



RESEARCH ARTICLE

10.1029/2022SW003388

Validation of Ionospheric Specifications During Geomagnetic Storms: TEC and foF2 During the 2013 March Storm Event-II

Key Points:

- F2-layer critical frequency/Total Electron Content (foF2/TEC) and their changes during a storm predicted by ionosphere-thermosphere coupled models are evaluated against Global Ionosphere Radio Observatory foF2 and GPS TEC measurements
- Model simulations tend to underestimate the storm-time enhancements of foF2 and TEC and to predict them better in the northern hemisphere
- Ensemble of all simulations for TEC is comparable to the data assimilation model (USU-GAIM)

Supporting Information:

Supporting Information may be found in the online version of this article.





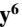








Correspondence to:

I.-S. Song,
songi@yonsei.ac.kr

Citation:

Shim, J. S., Song, I.-S., Jee, G., Kwak, Y.-S., Tsagouri, I., Goncharenko, L., et al. (2023). Validation of ionospheric specifications during geomagnetic storms: TEC and foF2 during the 2013 March storm event-II. *Space Weather*, 21, e2022SW003388. <https://doi.org/10.1029/2022SW003388>

Received 13 DEC 2022
Accepted 16 MAR 2023

J. S. Shim¹, I.-S. Song¹ , G. Jee² , Y.-S. Kwak³ , I. Tsagouri⁴ , L. Goncharenko⁵, J. McInerney⁶ , A. Vitt⁶ , L. Rastaetter⁷ , J. Yue^{7,8} , M. Chou^{7,8} , M. Codrescu⁹, A. J. Coster⁵, M. Fedrizzi⁹ , T. J. Fuller-Rowell⁹, A. J. Ridley¹⁰ , S. C. Solomon⁶ , and J. B. Habarulema¹¹ 

¹Department of Atmospheric Sciences, Yonsei University, Seoul, South Korea, ²Division of Atmospheric Sciences, Korea Polar Research Institute, Incheon, South Korea, ³Space Science Division, Korea Astronomy and Space Science Institute, Daejeon, South Korea, ⁴National Observatory of Athens, Penteli, Greece, ⁵Haystack Observatory, Westford, MA, USA, ⁶High Altitude Observatory, NCAR, Boulder, CO, USA, ⁷NASA GSFC, Greenbelt, MD, USA, ⁸Catholic University of America, Washington, DC, USA, ⁹NOAA SWPC, Boulder, CO, USA, ¹⁰Space Physics Research Laboratory, University of Michigan, Ann Arbor, MI, USA, ¹¹South African National Space Agency (SANSA) Space Science, Hermanus, South Africa

Abstract Assessing space weather modeling capability is a key element in improving existing models and developing new ones. In order to track improvement of the models and investigate impacts of forcing, from the lower atmosphere below and from the magnetosphere above, on the performance of ionosphere-thermosphere models, we expand our previous assessment for 2013 March storm event (Shim et al., 2018, <https://doi.org/10.1029/2018SW002034>). In this study, we evaluate new simulations from upgraded models (the Coupled Thermosphere Ionosphere Plasmasphere Electrodynamics (CTIPE) model version 4.1 and the Global Ionosphere Thermosphere Model (GITM) version 21.11) and from the NCAR Whole Atmosphere Community Climate Model with thermosphere and ionosphere extension (WACCM-X) version 2.2 including eight simulations in the previous study. A simulation from the NCAR Thermosphere-Ionosphere-Electrodynamics General Circulation Model version 2 (TIE-GCM 2.0) is also included for comparison with WACCM-X. TEC and foF2 changes from quiet-time background are considered to evaluate the model performance on the storm impacts. For evaluation, we employ four skill scores: Correlation coefficient (CC), root-mean square error (RMSE), ratio of the modeled to observed maximum percentage changes (Yield), and timing error (TE). It is found that the models tend to underestimate the storm-time enhancements of foF2 (F2-layer critical frequency) and TEC (Total Electron Content) and to predict foF2 and/or TEC better in North America but worse in the Southern Hemisphere. The ensemble simulation for TEC is comparable to results from a data assimilation model (Utah State University-Global Assimilation of Ionospheric Measurements (USU-GAIM)) with differences in skill score less than 3% and 6% for CC and RMSE, respectively.

Plain Language Summary The Earth's ionosphere-thermosphere (IT) system, which is present between the lower atmosphere and the magnetosphere, is highly variable due to external forcings from below and above as well as internal forcings mainly associated with ion-neutral coupling processes. The variabilities of the IT system can adversely affect our daily lives, therefore, there is a need for both accurate and reliable weather forecasts to mitigate harmful effects of space weather events. In order to track the improvement of predictive capabilities of space weather models for the IT system, and to investigate the impacts of the forcings on the performance of IT models, we evaluate new simulations from upgraded models (Coupled Thermosphere Ionosphere Plasmasphere Electrodynamics model version 4.1 and Global Ionosphere Thermosphere Model version 21.11) and from NCAR Whole Atmosphere Community Climate Model with thermosphere and ionosphere extension (WACCM-X) version 2.2 together with 8 simulations in the previous study. A simulation of NCAR Thermosphere-Ionosphere-Electrodynamics General Circulation Model version 2 is also included for the comparison with WACCM-X. Quantitative evaluation is performed by using four skill scores including Correlation coefficient, root-mean square error, ratio of the modeled to observed maximum percentage changes (Yield), and timing error. The findings of this study will provide a baseline for future validation studies of new and improved models.

© 2023. The Authors. Space Weather published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

1. Introduction

Variabilities of the Earth's ionosphere-thermosphere (IT) system, caused by charged particles and electromagnetic radiation emitted from the sun, can adversely affect our daily lives, which are highly dependent on space-based technological infrastructures such as Low-Earth Orbit (LEO) satellites and the Global Navigation Satellite System (GNSS). To mitigate harmful effects of space weather events, modeling plays a critical role in our quest to understand the connection between solar eruptive phenomena and their impacts in interplanetary space and near-Earth space environment. In particular, the Earth's upper atmosphere including the IT system is the space environment closest to human society. Thus, during the past few decades, first-principles physics-based (PB) IT models have been developed for specifications and forecasts of the near-Earth space environment. In addition, there have been recent developments of whole atmosphere models with a thermospheric and ionospheric extension to fully understand variabilities of the IT system by considering coupling between the IT system and the lower atmosphere (e.g., Akmaev, 2011; T. Fuller-Rowell et al., 2010; Jin et al., 2011; Liu et al., 2018).

For more accurate space weather forecasting, assessing space weather modeling capability is a key element to improve existing models and to develop new models. Over the last decade, in an effort to address the needs and challenges of the assessment of our current knowledge about space weather effects on the IT system and the current state of IT modeling capabilities, the NASA GSFC Community Coordinated Modeling Center (CCMC) has been supporting community-wide model validation projects, including Coupling, Energetics and Dynamics of Atmospheric Regions (CEDAR) (Shim et al., 2011, 2012, 2014) and Geospace Environment Modeling (GEM)-CEDAR modeling challenges (Rastätter et al., 2016; Shim, Rastätter, et al., 2017).

Furthermore, in 2018, the CCMC established an international effort, the “International Forum for Space Weather Modeling Capabilities Assessment,” to evaluate and assess the predictive capabilities of space weather models (<https://ccmc.gsfc.nasa.gov/iswat/IFSWCA/>). As a result of this international effort, four ionosphere/thermosphere working groups were established with an overarching goal to devise a standardized quantitative validation procedure for IT models (Scherliess et al., 2019).

The working group, focusing on neutral density and orbit determination in LEO, reported their initial results for specific metrics for thermosphere model assessment over the selected three full years and two geomagnetic storms in 2005 (Bruinsma et al., 2018). They reported that the tested models in general performed reasonably well, although seasonal errors were sometimes observed and impulsive geomagnetic events remain a challenge. Kalafatoglu Eyiguler et al. (2019) compared the neutral density estimates from two empirical and three PB models with those obtained from the CHAMP satellite. They suggested that several metrics that provide different aspects of the errors should be considered together for a proper performance evaluation.

Another working group, the “Ionosphere Plasmasphere Density Working Team,” performed the assessment of present modeling capabilities in predicting the ionospheric climatology of foF₂ and hmF₂ for the entire year of 2012 (Tsgouri et al., 2018). Tsgouri et al. (2018) identified a strong seasonal and local time dependence of the model performances, especially for PB models, which could provide useful insight for future model improvements. Tsgouri et al. (2018) cautioned that the quality of the ground truth data may play a key role in testing the model performance. Shim et al. (2018) assessed how well the ionospheric models predict storm time foF₂ and TEC by considering quantities, such as TEC and foF₂ changes and percentage changes compared to quiet time background, at 12 selected midlatitude locations in the American and European-African longitude sectors. They found that the performance of the model varies with location, even within a localized region like Europe, as well as with the metrics considered.

In this paper, we expand our previous assessment of modeled foF₂ and TEC during 2013 March storm event (17 March 2013) (Shim et al., 2018) to track improvement of the models and to investigate impacts of forcings from the lower atmosphere below and from the magnetosphere above on the performance of IT models. For this study, we evaluate the updated version of the coupled IT models available at the CCMC (Webb et al., 2009) since our previous study (Shim et al., 2018): Coupled Thermosphere Ionosphere Plasmasphere Electrodynamics (CTIPE) version 4.1 and Global Ionosphere Thermosphere Model (GITM) version 21.11. However, the other types of models such as empirical models, stand-alone ionospheric models, and data assimilation models are not included. In addition, for the first time, simulations from the NCAR Whole Atmosphere Community Climate Model with thermosphere and ionosphere extension (WACCM-X) 2.2 are included in our assessment. We also include a simulation from the NCAR TIE-GCM 2.0 to compare with results from WACCM-X 2.2. For TEC prediction, we compare a weighted mean of the ensemble of all 13 simulations (ensemble average), including eight simulations from our previous study with individual simulations to assess ensemble forecast capability. In Section 2, we

Table 1
Quantities and Skill Scores for Model-Data Comparison

Quantities and skill scores for model-data comparison	
Quiet time references	30-Day median value at a given time: TEC _{quiet} (UT) 30 Day consist of 15 days before (03/01–03/15/2013) and 15 days after (03/22–04/05/2013) the storm
Shifted TEC/foF2	For example, TEC*(doy, UT) = TEC(doy, UT) – minimum of TEC _{quiet} (UT)
TEC/foF2 changes w.r.t. the quiet time	For example, dTEC(doy, UT) = TEC(doy, UT) – TEC _{quiet} (UT)
TEC/foF2 percentage changes w.r.t. the quiet time	For example, dTEC[%](doy, UT) = 100* dTEC(doy, UT)/TEC _{quiet} (UT)
Normalized Percentage changes of TEC	dTEC[%] _{norm} = (dTEC[%] – ave_dTEC[%])/std_dTEC[%] ave_dTEC[%] is the average of dTEC[%] at a given time and at a given location over the quiet 30 days std_dTEC[%] is the standard deviation of the average percentage change
Skill Scores	
CC	Correlation Coefficient
RMSE	Root-Mean-Square Error = $\left(\sqrt{\frac{\sum (x_{\text{obs}} - x_{\text{mod}})^2}{N}} \right)$, where x_{obs} and x_{mod} are observed and modeled values
Yield	Ratio of the peak of modeled percentage change to that of the observed one = $\left(\frac{(x_{\text{mod}})_{\text{max}}}{(x_{\text{obs}})_{\text{max}}} \right)$
Timing Error (TE)	Difference between the modeled peak time and observed peak time: TE = t _{peak_model} – t _{peak_obs}

briefly describe observations, models, and metrics used for this study. Section 3 presents the results of model-data comparisons and performance of the models are presented. Section 4 shows comparisons of ensemble of TEC predictions with the individual simulations based on the skill scores used in this study. In Section 5, we summarize and discuss our results. Finally, we conclude in Section 6.

2. Methodology

2.1. Observations and Metrics

We use the foF2 and TEC measurements at 12 ionosonde stations selected in middle latitudes: eight northern hemisphere (NH) stations in the US (Millstone Hill, Idaho National Laboratory, Boulder, and Eglin AFB) and Europe (Chilton, Pruhonice, Ebre, and Athens) and 4 southern hemisphere (SH) stations in South America (SAM) (Port Stanley) and South Africa (SAF) (Louisvale, Hermanus, and Grahamstown) (Figure 1 and Table 1 in Shim et al. (2018) for details). The foF2 and GNSS vertical TEC (vTEC) data are provided by the Global Ionosphere Radio Observatory (GIRO) (<http://giro.uml.edu/>) (Reinisch & Galkin, 2011) and by the MIT Haystack Observatory (<http://cedar.openmadrigal.org/>) (Rideout & Coster, 2006), respectively.

