

Quality Assessment Report:

Graphical Turbulence Guidance (GTG) version 2.3

Michael P. Kay^{*12}, Judy K. Henderson¹, Stacey A. Krieger¹², Jennifer L. Mahoney¹,
Lacey D. Holland³, and Barbara G. Brown³

¹NOAA Earth System Research Laboratory, Boulder, CO

²Cooperative Institute for Research in Environmental
Sciences/University of Colorado, Boulder, CO

³National Center for Atmospheric Research, Boulder, CO

Quality Assessment Product Development Team

30 August 2006

Table of Contents

1. Introduction.....	8
2. Approach.....	8
3. Algorithms and Forecasts.....	9
4. Data.....	10
5. Methods.....	11
5.1 Creation of forecast/observation pairs.....	11
5.2 Stratifications.....	15
6. Results.....	17
6.1 Comparison of GTG2.3P with GTG2.3E.....	17
6.2 GTG2.3E performance.....	22
6.3 GTG2.3E vs. GTG Comparison.....	36
6.4 GTG2.3E vs. AIRMETs Comparison.....	45
6.5 GTG2.3E vs. SIGMETs Comparison.....	53
7. Conclusions.....	54
8. References.....	56
Appendix: Categorical distributions of forecasts and observations for GTG.....	58

List of Figures

Fig. 1: RTVS display of GTG version 2.3E 6-h forecast on 01 Dec. 2005. Top panel shows plan view of forecast using a threshold of 0.375, observation locations, and lightning data. Bottom panels show vertical distribution of forecasts and observations at all observation locations.....	13
Fig. 2: AWC forecast regions.....	16
Fig. 3: Climatological regions used for subcontinental characterization of GTG performance.....	17
Fig. 4: ROC diagrams for upper levels (20,000 to 40,000 ft) for the (a) 3-h, (b) 6-h, (c) 9-h, and (d) 12-h forecasts for GTG2.3E and GTG2.3P.....	19
Fig. 5: As in Fig. 4 except for midlevels (10,000 to 20,000 ft).....	20
Fig. 6: Height series of (a) MOG POD _y and (b) MOG POD _n for GTG2.3E and GTG2.3P. Data are plotted at the bottom of each 5,000 ft vertical layer. Forecast threshold is 0.475... 21	21
Fig. 7: ROC diagram for GTG2.3E for the 3-, 6-, 9-, and 12-h lead times for upper levels (20,000 to 40,000 ft).....	23
Fig. 8: Threshold plots of (a) MOG POD _y and (b) % Volume for GTG2.3E at upper levels for the 3-, 6-, 9- and 12-h lead times.	24
Fig. 9: As in Fig. 7 except for midlevels (10,000 to 20,000 ft).....	25
Fig. 10: MOG POD _y and MOG POD _n height series for GTG2.3E. Forecast threshold is 0.475.....	26
Fig. 11: MOG POD _y height series for GTG2.3E for 3-, 6-, 9-, and 12-h lead times. Forecast threshold is 0.475.....	27
Fig. 12: As in Fig. 11 except for MOG POD _n	28
Fig. 13: ROC diagram for GTG2.3E stratified by AWC forecast region for upper levels (20,000 to 40,000 ft).....	29
Fig. 14: As in Fig. 13 except for midlevels (10,000 to 20,000 ft).....	30
Fig. 15: Threshold plots of (a) MOG POD _y and (b) % Volume for GTG2.3E at midlevels for the stratified by AWC forecast region.....	31
Fig. 16: Overall (a) POD _y , (b) POD _n , (c) % Volume, and (d) TSS values for GTG2.3E at upper levels for the climatological regions. Forecast threshold is 0.475.....	32
Fig. 17: As in Fig. 16 except for midlevels.....	33
Fig. 18: ROC diagram for GTG2.3E and operational GTG at upper levels.....	37
Fig. 19: Threshold plots for (a) MOG POD _y and (b) % Volume for GTG2.3E and operational GTG at upper levels.....	38
Fig. 20: Height series of MOG POD _y for GTG2.3E and operational GTG at upper levels. GTG2.3E threshold is 0.475 and GTG threshold is 0.375.....	40
Fig. 21: As in Fig. 20 except for MOG POD _n	41
Fig. 22: ROC diagrams for GTG2.3E and operational GTG for the (a) West, (b) Central, and (c) East AWC forecast regions at upper levels.....	42
Fig. 23: Overall (a) POD _y , (b) POD _n , (c) % Volume, and (d) TSS values for GTG at	

upper levels for the climatological regions. Forecast threshold is 0.375.....44

Fig. 24: ROC diagram for GTG2.3E and AIRMETs at upper levels. GTG2.3 points corresponding to the light and moderate thresholds are highlighted.....46

Fig. 25: Boxplot of % Volume values for GTG2.3E with thresholds 0.475, 0.300, 0.250, and AIRMETs at upper levels.....48

Fig. 26: As in Fig. 24 except for midlevels.....49

Fig. 27: Height series of MOG PODy for GTG2.3E and AIRMETs. GTG2.3E threshold is 0.475.....51

Fig. 28: Boxplot of % Volume values for GTG2.3E and AIRMETs in the mid- and upper levels. GTG2.3E threshold is 0.475.....52

Fig. 29: Height series of MOG PODn for GTG2.3E and AIRMETs. GTG2.3E threshold is 0.475.....53

List of Tables

Table 1: The set of issue times and lead times for GTG2.3 and operational GTG used in this report.....	10
Table 2: Contingency table for a dichotomous forecast situation.....	11
Table 3: Dichotomous summary statistics used in this report. Terms in definitions column are linked to the contingency table presented in Table 2.....	14
Table 4: Dichotomous thresholds applied to GTG and GTG2.3 for POD and PODn computations used to resolve ROC curves.....	15
Table 5: Mapping of forecast and observed values to categories used in this report. Values for GTG and GTG2.3 represent lower bounds of the ranges of data associated with each category.....	15
Table 6: List of the climatologically-defined regions within the Continental U.S. used in this study and their abbreviations.....	16
Table 7: Joint distribution of GTG2.3E forecasts and PIREPs at upper levels.....	34
Table 8: Conditional probability of a forecast for each observation category, $p(\text{flx})$, for upper levels.....	34
Table 9: Joint distribution of GTG2.3E forecasts and PIREPS at midlevels.....	35
Table 10: Conditional probability of a forecast for each observation category, $p(\text{flx})$, for midlevels.....	36
Table 11: ROC area under the curve values for GTG2.3E and GTG for each AWC region.	43
Table 12: Difference table (GTG2.3E-GTG) for the conditional probability of a forecast given an observation, $p(\text{flx})$, between GTG2.3E and GTG.....	45
Table 13: MOG PODy, MOG PODn, and median % Volume values for GTG2.3E and AIRMETs at upper levels.....	47
Table 14: MOG PODy, MOG PODn, and median % Volume for GTG2.3E and AIRMETs in midlevels.....	49
Table 15: Regional AUC values.....	50
Table 16: Distribution of forecast/observation pairs and mean % Volume for GTG2.3E with thresholds 0.475, 0.625, and 0.8 and SIGMETs when severe intensity PIREPs were reported for both mid- and upper levels.....	54

Quality Assessment Report:

Graphical Turbulence Guidance (GTG) version 2.3

Michael P. Kay^{*12}, Judy K. Henderson¹, Stacey A. Krieger¹², Jennifer L. Mahoney¹,
Lacey D. Holland³, and Barbara G. Brown³

¹NOAA Earth System Research Laboratory, Boulder, CO

²Cooperative Institute for Research in Environmental
Sciences/University of Colorado, Boulder, CO

³National Center for Atmospheric Research, Boulder, CO

Quality Assessment Product Development Team

SUMMARY

Results of an evaluation of GTG version 2.3 (GTG2.3) are presented in this report. The algorithm was analyzed from 1 November 2005 through 31 January 2006. Additionally, GTG2.3 was compared to several existing operational turbulence forecasts including GTG version 1.0 (GTG). Forecasts were verified with pilot reports of turbulence.

The primary findings are:

- The two versions of GTG studied here (GTG2.3E and GTG2.3P) had nearly identical performance. The introduction of eddy dissipation rate information into the algorithm did not have any adverse effects on the results.
- GTG2.3E performed well in the mid- and upper levels for forecasts of moderate or greater turbulence.
- GTG2.3E showed limited ability to forecast the correct intensity of turbulence. It performs best for the None and Moderate categories.
- When compared to Airman's Meteorological Advisories (AIRMETs), GTG2.3E performed well in both mid- and upper levels. GTG2.3E forecast volumes were much

smaller than the volumes associated with the AIRMETS.

- GTG2.3E and Significant Meteorological Advisories (SIGMETs) did a poor job forecasting severe turbulence as was indicated by the statistics that were derived from the limited number of severe turbulence reports. GTG2.3E forecast volumes were several orders of magnitude smaller than those produced by SIGMETs. This result may have been due in part to the small numbers of PIREPs reporting severe turbulence.

1. INTRODUCTION

This report summarizes the quality of mid- and upper-level turbulence forecasts produced by the second generation (version 2.3) Graphical Turbulence Guidance (GTG) product (denoted as GTG2.3), which is under consideration for transition from experimental to operational status within the Aviation Weather Technology Transfer (AWTT) process. Takacs et al. (2004) evaluated the quality of the previous GTG product (version 2.0), which was accepted by the AWTT Board as an experimental product in 2004.

The GTG2.3 algorithm combines input from numerous data sources to provide forecasts of clear-air turbulence over the continental United States (CONUS) at altitudes greater than 10,000 ft (Sharman et al. 2004). GTG2.3 has been developed by the Federal Aviation Administration Weather Research Program's (AWRP) Turbulence Product Development Team (TPDT). The AWRP's Quality Assessment Team (QA PDT) evaluated GTG2.3 through specific algorithm comparison studies. The studies were conducted using the Real-Time Verification System (RTVS; Mahoney et al. 2002) developed by staff at the National Oceanic and Atmospheric Administration's Earth System Research Laboratory Global Systems Division.

The report is organized in the following manner. Section 2 provides an overview of the approach taken in evaluating GTG2.3. Section 3 describes the algorithms and forecasts that are assessed in this evaluation. The data used are described next in Section 4. Section 5 presents the verification methods that are employed for the evaluation while results are presented in Section 6. Finally, the report concludes with discussion and a summary of results in Section 7.

2. APPROACH

GTG2.3 was evaluated with respect to other operational turbulence forecasts, which included the operational version of GTG, Airman's Meteorological Advisories (AIRMETs) and Significant Meteorological Advisories (SIGMETs). It should be noted that this report is not intended as an evaluation of the turbulence AIRMETs and SIGMETs. The intercomparison is made in such a way as to treat all forecasts as equitably as possible. More explanation is provided in Section 5. Users of these statistics should keep these assumptions in mind when evaluating the strengths and weaknesses of each type of forecast.

Due to the emphasis placed on forecasting mid- and upper-level turbulence, the evaluation focused on the layers of the atmosphere from 10,000 to 20,000 ft, and 20,000 to 40,000 ft. In addition to the entire CONUS, forecasting performance across three large and fifteen small geographic subregions was also considered. Forecasts issued during the period 1 November 2005 through 31 January 2006 were analyzed. The verification approach applied in this evaluation is identical to the approach taken in previous studies (e.g., Takacs et al. 2004). Additional analyses that focus on the qualitative trend between

the forecast and observed turbulence intensity categories are also included in this report.

3. ALGORITHMS AND FORECASTS

This report is focused on the evaluation of GTG2.3 and its transition to National Weather Service (NWS) operations. The turbulence forecasts used for intercomparison with GTG2.3 in this report represent the current operational guidance available to forecasters. The forecasts considered in this report are:

GTG: This algorithm, formerly known as the Integrated Turbulence Forecast Algorithm (ITFA; Sharman et al. 2002), is intended to forecast moderate or greater (MOG) clear-air turbulence at altitudes from 20,000 to 40,000 ft. GTG forecasts are created by dynamically combining, and optimally weighting, a series of turbulence diagnostics using a fuzzy logic system. The Rapid Update Cycle (RUC; Benjamin 1998) model is used to provide the background fields from which the diagnostics are computed. The GTG is produced operationally by the National Weather Service's Aviation Weather Center (NWS/AWC).

GTG2.3 (versions E and P): The GTG2.3 forecast system represents an incremental improvement to the current experimental version of GTG (version 2.0) and expands the capability of the operational version of GTG (denoted as GTG) by providing turbulence predictions at both mid- (10,000 to 20,000 ft) and upper levels (20,000 ft and above). Additional changes include new diagnostics and the use of the 13-km RUC model for the large-scale atmospheric processes. For information on the performance of the experimental version of GTG (GTG2) and support for its transition to experimental status, see Takacs et al. (2004).

Two versions of GTG2.3 are currently produced by the TPDT: GTG2.3P and GTG2.3E. The configurations of the algorithms are identical except for the turbulence observations used in the forecast tuning, and therefore, the internal weighting of the various diagnostics. GTG2.3P is tuned using pilot reports (PIREPs) whereas GTG2.3E utilizes *in situ* eddy dissipation rate (EDR) measurements that are available from numerous commercial aircraft in addition to PIREPs. EDR data are used to augment many of the shortcomings of PIREP data: they provide high frequency, objective, quantitative observations of turbulence that are independent of aircraft size (Cornman et al. 2004). The EDR observations also lead to increased numbers of “No” reports of turbulence. No PIREPs of turbulence are much less frequent than Yes turbulence PIREPs despite the fact that much of the atmosphere is turbulence-free. Understanding the variations in forecast quality that are due to the differences between the two versions of GTG2.3 is important, since situations could occur when the EDR data are unavailable in operations. In these situations, the GTG2.3E version of the algorithm will revert to version GTG2.3P.

AIRMETs: AIRMETs are advisories issued for en-route hazardous weather phenomena (NWS; 2003). In this study, only AIRMET forecasts issued for turbulence were

considered. Turbulence AIRMETs are issued for moderate or greater turbulence conditions. Forecasts are issued four times per day for periods up to six hours and may be amended as needed. In this study, only non-amended AIRMETs are considered. The temporal aspect of AIRMETs, as it relates to the intercomparison with GTG2.3, is discussed further in Section 6. Attributes used from these forecasts include the areal extent of the forecast and the vertical layer where turbulence is expected. While AIRMETs provide more detailed information that could potentially aid in the analysis, this information is not encoded in a standard way and therefore cannot be decoded and used systematically in verification studies.

SIGMETs: SIGMETs are in-flight advisories that warn of internationally specified weather phenomena of an intensity and/or extent that concerns pilots and operators of all aircraft (NWS, 2003). SIGMETs can be issued at any time and are valid for up to four hours. In the conterminous United States, SIGMETs have been separated into two types: convective (i.e., thunderstorm-related) and nonconvective. In this study, only nonconvective turbulence SIGMETs are considered. Hereafter, the term SIGMET will be used to represent nonconvective SIGMETs.

