661 studies. Journal of Geophysical Research: Space Physics, 117(A8).

662 https://doi.org/10.1029/2012JA017782

663 Chen, G., Xu, J., Wang, W., & Burns, A. G. (2014). A comparison of the effects of CIR- and

664 CME-induced geomagnetic activity on thermospheric densities and spacecraft orbits: Statistical

665 studies. Journal of Geophysical Research: Space Physics, 119(9), 7928–7939.

666 https://doi.org/10.1002/2014JA019831

667 Codrescu, M. V., Fuller-Rowell, T. J., Munteanu, V., Minter, C. F., & Millward, G. H. (2008).

668 Validation of the coupled thermosphere ionosphere plasmasphere electrodynamics model:

669 CTIPE-mass spectrometer incoherent scatter temperature comparison. Space Weather, 6(9), 1–

670 10. https://doi.org/10.1029/2007SW000364

671 Codrescu, M. V., Negrea, C., Fedrizzi, M., Fuller-Rowell, T. J., Dobin, A., Jakowsky, N., …

672 Maruyama, N. (2012). A real-time run of the Coupled Thermosphere Ionosphere Plasmasphere

673 Electrodynamics (CTIPe) model. Space Weather, 10(1), 1–10.

674 https://doi.org/10.1029/2011SW000736

675 Connor, H. K., Zesta, E., Fedrizzi, M., Shi, Y., Raeder, J., Codrescu, M. V., & Fuller-Rowell, T.

676 J. (2016). Modeling the ionosphere-thermosphere response to a geomagnetic storm using

677 physics-based magnetospheric energy input: OpenGGCM-CTIM results. Journal of Space

678 Weather and Space Climate, 6, A25. https://doi.org/10.1051/swsc/2016019

679 Deng, Y., Fuller-Rowell, T. J., Ridley, A. J., Knipp, D., & Lopez, R. E. (2013). Theoretical

680 study: Influence of different energy sources on the cusp neutral density enhancement. Journal of

681 Geophysical Research: Space Physics, 118(5), 2340–2349. https://doi.org/10.1002/jgra.50197

682   Elvidge, S., Angling, M. J., & Nava, B. (2014). On the use of modified Taylor diagrams to

683   compare ionospheric assimilation models. Radio Science, 49(9), 737–745.

684   https://doi.org/10.1002/2014RS005435

685   Elvidge, S., Godinez, H. C., & Angling, M. J. (2016). Improved forecasting of thermospheric

686   densities using multi-model ensembles. Geoscientific Model Development, 9(6), 2279–2292.

687   https://doi.org/10.5194/gmd-9-2279-2016

688   Emmert, J. T., Warren, H. P., Segerman, A. M., Byers, J. M., & Picone, J. M. (2017).

689   Propagation of atmospheric density errors to satellite orbits. Advances in Space Research, 59(1),

690   147–165. https://doi.org/10.1016/j.asr.2016.07.036

691   Fedrizzi, M., Fuller-Rowell, T. J., & Codrescu, M. V. (2012). Global Joule heating index

692   derived from thermospheric density physics-based modeling and observations. Space Weather,

693   10(3). http://doi.org/10.1029/2011SW000724

694   Hejduk, M. D., & Snow, D. E. (2018). The Effect of Neutral Density Estimation Errors on

695   Satellite Conjunction Serious Event Rates. Space Weather.

696   https://doi.org/10.1029/2017SW001720

697   Huang, C. Y., Su, Y.-J., Sutton, E. K., Weimer, D. R., & Davidson, R. L. (2014). Energy

698   coupling during the August 2011 magnetic storm. Journal of Geophysical Research: Space

699   Physics, 119(2), 1219–1232. https://doi.org/10.1002/2013JA019297

700   Kim, K. H., Moon, Y. J., Cho, K. S., Kim, H. D., & Park, J. Y. (2006). Atmospheric drag effects

701   on the KOMPSAT-1 satellite during geomagnetic superstorms. Earth, Planets and Space,

702   58(12), 25–28. https://doi.org/10.1186/BF03351968

703 Knipp, D., Kilcommons, L., Hunt, L., Mlynczak, M., Pilipenko, V., Bowman, B., & Deng, Y.

704 (2013). Thermospheric damping response to sheath-enhanced geospace storms. Geophysical

705 Research Letters, 40(7), 1263–1267. https://doi.org/10.1002/grl.50197

706 Knipp, D. J. (2016), Advances in Space Weather Ensemble Forecasting, Space Weather, 14, 52–

707 53, doi:10.1002/2016SW001366.