Table 1 shows the quantities and skill scores calculated for the model-data comparison. To remove potential systematic uncertainties in the models and observations and baseline differences among the models and between models and observations, we use the shifted values and changes from their own quiet-time background values (e.g., shifted TEC (TEC*) = TEC (UT) on a particular DOY – minimum of 30-day median). Furthermore, using these quantities likely reduce the impacts of differing upper boundaries for TEC calculations, since the plasmaspheric TEC variations with geomagnetic activity are negligible in middle latitudes (Shim, Jee, & Scherliess, 2017).

To measure how well the observed and modeled values are linearly correlated (in phase) with each other and how different the values are on average over the time interval considered, Correlation coefficient (CC) and root-mean square error (RMSE) are calculated, respectively, for the error values below 95th percentile. We also calculate Yield and timing error (TE) to measure the models' capability to capture peak disturbances during the storm. For more detailed information on the quantities and skill scores used for the study, refer to Section 2 in Shim et al. (2018).

2.2. Models and Simulations

The simulations used in this study are obtained from the updated and newly incorporated coupled ionosphere-thermosphere models available at the CCMC (Webb et al., 2009) since our previous study (Shim

et al., 2018): CTIPe 4.1, GITM 21.11, and WACCM-X 2.2. The WACCM-X 2.2 simulations are provided by NCAR HAO. The WACCM-X version 2 (Liu et al., 2018) is a comprehensive numerical model that extends the atmospheric component model of the NCAR Community Earth System Model (CESM) (Hurrell et al., 2013) into the thermosphere up to 500–700 km altitude. WACCM-X is uniquely capable of being run in a configuration where the atmosphere is coupled to active or prescribed ocean, sea ice, and land components, enabling studies of thermospheric and ionospheric weather and climate. WACCM-X version 2 is based upon WACCM version 6 (Gettelman et al., 2019) with a top boundary of ~ 130 km, which is built upon the Community Atmosphere Model (CAM) version 6 having a top boundary of ~ 40 km. WACCM-X 2.2 includes WACCM6 physics for middle atmosphere and lower thermosphere as well as CAM6 physics for the troposphere and the lower stratosphere, and it fully incorporates the electrodynamical processes related to low-to mid-latitude wind dynamo that is implemented in the NCAR TIE-GCM. For this study, two specified-dynamics (SD) WACCM-X 2.2 simulations with different high-latitude electrostatic potential models (Heelis et al., 1982; Weimer, 2005) are used. The SD simulations are carried out by constraining the model's lower atmospheric neutral dynamics using meteorological reanalysis data. The constraining process is achieved by nudging the model toward MERRA-2 (Modern Era Retrospective Analysis for Research and Applications, Version 2) data (Gelaro et al., 2017) below around the altitude of 50 km in a way presented by Brakebusch et al. (2013). SD-WACCM-X is nudged at every 5 min time step with horizontal winds, temperatures, and surface pressure from MERRA-2 data to prevent divergence from real dynamical conditions. Additionally, SD-WACCM-X is forced with surface wind stress and sensible as well as latent surface heat flux. As suggested by Brakebusch et al. (2013), the nudging coefficient is 0.01 s^{-1} below the altitude of 50 km, and linearly decreases and becomes zero above the altitude of 60 km.

The resulting WACCM-X simulations are compared with the simulations of TIE-GCM. The comparisons between WACCM-X and TIE-GCM simulations will show differences and similarities in modeling capabilities between whole atmosphere modeling and ionosphere-thermosphere modeling with a specified low-boundary forcing (e.g., Global Scale Wave Model (GSWM) (Hagan et al., 1999) used for this study).

Table 2 shows the version of the models, input data used for the simulations, and models used for lower boundary forcing and high latitude electrodynamics. We utilized unique model setting identifiers to distinguish the current simulations from those used in our previous studies (Shim, Rastätter, et al., 2017; Shim et al., 2011, 2012, 2014, 2018). Additional information for the models and model setting identifiers is available in Shim et al. (2011) (Refer to all references therein) and at https://ccmc.gsfc.nasa.gov/support/GEM_metrics_08/tags_list.php.

To investigate improvement in foF2 and TEC predictions of the updated versions of CTIPE (12_CTIPE) and GITM (7_GITM), the simulations of the old versions of the models (11_CTIPE and 6_GITM) from our previous study are included. The comparison will be focused on the comparison between the simulations obtained from the same model. As for TIE-GCM, 12_TIE-GCM (run at 2.5° resolution) is presented for this study, but the comparison between 11_TIE_GCM and 12_TIE-GCM was not included in this study because the only difference between the two is horizontal resolution ($5^\circ\text{lat.} \times 5^\circ\text{long.}$ vs. $2.5^\circ\text{lat.} \times 2.5^\circ\text{long.}$).

We should take note of the difference between the simulations obtained from the same model that influence foF2 and TEC responses to geomagnetic storms. For two CTIPE runs, different lower atmospheric tides were specified: 11_CTIPE was driven by the imposed migrating semidiurnal (2, 2), (2, 3), (2, 4), (2, 5), and diurnal (1, 1) tidal modes, while 12_CTIPE was run with monthly mean spectrum of tides obtained from WAM (Whole Atmosphere Model) (Akmaev, 2011; T. Fuller-Rowell et al., 2010). For two GITM simulations, 7_GITM used the T. J. Fuller-Rowell and Evans (1987) model, while 6_GITM used the Ovation model (Newell & Gjerloev, 2011; Newell et al., 2009) for specifying the patterns of auroral precipitation average energy and total energy flux. For energy deposition from energetic particle precipitation (EPP) into the atmosphere, results of Fang et al. (2010) and Sharber et al. (1996) were used for 7_GITM and 6_GITM, respectively. For two WACCM-X simulations, Heelis et al. (1982) and Weimer (2005) electric potential models were used for 3_WACCM-X and 4_WACCM-X, respectively. 12_TIE-GCM was driven by Weimer-2005 electric potential model and GSWM.

3. Performance of the Models in Predictions of foF2 and vTEC on 17 March 2013

Most simulations newly added for this study show similar behavior to those used in Shim et al. (2018), in predicting foF2 and TEC during the storm. For example, the simulations are not able to reproduce (a) the difference between eastern and western parts of the North American sector (e.g., TEC increases at Millstone Hill but decreases at Idaho and Boulder around 20 UT), and (b) different responses between foF2 (negligible changes)

Table 2
Models Used for This Study

Model setting ID	Model version	Input data	Drivers		Upper boundary for TEC calculation/Resolution
			Models used for thermosphere, tides from lower boundary, and high latitude electrodynamics		
Physics-based coupled ionosphere-thermosphere model					
11_CTIPE ^a	CTIPE3.2 (Codrescu et al., 2000; Millward et al., 2001)	F10.7, ACE IMF data and solar wind speed and density, NOAA POES Hemispheric Power data	Tides (2, 2), (2, 3), (2, 4), (2, 5), and (1, 1) propagating tidal modes	High latitude electrodynamics Weimer-2005 high latitude electric potential (Weimer, 2005), T. J. Fuller-Rowell and Evans (1987) auroral precipitation	~2,000 km, 2° lat. × 18° long.
12_CTIPE ^a	CTIPE4.1		WAM (Akmaev, 2011; T. Fuller-Rowell et al., 2010) tides		
6_GITM ^a	GITM2.5 (Ridley et al., 2006)	FISM solar EUV irradiance (Chamberlin et al., 2007), ACE IMF data and solar wind speed and density	MSIS (Hedin, 1991) migrating diurnal and semidiurnal tides	Weimer-2005 high latitude electric potential, Ovation auroral precipitation (Newell & Gjerloev, 2011; Newell et al., 2009), Fang's EPP energy deposition (Fang et al., 2010)	~600 km, 2.5° lat. × 5° long.
7_GITM	GITM21.11			Weimer-2005 high latitude electric potential, T. J. Fuller-Rowell and Evans (1987) auroral precipitation, Sharber's EPP energy deposition (Sharber's et al., 1996)	
12_TIE-GCM ^a	TIE-GCM2.0 (Richmond et al., 1992; Roble et al., 1988; Solomon et al., 2012)	F10.7, Kp, OMNI IMF data and solar wind speed and density	GSWM (Hagan et al., 1999) migrating diurnal and semidiurnal tides	Weimer-2005 high latitude electric potential, Roble and Ridley (1987) auroral precipitation	~600 km, 2.5° lat. × 2.5° long.
Whole atmosphere model					
3_WACCM-X	CESM2.2 (Gottelman et al., 2019; Liu et al., 2018)	F10.7, Kp, OMNI IMF data and solar wind speed and density	Heelis high latitude electric potential (Heelis et al., 1982), Roble and Ridley (1987) auroral precipitation		~600 km, 1.9° lat. × 2.5° long.
4_WACCM-X			Weimer-2005 high latitude electric potential, Roble and Ridley (1987) auroral precipitation		

^aThe model results are submitted by the Community Coordinated Modeling Center (CCMC) using the models hosted at the CCMC.

and TEC (noticeable increase) found in European (Chilton) and South-African (Grahamstown) stations (See Figure 4 of Shim et al. (2018) for reference). However, compared to other simulations, 4_WACCM-X driven by Weimer-2005 high latitude electric potential model captures relatively well the two differences in TEC and foF2 described above (Figure S1).

Scatter plots of the observed (*x* axis) and modeled (*y* axis) shifted foF2 and TEC, and percentage change of foF2 and TEC during the storm (17 March 2013) are shown in Figure 1 for CTIPE, in Figure 2 for GITM, and in Figure 3 for TIE-GCM and WACCM-X. Figures 1–3 display the values of all 12 locations grouped into 4 sectors: North America (NA, green), Europe (EU, blue), SAF (red), and SAM (black). The modeled foF2 was calculated from the maximum electron density of the F2 layer, NmF2, by using the relation, $NmF2 = 1.24 \times 10^{10} \times (foF2)^2$, where NmF2 is in electrons/m³ and foF2 is in MHz. First, the qualitative comparison between the simulations from the same model can be summarized as follows. 11_CTIPE/12_CTIPE tends to underestimate/overestimate foF2 for both quiet and disturbed conditions, but 12_CTIPE predicts much better both foF2 and TEC during the storm than 11_CTIPE (Figure 1). 6_GITM and 7_GITM underestimate foF2 and TEC for all cases and show relatively small response to the storm compared to the other simulations (Figure 2). 12_TIE-GCM and WACCM-Xs produce similar foF2 and TEC changes during the storm. All three simulations give substantial underestimation of TEC in SAF. 12_TIE-GCM and 3_WACCM-X produce larger overestimation of foF2 and TEC in the NA

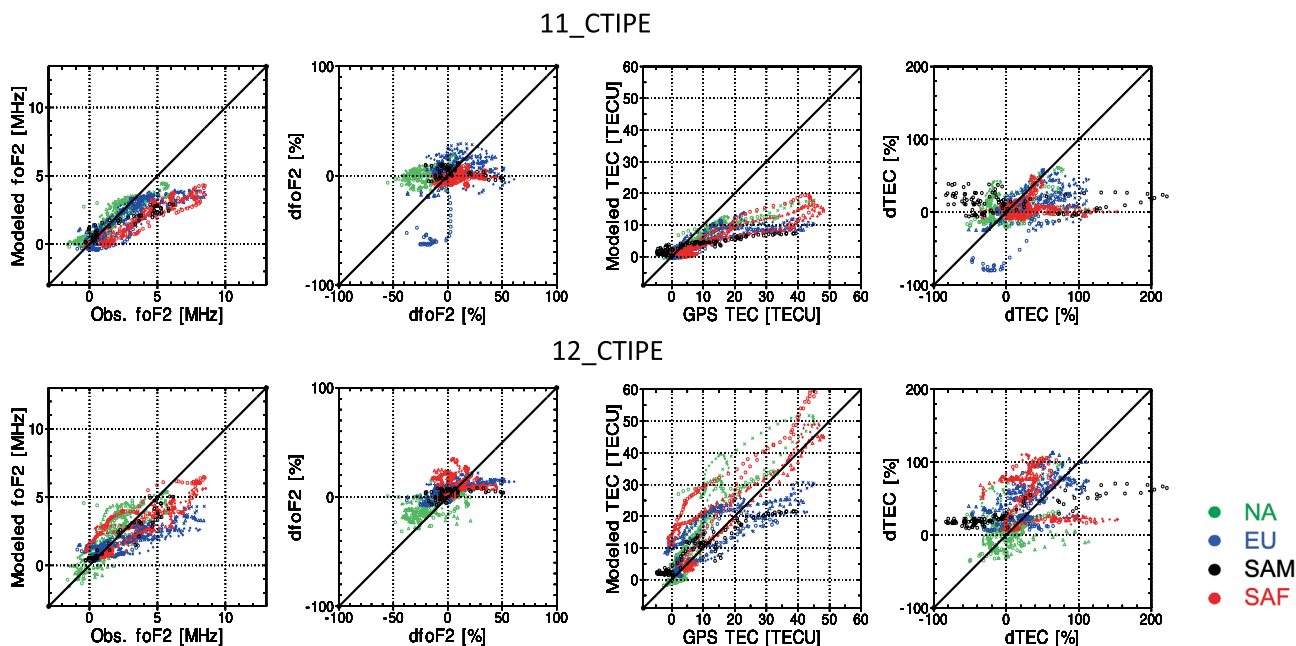


Figure 1. Scatter plots of the observed (x axis) and modeled (y axis) shifted foF2 and TEC (foF2* in the first, TEC* in the third columns), and percentage change of foF2 and TEC (dfoF2[%] in the second, dTEC[%] in the fourth columns) during the storm (17 March 2013) for 11_CTIPE and 12_CTIPE. The displayed values are for all 12 locations grouped into North America (NA, green), Europe (EU, blue), South Africa (SAF, red), and South America (SAM, black).

sector than 4_WACCM-X. 4_WACCM-X shows substantial improvement in the TEC overestimation in NA. 3_WACCM-X, of which the high latitude electric potential is specified by Heelis et al. (1982), tends to overestimate foF2 and TEC compared with 4_WACCM-X (Figure 3). 3_WACCM-X and 4_WACCM-X produce better quiet time foF2 and TEC than 12_TIE-GCM does and capture wave-like small increases in foF2 and TEC at Idaho National Lab around 10–11UT (2–3 LT) (Figure S1).