4. DATA

Data were collected for analysis from 1 November 2005 - 31 January 2006. A subset of all possible GTG2.3E and GTG2.3P issuance and lead times were used in this study (Table 1). The study focuses on the valid time period between 1500 and 0000 UTC in order to maximize the number of pilot reports available for verification. GTG2.3 forecasts with issuance times of 1800 UTC were not used in this study since they were used by the TPDT to alter weighting parameters within GTG2.3. GTG2.3 algorithms were applied to data from the 13-km RUC model output obtained from the NWS.

Table 1: The set of issue times and lead times for GTG2.3 and operational GTG used in this report.

<i>Issue Time (UTC)</i>	<i>Lead Time (h)</i>
1200	3, 6, 9, 12
1500	3, 6, 9
2100	3

PIREPs of turbulence were used as the observational dataset for this evaluation. PIREPs are subjective, non-systematic reports of aircraft encounters with weather hazards such as turbulence and icing. They can also be issued by pilots when a hazard is expected but none is observed (these are referred to as null PIREPs in this study). PIREPs represent the best available operational source of turbulence observations available today. The attributes of PIREPs that were considered include the report location (latitude,

longitude, and altitude) and the intensity of the turbulence encountered. Because a PIREP may be issued for a hazard over a vertical range instead of a single level, they are broken down to create a series of one or more reports for each PIREP. For instance, a PIREP for turbulence between 34,000 ft and 37,000 ft would be split into a series of four PIREPs having the same latitude, longitude, and time but with altitudes 34,000, 35,000, 36,000 and 37,000 ft, respectively. These reports are then used for verification. Throughout this report the term PIREP will refer to these post-processed reports instead of the original PIREPs unless otherwise noted. For more information on PIREPs and their characteristics, see Schwartz (1996). No attempt to stratify the PIREPs by origin (i.e., mountain wave, convection, and clear-air) was pursued for this evaluation.

5. METHODS

This section describes the verification methodology and statistics employed to assess GTG2.3. The methodology is similar to past evaluations of turbulence by the QAPDT. More detail and background can be found in Brown and Mahoney (1998). Verification results were obtained from the Real-Time Verification System (RTVS) (Mahoney et al. 2002).

5.1 Creation of forecast/observation pairs

In order to assess the accuracy of GTG2.3, the forecast values must be matched in space and time to the PIREP observations. Because GTG2.3 is a gridded product, and PIREPs are point observations, the product is only assessed at observation locations. The gridded forecast values are bilinearly interpolated to the PIREP positions (latitude, longitude, and altitude) using the four surrounding grid points representing the bounding volume for a PIREP observation. In order to allow for timing differences between reports and forecast valid times, a temporal window of ± 60 min. is used to collect PIREPs for each forecast valid time. This window also increases the set of observations available for verification.

Table 2: Contingency table for a dichotomous forecast situation.

		<i>Observed</i>	
		<i>Yes</i>	<i>No</i>
<i>Forecast</i>	<i>Yes</i>	YY	YN
	<i>No</i>	NY	NN

Once the forecast/observation pairs have been generated, verification is performed in one of two fundamental ways. The first approach is to treat the forecast dichotomously (i.e., Yes/No) by thresholding the forecast values to derive a 2x2 contingency table (Table 2). The choice of forecast thresholds for both GTG and GTG2.3 and their associated turbulence categories (as used in PIREPs) are discussed later in this section. For all

analyses, unless otherwise noted, PIREPs representing moderate or greater intensities are treated as “Yes” observations of turbulence. “No” observations are represented by PIREPs with reported intensities less than moderate. Moderate intensity is used as the threshold for Yes and No events because moderate-or-greater (MOG) turbulence represents a greater hazard to aviation than do lesser intensities. Is the most often observed level of turbulence intensity (aside from no turbulence).

Because of the limitations of PIREP data, which do not sample the airspace systematically, not all scalar dichotomous summary statistics can be computed (Brown and Young 2000). The three statistics that will be the focus of this report are the probability of detecting an event (POD_y), the probability of detection of a non-event (POD_n), and the True Skill Statistic (TSS), which can be represented as $POD_y + POD_n - 1$. TSS is a measure of a forecast's ability to distinguish between Yes (turbulence) and No (no turbulence) events. Due to the nature of the non-systematic observations, POD_y and POD_n must not be considered true probabilities but instead as proportions of the observed set of Yes and No PIREPs that are correctly categorized by the forecasts. An additional summary statistic that is used within the report is Percent Volume (or % Volume). This statistic measures the percent of the total air space volume where turbulence is forecast. Possible values range from 0 to 100 %. The value is not itself a measure of accuracy but should instead be used in conjunction with other scores such as POD_y to gain greater understanding of the forecast quality as a function of areal forecast coverage. A description of these statistics is given in Table 3.

A sample GTG2.3E forecast obtained from RTVS is presented in Fig. 1 to illustrate the verification mechanics described above¹. The top panel of the display shows the 6-h lead time forecast issued at 1500 UTC on 01 December 2005. A threshold of 0.375 was used to create a dichotomous forecast situation. The spatial distribution of PIREPs is illustrated by the numbers on the map. The lower panels show the vertical profile of the GTG2.3E forecast at each of the PIREP locations. Within each column, multiple PIREP values may be noted depending on the depth of the layer reported within the original PIREP. Recall that PIREPs are broken into multiple reports each having a vertical depth of 1000 ft. The set of PIREPs and the corresponding forecast values at the PIREP locations make up the cells of the 2x2 contingency table (Table 2) from which the relevant statistics are computed (Table 3).

¹ GTG2.3 forecasts on the Experimental Aviation Digital Data Service are available from the following URL: <http://www.weather.aero/turbulence/>

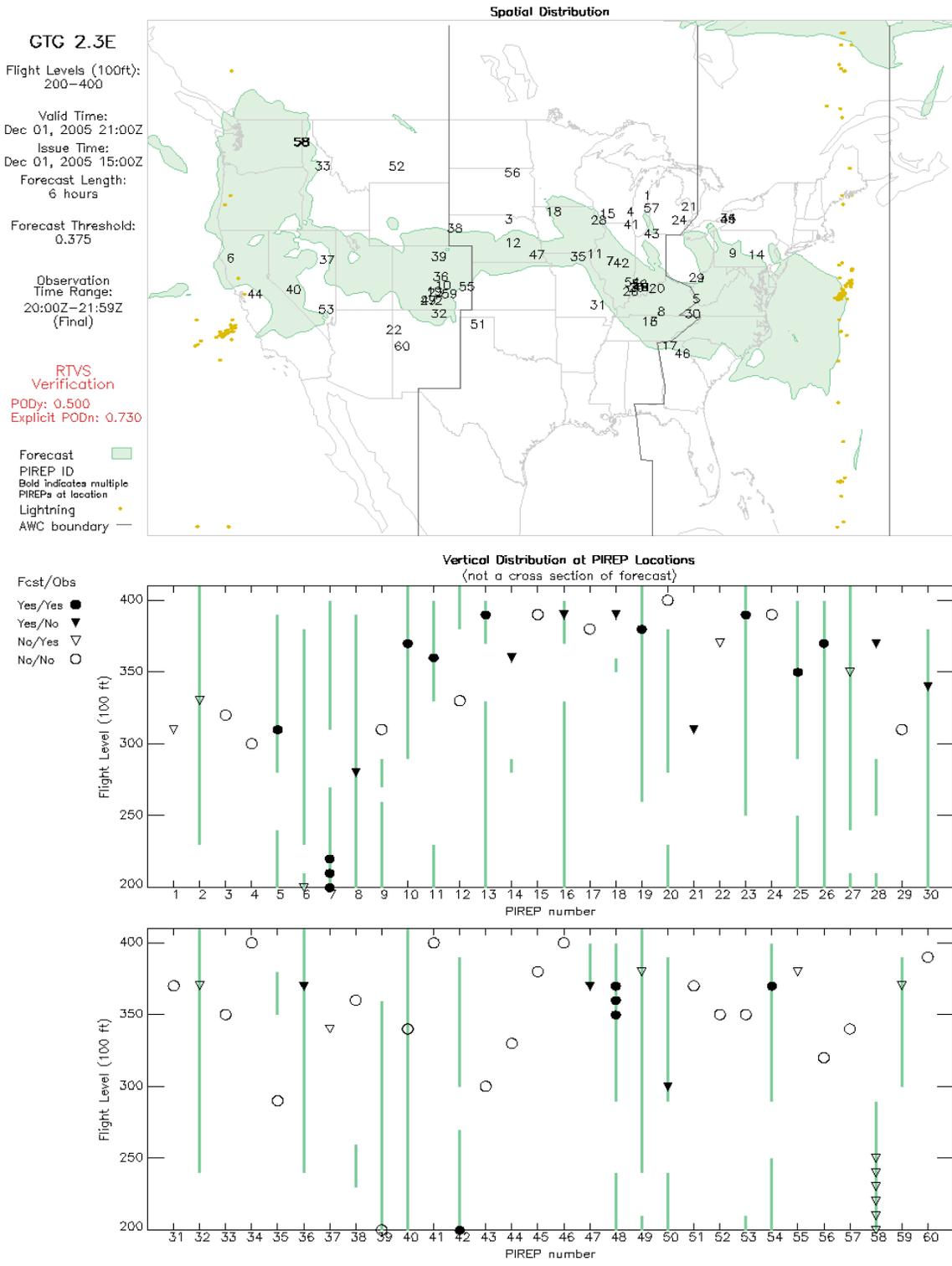


Fig. 1: RTVS display of GTG version 2.3E 6-h forecast on 01 Dec. 2005. Top panel shows plan view of forecast using a threshold of 0.375, observation locations, and lightning data. Bottom panels show vertical distribution of forecasts and observations at all observation locations.

Table 3: Dichotomous summary statistics used in this report. Terms in definitions column are linked to the contingency table presented in Table 2.

Statistic	Description	Definition
POD _y	Proportion of events detected correctly	$YY/(YY+NY)$
POD _n	Proportion of non-events detected correctly	$NN/(NN+YN)$
TSS	True Skill Statistic	$POD_y + POD_n - 1$
% Volume	Percent of the possible volume covered by the forecast	$100 * \text{Volume}_{\text{forecast}} / \text{Volume}_{\text{possible}}$

The second analysis approach involves the use of signal detection theory and, more specifically, the Relative Operating Characteristic (ROC) diagram (Mason 1982). Rather than choosing a single decision threshold (such as 0.25) from the forecast values to compute the dichotomous statistics, a set of thresholds is chosen and for each threshold the dichotomous statistics POD_y and POD_n are computed. Each of these pairs of points is then plotted on a diagram called a ROC diagram where the x-axis is 1-POD_n and the y-axis is POD_y. The line connecting these points is the Relative Operating Characteristic curve. If a forecast shows no ability to distinguish between Yes and No events, the POD_y and POD_n values will be identical and values will lie on the diagonal of the diagram for all decision thresholds. Forecasts for which POD_y exceeds 1-POD_n have skill in separating the events from nonevents; for these forecasts the points will lie above the diagonal on the diagram. Perfect forecasts will have points near the upper left-hand corner of the diagram where correct forecasts are maximized and false alarms are minimized. The area under the ROC curve, commonly referred to as the AUC, is used as a summary measure of performance. Possible values for the AUC range from 0 to 1, with values of 0.5 indicating no skill. For the ROC computations, additional thresholds were used to create sufficient data points to resolve the curves. Slight differences in thresholds were necessary for GTG and GTG2.3 owing to differences in thresholds used to define the categories None, Light, Moderate, and Severe. The categories are discussed later in this section. The sets of thresholds used for each algorithm are presented in Table 4.

Table 4: Dichotomous thresholds applied to GTG and GTG2.3 for POD and PODn computations used to resolve ROC curves.

Algorithm	Thresholds
GTG	0.06, 0.125, 0.15, 0.20, 0.25, 0.312, 0.375, 0.437, 0.50, 0.562, 0.625, 0.75
GTG2.3	0.06, 0.125, 0.15, 0.20, 0.25, 0.312, 0.375, 0.475, 0.50, 0.562, 0.625, 0.75, 0.80

The focus for the overall evaluation of GTG2.3 is the accuracy of its predictions of moderate or greater (MOG) intensity turbulence. This aspect of GTG2.3 performance is well captured by the techniques described above. However, when GTG2.3 becomes an operational forecast product and its output is displayed to end users through the operational Aviation Digital Data Service (ADDS), forecasts of specific turbulence intensity categories of None, Light, Moderate, and Severe will be provided. Therefore, it is imperative that the ability of GTG2.3 to predict the correct category of turbulence intensity also be evaluated. An analysis will be performed that focuses on the qualitative trend between the forecast and observed categories. Table 5 shows the mapping between the categorical labels and the associated forecast thresholds that define the lower bound of the range of values tied to each label. For instance, GTG2.3 forecasts of moderate intensity turbulence are associated with forecast values greater than or equal to 0.475 and less than 0.8.

Table 5: Mapping of forecast and observed values to categories used in this report. Values for GTG and GTG2.3 represent lower bounds of the ranges of data associated with each category.

Category	GTG	GTG2.3	PIREP Intensities
None	0.0	0.0	Smooth/None, Smooth to occasional light
Light	0.125	0.3	Light to occasional moderate
Moderate	0.375	0.475	Moderate, Moderate to occasional severe
Severe	0.625	0.8	Severe, Severe to occasional extreme, Extreme

5.2 Stratifications

Data were analyzed over the CONUS and nearby oceanic regions according to the east,

west, and central forecast regions used by the AWC (Fig. 2). The national domain is simply the aggregate of the three regions. Additionally, verification results for 15 smaller regions that are based upon differing climatological attributes (Table 6 and Fig. 3) are considered in Section 6 with results presented for each region. For mid- and upper-level results, forecast/observation counts were aggregated vertically through the 10,000 to 20,000 ft and 20,000 to 40,000 ft layers, respectively. For evaluation of the forecast performance across vertical profiles, data were aggregated into 5,000 ft layers.