708 Kwak, Y. S., Richmond, A. D., Deng, Y., Forbes, J. M., & Kim, K. H. (2009). Dependence of

709 the high-latitude thermospheric densities on the interplanetary magnetic field. Journal of

710 Geophysical Research: Space Physics, 114(5), 1–7. https://doi.org/10.1029/2008JA013882

711 Lathuillère, C., Menvielle, M., Marchaudon, A., & Bruinsma, S. (2008). A statistical study of

712 the observed and modeled global thermosphere response to magnetic activity at middle and low

713 latitudes. Journal of Geophysical Research: Space Physics, 113(7), 1–9.

714 https://doi.org/10.1029/2007JA012991

715 Lei, J., Thayer, J. P., Lu, G., Burns, A. G., Wang, W., Sutton, E. K., & Emery, B. A. (2011).

716 Rapid recovery of thermosphere density during the October 2003 geomagnetic storms. Journal

717 of Geophysical Research: Space Physics, 116(3), 1–10. https://doi.org/10.1029/2010JA016164

718 Liu, H., H. Lühr, V. Henize, and W. Köhler (2005), Global distribution of the thermospheric

719 total mass density derived from CHAMP, J. Geophys. Res.,110, A04301,

720 doi:10.1029/2004JA010741.

721 Lühr, H., Rother, M., Köhler, W., Ritter, P., & Grunwaldt, L. (2004). Thermospheric up-welling

722 in the cusp region: Evidence from CHAMP observations. Geophysical Research Letters, 31(6).

723 https://doi.org/10.1029/2003GL019314

724 Lühr, H., Liu, H., Park, J., & Müller, S. (2011). New Aspects of the Coupling Between

725 Thermosphere and Ionosphere, with Special regards to CHAMP Mission Results. In M. A.

Author Manuscript

726     Abdu & D. Pancheva (Eds.), Aeronomy of the Earth's Atmosphere and Ionosphere (pp. 303–

727     316). Dordrecht: Springer Netherlands. https://doi.org/10.1007/978-94-007-0326-1

728     McGranaghan, R., Knipp, D. J., McPherron, R. L., & Hunt, L. A. (2014). Impact of equinoctial

729     high-speed stream structures on thermospheric responses. Space Weather, 12(4), 277–297.

730     https://doi.org/10.1002/2014SW001045

731     Mehta, P. M., Walker, A. C., Sutton, E. K., & Godinez, H. C. (2017). New density estimates

732     derived using accelerometers on board the CHAMP and GRACE satellites. Space Weather,

733     15(4), 558–576. https://doi.org/10.1002/2016SW001562

734     Millward, G. H., Müller-Wodarg, I. C. F., Aylward, A. D., Fuller-Rowell, T. J., Richmond, A.

735     D., & Moffett, R. J. (2001). An investigation into the influence of tidal forcing on F region

736     equatorial vertical ion drift using a global ionosphere-thermosphere model with coupled

737     electrodynamics. Journal of Geophysical Research, 106(A11), 24733.

738     https://doi.org/10.1029/2000JA000342

739     Murray, S. A. (2018). The importance of ensemble techniques for operational space weather

740     forecasting. Space Weather, 1–10. https://doi.org/10.1029/2018SW001861

741     Pardini, C., Moe, K., & Anselmo, L. (2012). Thermospheric density model biases at the 23rd

742     sunspot maximum. Planetary and Space Science, 67(1), 130–146.

743     https://doi.org/10.1016/j.pss.2012.03.004

744     Picone, J. M., Hedin, A. E., Drob, D. P., & Aikin, A. C. (2002). NRLMSISE-00 empirical

745     model of the atmosphere: Statistical comparisons and scientific issues. Journal of Geophysical

746     Research: Space Physics, 107(A12), SIA 15-1-SIA 15-16.

747     https://doi.org/10.1029/2002JA009430

748     Prölss, G. W. (2011). Density Perturbations in the Upper Atmosphere Caused by the Dissipation

749     of Solar Wind Energy. Surveys in Geophysics, 32(2), 101–195. https://doi.org/10.1007/s10712-

750     010-9104-0

751     Qian, L., Solomon, S. C., Roble, R. G., Bowman, B. R., & Marcos, F. A. (2008). Thermospheric

752     neutral density response to solar forcing. Advances in Space Research, 42(5), 926–932.

753     https://doi.org/10.1016/j.asr.2007.10.019

754     Qian, L., & Solomon, S. C. (2012). Thermospheric Density: An Overview of Temporal and

755     Spatial Variations. Space Science Reviews, 168(1–4), 147–173. https://doi.org/10.1007/s11214-

756     011-9810-z.