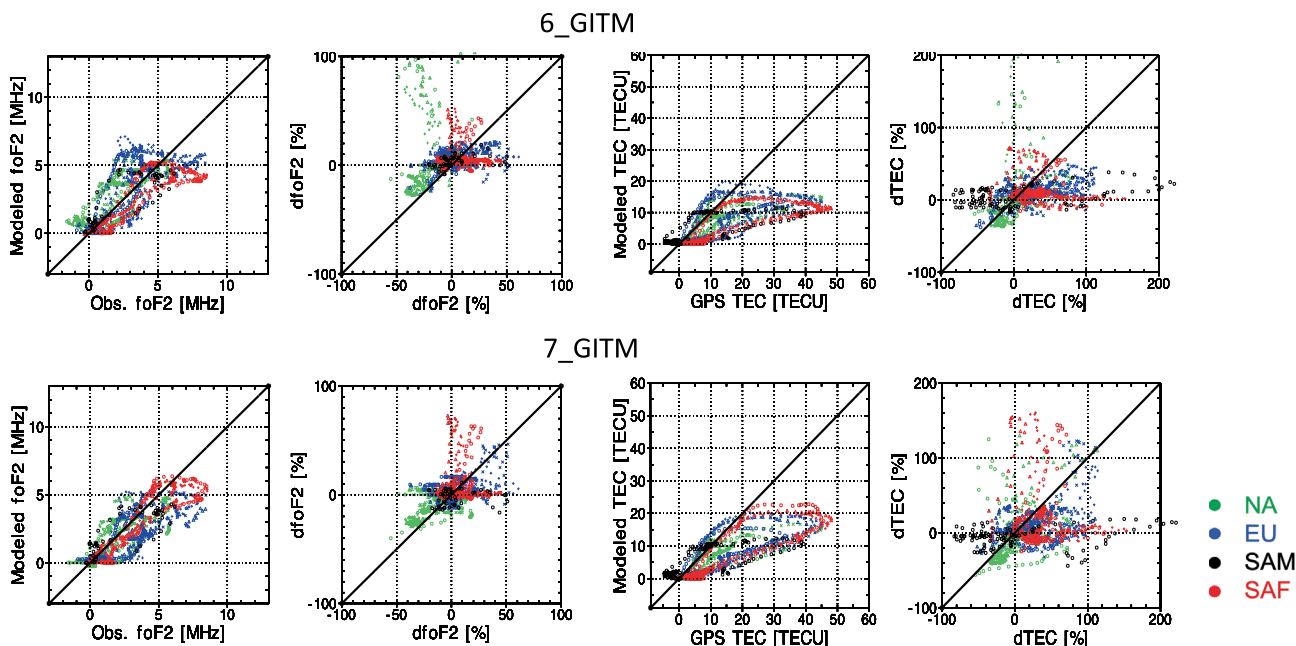


Figure 2. Same as Figure 1 but for 6_GITM and 7_GITM.

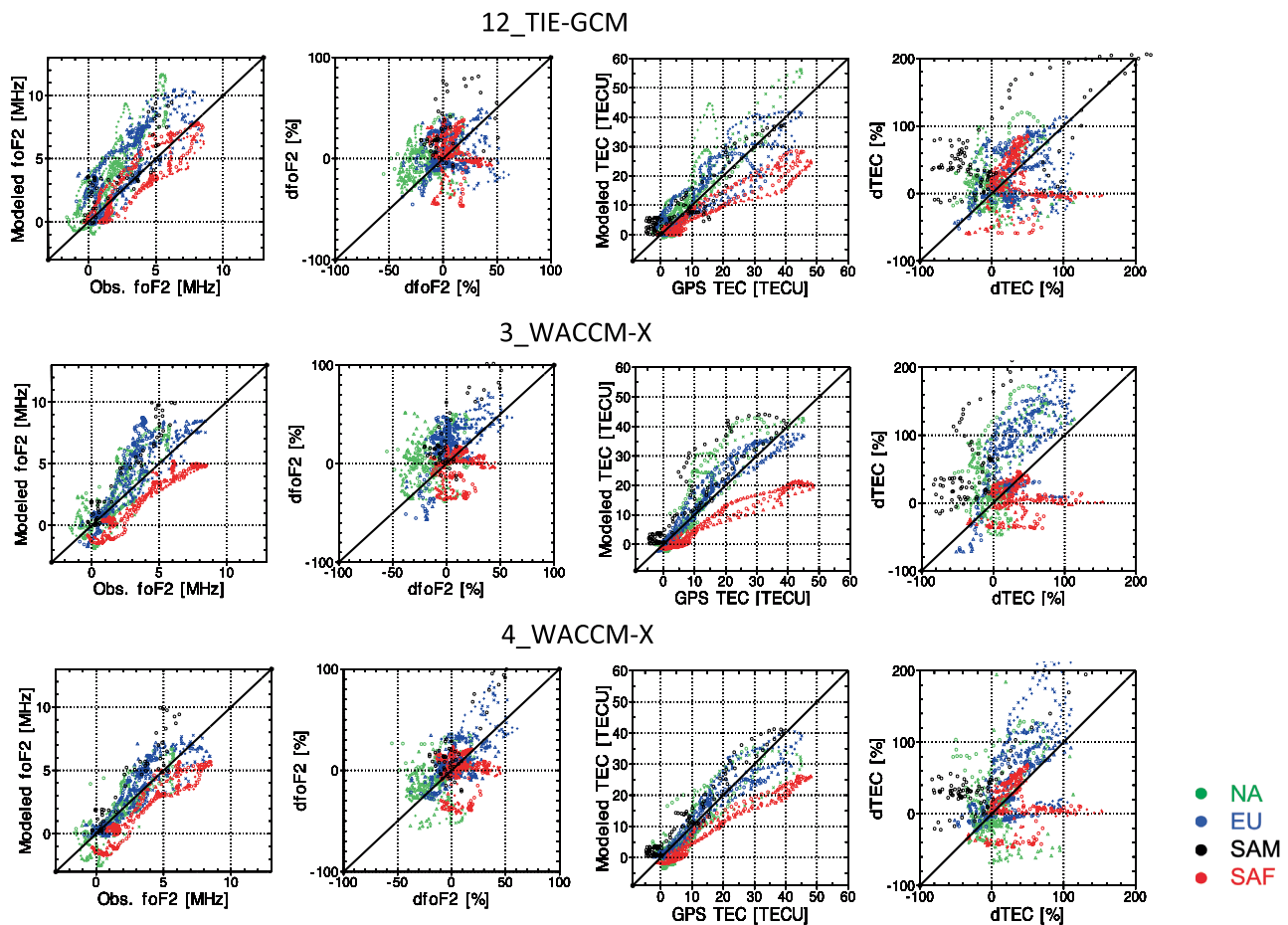


Figure 3. Same as Figure 1 but for 12_TIE-GCM, 3_WACCM-X, and 4_WACCM-X.

As shown for 6_GITM and 11_CTIPE in Shim et al. (2018), the modeled foF2 values from 7_GITM and 12_CTIPE better agree with the observed ones when they are shifted by subtracting the minimum of the 30-day median (see Figure S2 in Supporting Information in Shim et al. (2018)). Most foF2 and TEC data points from 7_GITM and 12_CTIPE before shifting are below and above the line with slope 1 (black solid line), respectively. This indicates that 7_GITM underestimates foF2 and TEC like 6_GITM, while 12_CTIPE overestimates them. The models that tend to underestimate foF2, such as 6_GITM, 7_GITM, and 11_CTIPE, seem to be unable to produce foF2* larger than about 7 MHz, and underestimate TEC* being less than about 20 TECU during the storm as reported in Shim et al. (2018). 12_TIE-GCM and WACCM-Xs show similar distribution of the data points after shifting foF2 and TEC with a tendency to underestimate foF2 and TEC in the SAF region. This shifting procedure by the minimum of the 30-day median (i.e., quiet-time minimum) for each model simulation and observation should effectively remove any differences among the models and observations that may be associated with potential biases of the models and observations. Note that this comparative study focuses on the storm-time variations of the models from their quiet-time values.

The modeled dfoF2[%] and dTEC[%] show less agreement with the observed values than the modeled foF2* and TEC* do. The data points in the second quadrant (top left) and the fourth quadrant (bottom right) indicate that the modeled and observed percentage changes are in opposite sign. 7_GITM and 3_WACCM-X have more data points in the second quadrant for the dfoF2[%] prediction than 6_GITM and 4_WACCM-X, respectively. Like most simulations used in our previous evaluation (Shim et al., 2018), 12_CTIPE and 7_GITM do not appear to reproduce the large dTEC[%] (about 200%) at Port Stanley in SAM. However, 12_TIE-GCM and WACCM-Xs better produce the enhancement in TEC percentage change. Compared to 4_WACCM-X and 12_TIE-GCM, 3_WACCM-X overestimates dTEC[%] especially in the NA and EU regions. 12_CTIPE and 6_GITM have more data points of overestimated dTEC[%] in SAF than 11_CTIPE and 7_GITM, respectively.

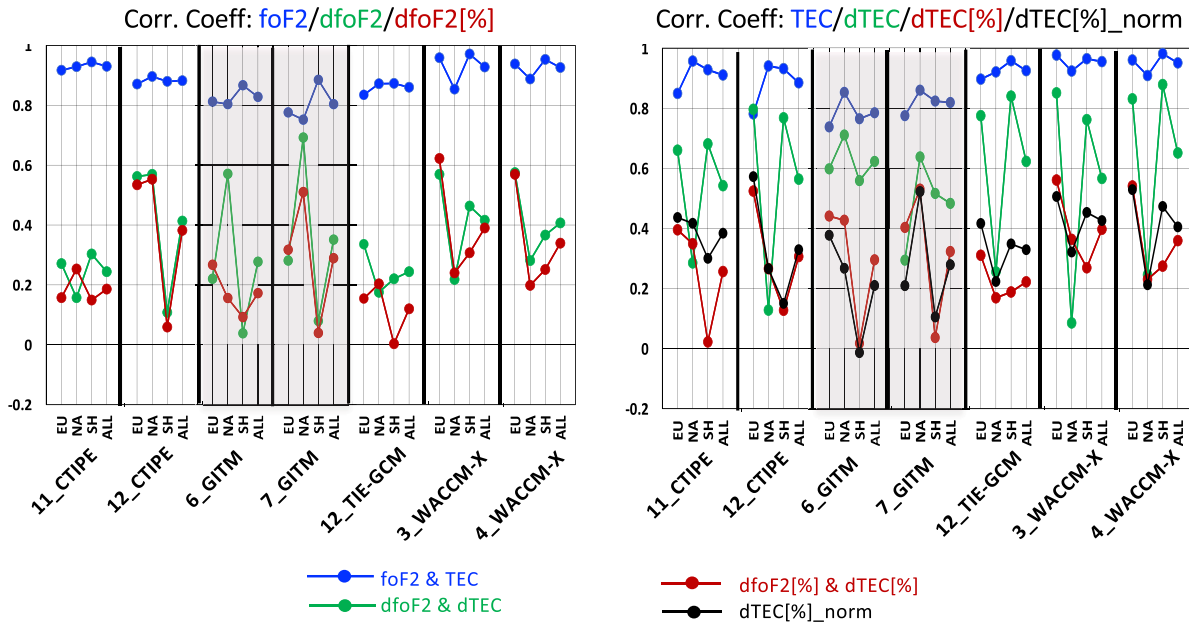


Figure 4. Correlation Coefficients (CC) between modeled and observed F2-layer critical frequency (foF2) (left panel) and Total Electron Content (TEC) (right panel). Four CCs are displayed for each simulation: CC averaged over Europe (EU), North America (NA), Southern Hemisphere (SH refers to South Africa and South America combined), and all 12 locations, from left to right. Different colors denote different quantities. Blue denotes shifted foF2 and TEC, green and red the change and percentage changes, and black normalized percentage change. The closer the circles are to the horizontal line of 1, the better the model performances are.