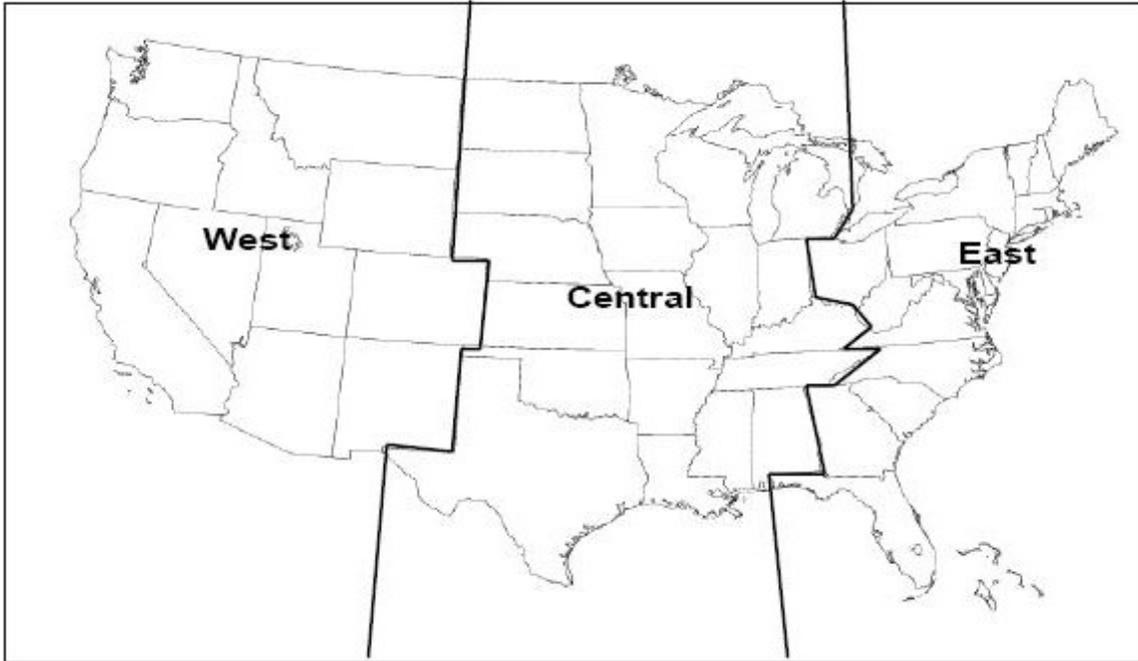


Fig. 2: AWC forecast regions.

Table 6: List of the climatologically-defined regions within the Continental U.S. used in this study and their abbreviations.

<i>Abbreviation</i>	<i>Region</i>
WCN	West Coast North
WCS	West Coast South
IMN	Intermountain North
IMS	Intermountain South
RMN	Rocky Mountain
HPN	High Plains North
HPS	High Plains South
GPN	Great Plains North

<i>Abbreviation</i>	<i>Region</i>
GPS	Great Plains South
GLA	Great Lakes
OMV	Ohio and Mississippi Valley
GCO	Gulf Coast
APP	Appalachians
ECN	East Coast North
ECS	East Coast South

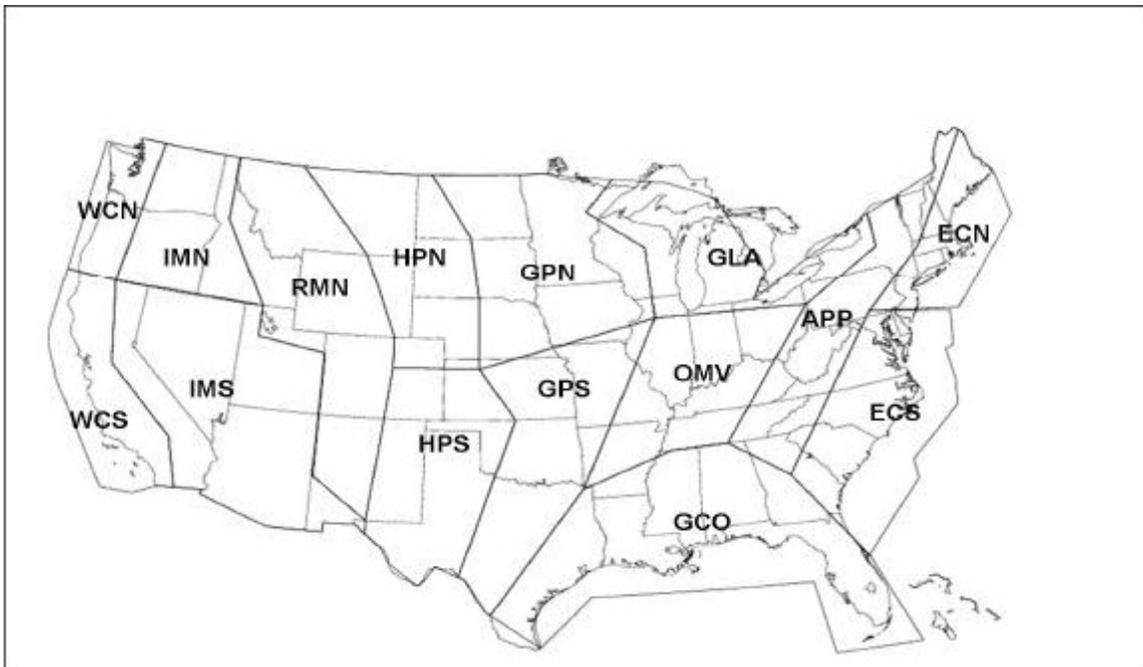


Fig. 3: Climatological regions used for subcontinental characterization of GTG performance.

6. RESULTS

6.1 Comparison of GTG2.3P with GTG2.3E

The purpose of this section is to compare the performance of the two versions of GTG2.3 that are being considered for transition to operations, GTG2.3E and GTG2.3P. The two algorithms are identical, but use differing observations at initialization: GTG2.3P uses PIREPs alone whereas GTG2.3E incorporates EDR measurements in addition to PIREPs. Recall the importance of this evaluation. If EDR observations are

unavailable in operations, GTG2.3E will revert to GTG2.3P. This intercomparison should illuminate any differences that may arise owing to the incorporation of EDR data into GTG2.3E.

Overall performance, depicted through ROC curves, is shown in Fig. 4 for the 3-, 6-, 9-, and 12-h lead times for the upper levels (20,000 to 40,000 ft) and in Fig. 5 for midlevels (10,000 to 20,000 ft). Both forecasts show convex curves indicating significant skill at discriminating between Yes and No observations of turbulence throughout the airspace. Minor differences appear between the MOG POD_y and 1-MOG POD_n values for the two forms of the algorithm, particularly at the lower thresholds, but these differences do not appear to be significant. Areas under the ROC curves are identical for each forecast system at each lead time, with values of 0.87, 0.85, 0.84, and 0.84 for the 3-, 6-, 9-, and 12-h lead times, respectively. The two algorithms appear identical in the overall results for both upper- and midlevels.

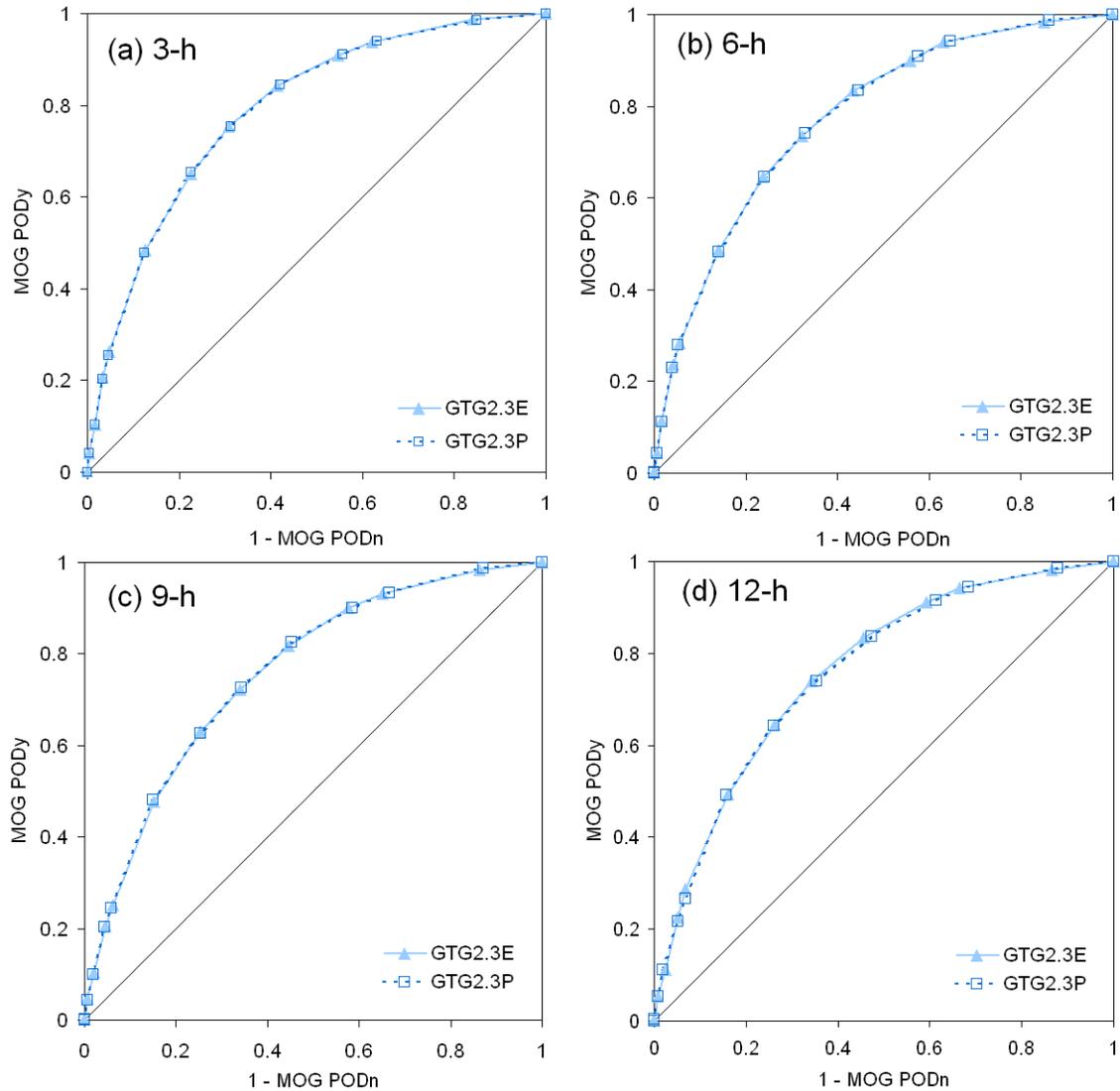


Fig. 4: ROC diagrams for upper levels (20,000 to 40,000 ft) for the (a) 3-h, (b) 6-h, (c) 9-h, and (d) 12-h forecasts for GTG2.3E and GTG2.3P.

Height series of both MOG PODy and MOG PODn, using a threshold of 0.475 and shown in Fig. 6, illustrate good agreement between the two algorithms from 10,000 to 40,000 ft. Minor differences are apparent in the layer from 20,000 ft to 35,000 ft. GTG2.3E has slightly larger MOG PODy values while having slightly smaller MOG PODn values than GTG2.3P.

The results presented above suggest that the introduction of EDR measurements into GTG2.3 does not decrease the skill of the algorithm, nor do they appear to significantly enhance the algorithm. This may be due to the fact that EDR measurements are still not particularly widespread. Moreover, the *in situ* EDR observations were not used in the verification analyses; inclusion of these observations likely would have some impacts on the verification results. Given the known benefits of EDR data at detecting and

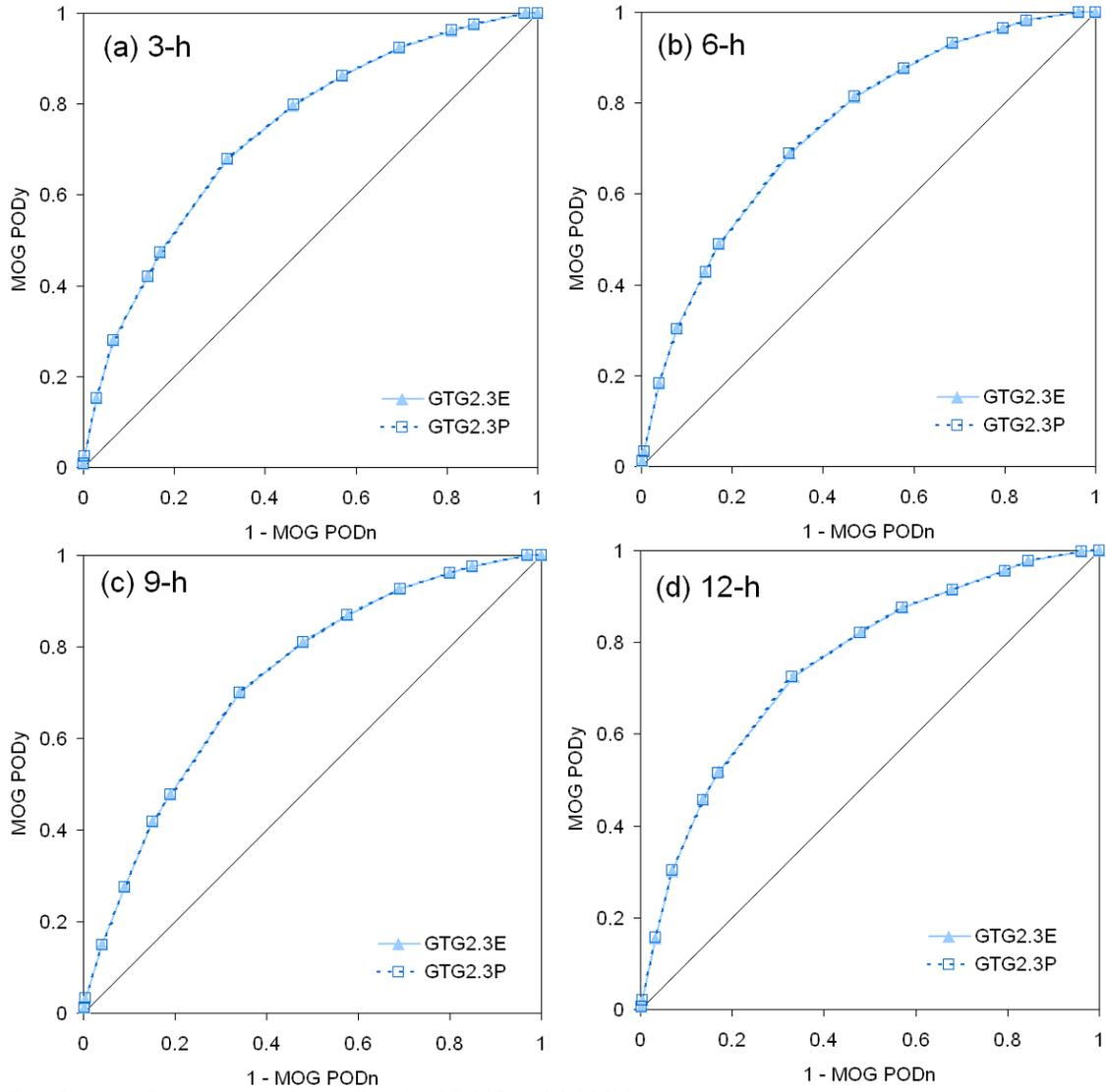


Fig. 5: As in Fig. 4 except for midlevels (10,000 to 20,000 ft).

quantifying turbulence in the free atmosphere, the rest of this report will focus solely on the quality of GTG2.3E.