757     Rhoden, E. A., Forbes, J. M., & Marcos, F. A. (2000). The influence of geomagnetic and solar

758     variabilities on lower thermosphere density. Journal of Atmospheric and Solar-Terrestrial

759     Physics, 62(11), 999–1013. https://doi.org/10.1016/S1364-6826(00)00066-3.

760     Richmond, A. D., Ridley, E. C., & Roble, R. G. (1992). A thermosphere/ionosphere general

761     circulation model with coupled electrodynamics. Geophysical Research Letters, 19(6), 601–604.

762     https://doi.org/10.1029/92GL00401

763     Ridley, A. J., Deng, Y., & Tóth, G. (2006). The global ionosphere–thermosphere model. Journal

764     of Atmospheric and Solar-Terrestrial Physics, 68(8), 839–864.

765     https://doi.org/10.1016/j.jastp.2006.01.008

766     Shim, J. S., et al. (2011), CEDAR Electrodynamics Thermosphere Ionosphere (ETI) Challenge

767     for systematic assessment of ionosphere/thermosphere models: NmF2, hmF2, and vertical drift

768     using ground-based observations, Space Weather, 9, S12003, doi:10.1029/2011SW000727.

769     Shim, J. S., Kuznetsova, M., Rastätter, L., Bilitza, D., Butala, M., Codrescu, M., … Sutton, E.

770     (2012). CEDAR Electrodynamics Thermosphere Ionosphere (ETI) Challenge for systematic

771    assessment of ionosphere/thermosphere models: Electron density, neutral density, NmF2, and

772    hmF2 using space-based observations. Space Weather, 10(10).

773    https://doi.org/10.1029/2012SW000851

774    Shim, J. S., Kuznetsova, M., Rastätter, L., Bilitza, D., Butala, M., Codrescu, M., … Sutton, E.

775    (2014). Systematic Evaluation of Ionosphere/Thermosphere (IT) Models (pp. 145–160).

776    American Geophysical Union (AGU). https://doi.org/10.1002/9781118704417.ch13

777    Shim, J. S., L. Rastaetter, K. M. Kuznetsova, E. C. Kalafatoglu, and Y. Zheng (2015).

778    Assessment of the predictive capability of IT models at the Community Coordinated Modeling

779    Center. Presented at Ionospheric Effect Symposium, Alexandria VA.

780    Solomon, S. C., Qian, L., Didkovsky, L. V., Viereck, R. A., & Woods, T. N. (2011). Causes of

781    low thermospheric density during the 2007-2009 solar minimum. Journal of Geophysical

782    Research: Space Physics, 116(7), 1–14. https://doi.org/10.1029/2011JA016508

783    Storz, M. F., Bowman, B. R., Branson,M. J. I., Casali, S. J., & Tobiska,W. K. (2005). High

784    accuracy satellite drag model (HASDSM). Advances in Space Research, 36(12), 2497–2505.

785    Sutton, E. K. (2008), Effects of solar disturbances on the thermosphere densities and winds from

786    CHAMP and GRACE satellite accelerometer data, Doctoral Dissertation, Dept. of Aerosp. Eng.

787    Sci., Univ. of Colorado, Boulder, Colo.

788    Sutton, E. K. (2009), Normalized force coefficients for satellites with elongated shapes, J.

789    Spacecraft and Rockets, 46(1), doi:10.2514/1.40940.

790    Sutton, E. K. (2011). Accelerometer-Derived Atmospheric Density from the CHAMP and GRACE

791    Satellites. Version 2.3. AIR FORCE RESEARCH LAB KIRTLAND AFB NM.

Sutton, E. K. (2018). A New Method of Physics-Based Data Assimilation for the Quiet and

Disturbed Thermosphere. Space Weather, 16(6), 736–753.

https://doi.org/10.1002/2017SW001785

Sutton, E. K., Forbes, J. M., & Nerem, R. S. (2005). Global thermospheric neutral density and

wind response to the severe 2003 geomagnetic storms from CHAMP accelerometer data.