From now on, foF2 and TEC will represent shifted foF2 (foF2*) and shifted TEC (TEC*), respectively.

3.1. Correlation Coefficient (CC)

We first calculate CC between the modeled and observed foF2 and TEC for DOY 076 (17 March 2013) for quantitative assessment of the model performance of TEC and foF2 predictions. In Figure 4, the CCs for each simulation are presented for foF2 in the left panel and for TEC in the right panel. For each simulation, four CC values are displayed. The first three of the values correspond to the average CC over Europe (EU), NA, Southern Hemisphere (SH refers to SAF and SAM combined), and the last one is the average of all 12 locations. The modeled foF2 and TEC (blue dots) are highly correlated with the observed values. The average CC values over all 12 locations for both foF2 and TEC are about 0.8–0.95, but the average CCs for their changes are much smaller. For example, the CCs for TEC changes (dTEC) are 0.5–0.6 and even smaller for foF2. The modeled foF2 changes (green), percentage changes (red) and normalized percentage changes (black only applicable for TEC) are much less correlated (closer to uncorrelated) with the observed values (about $0.1 < \text{average CC} < 0.4$). There is no big difference between dTEC[%] and dTEC[%]_norm based on the average values for each simulation as reported in Shim et al. (2018).

Note that the CC values for the changes and percentage changes of foF2 and TEC are highly dependent on location. Most simulations, except for 12_CTIPE and GITMs, show lower CC for dfoF2 and dTEC in NA. It seems to be caused by the decreases of foF2 and TEC during the storm (negative phase) in the western parts of NA that are not captured well. GITMs show the negative phase well although it underestimated the magnitude of the change. The CCs for the percentage changes of foF2 and TEC are particularly small for CTIPEs and GITMs.

11_CTIPE's foF2 and TEC averaged over 12 locations are slightly better correlated with the observed values than 12_CTIPE. However, the changes and percentage changes of foF2 and TEC from 12_CTIPE are better correlated with the observed values than 11_CTIPE's values in most regions. Although the two GITMs produce similar CCs, 7_GITM shows better CC in NA regions for dfoF2, dfoF2[%], dTEC[%], and n_dTEC[%], while 6_GITM shows better CC for foF2 and dTEC. WACCM-Xs perform better than 12_TIE_GCM for all the considered quantities based on the average except for dTEC. WACCM-Xs perform similar to each other.

Close inspection of Figures 1 and 4 indicates that a linearity between CTIPEs and observations is improved in the newer version of CTIPE (12_CTIPE), but 12_CTIPE gives more scattered distribution around a linear

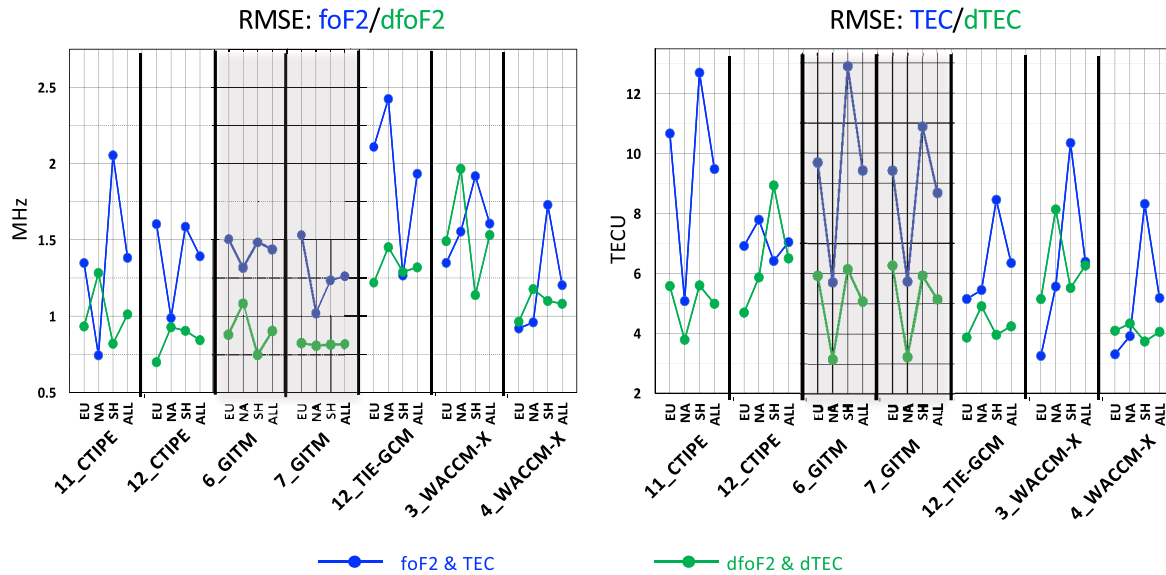


Figure 5. Same as Figure 4 but for root-mean square error of shifted F2-layer critical frequency (foF2) and Total Electron Content (TEC), and changes of foF2 and TEC.

relation (Figure 1), which seems to lead to the lower CC in 12_CTIPE than in 11_CTIPE. 7_GITM exhibits a slight improvement in a linearity between the model and observations (Figure 2), but this improvement is not clearly seen in the correlation analysis (Figure 4). For 12_TIE-GCM and WACCM-Xs, both a linearity between the models and observations (Figure 3) and CCs (Figure 4) demonstrate that the model performances are overall improved in WACCM-Xs compared with TIE-GCM. In terms of the model-observation linearity, 4_WACCM-X is somewhat better than 3_WACCM-X (Figure 3), but their CCs seems comparable to each other (Figure 4).

3.2. Root Mean Square Error (RMSE)

Figure 5 shows RMSE of foF2 and dfFoF2 in the left panel, and TEC and dTEC in the right panel. For foF2 (blue) and dfFoF2 (green) predictions, based on the average RMSE values, the RMSEs from the updated version (12_CTIPE and 7_GITM) are about 1.5 MHz for foF2 and about 1 MHz for dfFoF2, and they are slightly lower than RMSEs in their old versions. 12_CTIPE shows improvement in foF2 in SH and dfFoF2 in NA and EU compared to 11_CTIPE. 7_GITM performs better in foF2 and dfFoF2 in EU and SH than 6_GITM. 4_WACCM-X has smaller RMSE (~1 MHz) than 3_WACCM-X and 12_TIE-GCM (~1.3 MHz for dfFoF2 and ~2 MHz for foF2).

12_CTIPE is better in TEC prediction than 11_CTIPE, while the opposite holds true for dTEC prediction. The two GITMs' average RMSE values for TEC and dTEC predictions are similar to each other, about 9 TECU for TEC and 5 TECU for dTEC. Like foF2 and dfFoF2 prediction, 4_WACCM-X has smaller RMSE (~5 TECU for TEC and 4 TECU for dTEC) than 12_TIE-GCM and 3_WACCM-X (~6 TECU).

As seen in Shim et al. (2018), RMSE is highly variable with location. Most simulations appear to predict foF2 and/or TEC better in NA and worse in SH (except for 12_TIE-GCM for foF2 and 12_CTIPE for TEC). This hemispheric asymmetry in the performance of the models may readily be expected from the fact that the ionospheric density structures in SH are typically more complex and therefore relatively less understood compared with the density structures in NH, mainly due to more complex structure of the geomagnetic field, for example, larger declination and larger offset between geographic and magnetic poles in SH (e.g., Jee et al., 2009; Laundal et al., 2017; Kim et al., 2023) and resulting hemispheric asymmetry in thermospheric O/N₂ ratio (Qian et al., 2022). Shim et al. (2018) also suggested that this hemispheric asymmetry is possibly partly attributed to the fact that the models do not include the energy input from the inner magnetosphere that affects the ionosphere (e.g., foF2 and TEC enhancements) in the South Atlantic Anomaly (SAA) region (Dmitriev et al., 2017; Zhao et al., 2016) where the 4 stations in SH are situated nearby. Both 11_CTIPE and GITMs tend to perform better in NA for dTEC, while WACCM-Xs show the opposite tendency for dfFoF2 and dTEC. 7_GITM and 4_WACCM-X show the least RMSE dependence on location for dfFoF2 and for dTEC, respectively, among seven simulations.

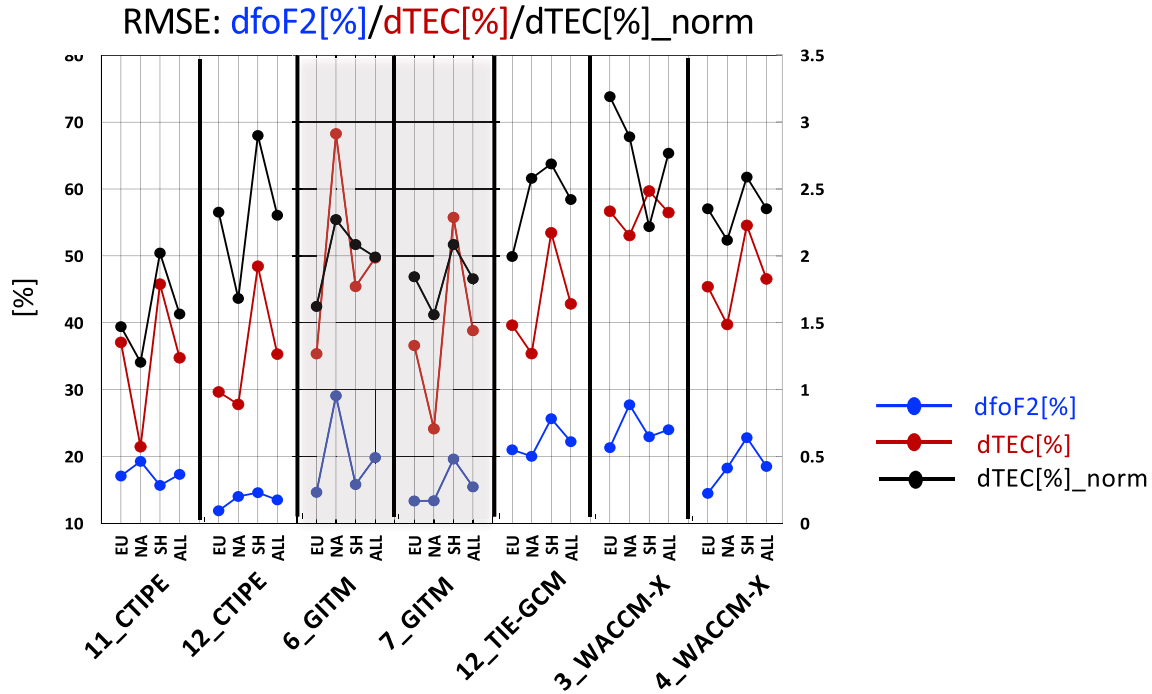


Figure 6. Same as Figure 4 but for root-mean square error of percentage change of F2-layer critical frequency (foF2) and Total Electron Content (TEC), and normalized percentage change. Blue denotes dfoF2[%], red and black dTEC[%] and dTEC[%]_norm.

Figure 6 shows the RMSE of percentage changes of foF2 (blue) and TEC (red) and normalized percentage changes of TEC (black). The two CTIPEs produce similar RMSE for dTEC[%], but 12_CTIPE and 11_CTIPE produce lower RMSE for dfoF2[%] and dTEC[%]_norm, respectively. For all three percentage changes of dfoF2[%], dTEC[%], and dTEC[%]_norm, 7_GITM seems to perform better than 6_GITM based on the average RMSEs over the 12 locations. 4_WACCM-X and 12_TIE-GCM perform very similarly for dfoF2[%] and dTEC[%] and better than 3_WACCM-X.