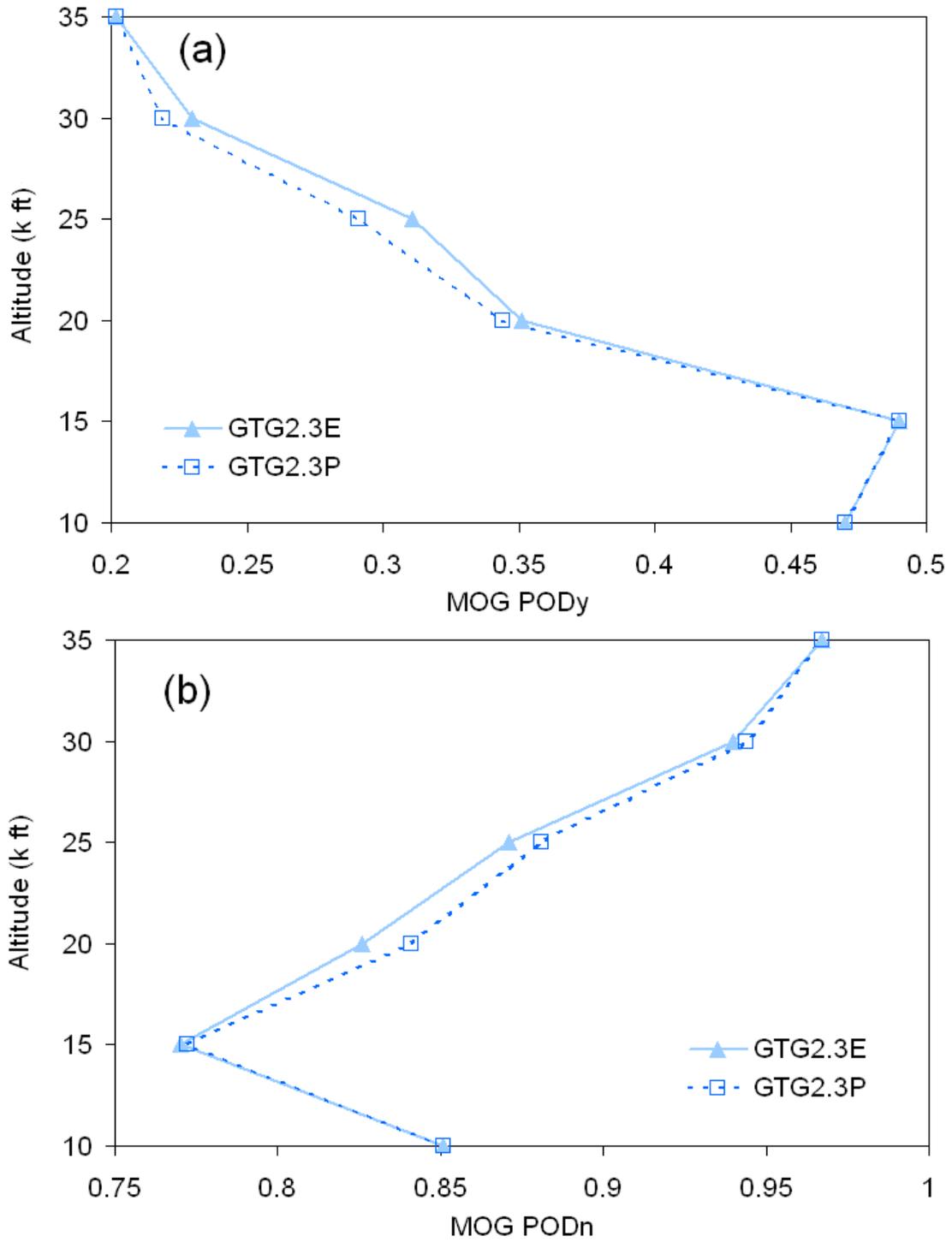


Fig. 6: Height series of (a) MOG PODy and (b) MOG PODn for GTG2.3E and GTG2.3P. Data are plotted at the bottom of each 5,000 ft vertical layer. Forecast threshold is 0.475.

6.2 GTG2.3E performance

The performance of GTG2.3E is presented in this section. Results are summarized by forecast lead time, height, and region.

GTG2.3E performance at upper levels (20,000 to 40,000 ft) is very similar for all lead times (Fig. 7). The 3-h forecasts perform best while the 9-h and 12-h forecasts perform slightly worse than the 6-h forecasts. Areas under the curves for the 3-, 6-, 9-, and 12-h forecasts are 0.789, 0.776, 0.757, and 0.761, respectively, indicating positive skill at all lead times. An identical pattern of performance is seen when one considers how probability of detection varies along with the volume of the airspace where turbulence is forecast (Fig. 8). For lower thresholds, such as the interval between 0.125 through 0.25, the forecasts at all lead times achieve roughly comparable MOG PODy values along with decreased volumes of impacted airspace.

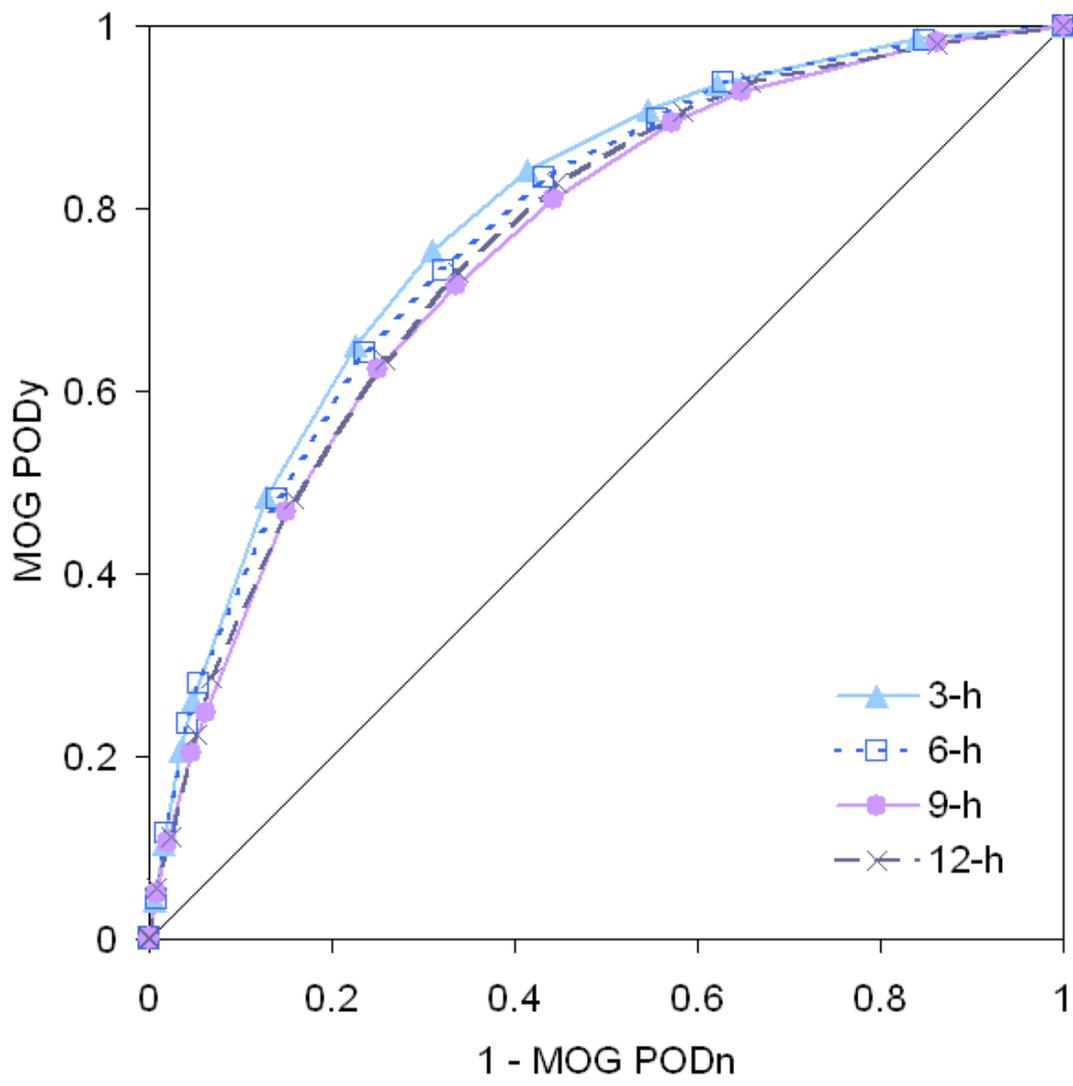


Fig. 7: ROC diagram for GTG2.3E for the 3-, 6-, 9-, and 12-h lead times for upper levels (20,000 to 40,000 ft).

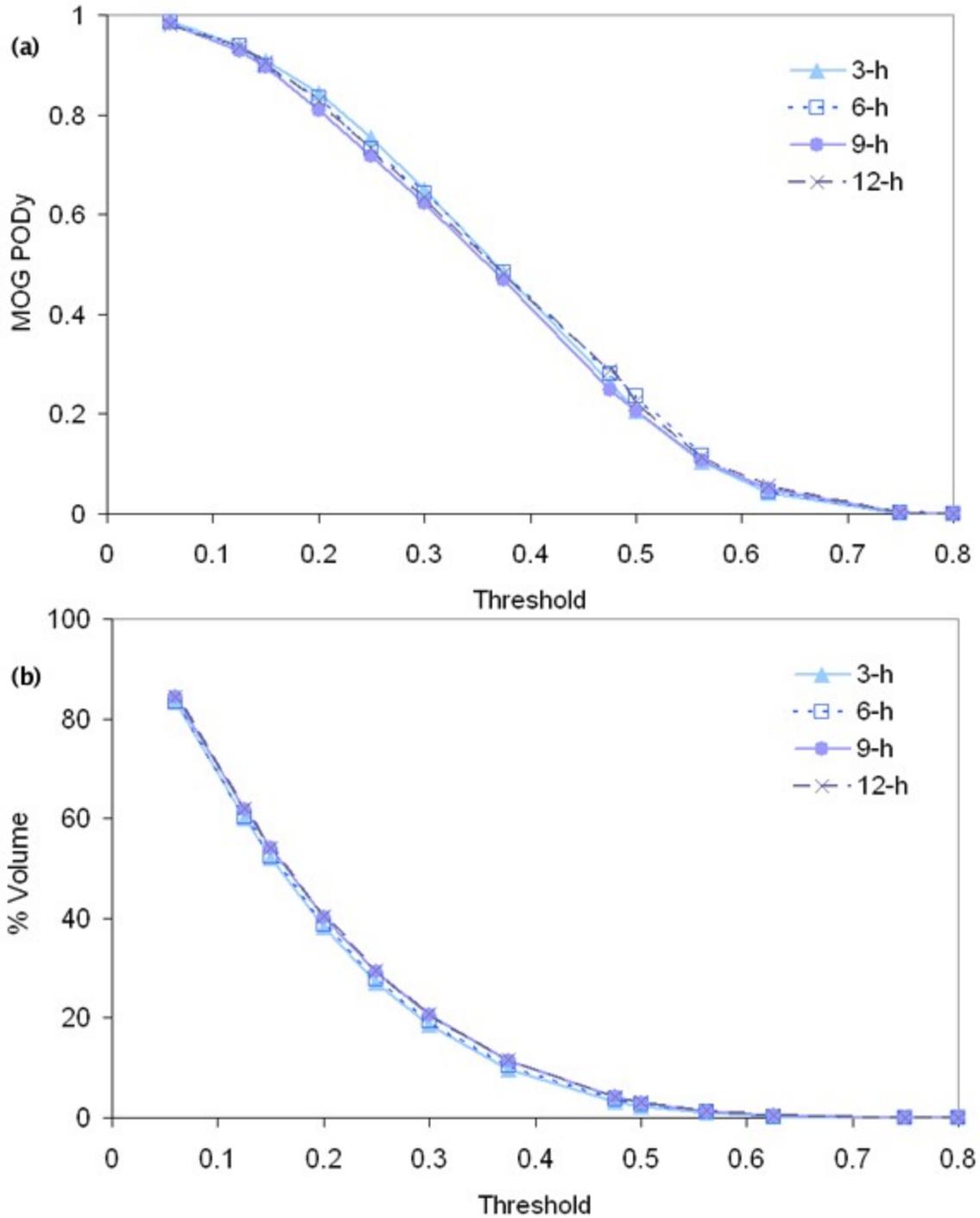


Fig. 8: Threshold plots of (a) MOG PODy and (b) % Volume for GTG2.3E at upper levels for the 3-, 6-, 9- and 12-h lead times.

The ROC diagram for midlevel (10,000 to 20,000 ft) forecasts from GTG2.3E is shown in Fig. 9. These results vary somewhat from those shown for upper levels (Fig. 7). The

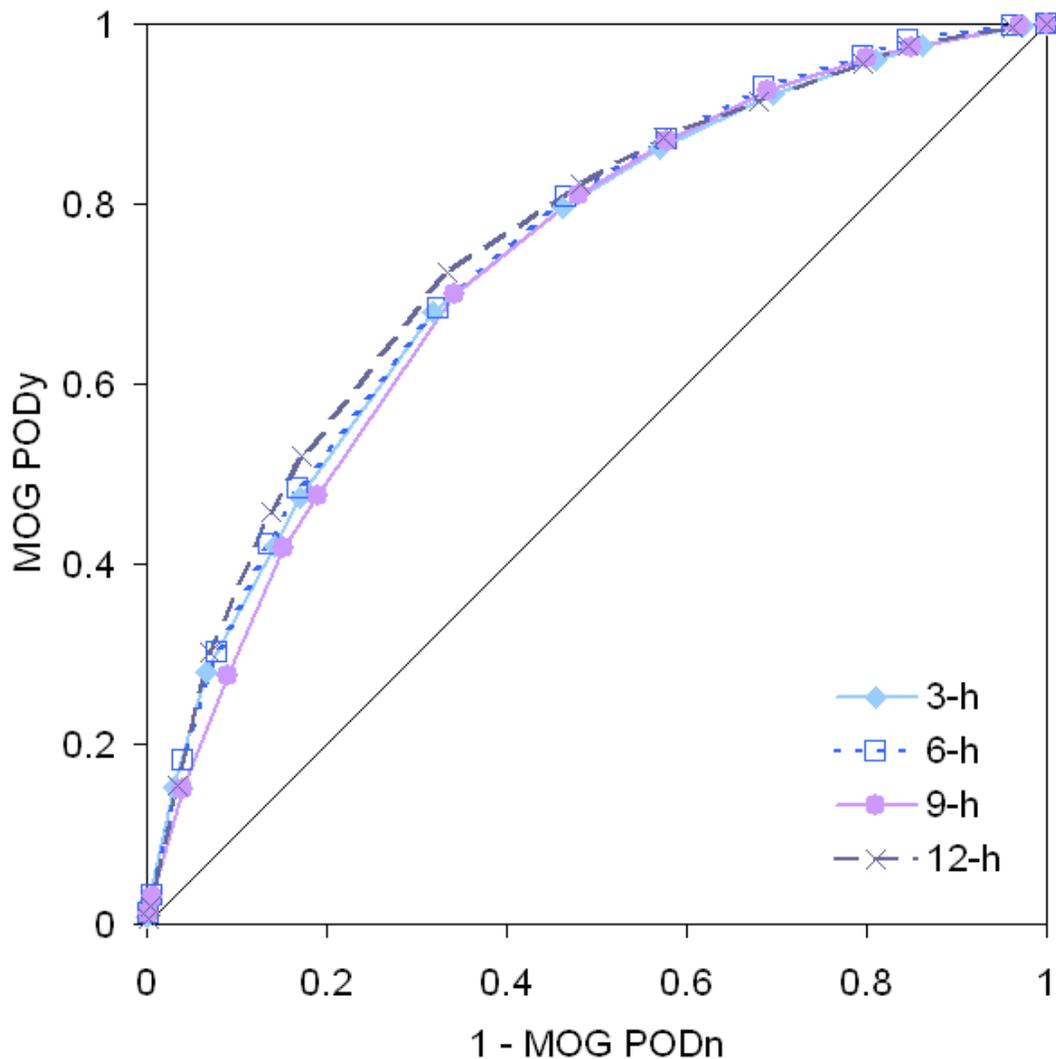


Fig. 9: As in Fig. 7 except for midlevels (10,000 to 20,000 ft).