Journal of Geophysical Research: Space Physics, 110(A9), 1–10.

https://doi.org/10.1029/2004JA010985

Sutton, E. K., Forbes, J. M., Nerem, R. S., & Woods, T. N. (2006). Neutral density response to

the solar flares of October and November, 2003. Geophysical Research Letters, 33(22), 1–5.

https://doi.org/10.1029/2006GL027737

Thayer, J. P., Lei, J., Forbes, J. M., Sutton, E. K., & Nerem, R. S. (2008). Thermospheric

density oscillations due to periodic solar wind high speed streams. Journal of Geophysical

Research: Space Physics, 113(6), A06307. https://doi.org/10.1029/2008JA013190

Xu, J., Wang, W., Lei, J., Sutton, E. K., & Chen, G. (2011). The effect of periodic variations of

thermospheric density on CHAMP and GRACE orbits. Journal of Geophysical Research: Space

Physics, 116(2), A02315. https://doi.org/10.1029/2010JA015995

Weimer, D. R. (2005). Improved ionospheric electrodynamic models and application to

calculating Joule heating rates. Journal of Geophysical Research, 110(A5), A05306.

https://doi.org/10.1029/2004JA010884

Willmott, C. J., & Matsuura, K. (2005). Advantages of the mean absolute error (MAE) over the

root mean square error (RMSE) in assessing average model performance. Climate Research, 30,

79–82.

814  Zesta, E., & Huang, C. Y. (2016). Satellite orbital drag. In G. V. Khazanov (Ed.), Space weather

815  fundamentals (pp. 329–351). Boca Raton, FL: CRC Press.

816  **Table 1.** GEM-CEDAR Challenge events. Table shows the maximum values of geomagnetic and

817  solar indices ($Kp_{max}$, F10.7, $Dst_{max}$, $HP_{max}$) and solar wind drivers of the events.

| Event | $Kp_{max}$ | F10.7 | $Dst_{min}$ (nT) | $HP_{max}$ (GW) | Driver |
|---|---|---|---|---|---|
| 2005-135 | 8+ | 103 | -247 | 1225 | CME |
| 2006-348 | 8+ | 93.6 | -162 | 504 | CME |
| 2005-243 | 7 | 84 | -122 | 260 | HSS |
| 2005-190 | 6+ | 106.6 | -92 | 238 | HSS |
| 2007-142 | 5+ | 72 | -58 | 197 | HSS |
| 2007-091 | 5 | 71.7 | -63 | 286 | HSS |

818

819

820  **Table 2.** Baseline shifts. $\rho_{old}$ is the original orbit-averaged time series whereas $\rho_{new}$ is the

821  baseline shifted time series. Subscript index "n" represents the orbit numbers for the entire event

822  (quiet+storm) interval, "i" stands for the orbit number during the selected quiet time interval of

823  the event. Overbars denote the mean.

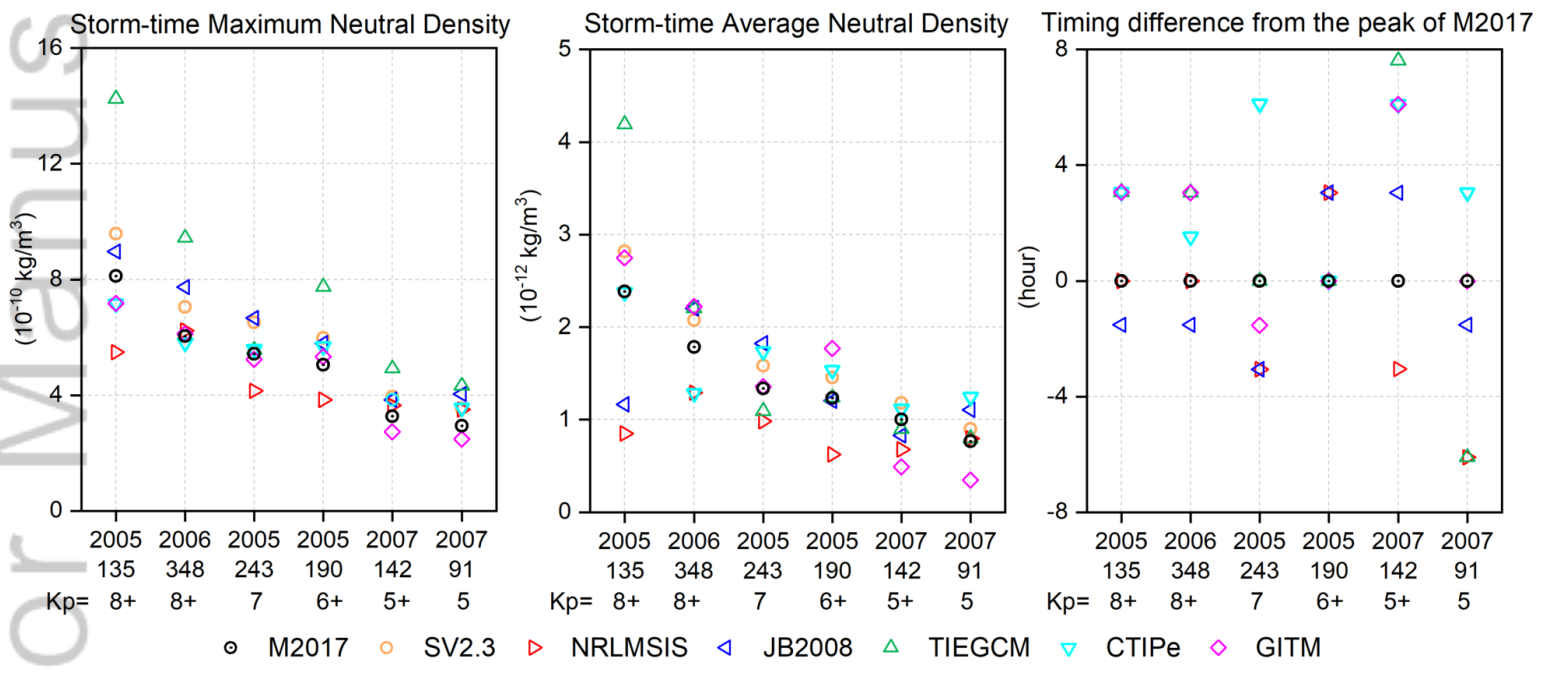| Shifts | Shifting Parameter | Shifted Series | Reference Level |
|---|---|---|---|
| Shift1 (SH1) | $S_1 = \overline{\rho_{champ,i}} - \overline{\rho_{model,i}}$ | $\rho_{new,n} = \rho_{old,n} - S_1$ | CHAMP |
| Shift2 (SH2) | $S_2 = \overline{\rho_{champ,i}}$ for CHAMP $S_2 = \overline{\rho_{models,i}}$ for models | $\rho_{new,n} = \rho_{old,n} - S_2$ | Zero |
| Shift3 (SH3) | $S_3 = \overline{\rho_{champ,i} / \rho_{model,i}}$ | $\rho_{new,n} = \rho_{old,n} \times S_3$ | CHAMP |