Difference in the performance among locations is more noticeable in dTEC[%] and dTEC[%]_norm than in dfoF2[%] as found in Shim et al. (2018). All simulations, except 6_GITM, produce lower RMSE of dTEC[%] in NA and higher in SH region. This tendency remains the same for dTEC[%]_norm with the exception of 3_WACCM-X, which has lower RMSE for dTEC[%]_norm in SH. For 3_WACCM-X, the higher RMSE for dTEC[%] and the lower RMSE for dTEC[%]_norm in SH than in NA are probably due to the normalization factor, standard deviation of dTEC[%] in the locations.

3.3. Yield and Timing Error (TE)

To measure how well the models capture the degree of TEC and foF2 disturbances during the main phase, Yield and TE of dfoF2[%], dTEC[%], and dTEC[%]_norm are calculated. Shim et al. (2018) considered two time intervals, 06–15UT and 15–22UT, when peaks are observed in most of 12 locations. In each time interval, we calculate one Yield value and one TE value. Definitions of Yield and TE are presented in Table 1.

In each sector, average Yield and TE are calculated over the number of stations where the model correctly predicts the storm phase, that is, Yield is positive. Table 3 shows the total number of stations where the models show correct storm phase, either positive or negative. The numbers in bold are the higher values between the simulations compared. 12_CTIPE predicts the storm phase better for dTEC[%] than 11_CTIPE, but 11_CTIPE predicts better for dfoF2[%] than 12_CTIPE. 7_GITM is improved in predicting the storm phase of dfoF2[%], while 6_GITM predicts better the storm phase of dTEC[%]. 4_WACCM-X, compared to 12_TIE-GCM and 3_WACCM-X, is better for predicting the phase of dfoF2[%] and worse for predicting that of dTEC[%].

Figure 7 shows average Yield (left) and average of absolute values of TE (right) over the two time intervals: dfoF2[%] in blue, dTEC[%] in red, and dTEC[%]_norm in black. Concerning the average of all 12 locations,

Table 3
Number of Locations Where the Models Correctly Predict Negative or Positive Phase

	Time interval	11_CTIPE	12_CTIPE	6_GITM	7_GITM	12_TIE-GCM	3_WACCM-X	4_WACCM-X
dfoF2[%]	06–15 UT	8	7	5	9	9	6	10
	15–22 UT	10	6	7	8	7	7	10
dTEC[%]	06–15 UT	9	10	10	10	7	10	9
	15–22 UT	7	10	12	11	10	7	8

12_CTIPE appears to overestimate peak values of dTEC[%] and dTEC[%]_norm with larger variation with location (e.g., $\sim 1 < \text{Yield of dTEC}[\%]_{\text{norm}} < \sim 2.5$) than 11_CTIPE, of which Yield is less than 1 for all three quantities of percentage changes (e.g., $0.7 < \text{Yield of dTEC}[\%]_{\text{norm}} < 0.9$). Yields of 12_CTIPE for dTEC[%] and dTEC[%]_norm are closer to 1 in NA. GITMs produce similar ratios based on the average over all locations, but 7_GITM shows smaller differences in Yield among locations (e.g., $\sim 0.5 < \text{Yield of dTEC}[\%]_{\text{norm}} < \sim 1$) than 6_GITM (e.g., $0.5 < \text{Yield of dTEC}[\%]_{\text{norm}} < \sim 2.5$). In terms of average Yield, 12 TIE-GCM and two WACCM-Xs tend to overestimate the peak values and show similar performance, although 12_TIE-GCM's ratios are closer to 1 than those of WACCM-Xs. 3_WACCM-X shows larger variation in Yield among locations (e.g., $\sim 0.9 < \text{Yield of dTEC}[\%]_{\text{norm}} < \sim 2.7$) than 12_TIE-GCM and 4_WACCM-X (e.g., $\sim 1.7 < \text{Yield of dTEC}[\%]_{\text{norm}} < \sim 2.3$).

Average Timing Errors of dfoF2[%] and dTEC[%]_norm are between 1 and 2 hr, and TE of dTEC[%] are about 0.8–1.5 hr. With respect to the average TE, 12_CTIPE has smaller TE (~ 1 hr) than 11_CTIPE (about 1.5 hr) for all three percentage changes with less location dependence as well. 7_GITM's three TEs are about 1.5 hr, while 6_GITM's TEs of dfoF2[%], dTEC[%], and dTEC[%]_norm are ~ 1 , ~ 1.4 , and ~ 2 hr, respectively. 12 TIE-GCM has smaller TE for dfoF2[%] and 3_WACCM-X has smaller TE for dTEC[%] and dTEC[%]_norm, however 3_WACCM-X show larger location dependence of TE for dTEC[%]_norm and dfoF2[%].

4. Ensemble of TEC Obtained From 13 Simulations

The linearity check, RMSE, and CC between model results and observations for shifted foF2 and TEC and their relative changes indicate that the newer versions of the models (i.e., 12_CTIPE, 7_GITM, and 4_WACCM-X)

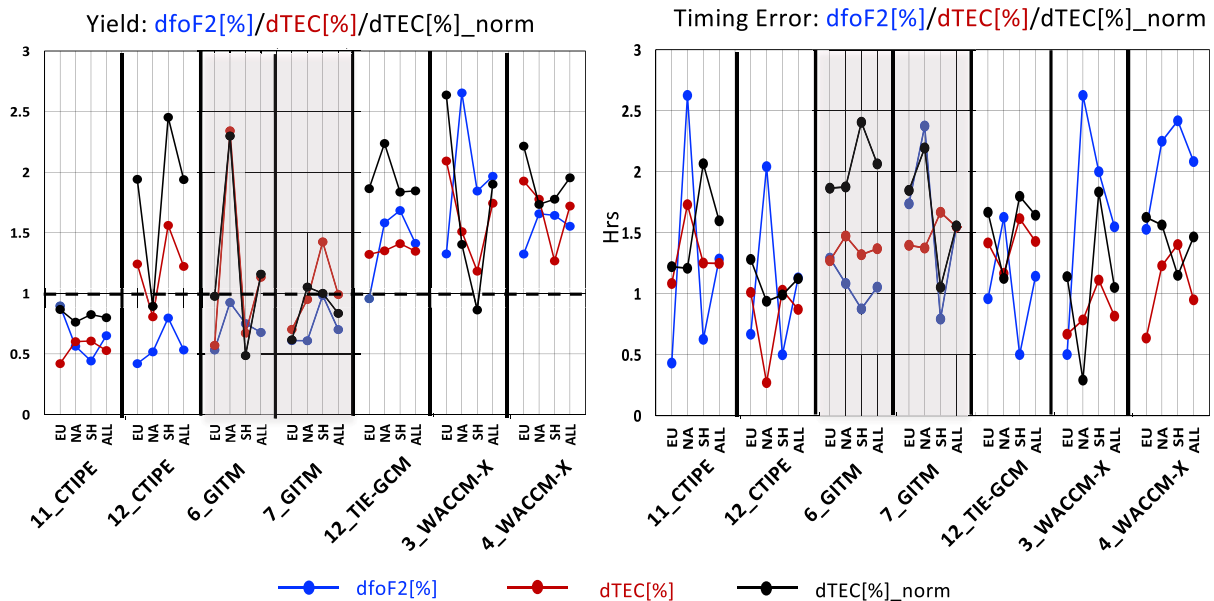


Figure 7. Same as Figure 4 but for Yield (ratio) and absolute of Timing Error ($\text{ITEI} = |t_{\text{peak_model}} - t_{\text{peak_obs}}|$).

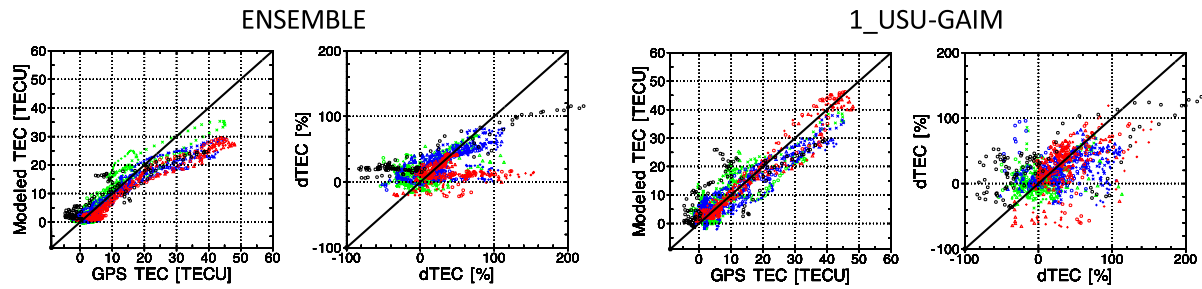


Figure 8. Same as Figure 1 but for only Total Electron Content (TEC) and dTEC[%] from the ensemble of the simulations (ENSEMBLE) and 1_USU-GAIM.

produce the better results. From the viewpoints of correct prediction of storm phases (Table 3), Yields, and TES (Figure 7), however, there is no one best simulation for all locations, and the performance of the models varies with location as well as the Yields and TE.

The differences in performance among the simulations could be caused by inherent differences among the models, for example, different methods to solve for chemistry and advection, and different ways to treat eddy diffusion and vertical transport (T. J. Fuller-Rowell et al., 1996; Liu et al., 2018; Perlongo et al., 2018; Ridley et al., 2006; Solomon et al., 2012), or by a combination of different input data and different models used for lower boundary forcing and high-latitude electrodynamics (Shim et al., 2018). Even different data assimilation models for the same weather condition can yield different results, due to numerous reasons (e.g., the use of different background weather models, spatial/temporal resolutions, assimilation methods, and data error analyses), even if the same data are assimilated (Schunk et al., 2021). The common way to handle these differences is to use model ensembles and the use of ensembles enables estimations of the certainty of results. Thus, we used a weighted mean of the ensemble of all 13 simulations including eight simulations from our previous study (Shim et al., 2018) for TEC, dTEC and dTEC[%] to compare the ensemble average with the individual simulations. To get the weighted mean ($\bar{x} = \sum w_i x_i / \sum w_i$), we used the RMSE of shifted TEC ($w_i = 1/\text{RMSE}$).

Figure 8 is the same as Figure 1 but for the ensemble of the simulations (ENSEMBLE will be used as model setting ID) and a simulation (1_USU-GAIM) from a data assimilation model (DA), Utah State University-Global Assimilation of Ionospheric Measurements (USU-GAIM). For TEC less than about 20 TECU, ENSEMBLE shows better agreement with GPS TEC than the individual simulations, including 1_USU-GAIM. However, as we can expect, ENSEMBLE underestimates TEC larger than about 30 TECU due to the tendency to underestimate TEC of many simulations as pointed out in Section 3 and Shim et al. (2018). For dTEC[%], ENSEMBLE appears to be correlated better with GPS dTEC[%] than the other simulations, although there are some underestimations in SAF, as well as in SAM with opposite prediction of the storm phase.

Figure 9 shows averaged CC and RMSE values over all 12 locations of 13 simulations, the ensemble of them, and the ensemble of 12 simulations excluding 1_USU-GAIM (ENSEMBLE_wo_DA). The detailed settings of the simulations that are used in Shim et al. (2018) but not listed in Table 2, such as 4_IRI, 1_IFM, 1_SAMI3, are presented in Table 2 in Shim et al. (2018). The simulations in Figure 9a were arranged by the average of the three averaged CC values for TEC, dTEC and dTEC[%] from the smallest to the largest (closer to 1). In Figure 9b, the simulations were arranged by the average of the two averaged RMSEs for TEC and dTEC from the largest to the smallest. Based on the averaged CC and RMSE, ENSEMBLES (ENSEMBLE and ENSEMBLE_wo_DA) of the simulations perform very similarly and outperform all 12 simulations but a data assimilation model, 1_USU-GAIM, which assimilated GNSS TEC data and shows the best performance for TEC prediction in most cases with the least location dependence of RMSE in our former study (Shim et al., 2018). However, ENSEMBLES and 1_USU-GAIM do not show big difference in their performance. The differences in RMSE of TEC and dTEC between ENSEMBLE and 1_USU-GAIM are less than 0.5 and 0.1 TECU, respectively. For dTEC[%], ENSEMBLE performs slightly better than 1_USU-GAIM with about 1.5% lower RMSE. The fact that ENSEMBLES are comparable to the data assimilation model 1_USU-GAIM indicates that the multi-model ensemble can be useful in forecasting the IT system, although this result is obtained from a single geomagnetic storm event.