AUCs are 0.738, 0.744, 0.730, and 0.752 for the 3-, 6-, 9-, and 12-h forecasts, respectively. These values are slightly smaller than those found for the upper level forecasts. This difference may be an indication that the midlevel forecasts from GTG2.3E are not quite as mature as those from the upper levels where the algorithm was initially developed. In addition, fewer PIREPs and EDR observations are typically found at midlevels, which may have an impact on the algorithm performance since these observations are one of the important inputs to GTG2.3E. The 12-h forecasts are the most skillful forecasts at midlevels according to the AUC calculations. MOG PODy and % Volume statistics as a function of threshold (not shown) are similar to the results found for upper levels (Fig. 8).

The performance of GTG2.3E at different vertical levels is somewhat variable (Fig.

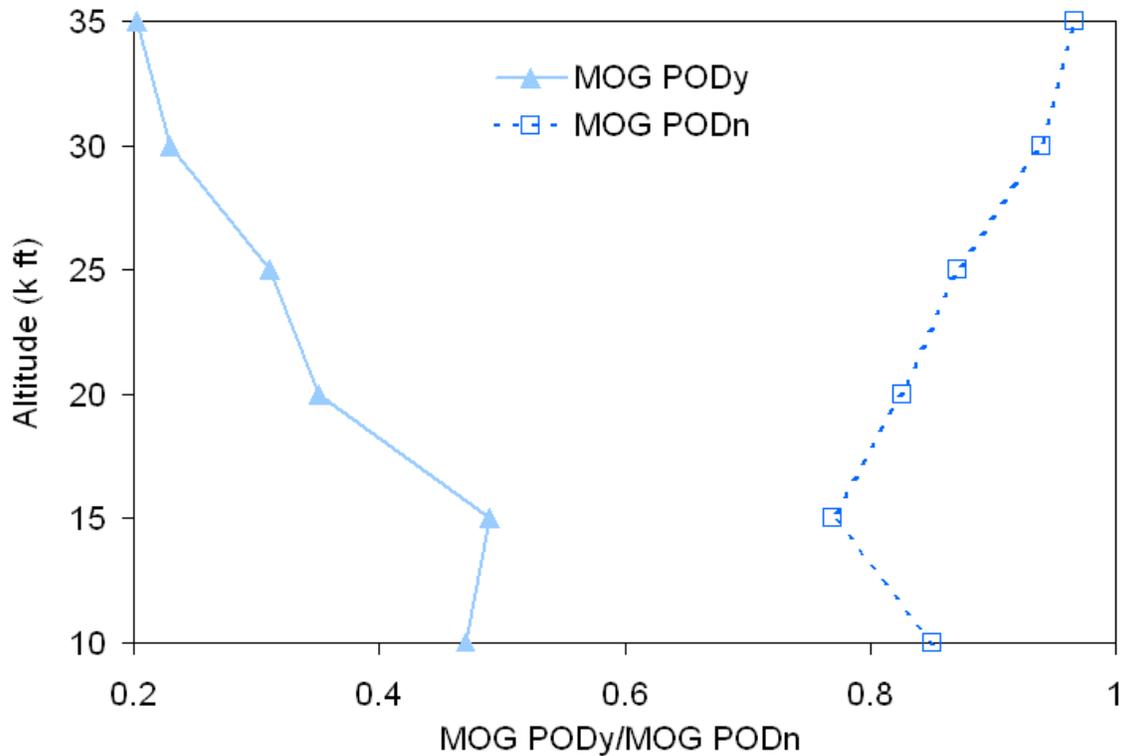


Fig. 10: MOG PODy and MOG PODn height series for GTG2.3E. Forecast threshold is 0.475.

10). MOG PODy values are maximized in the 15,000 to 20,000-ft layer. Approximately 51% of all MOG PIREPs are captured correctly in that layer with slightly more than 20% correct in the 35,000 to 40,000-ft layer (Fig. 10). MOG PODn values change in an opposite sense as PODy with the lowest value occurring for the 15,000-20,000 ft layer and a maximum of 97% of PIREPs with intensities less than moderate having correct forecasts in the 35,000 to 40,000-ft layer. The vertical profile of MOG PODn is also less variable than the MOG PODy profile.

Vertical profiles of MOG PODy and MOG PODn as a function of lead time are shown in Figs. 11 and 12, respectively. Few differences are seen in the MOG PODy values among the 3-, 6-, and the 9-h forecasts from GTG2.3E. However, the 12-h forecasts show a statistically significant departure from the other lead times for the 25,000 to 30,000-ft layer. MOG PODn values cluster tightly with several minor differences between the profiles. The only statistically significant differences are the 3- and 9-h values for the 15,000 to 20,000-ft layer and the 3- and 12-h values for the 20,000 to 25,000-ft layer. No statistically significant differences exist for any other MOG PODn values given the number of PIREPs available during the evaluation period.

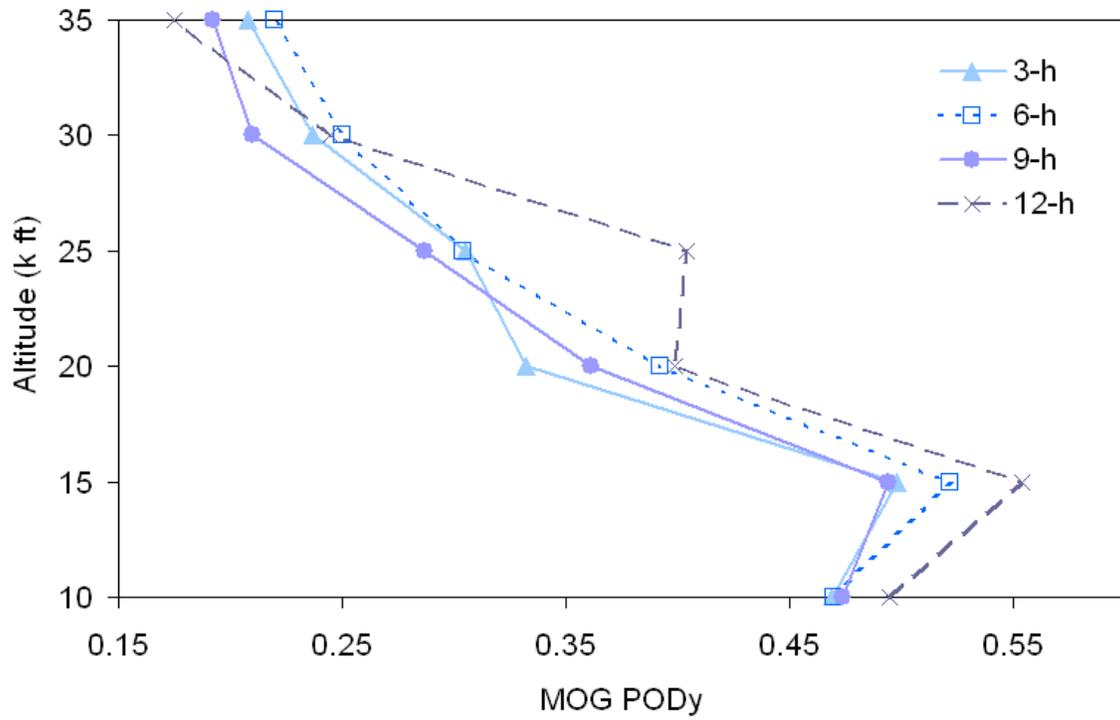


Fig. 11: MOG PODy height series for GTG2.3E for 3-, 6-, 9-, and 12-h lead times. Forecast threshold is 0.475.

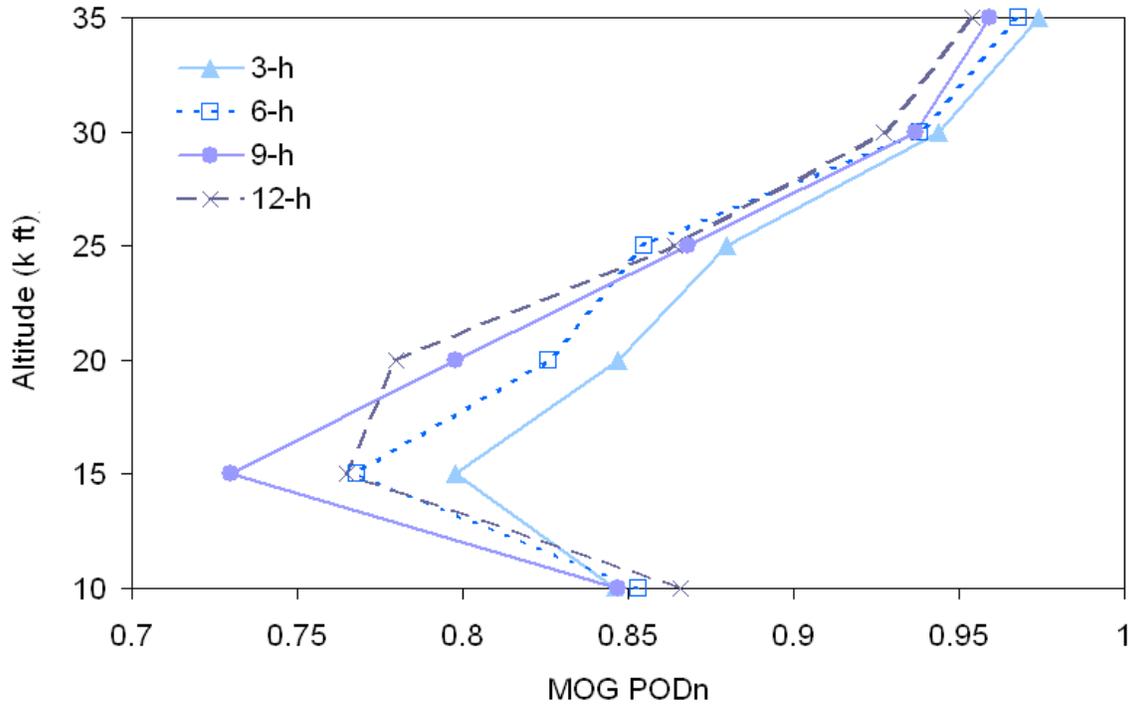


Fig. 12: As in Fig. 11 except for MOG PODn.

Regional analysis is important as GTG2.3 is meant to be used as guidance to support forecasting efforts, such as those performed at the AWC where regional forecasting is required. Additionally, the varied topography of the CONUS is such that terrain-induced turbulence may cause regional differences in the performance of GTG2.3E to become apparent. This knowledge is beneficial to forecasters as well as forecast users.

At upper levels, GTG2.3E performs the worst in the western U.S. while performing best in the central U.S. (Fig. 13). AUC values for the West, Central, and East regions are 0.724, 0.810, and 0.792, respectively. It is possible that mountain wave activity is influencing the results in the West region. Sharman et al. (2006) note that GTG attempts only to forecast clear-air turbulence. Even though some of the diagnostics used to create the GTG product may capture mountain wave conditions, the current algorithm is not expected to forecast them well. Additionally, the PIREP dataset is not filtered to remove any PIREPs that may be due to mountain waves.

For the midlevels, overall performance is degraded somewhat from that found for upper levels in the East and Central regions (Fig. 14). In the West region, performance is slightly better. AUC values are 0.763, 0.729, and 0.699 for the West, Central, and East regions, respectively. The West region has noticeably larger MOG PODy values compared to the East and Central regions. The larger MOG PODy appears to be directly related to the fact that MOG turbulence was forecast over a much larger volume in the West (Fig. 15) than in the other regions.

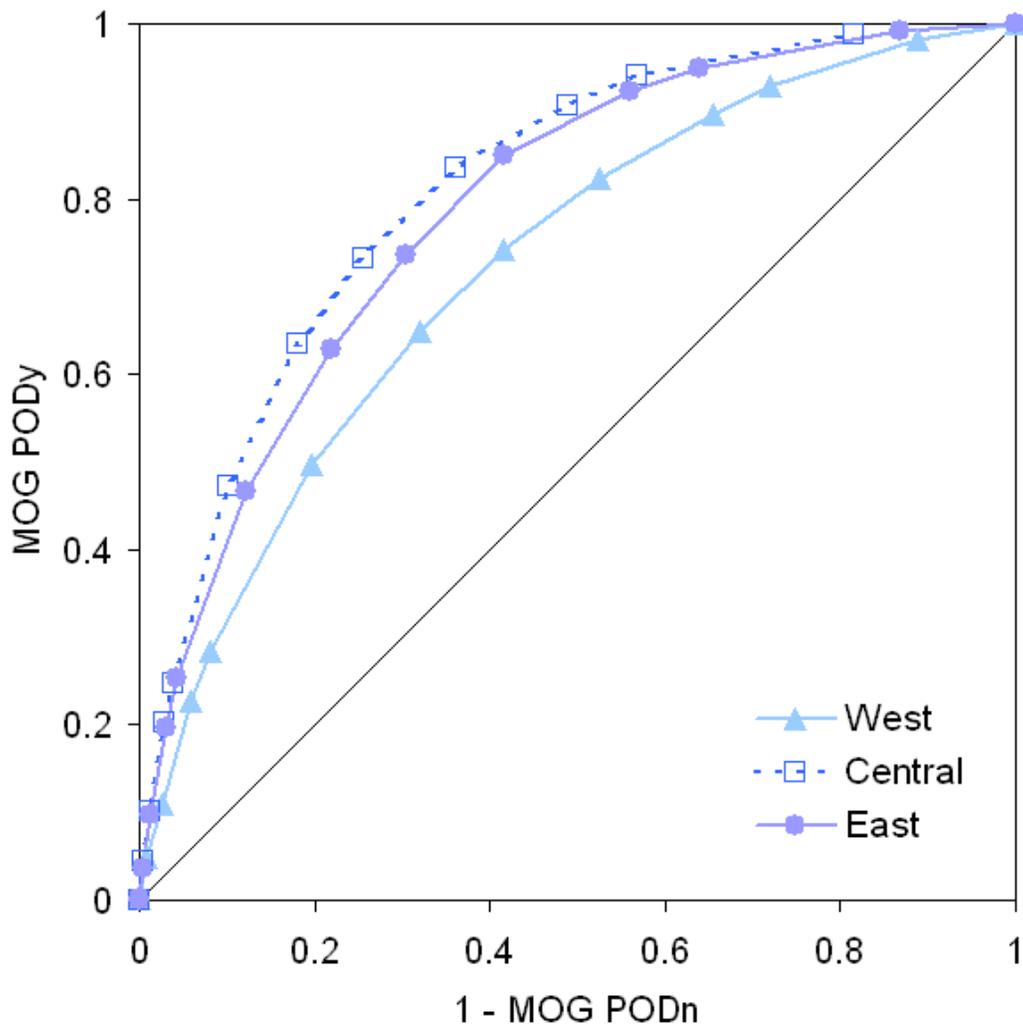


Fig. 13: ROC diagram for GTG2.3E stratified by AWC forecast region for upper levels (20,000 to 40,000 ft).