824

825 **Figures:**

826 **Figure 1.** From left to right: storm-time maximum in neutral density, storm-time average neutral

827 density, timing difference between the peak of models and M2017. The circles denote neutral

828 density estimations based on accelerometers on CHAMP: orange, SV2.3 and dot-centered black,

829 M2017. The triangles and the diamond show the model estimations: red, right-triangle: MSIS;

830 blue, left-triangle: JB2008; green, up-triangle: TIEGCM; cyan, down-triangle: CTIPe; pink,

831 diamond: GITM. X-label is the events listed from severe (Kp>8) to weak (Kp=5) starting from

832 left to right, according to the NOAA classification based on Kp values.

833 **Figure 2:** An example event: 2006-348. First row, from left to right: a) top: Neutral density from

834 the model and observations without shift; below: Kp and Dst indices, neutral density estimations

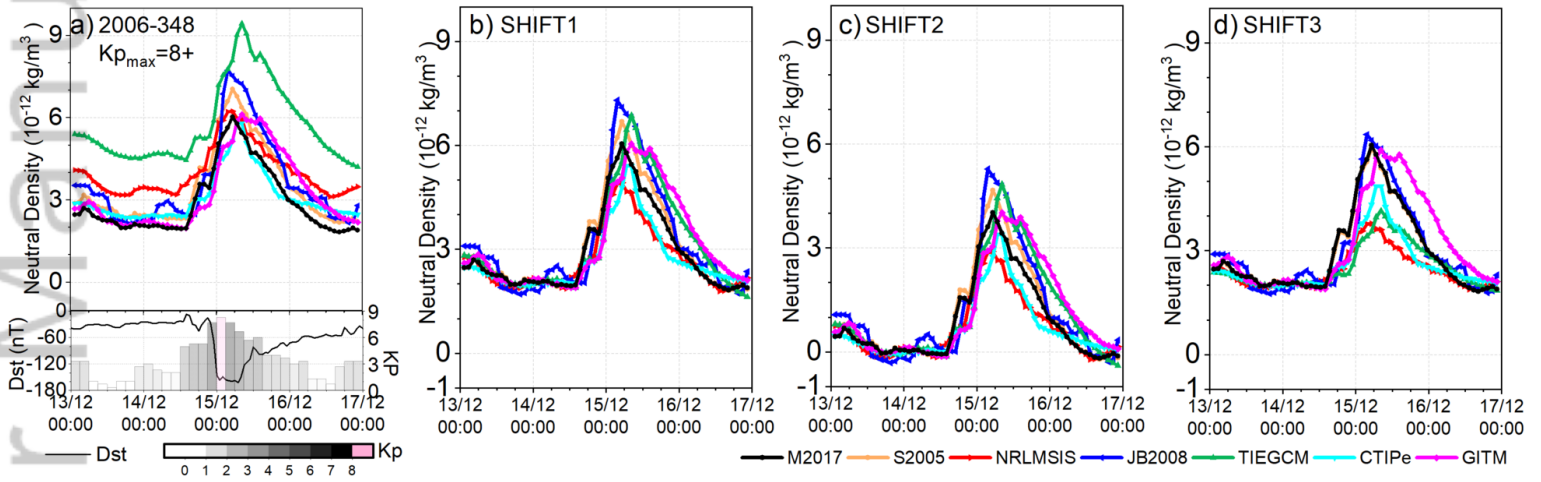835 from the models and M2017 after b) SH1, c) SH2, d) SH3.