Figure 10 shows Yield and TE of dTEC[%] for all 13 simulations along with ENSEMBLE. The values correspond to the average over all 12 locations. Unlike CC and RMSE, ENSEMBLE does not outperform all physic-based

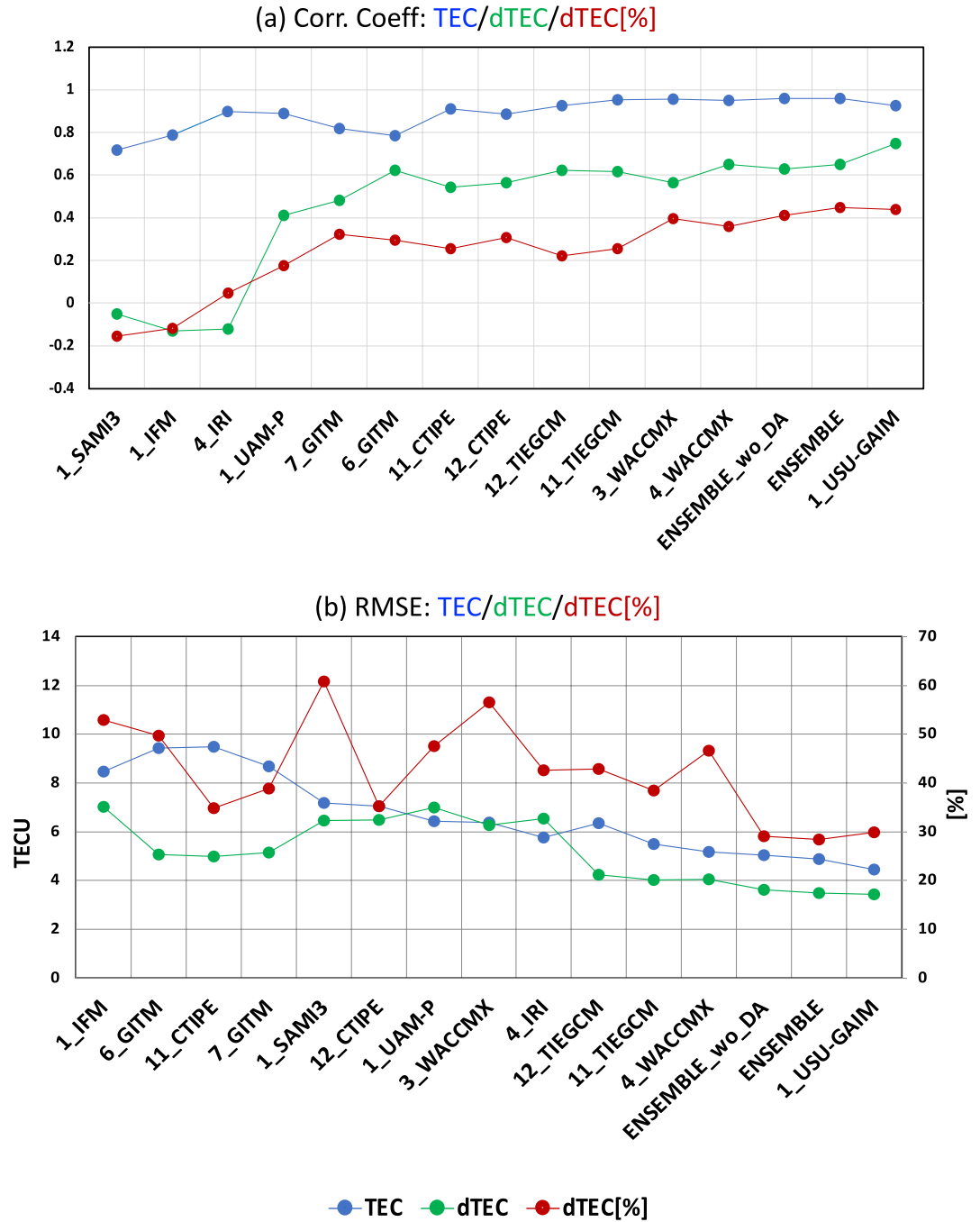


Figure 9. Averaged Correlation coefficient (CC) (a) and root-mean square error (RMSE) (b) over all 12 locations of 13 simulations, the ensemble of them (ENSEMBLE), and the ensemble of 12 simulations excluding 1_USU-GAIM (ENSEMBLE_wo_DA). Blue denotes shifted Total Electron Content (TEC), green and red the change and percentage changes of TEC. CCs are plotted from the smallest to the largest (closer to 1) according to the average of the three averaged CC values of TEC, dTEC, and dTEC[%]. RMSEs are plotted from the largest to the smallest according to the average RMSE for TEC and dTEC.

coupled models in terms of Yield and TE, although the difference is small. ENSEMBLE underestimates Yield, while most of the simulations overestimate it, except 4_IRI and 11_CTIPE. 7 simulations from PB coupled IT models and 1_USU-GAIM produce Yield closer to 1 than ENSEMBLE does.

Timing Error of dTEC[%] from ENSEMBLE is about 1 hr, which is slightly larger than TE from 4 simulations from CTIPE and WACCM-X, but the difference from the smallest TE is less than 0.5 hr.

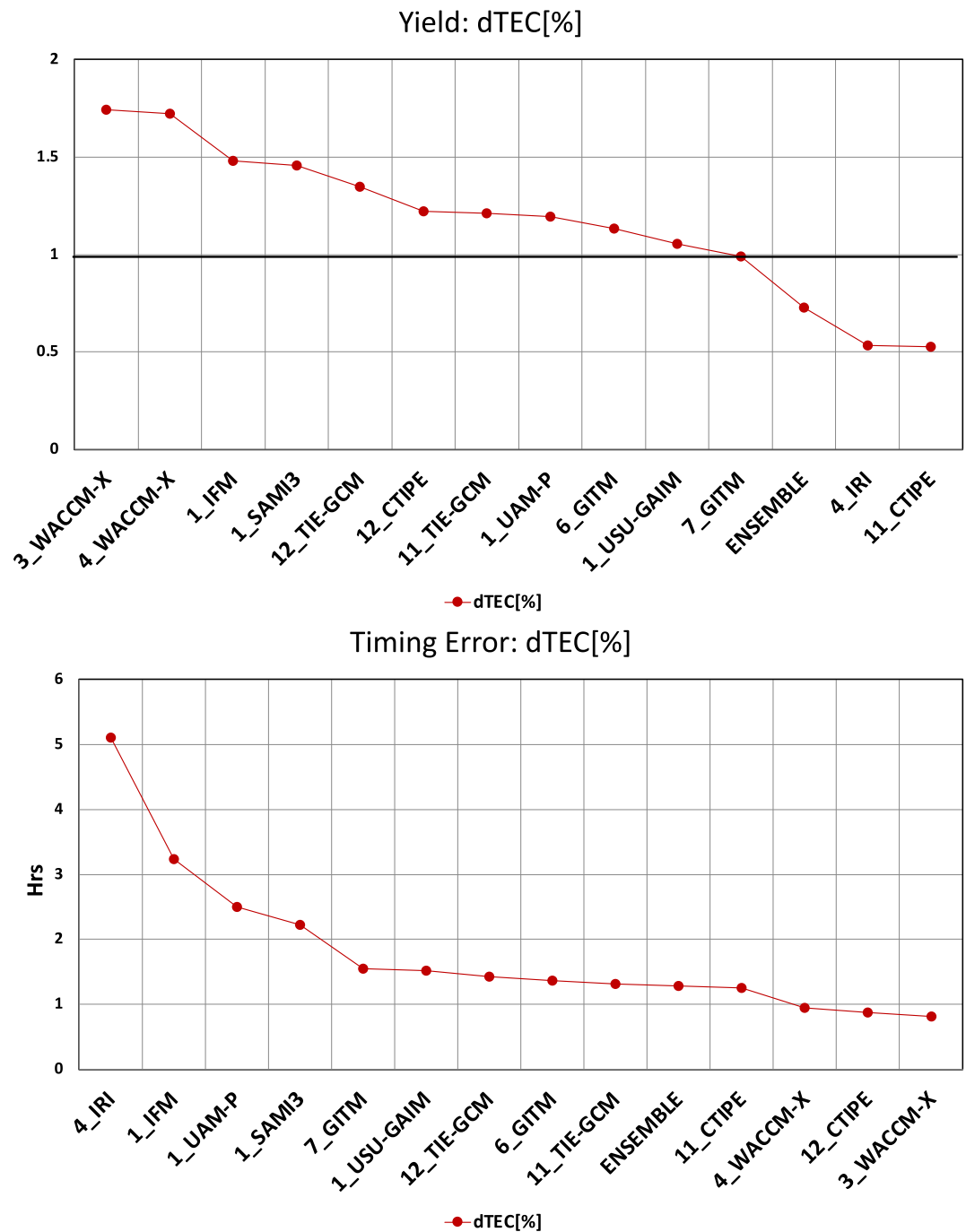


Figure 10. Yield and Timing Error of dTEC[%] for all 13 simulations and ENSEMBLE.

Regarding the averaged skill scores for all 12 locations, the five newly added simulations in this study produce comparable TEC and TEC changes to the simulations from PB IT models used in our previous study. The simulations of newer versions of the models (12_CTIPE, 7_GITM and 4_WACCM-X) are found to give overall improved forecast results. Based on the average RMSE, the ensemble of simulations of the models' newer versions is comparable to 1_USU-GAIM and performs better than the ensemble of the simulations of older versions of the models (11_CTIPE, 6_GITM and 12_TIE-GCM) (Table 4).

Table 4

Averaged Root-Mean Square Error Over All 12 Locations of the Ensemble of Newer Versions (ENSEMBLE_new) of Models (12_CTIPE, 7_GITM and 4_WACCM-X) Driven by Weimer-2005 Electric Potential Model, the Ensemble of Older Versions (ENSEMBLE_old) of Models (11_CTIPE, 6_GITM and 12_TIE-GCM), and 1_USU-GAIM

	TEC (TECU)	dTEC (TECU)	dTEC[%]
ENSEMBLE_old	6.6	4.1	33.4
ENSEMBLE_new	4.6	3.2	29.8
1_USU-GAIM	4.5	3.4	29.9

5. Summary and Discussion

We expanded on our previous systematic assessment of modeled foF2 and TEC during the 2013 March storm event (17 March 2013) to track the improvement of the models and investigate impacts of forcings from the lower atmosphere and the magnetosphere, on the performance of ionosphere-thermosphere coupled models.

We evaluated simulations from upgraded models (CTIPE4.1 and GITM21.11) since our previous assessment and a whole atmosphere model (WACCM-X2.2). To compare with results from WACCM-X2.2, we also included a simulation of TIE-GCM2.0, of which the electrodynamic processes are implemented in WACCM-X 2.2. Furthermore, to evaluate TEC prediction of the simulations, we used a weighted mean of the ensemble of

all 13 simulations including eight simulations from our previous study to compare the ensemble average with the individual simulations.

For evaluation of the simulations, we used the exact same procedure with the same data set, same physical quantities, and same skill scores as our previous study (Shim et al., 2018). The skill scores were calculated for the three sectors, EU (Europe), NA, and SH to investigate the longitudinal and hemispheric dependence of the performance of the models.

From the five simulations used in the study, we also found the general behaviors of most simulations identified in Shim et al. (2018): (a) tendency to underestimate storm-time enhancements of foF2 and TEC and not to reproduce large enhancements of dTEC[%] (e.g., about 200% TEC increase at Port Stanley in the SAA region), (b) being unable to capture opposite responses to the storm in the eastern and western parts of NA, especially the negative phase (except for GITM), which is what in part causes lower CC in NA, (c) tendency to predict foF2 and/or TEC better in NA and worse in SH with respect to RMSE. However, it was found that 12_TIE-GCM and WACCM-Xs better produce the large TEC percentage changes at Port Stanley in SAM. Based on the averaged skill scores for all 12 locations, the five simulations used in this study show skill scores better or comparable to those of the simulations from PB IT models used in our previous study.

Compared to 11_CTIPE (obtained from CTIPE3.2), 12_CTIPE (from CTIPE4.1) driven by tides from WAM tends to overestimate foF2 and TEC for both quiet and disturbed conditions and predicts better TEC peaks during the storm. For more cases, 12_CTIPE performs largely better than 11_CTIPE based on the average scores. 12_CTIPE predicts the storm phase better for dTEC[%], but 11_CTIPE does better for dfoF2[%]. 12_CTIPE appears to overestimate peak values of dTEC[%] and dTEC[%]_norm, while 11_CTIPE produces Yield less than 1.