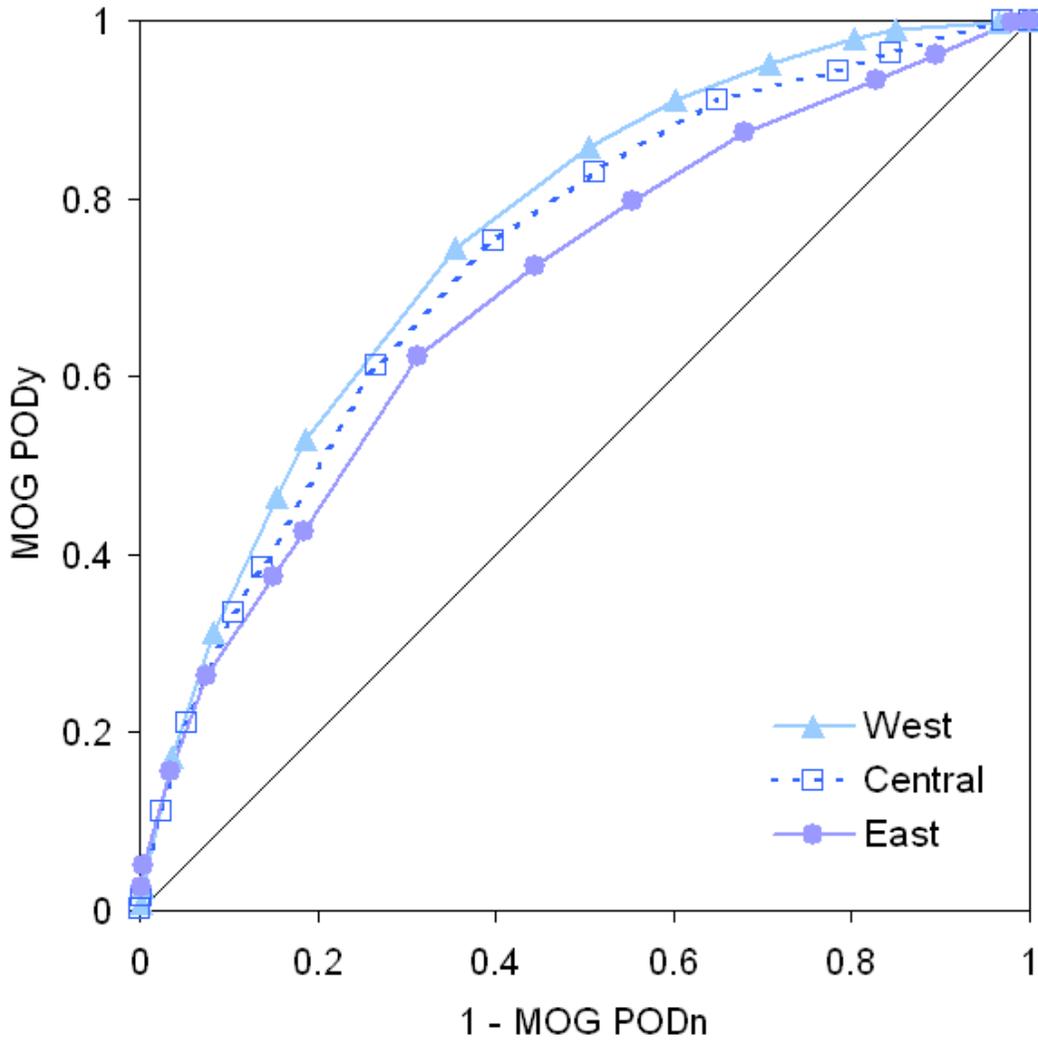


Fig. 14: As in Fig. 13 except for midlevels (10,000 to 20,000 ft).

GTG2.3E performance in each of the climatological regions (Fig. 3; Table 6) are discussed next. For this analysis, all issue times and lead times are combined to increase sample sizes within each region. Results are presented for moderate or greater turbulence only (GTG2.3E threshold of 0.475).

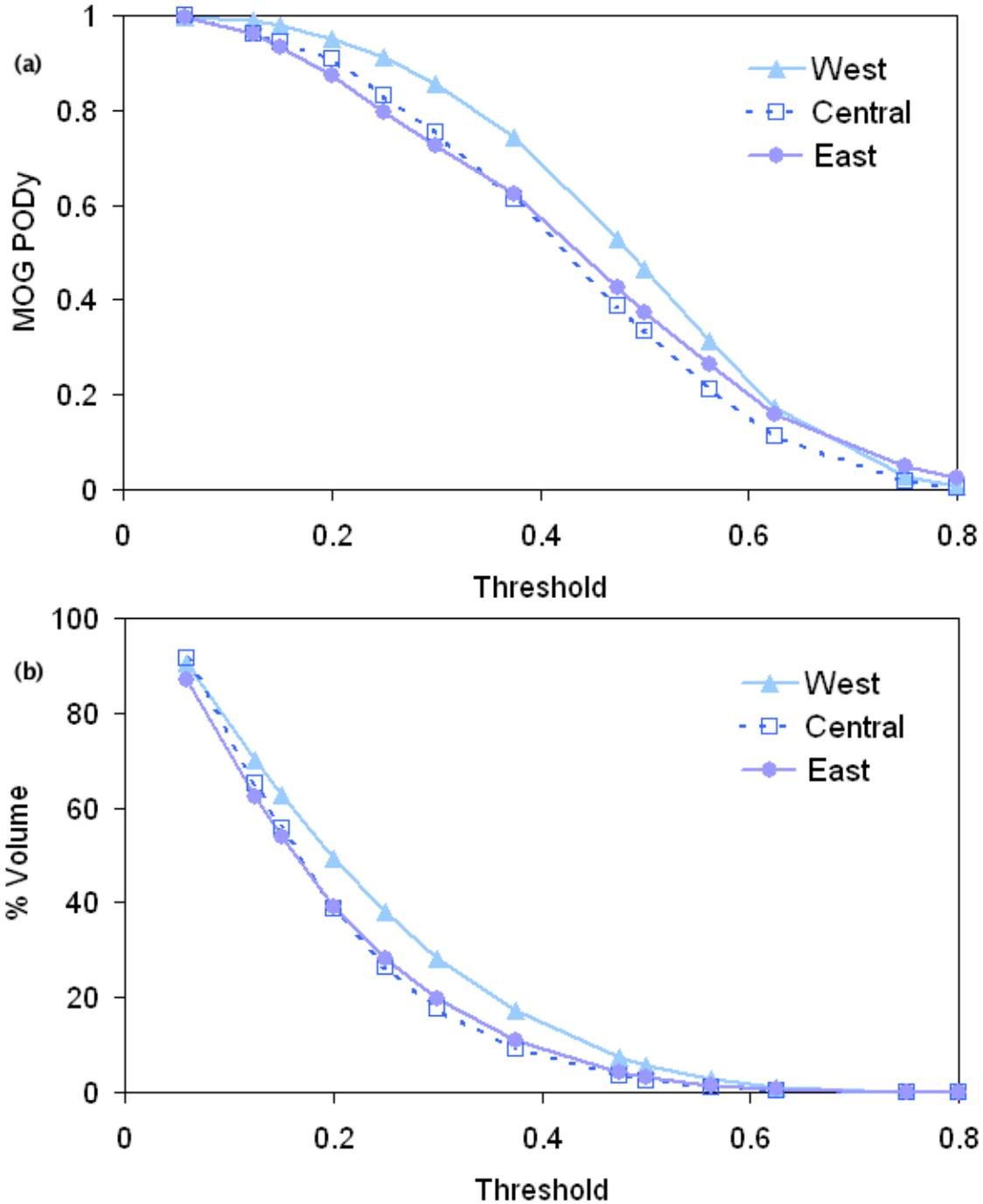


Fig. 15: Threshold plots of (a) MOG PODy and (b) % Volume for GTG2.3E at midlevels for the stratified by AWC forecast region.

For upper levels, GTG2.3E has the largest MOG PODy in the Rocky Mountains North (RMN) region (MOG PODy of 0.405; Fig. 16a.). The RMN region is also associated with the largest percent volume of any region (5.8%; nearly double the largest value of any other region). The MOG TSS values are less than 0.300 for most regions (Fig. 16d.).

The statistics for the WCN, WCS, GCO, and APP regions indicate that GTG2.3E has relatively little skill in these regions compared to other regions. The algorithm has the greatest skill in distinguishing between MOG PIREPs and PIREPs of lesser intensities in the GLA and RMN regions (MOG TSS values of 0.327 and 0.306, respectively).

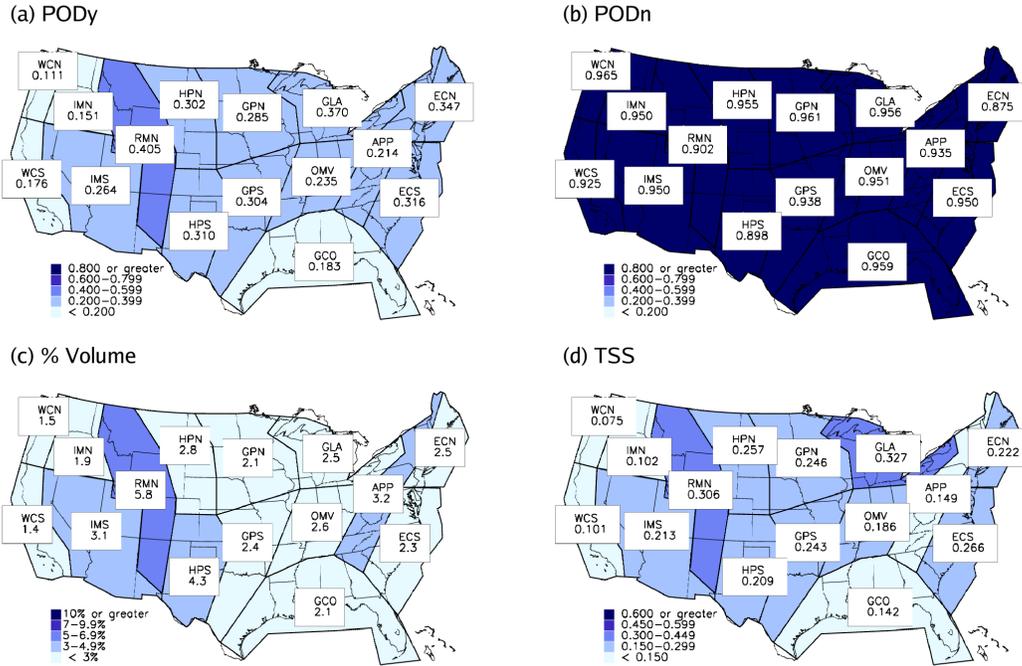


Fig. 16: Overall (a) PODy, (b) PODn, (c) % Volume, and (d) TSS values for GTG2.3E at upper levels for the climatological regions. Forecast threshold is 0.475.

For the midlevels, the performance in the climatological regions is much more variable than at upper levels (Fig. 17). While some regions, such as HPN, have a large PODy value for MOG turbulence (MOG PODy of 0.873), other regions like GPN have an extremely small MOG PODy value (0.130). Furthermore, the MOG PODy value in the GPN region is quite different from the values for adjacent regions. The largest TSS value for GTG2.3E occurs in the HPS region (Fig. 17d); the algorithm also has relatively good skill in the RMN region. The MOG TSS values for most regions are greater than 0.200, but the TSS score for the GPN region is negative, indicating negative skill for this region. The negative skill is due to the poor MOG PODy value here (0.130). GTG2.3E performs well over the mountains and high plains, while having smaller TSS values over the coastal areas (except for the ECS region).

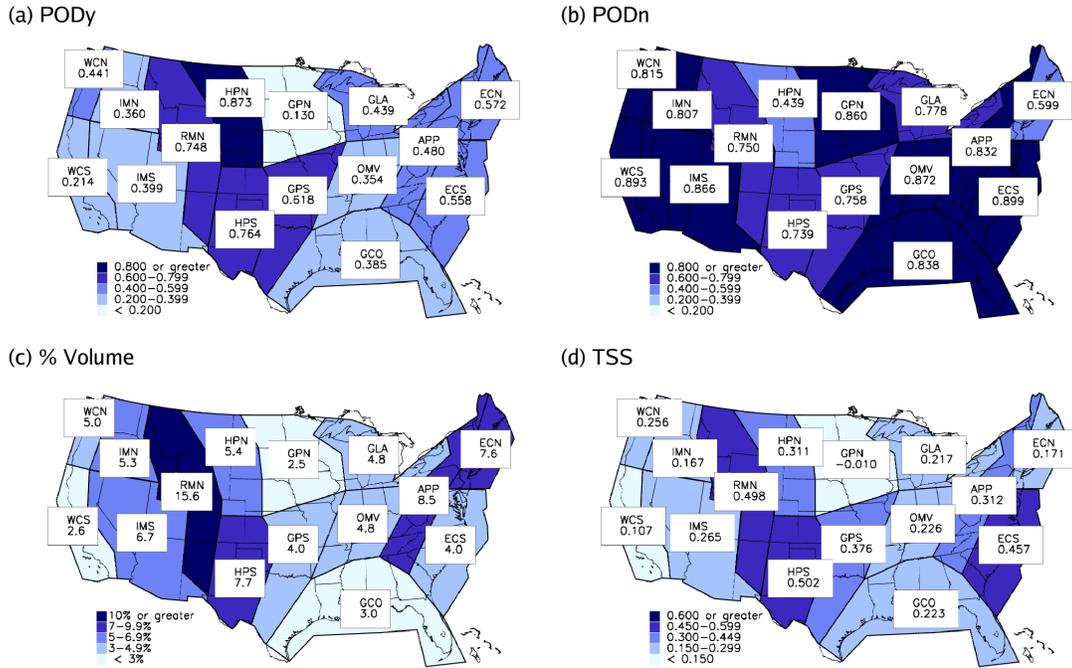


Fig. 17: As in Fig. 16 except for midlevels.

GTG2.3E is displayed to users through ADDS using four categories: None, Light, Moderate and Severe. The forecast and observation values used to define the categories were listed previously in Table 5. Because the forecasts will be provided in this way to users and decision makers, it is important to provide an assessment of its ability to discern the correct relationship with PIREP observations. The analysis will be based upon the joint distribution of forecasts and observations. This distribution provides all nontime-dependent information about the forecasts and observations (Murphy and Winkler, 1987). From the joint distribution one can derive two additional distributions that provide important information about the forecast performance. For this report, only results for the distribution of forecast intensities given the observation intensities, denoted $p(\text{flx})$, will be presented. The conditioning variable provides the frame of reference within which the results are interpreted.