836 **Figure 3:** From top to bottom: storm-time ratio of maximum neutral density of the models to

837 M2017, storm-time ratio of average neutral density from the models to M2017, timing difference

838 between the peak of models and M2017. From left to right: SV2.3, MSIS, JB2008, CTIPe,

839 GITM and TIEGCM. O denotes the results for the original, unshifted time series whereas SH1 to

840 SH3 represents the shifts from Shift1 to Shift3. Red symbols represent the severe events with

841 high Kp; cyan denotes strong event with Kp=7; black is for 7>Kp>6; and green color is for weak

842 events with Kp around 5. Circle represents the event 2005-135; square, 2006-348; up-triangle,

843 2005-243; down-triangle, 2005-190; cross, 2007-142; plus, 2007-91.

844 **Figure 4:** From left to right: storm-time orbit and time integrated neutral density, storm-time

845 change in maximum neutral density, storm-time change in mean neutral density. The symbol and
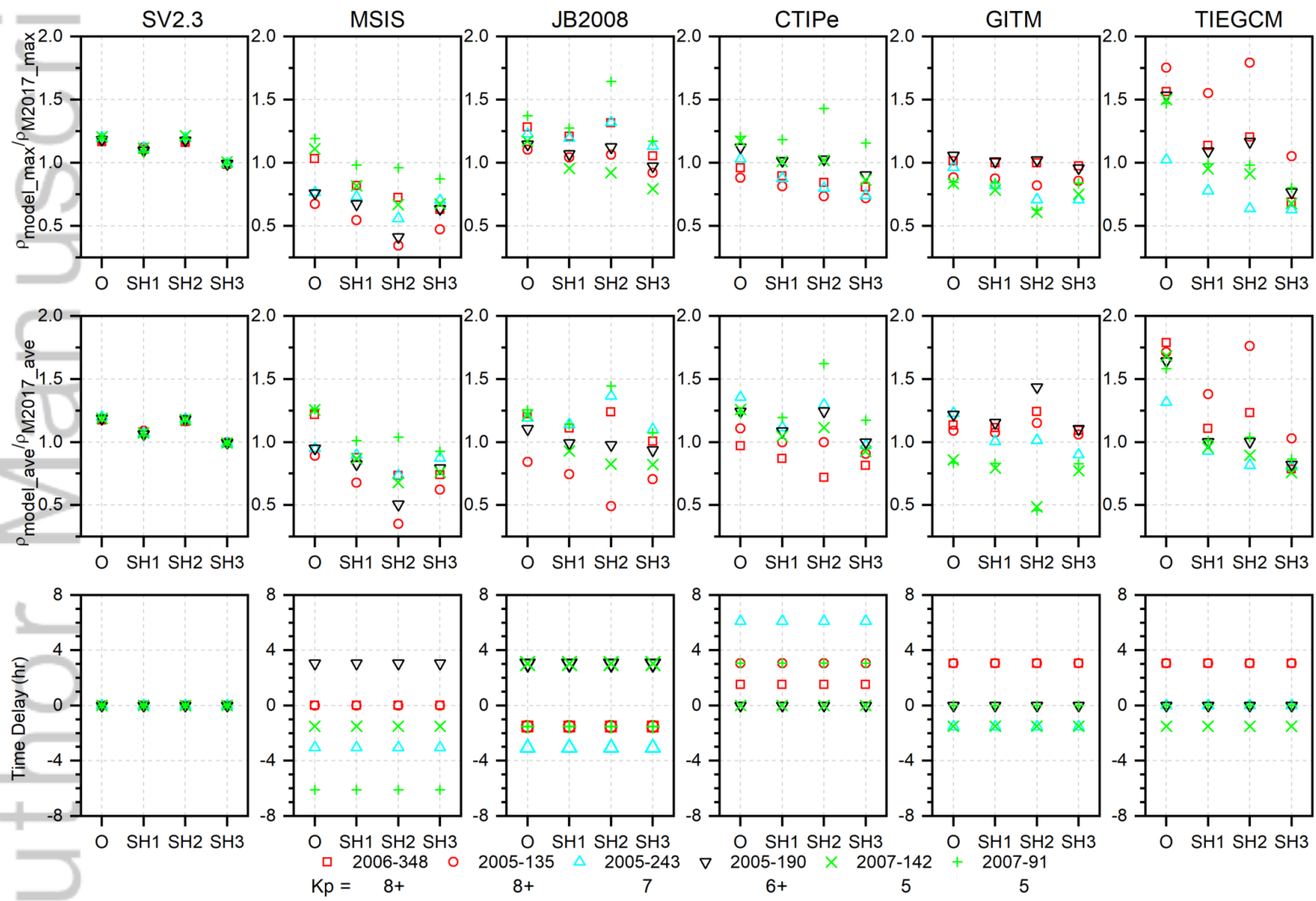
846 colors are the same as Figure 1.