The two GITMs, 7_GITM (with Fuller-Rowell and Evans auroral model and Fang's EPP energy deposition) and 6_GITM (with Ovation model and Sharber's energy deposition), underestimate foF2 and TEC for all cases and show relatively small response to the storm compared to the other simulations that do not appear to reproduce the large dTEC[%] (about 200% increase at Port Stanley in SAM). 7_GITM and 6_GITM perform very similarly for most cases with similar skill scores. However, 7_GITM shows better CC for most quantities except for dTEC, and lower RMSEs and Yield closer to 1 for most regions and quantities considered. 7_GITM shows the least RMSE dependence on location for dfoF2 among all simulations.

Comparing the two WACCM-Xs and 12_TIE-GCM, the two WACCM-Xs, 3_WACCM-X with Heelis high latitude electric potential model and 4_WACCM-X with Weimer-2005, predict quiet time foF2 and TEC better than 12_TIE-GCM. During the storm, 12_TIE-GCM and 4_WACCM-X produce similar foF2 and TEC in the NA sector, while 3_WACCM-X tends to overestimate these variables, producing larger changes in foF2 and TEC. In most cases, the WACCM-Xs and 12_TIE_GCM perform similarly in terms of average values of skill scores, but 3_WACCM-X and/or 4_WACCM-X perform better than 12_TIE-GCM except for Yield of percentage changes. 4_WACCM-X slightly outperforms 3_WACCM-X for all cases but not for TE for percentage changes.

Our findings suggest that the newer versions of the models (12_CTIPE, 7_GITM, and 4_WACCM-X) with Weimer-2005 electric potential model give overall improved forecast, and the performance of the models depends on forcing from the magnetosphere and also forcing from the lower atmosphere even during storms. Differences in upward-propagating tides generate differences in foF2/TEC responses to the storm by E-region wind dynamo

and tidal mixing effects (Yamazaki and Richmond, 2013). The tidal differences between the two CTIpe simulations produce differences in O/N_2 column density ratio (not shown here), and better prediction of TEC peaks of 12_CTipe with the tendency of overestimation during the storm is possibly caused by larger O/N_2 ratio. The differences in the performance between the two GITM simulations and between the two WACCM-X simulations may partially be caused by different O/N_2 ratios affected by different auroral particle heating and Joule heating that cause expansion of the upper atmosphere and the resulting thermospheric composition changes (Richmond, 2021 and references therein). Furthermore, the disturbed neutral composition in the high-latitude region is transferred to the lower latitude region by the disturbed vertical wind and equatorward thermospheric circulation. The investigation of the actual causes of the differences in the simulations will require systematic modeling studies, which are beyond the scope of this paper.

For TEC, dTEC and dTEC[%], our results indicate that the ensemble of all 13 simulations (ENSEMBLE), including eight simulations from our previous study (Shim et al., 2018) is comparable to the data assimilation model (1_USU-GAIM) with differences in skill score less than 3% and 6% for CC and RMSE, respectively. However, ENSEMBLE underestimates Yield (0.73) while 7 simulations from PB coupled IT models and 1_USU-GAIM produce Yield closer to 1. Timing Error of dTEC[%] from ENSEMBLE is about 1 hr, but the difference from the smallest TE of the simulations is less than 0.5 hr. In addition, based on RMSE, the ensemble of the newer versions of the models (12_CTipe, 7_GITM and 4_WACCM-X) is comparable to 1_USU-GAIM.

To advance our understanding of the ionosphere-thermosphere system requires significant efforts to improve the capability of numerical models along with expanding the scope of observations (Heelis & Maute, 2020). There have been recent new developments of theoretical models, including AMGeO (Assimilative Mapping of Geospace Observations) for High-Latitude Ionospheric Electrodynamics (Matsuo, 2020) and MAGE geospace model that couples the Grid Agnostic MHD for Extended Research Applications (GAMERA) global MHD model of the magnetosphere (Sorathia et al., 2020; Zhang et al., 2019), the Rice Convection Model (RCM) model of the ring current (Toffoletto et al., 2003), TIE-GCM of the upper atmosphere and the RE-developed Magnetosphere-Ionosphere Coupler/Solver (REMIX) (Merkin & Lyon, 2010). These models will be available soon to the public through CCMC, and then the modeling capability will help us better understand the processes responsible for the observed characteristics and features during disturbed conditions. In addition, CCMC will also provide users with the capability to run PB IT models with various combination of models for lower atmospheric forcing and for magnetosphere forcing, which enable us to research further the impacts of the forcings on the IT system.

The findings of this study will provide a baseline for future validation studies using new models and improved models, along with earlier results (Shim, Rastätter, et al., 2017; Shim et al., 2011, 2012, 2014, 2018) obtained through CEDAR ETI, GEM-CEDAR Modeling Challenges, and the international effort, “International Forum for Space Weather Modeling Capabilities Assessment.” We will extend our study to include more geomagnetic storm events and also geomagnetically quiet times to investigate differences and similarities in the performance of the models. In addition, we will also include foF2 and TEC predictions for the high- and low-latitude regions.

6. Conclusion

As an expansion of the model assessment study for 2013 March storm event (Shim et al., 2018), new simulations from the upgraded models including CTIpe model version 4.1, GITM version 21.11, WACCM-X version 2.2, and TIE-GCM 2.0 were evaluated to track the status of model improvement and to investigate the impacts of lower atmospheric and magnetospheric forcings on the performance of the ionosphere-thermosphere models. Here are the main results of the study.

- Model simulations tend to underestimate the storm-time enhancements of foF2 and TEC and to predict them better in the NH (specifically in the NA) but worse in the SH. It seems to be associated with more complex structure of the geomagnetic field in the SH such as larger declination and offset between geographic and magnetic poles. Furthermore, the models do not include the energy input from the inner magnetosphere that affects the ionosphere (e.g., foF2 and TEC enhancements) in the South Atlantic Anomaly (SAA) region.
- The performance of the models is strongly dependent on forcings from the magnetosphere and also from the lower atmosphere even during storms. The newer versions of the models (12_CTipe, 7_GITM, and 4_WACCM-X) with Weimer-2005 electric potential model provide overall improved forecast.

- Ensemble of all simulations for TEC is comparable to the data assimilation model (USU-GAIM) that showed best performance for TEC prediction in most cases, by assimilating GNSS TEC data, in our former study (Shim et al., 2018).
- The performance of the models substantially varies with the quantity and location considered, and the type of metrics used.
- New developments of theoretical models have recently been performed to improve the capability of numerical models along with expanding the scope of observations, including AMGeO for high-latitude ionospheric electrodynamics and MAGE geospace model, which will be available soon to the public through CCMC.
- Results of this study will provide a baseline for future validation studies using new/improved models.

Data Availability Statement

The foF2 and vertical TEC data are available on the Global Ionosphere Radio Observatory (GIRO) database (<https://giro.uml.edu/didbase/scaled.php>) and the CEDAR Madrigal database (<http://cedar.openmadrigal.org/>), respectively. Data from the South African Ionosonde network is made available through the South African National Space Agency (SANSA) (https://sandims.sansa.org.za/user/login?_next=/portal/searchBySite).

Acknowledgments

This work is supported by Korea Polar Research Institute (KOPRI) grant funded by the Ministry of Oceans and Fisheries (KOPRI PE23020) and basic research funding from the Korea Astronomy and Space Science Institute (KASI) (KASI2022185009). This work is supported by grants from the National Science Foundation (NSF) Space Weather Program. This model validation study is supported by the Community Coordinated Modeling Center (CCMC) at the Goddard Space Flight Center. Data processing and research at MIT Haystack Observatory are supported by cooperative agreement AGS-1242204 between the U.S. National Science Foundation and the Massachusetts Institute of Technology. The National Center for Atmospheric Research is sponsored by the National Science Foundation. WACCM-X source code is publicly available at the NCAR Community Earth System Model web site. WACCM-X simulations were performed using computational resources at the NCAR-Wyoming Supercomputing Center (<https://doi.org/10.5065/D6RX99HX>). Model output and observational data used for the study will be permanently posted at the CCMC website (<http://ccmc.gsfc.nasa.gov>) and provided as a resource for the space science community to use in the future.

References

- Akmaev, R. A. (2011). Whole atmosphere modeling: Connecting terrestrial and space weather. *Reviews of Geophysics*, 49(4), 390. <https://doi.org/10.1029/2011RG000364>
- Brakebusch, M., Randall, C. E., Kinnison, D. E., Tilmes, S., Santee, M. L., & Manney, G. L. (2013). Evaluation of whole atmosphere community climate model simulations of ozone during Arctic winter 2004–2005. *Journal of Geophysical Research*, 118(6), 2673–2688. <https://doi.org/10.1002/jgrd.50226>
- Bruinsma, S., Sutton, E., Solomon, S. C., Fuller-Rowell, T., & Fedrizzi, M. (2018). Space weather modeling capabilities assessment: Neutral density for orbit determination at low Earth orbit. *Space Weather*, 16(11), 1806–1816. <https://doi.org/10.1029/2018SW002027>
- Chamberlin, P. C., Woods, T. N., & Eparvier, F. G. (2007). Flare irradiance spectral model (FISM): Daily component algorithms and results. *Space Weather*, 5(7), S07005. <https://doi.org/10.1029/2007SW000316>
- Codrescu, M. V., Fuller-Rowell, T. J., Foster, J. C., Holt, J. M., & Cariglia, S. J. (2000). Electric field variability associated with the Millstone Hill electric field model. *Journal of Geophysical Research*, 105(A3), 5265–5273. <https://doi.org/10.1029/1999JA900463>
- Dmitriev, A. V., Suvorova, V., Klimenko, M. V., Klimenko, V. V., Ratovsky, K. G., Rakhmatulin, R. A., & Parkhomov, V. A. (2017). Predictable and unpredictable ionospheric disturbances during St. Patrick's Day magnetic storms of 2013 and 2015 and on 8–9 March 2008. *Journal of Geophysical Research: Space Physics*, 122(2), 2398–2423. <https://doi.org/10.1002/2016JA0232>
- Fang, X., Randall, C. E., Lummerzheim, D., Wang, W., Lu, G., Solomon, S. C., & Frahm, R. A. (2010). Parameterization of monoenergetic electron impact ionization. *Geophysical Research Letters*, 37(22), L22106. <https://doi.org/10.1029/2010GL045406>
- Fuller-Rowell, T., Wu, F., Akmaev, R., Fang, T.-W., & Araujo-Pradere, E. (2010). A whole atmosphere model simulation of the impact of a sudden stratospheric warming on thermosphere dynamics and electrodynamics. *Journal of Geophysical Research*, 115(A10), A00G08. <https://doi.org/10.1029/2010JA015524>
- Fuller-Rowell, T. J., Codrescu, M. V., Rishbeth, H., Moffett, R. J., & Quegan, S. (1996). On the seasonal response of the thermosphere and ionosphere to geomagnetic storms. *Journal of Geophysical Research*, 101(A2), 2343–2353. <https://doi.org/10.1029/95ja01614>
- Fuller-Rowell, T. J., & Evans, D. S. (1987). Height-integrated Pedersen and Hall conductivity patterns inferred from the TIROS-NOAA satellite data. *Journal of Geophysical Research*, 92(A7), 7606–7618. <https://doi.org/10.1029/ja092ia07p07606>
- Gelaro, R., McCarty, W., Suárez, M. J., Todling, R., Molod, A., Takacs, L., et al. (2017). The modern-era retrospective analysis for research and applications, version 2 (MERRA-2). *Journal of Climate*, 30(14), 5419–5454. <https://doi.org/10.1175/JCLI-D-16-0758.1>
- Gottelman, A., Mills, M. J., Kinnison, D. E., Garcia, R. R., Smith, A. K., Marsh, D. R., et al. (2019). The whole atmosphere community climate model version 6 (WACCM6). *Journal of Geophysical Research: Atmospheres*, 124(23), 12380–12403. <https://doi.org/10.1029/2019JD030943>
- Hagan, M. E., Burrage, M. D., Forbes, J. M., Hackney, J., Randel, W. J., & Zhang, X. (1999). GSWM-98: Results for migrating solar tides. *Journal of Geophysical Research*, 104(A4), 6813–6828. <https://doi.org/10.1029/1998ja900125>
- Hedin, A. E. (1991). Extension of the MSIS thermospheric model into the middle and lower atmosphere. *Journal of Geophysical Research*, 96(A2), 1159–1172. <https://doi.org/10.1029/90ja02125>
- Heelis, R. A., Lowell, J. K., & Spiro, R. W. (1982). A model of the high-latitude ionospheric convection pattern. *Journal of Geophysical Research*, 87(A8), 6339. <https://doi.org/10.1029/ja087ia08p06339>
- Heelis, R. A., & Maute, A. (2020). Challenges to understanding the Earth's ionosphere and thermosphere. *JGR: Space Physics*, 125(7), e2019JA027497. <https://doi.org/10.1029/2019JA027497>
- Hurrell, J. W., Holland, M. M., Gent, P. R., Ghan, S., Kay, J. E., Kushner, P. J., et al. (2013). The community Earth system model: A framework for collaborative research. *Bulletin of the American Meteorological Society*, 94(9), 1339–1360. <https://doi.org/10.1175/BAMS-D-12-00121.1>
- Jee, G., Burns, A. G., Kim, Y. H., & Wang, W. (2009). Seasonal and solar activity variations of the Weddell Sea Anomaly observed in the TOPEX total electron content measurements. *Journal of Geophysical Research - Atmospheres*, 114(A4), A04307. <https://doi.org/10.1029/2008ja013801>
- Jin, H., Miyoshi, Y., Fujiwara, H., Shinagawa, H., Terada, K., Terada, N., et al. (2011). Vertical connection from the tropospheric activities to the ionospheric longitudinal structure simulated by a new Earth's whole atmosphere-ionosphere coupled model. *Journal of Geophysical Research*, 116(A1), A01316. <https://doi.org/10.1029/2010JA015925>
- Kalafatoglu Eyiguler, E. C., Shim, J. S., Kuznetsova, M. M., Kaymaz, Z., Bowman, B. R., Codrescu, M. V., et al. (2019). Quantifying the storm time thermospheric neutral density variations using model and observations. *Space Weather*, 17(2), 269–284. <https://doi.org/10.1029/2018SW002033>
- Kim, E., Jee, G., Wang, W., Kwak, Y.-S., Shim, J.-S., Ham, Y.-B., & Kim, Y. H. (2023). Hemispheric asymmetry of the polar ionospheric density investigated by ESR and JVD radar observations and TIEGCM simulations for the solar minimum period. *Journal of Geophysical Research: Space Physics*, 128(2), e2022JA031126. <https://doi.org/10.1029/2022ja031126>