Ideally, the forecasts and observations match identically. Realistically, they do not match perfectly and the conditional distributions allow for greater understanding about the behavior of the forecasts and observations. It is important to reiterate that in the results that follow, the data are only valid wherever there was a PIREP. Since the true distribution of turbulence in the atmosphere is unknown, all forecasts points cannot be considered.

The joint distribution for upper levels is shown in Table 7. The observations are not uniform and are instead dominated by reports of None or Moderate intensity. Forecasts

of None dominate the data with values decreasing as turbulence intensity increases. Of more interest are the conditional probabilities associated with this joint distribution. The results of conditioning on the observations [i.e., $p(f|x)$] are shown in Table 8. For observations of None, the most likely forecast is None (75% of the time) followed by Light (18% of the time). When Light is observed, None and Light account for approximately 70% of the forecasts. These results suggest that when the PIREPs indicate non-threatening turbulence, GTG2.3E typically agrees. However, GTG2.3E also displays a tendency to overforecast the severity of turbulence as indicated by 30% of the Light PIREPs that occurred where the forecast was for Moderate turbulence. Moderate forecasts are the most prevalent forecast category (38% of the time) when Moderate PIREPs are received; forecasts of Light are associated with these PIREPs slightly less frequently (32% of the time) followed by forecasts of None (29% of the time). For observations of Severe turbulence, more than 50% of the time GTG2.3E indicated an intensity of at least Moderate, suggesting a trend towards increasing forecast category with increasing observed turbulence intensity.

Table 7: Joint distribution of GTG2.3E forecasts and PIREPs at upper levels.

		<i>Observed</i>				Total
		None	Light	Moderate	Severe	
<i>Forecast</i>	None	29734	1697	3288	56	34775
	Light	7081	1677	3604	76	12438
	Moderate	3006	1493	4268	173	8940
	Severe	33	14	111	4	162
	total	39854	4881	11271	309	56315

Table 8: Conditional probability of a forecast for each observation category, $p(f|x)$, for upper levels.

		<i>Observed</i>			
		None	Light	Moderate	Severe
<i>Forecast</i>	None	0.746	0.348	0.292	0.181
	Light	0.178	0.346	0.320	0.246
	Moderate	0.075	0.306	0.379	0.560
	Severe	0.001	0.003	0.010	0.013
	$p(x)$	0.708	0.087	0.200	0.005

For midlevels, reports of None and Moderate are still the dominant categories, accounting for over 80% of the PIREPs during the evaluation period (Table 9). For observations of None, 87% of the time forecasts of None or Light occurred. For observations of Moderate, 45% of the time forecasts of Moderate occurred. More Severe PIREPs were reported at midlevels than at upper levels with a total of 401 being available for verification. The total number of PIREPs for available in midlevels was 18,279 compared to 56,315 at upper levels. The largest part of the decrease is due to the smaller number of None reports, decreasing from 39,854 reports at upper levels to 7,457 at midlevels.

The conditional probability of a forecast given an observation (p(flx)) for midlevels, shown in Table 10, illustrates several interesting features. Except for None, GTG2.3E is often associated with forecast intensities other than the expected values. For Moderate PIREPs, the probabilities of None, Light, Moderate, and Severe forecasts are 0.203, 0.335, 0.452, and 0.01, respectively. While Moderate is the most likely forecast when Moderate PIREPs are observed, nearly 50% of all forecasts fall into either the None or Light categories. Of the 401 Severe observations, 62% were associated with forecasts of intensities that were at least Moderate while 25% were associated with forecasts of Light turbulence.

Table 9: Joint distribution of GTG2.3E forecasts and PIREPS at midlevels.

		<i>Observed</i>				Total
		None	Light	Moderate	Severe	
<i>Forecast</i>	None	4541	810	1572	51	6974
	Light	1982	1009	2595	100	5686
	Moderate	922	852	3501	238	5513
	Severe	12	4	78	12	106
	total	7457	2675	7746	401	18279

Table 10: Conditional probability of a forecast for each observation category, $p(f|x)$, for midlevels.

		<i>Observed</i>			
		None	Light	Moderate	Severe
<i>Forecast</i>	None	0.609	0.303	0.203	0.127
	Light	0.266	0.377	0.335	0.249
	Moderate	0.124	0.319	0.452	0.594
	Severe	0.002	0.001	0.010	0.030
	p(x)	0.408	0.146	0.424	0.022

6.3 GTG2.3E vs. GTG Comparison

In this section, GTG2.3E is compared with the operational version of GTG that is produced at the AWC. Because operational GTG only provides upper-level forecasts, the comparisons are limited to upper levels (20,000 to 40,000 ft). Data for all issue- and lead times from Table 1 are evaluated in this section.

An evaluation of overall performance, as depicted by the ROC curves for the two forecasting systems, indicates that GTG2.3E has greater skill than GTG (Fig. 18). The AUC values for GTG2.3E and GTG are 0.775 and 0.677, respectively. The larger AUC value for GTG2.3E shows that it is better able to discriminate between events (i.e., moderate or greater turbulence PIREPs) and nonevents (PIREPs with intensities less than moderate) than GTG. For most thresholds, GTG2.3E provides larger MOG POD_y values and smaller MOG POD_n values than GTG. The larger MOG POD_y values might be related to increased volumes produced by GTG2.3E as compared to the operational GTG (Fig. 19).

Since specific thresholds of GTG2.3E will be used on ADDS displays to create the categorical forecasts, changes in forecast performance between GTG2.3E and GTG are noted here. For instance, the direct comparison of GTG2.3E and GTG for the MOG threshold (Fig. 19) indicates a considerable decrease in the % Volume at the expense of a slight decrease in forecast performance for GTG2.3E over the current GTG product. Similar trends are also identified for the Light and Severe thresholds. Therefore, as highlighted here, the skill of GTG2.3E is highly dependent upon the selected threshold(s), which are chosen to optimize the forecast performance while decreasing the % Volume. Kay et al. (2006) address specific changes to the GTG thresholds and resulting changes in forecast performance.

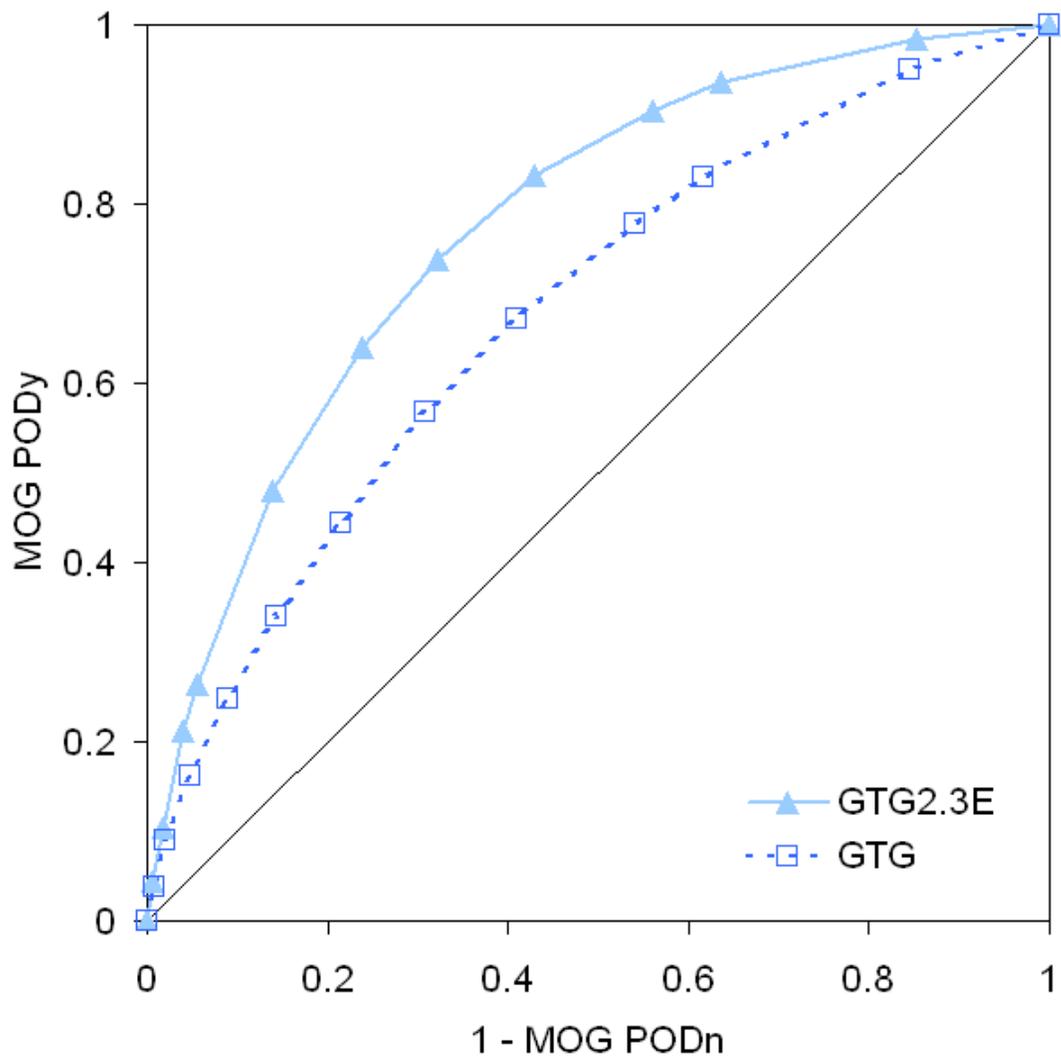


Fig. 18: ROC diagram for GTG2.3E and operational GTG at upper levels.

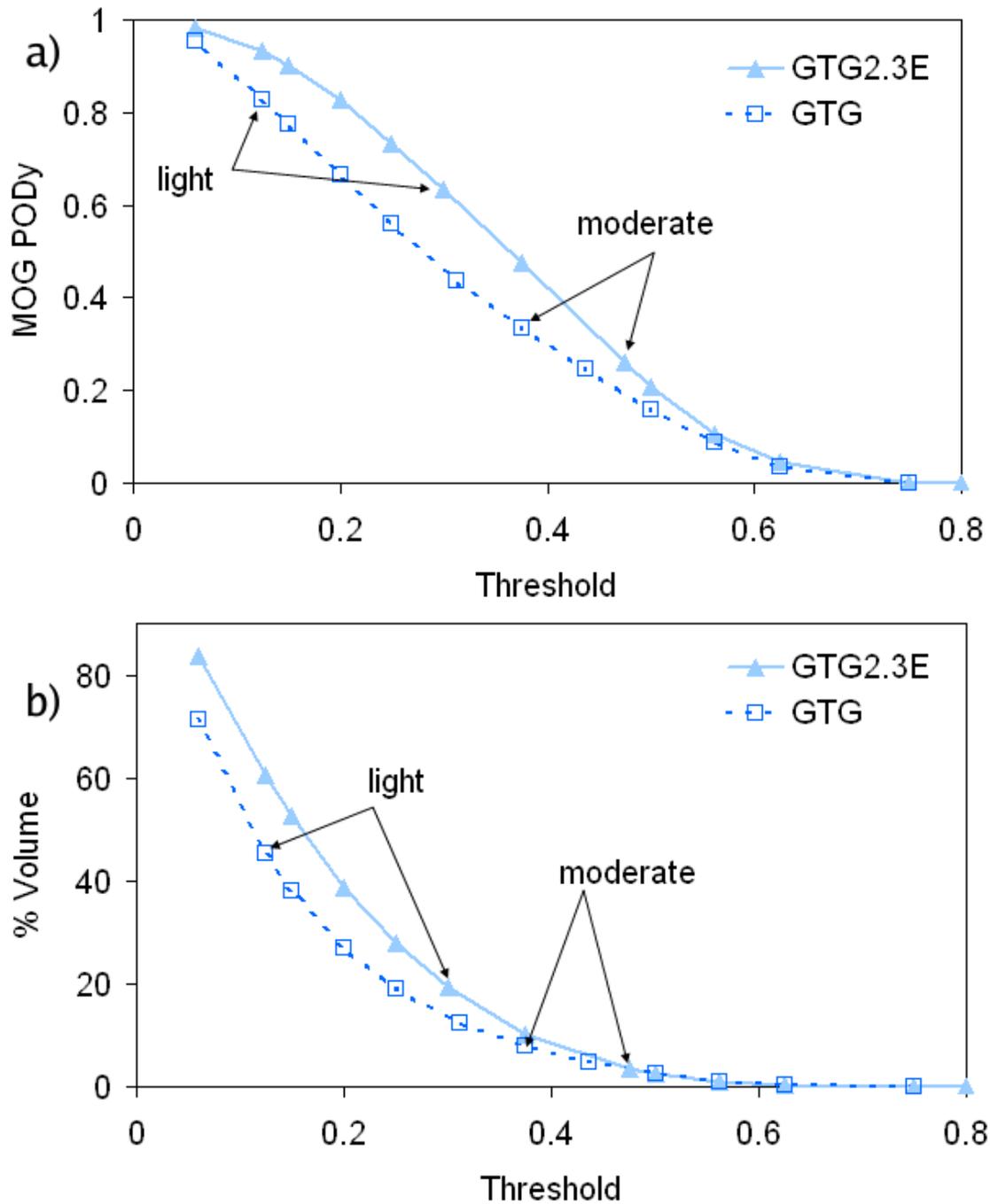


Fig. 19: Threshold plots for (a) MOG PODy and (b) % Volume for GTG2.3E and operational GTG at upper levels.

In the vertical, GTG2.3E and GTG are nearly identical in detection of MOG PIREPs from 20,000 to 30,000 ft (Fig. 20). MOG PODy values in this layer average around 0.30 indicating that approximately 30% of all MOG PIREPs were correctly forecast by both

algorithms. For the 30,000 to 40,000 ft layer the performance of the two algorithms diverges with GTG having larger MOG PODy values than GTG2.3E. The MOG PODy value for GTG2.3E decreases to 0.21 for the 35,000 to 40,000-ft layer, while the MOG PODy value for GTG is 0.42 for this layer. The MOG PODy differences for both of the 5,000 ft layers between 30,000 to 40,000 ft are statistically significant at the 0.05 level.

Both algorithms perform nearly identically with regards to correct detection of less-than-MOG PIREPs throughout the 20,000 to 30,000-ft layer (Fig. 21). The MOG PODn values are at a minimum in the 20,000 to 25,000-ft layer and generally increase with height. Maximum MOG PODn values are 0.90 for GTG in the 30,000 to 35,000-ft layer and 0.97 for GTG2.3E in the 35,000 to 40,000 ft layer. Similar to the MOG PODy values (Fig. 20), the differences from 30,000 to 40,000 ft are statistically significantly different from one another at the 95% level. GTG2.3E clearly performs better with respect to MOG PODn than GTG at the upper levels (Fig. 21), but at the expense of smaller MOG PODy values (Fig. 20).