847    **Figure 5:** From top to bottom: MAE, NRMSE, PE. From left to right: SV2.3, MSIS, JB2008,

848    CTIPe, GITM and TIEGCM. KP scales, axis labels, colors and symbols are the same as Figure 3.

849    Please note that the y-axis scales for TIEGCM is different than the other panels for the three

850    parameters. Additionally, for TIEGCM, PE results after the shifts SH1 to SH3 are shown in

851    another frame inside the PE panel with scaled y-axis. The inside frame has the same y-axis scale
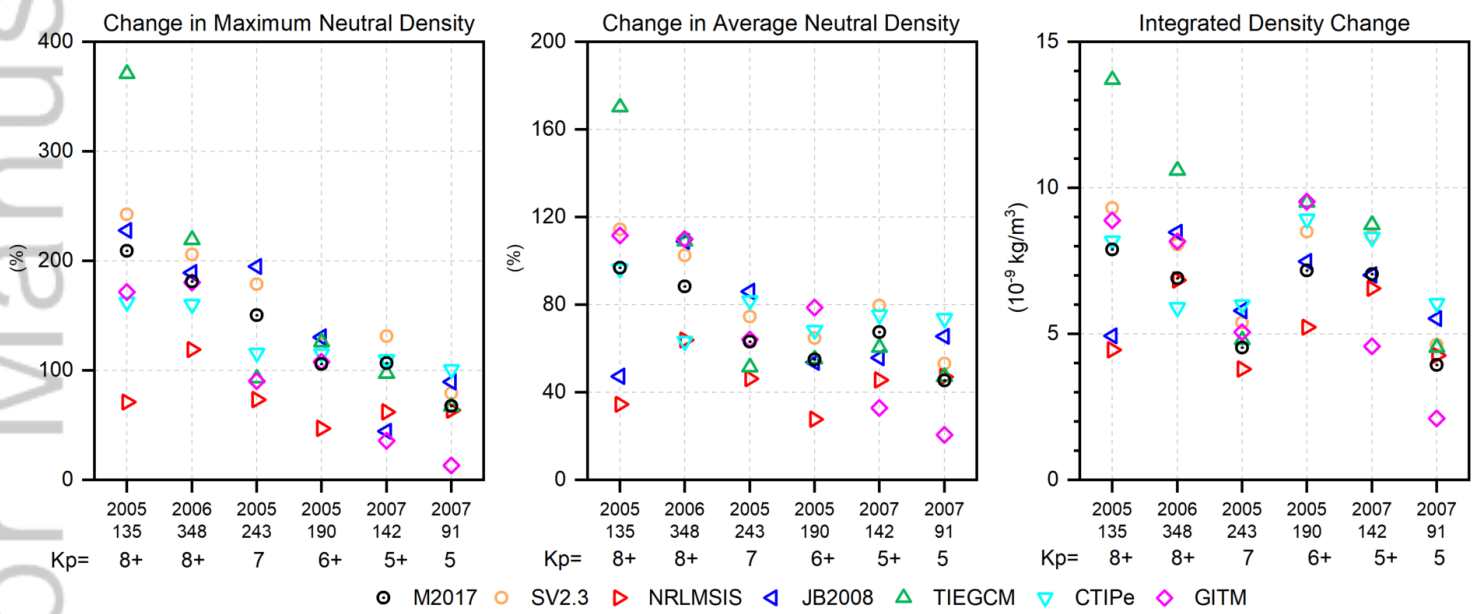
852    as the other panel for PEs.