- Laundal, K. M., Cnossen, I., Milan, S. E., Haaland, S. E., Coxon, J., Pedatella, N. M., et al. (2017). North–South asymmetries in Earth’s magnetic field. *Space Science Reviews*, 206(1–4), 225–257. <https://doi.org/10.1007/s11214-016-0273-0>
- Liu, H.-L., Bardeen, C. G., Foster, B. T., Lauritzen, P., Liu, J., Lu, G., et al. (2018). Development and validation of the whole atmosphere community climate model with thermosphere and ionosphere extension (WACCM-X 2.0). *Journal of Advances in Modeling Earth Systems*, 10(2), 381–402. <https://doi.org/10.1002/2017MS001232>
- Matsuo, T. (2020). Recent progress on inverse and data assimilation procedure for high-latitude ionospheric electrodynamics. In M. Dunlop & H. Lühr (Eds.), *Ionospheric multi-spacecraft analysis tools. ISSI scientific report series* (Vol. 17). Springer. https://doi.org/10.1007/978-3-030-26732-2_10
- Merkin, V., & Lyon, J. (2010). Effects of the low-latitude ionospheric boundary condition on the global magnetosphere. *Journal of Geophysical Research*, 115(A10), A10202. <https://doi.org/10.1029/2010JA015461>
- Millward, G. H., Müller-Wodrag, I. C. F., Aylward, A. D., Fuller-Rowell, T. J., Richmond, A. D., & Moffett, R. J. (2001). An investigation into the influence of tidal forcing on F region equatorial vertical ion drift using a global ionosphere-thermosphere model with coupled electrodynamics. *Journal of Geophysical Research*, 106(A11), 24733–24744. <https://doi.org/10.1029/2000JA000342>
- Newell, P. T., & Gjerloev, J. W. (2011). Substorm and magnetosphere characteristic scales inferred from the SuperMAG auroral electrojet indices. *Journal of Geophysical Research*, 116(A12), A12232. <https://doi.org/10.1029/2011JA016936>
- Newell, P. T., Sotirelis, T., & Wing, S. (2009). Diffuse, monoenergetic, and broadband aurora: The global precipitation budget. *Journal of Geophysical Research*, 114(A9), A09207. <https://doi.org/10.1029/2009JA014326>
- Perlongo, N. J., Ridley, A. J., Cnossen, I., & Wu, C. (2018). A year-long comparison of GPS TEC and global ionosphere-thermosphere models. *Journal of Geophysical Research: Space Physics*, 123(2), 1410–1428. <https://doi.org/10.1002/2017JA024411>
- Qian, L., Gan, Q., Wang, W., Cai, X., Eastes, R., & Yue, J. (2022). Seasonal variation of thermospheric composition observed by NASA GOLD. *Journal of Geophysical Research: Space Physics*, 127(6), e2022JA030496. <https://doi.org/10.1029/2022JA030496>
- Rastätter, L., Shim, J. S., Kuznetsova, M. M., Kilcommons, L. M., Knipp, D. J., Codrescu, M., et al. (2016). GEM-CEDAR challenge: Poynting flux at DMSP and modeled Joule heat. *Space Weather*, 14(2), 113–135. <https://doi.org/10.1002/2015SW001238>
- Reinisch, B., & Galkin, I. (2011). Global ionospheric Radio observatory (GIRO). *Earth Planets and Space*, 63(4), 377–381. <https://doi.org/10.5047/eps.2011.03.001>
- Richmond, A. D. (2021). Joule heating in the thermosphere. In W. Wang & Y. Zhang (Eds.), *Upper atmosphere dynamics and energetics (AGU Geophysical Monograph 261)* (pp. 3–18). John Wiley & Sons. <https://doi.org/10.1002/9781119815631.ch1>
- Richmond, A. D., Ridley, E. C., & Roble, R. G. (1992). A thermosphere/ionosphere general circulation model with coupled electrodynamics. *Geophysical Research Letters*, 19(6), 601–604. <https://doi.org/10.1029/92gl00401>
- Rideout, W., & Coster, A. (2006). Automated GPS processing for global total electron content data. *GPS Solutions*, 10(3), 219–228. <https://doi.org/10.1007/s10291-006-0029-5>
- Ridley, A. J., Deng, Y., & Toth, G. (2006). The global ionosphere-thermosphere model. *Journal of Atmospheric and Solar-Terrestrial Physics*, 68(8), 839–864. <https://doi.org/10.1016/j.jastp.2006.01.008>
- Roble, R. G., & Ridley, E. C. (1987). An auroral model for the NCAR thermospheric general circulation model (TGCM). In *Annales Geophysicae Series A* (Vol. 5, pp. 369–382).
- Roble, R. G., Ridley, E. C., Richmond, A. D., & Dickinson, R. E. (1988). A coupled thermosphere/ionosphere general circulation model. *Geophysical Research Letters*, 15(12), 1325–1328. <https://doi.org/10.1029/GL015i012p01325>
- Scherliess, L., Tsagouri, I., Yizengaw, E., Bruinsma, S., Shim, J. S., Coster, A., & Retterer, J. M. (2019). The International Community Coordinated Modeling Center space weather modeling capabilities assessment: Overview of ionosphere/thermosphere activities. *Space Weather*, 17(4), 527–538. <https://doi.org/10.1029/2018SW002036>
- Schunk, R. W., Scherliess, L., Eccles, V., Gardner, L. C., Sojka, J. J., Zhu, L., et al. (2021). Challenges in specifying and predicting space weather. *Space Weather*, 19(2), e2019SW002404. <https://doi.org/10.1029/2019SW002404>
- Sharber, J. R., Link, R., Frahm, R. A., Winningham, J. D., Lummerzheim, D., Rees, M. H., et al. (1996). Validation of UARS PEM electron energy deposition. *Journal of Geophysical Research*, 101(D6), 9571–9582. <https://doi.org/10.1029/95jd02702>
- Shim, J. S., Jee, G., & Scherliess, L. (2017). Climatology of plasmaspheric total electron content obtained from Jason 1 satellite. *Journal of Geophysical Research: Space Physics*, 122(2), 1611–1623. <https://doi.org/10.1002/2016JA023444>
- Shim, J. S., Kuznetsova, M., Rastätter, L., Bilitza, D., Butala, M., Codrescu, M., et al. (2014). Systematic evaluation of ionosphere/thermosphere (IT) models: CEDAR electrodynamics thermosphere ionosphere (ETI) challenge (2009–2010). In *Modeling the ionosphere-thermosphere system, AGU Geophysical Monograph series*.
- Shim, J. S., Kuznetsova, M., Rastätter, L., Bilitza, D., Butala, M., Codrescu, M., et al. (2012). CEDAR Electrodynamics Thermosphere Ionosphere (ETI) Challenge for systematic assessment of ionosphere/thermosphere models: Electron density, neutral density, NmF2, and hmF2 using space based observations. *Space Weather*, 10, S10004. <https://doi.org/10.1029/2012SW000851>
- Shim, J. S., Kuznetsova, M., Rastätter, L., Hesse, M., Bilitza, D., Butala, M., et al. (2011). CEDAR electrodynamics thermosphere ionosphere (ETI) challenge for systematic assessment of ionosphere/thermosphere models: NmF2, hmF2, and vertical drift using ground-based observations. *Space Weather*, 9(12), S12003. <https://doi.org/10.1029/2011SW000727>
- Shim, J. S., Rastätter, L., Kuznetsova, M., Bilitza, D., Codrescu, M., Coster, A. J., et al. (2017). CEDAR-GEM challenge for systematic assessment of Ionosphere/thermosphere models in predicting TEC during the 2006 December storm event. *Space Weather*, 15(10), 1238–1256. <https://doi.org/10.1002/2017SW001649>
- Shim, J. S., Tsagouri, I., Goncharenko, L., Rastätter, L., Kuznetsova, M., Bilitza, D., et al. (2018). Validation of ionospheric specifications during geomagnetic storms: TEC and foF2 during the 2013 March storm event. *Space Weather*, 16(11), 1686–1701. <https://doi.org/10.1029/2018SW002034>
- Solomon, S. C., Burns, A. G., Emery, B. A., Mlynczak, M. G., Qian, L., Wang, W., et al. (2012). Modeling studies of the impact of high-speed streams and co-rotating interaction regions on the thermosphere-ionosphere. *Journal of Geophysical Research*, 117(A9), A00L11. <https://doi.org/10.1029/2011JA017417>
- Sorathia, K., Merkin, V., Panov, E., Zhang, B., Lyon, J., Garretson, J., et al. (2020). Ballooning-interchange instability in the near-Earth plasma sheet and auroral beads: Global magnetospheric modeling at the limit of the MHD approximation. *Geophysical Research Letters*, 47(14), e2020GL088227. <https://doi.org/10.1029/2020GL088227>
- Toffoletto, F., Sazykin, S., Spiro, R., & Wolf, R. (2003). Inner magnetospheric modeling with the rice convection model. *Space Science Reviews*, 107(1–2), 175–196. <https://doi.org/10.1023/A:1025532008047>
- Tsagouri, I., Goncharenko, L., Shim, J. S., Belehaki, A., Buresova, D., & Kuznetsova, M. M. (2018). Assessment of current capabilities in modeling the ionospheric climatology for space weather applications: FoF2 and hmF2. *Space Weather*, 16(12), 1930–1945. <https://doi.org/10.1029/2018SW002035>

- Webb, P. A., Kuznetsova, M. M., Hesse, M., Rastaetter, L., & Chulaki, A. (2009). Ionosphere-thermosphere models at the community coordinated modeling center. *Radio Science*, *44*(1), RS0A34. <https://doi.org/10.1029/2008RS004108>
- 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>
- Yamazaki, Y., & Richmond, A. D. (2013). A theory of ionospheric response to upward-propagating tides: Electrodynamic effects and tidal mixing effects. *Journal of Geophysical Research: Space Physics*, *118*, 5891–5905. <https://doi.org/10.1002/jgra.50487>
- Zhang, B., Sorathia, K. A., Lyon, J. G., Merkin, V. G., Garretson, J. S., & Wiltberger, M. (2019). GAMERA: A three-dimensional finite-volume MHD solver for non-orthogonal curvilinear geometries. *The Astrophysical Journal Supplement Series*, *244*(1), 20. <https://doi.org/10.3847/1538-4365/ab3a4c>
- Zhao, H., Li, X., Baker, D. N., Claudepierre, S. G., Fennell, J. F., Blake, J. B., et al. (2016). Ring current electron dynamics during geomagnetic storms based on the Van Allen Probes measurements. *Journal of Geophysical Research: Space Physics*, *121*(4), 3333–3346. <https://doi.org/10.1002/2016JA022358>