Storm-time Maximum Neutral Density    Storm-time Average Neutral Density    Timing difference from the peak of M2017

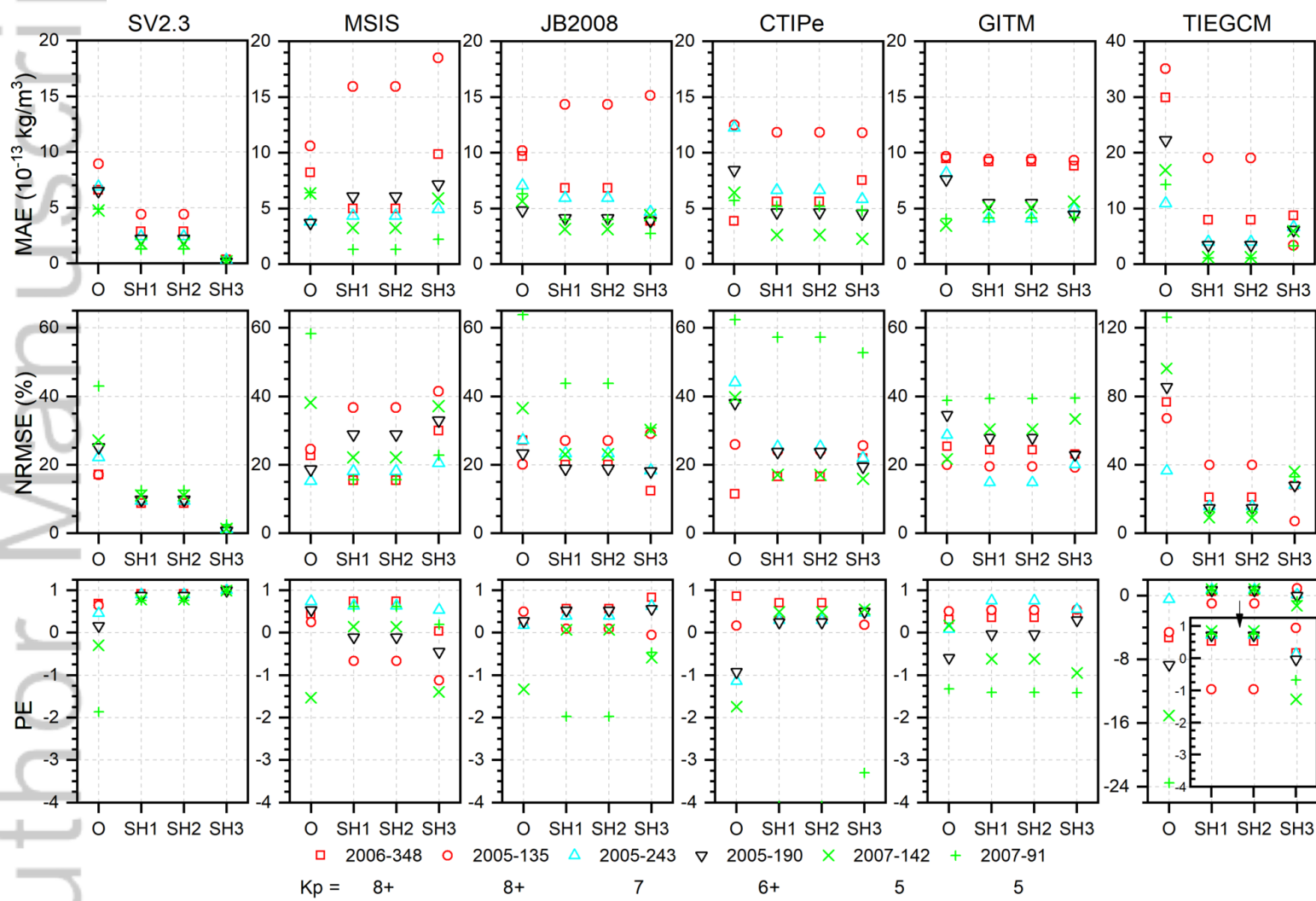Legend: M2017, SV2.3, NRLMSIS, JB2008, TIEGCM, CTIPe, GITM

2018SW002033-f01-z-.png

2018SW002033-f02-z-.png

2018SW002033-f03-z-.png

2018SW002033-f04-z-.png

| | | | | | |
|---|---|---|---|---|---|
| □ 2006-348 | ○ 2005-135 | △ 2005-243 | ▽ 2005-190 | ✕ 2007-142 | + 2007-91 |
| Kp = 8+ | 8+ | 7 | 6+ | 5 | 5 |

2018SW002033-f05-z-.png