GTG2.3E outperforms GTG within each of the AWC forecast regions. ROC diagrams for the West, Central, and East regions can be found in Fig. 22 a-c. The East and Central regions show the largest increase in performance for GTG2.3E. The West region, where terrain-induced turbulence may influence the results more strongly, is associated with the smallest increase in performance as compared to GTG. The West also has the lowest AUC value of the three AWC forecast regions. AUC values for GTG2.3E increase more than 0.1 compared to GTG for both the Central and Eastern regions with a much more modest increase in the Western region (Table 11).

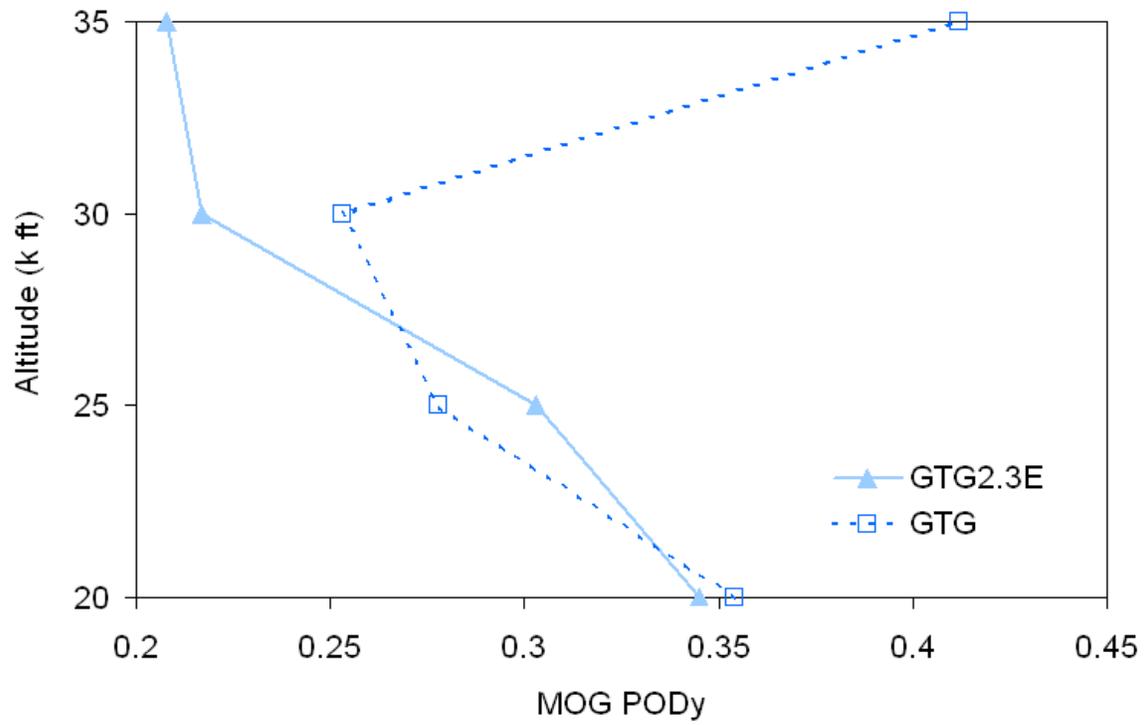


Fig. 20: Height series of MOG PODy for GTG2.3E and operational GTG at upper levels. GTG2.3E threshold is 0.475 and GTG threshold is 0.375.

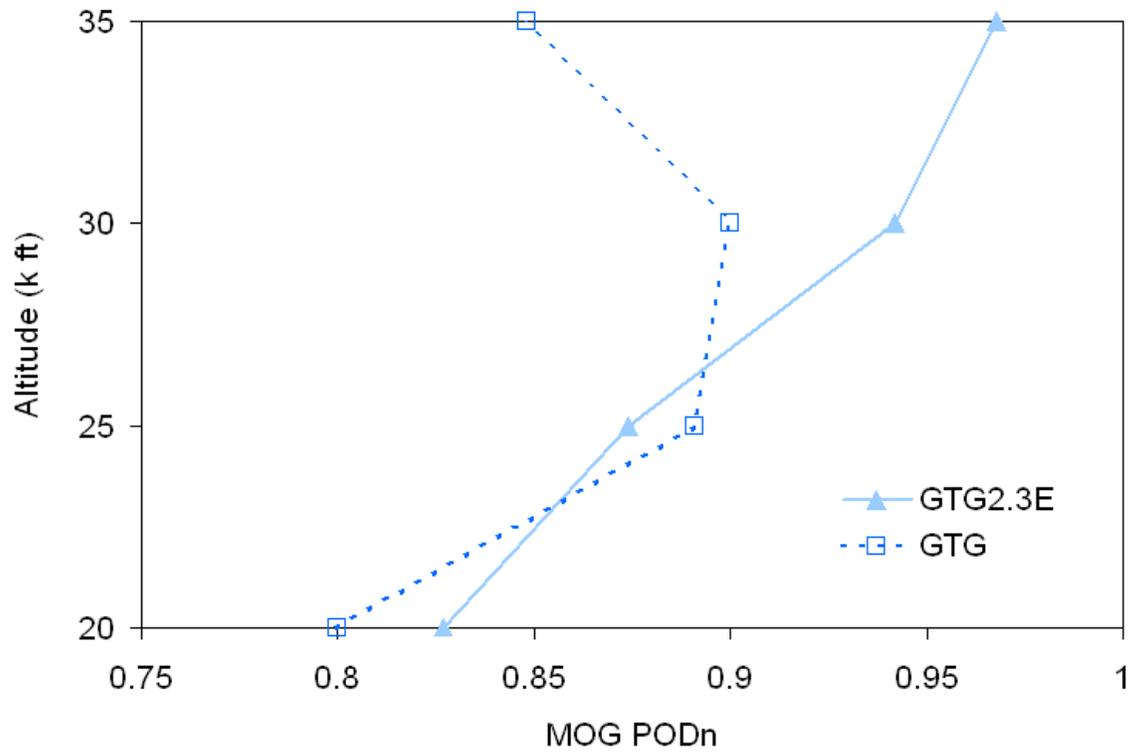


Fig. 21: As in Fig. 20 except for MOG PODn.

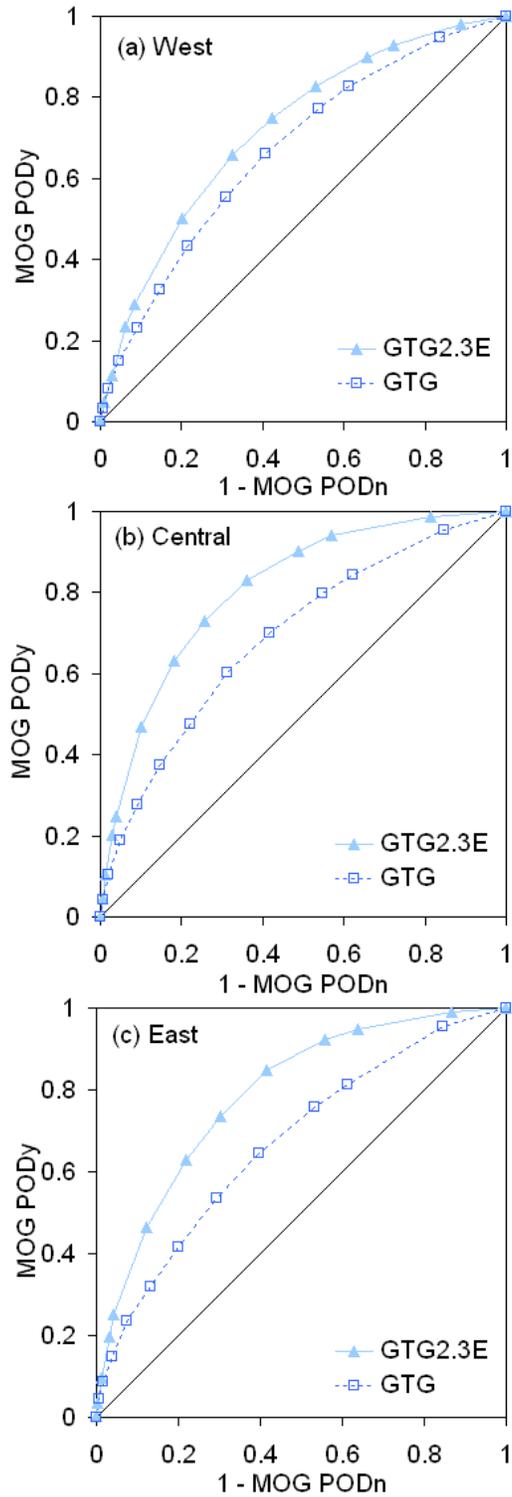


Fig. 22: ROC diagrams for GTG2.3E and operational GTG for the (a) West, (b) Central, and (c) East AWC forecast regions at upper levels.

Table 11: ROC area under the curve values for GTG2.3E and GTG for each AWC region.

Region	GTG2.3E	GTG	GTG2.3E - GTG
West	0.720	0.673	0.047
Central	0.813	0.690	0.123
East	0.791	0.669	0.122

For the comparison of GTG and GTG2.3E performance across the 15 climatological regions of the CONUS all issue and lead times were combined together in order to obtain adequate sample sizes within each region. The operational threshold of 0.375 was used to define areas of MOG turbulence for GTG while 0.475 was again used for GTG2.3E. Compared to GTG2.3E (Fig. 16), GTG generally has a larger MOG POD_y value and forecast volume over every region (Fig. 23a). However, it also has a smaller MOG POD_n for every region. This indicates that GTG2.3E does a better job discriminating between MOG PIREPs and less-than-MOG PIREPs. The TSS scores for GTG2.3E (Fig. 23d) are larger than those for GTG for nearly every region. However, the TSS scores for the WCN, IMN, and HPN regions were larger for GTG than for GTG2.3E. In general, GTG2.3E has improved TSS values for the eastern U.S. mountainous regions, Gulf Coast, Great Lakes, and the southwestern U.S., while GTG has larger TSS values over the midwestern U.S. Therefore, with the exception of a few regions, GTG2.3E generally shows improved skill over GTG in distinguishing between MOG and lesser intensity turbulence.

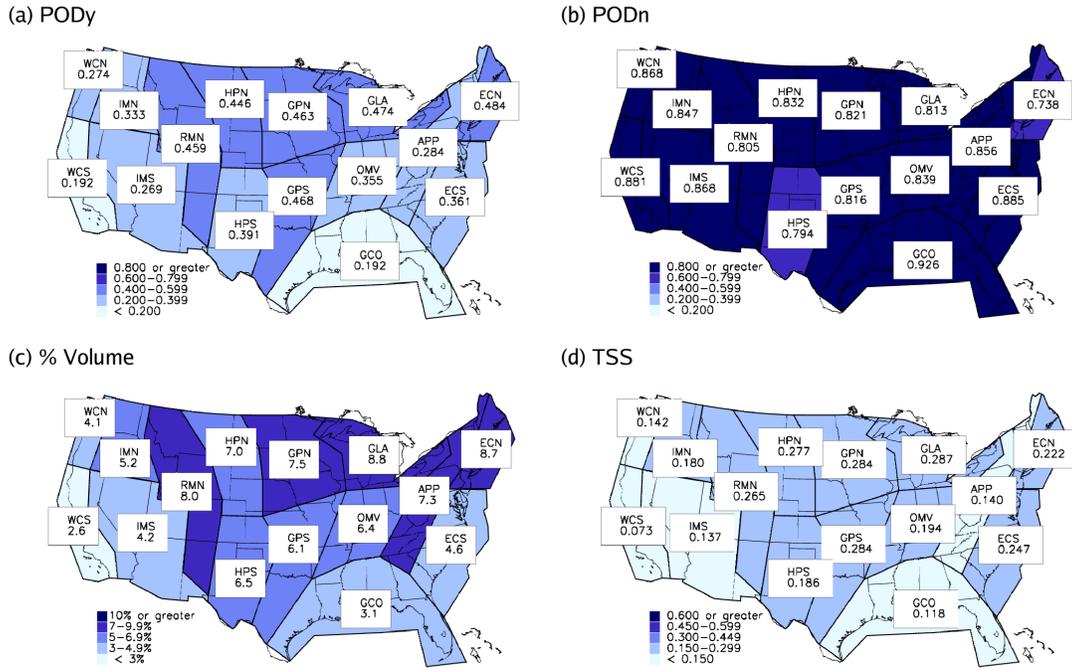


Fig. 23: Overall (a) POD_y , (b) POD_n , (c) % Volume, and (d) TSS values for GTG at upper levels for the climatological regions. Forecast threshold is 0.375.

The next focus of the evaluation is an examination of the impact of the change from GTG to GTG2.3E on the categorical forecasts displayed on ADDS. Recall that in Section 5.2, it was shown that GTG2.3E is modestly able to classify observations correctly, particularly for the None-to-Light and MOG categories. The following discussion compares the findings for GTG2.3E with similar results for GTG. Complete information on GTG categorical forecasting performance will not be displayed here but is available in the Appendix. For the raw GTG2.3E contingency tables, refer to Section 5.2. Table 12 provides information on the difference in conditional probabilities between GTG2.3E and GTG for $p(\text{flx})$. Positive (negative) values indicate that the GTG2.3E value is larger (smaller) than the GTG value.

Table 12: Difference table ($GTG2.3E - GTG$) for the conditional probability of a forecast given an observation, $p(\text{flx})$, between GTG2.3E and GTG.

		<i>Observed</i>			
		None	Light	Moderate	Severe
<i>Forecast</i>	None	+0.341	+0.147	+0.121	+0.052
	Light	-0.294	-0.170	-0.174	-0.259
	Moderate	-0.044	+0.047	+0.080	+0.233
	Severe	-0.004	-0.021	-0.026	-0.026

When PIREPs of type None are found, marked improvement is observed for GTG2.3E over GTG. For GTG2.3E, $p(\text{flx}=\text{None})$ is 0.341, greater than the value for GTG while at the same time forecasts of Light have decreased by 0.294 for GTG2.3E as compared to GTG. For observations of light turbulence, less correct forecasts were made by GTG2.3E than for GTG. More forecasts of None were made by GTG2.3E than with GTG when Light turbulence was observed. For observations of Moderate, the largest changes in GTG2.3E were an increase in the likelihood of None being forecast and a large decrease (of 0.174) in the likelihood of Light being forecast. For severe observations, a large decrease is noted in the likelihood of Light forecasts along with an increase of the likelihood of Moderate forecasts (an increase of 0.233). GTG2.3E appears to do a better job than GTG of discriminating the None and Severe observation categories and performs somewhat less well for the Light and Moderate observation categories.

6.4 GTG2.3E vs. AIRMETs Comparison

In order to perform the intercomparison of GTG2.3 with AIRMETs, the temporal attribute of AIRMETs was modified. AIRMETs are issued for 6-h periods and are intended to capture turbulence events within the spatial bounds of the forecast during the 6-h period rather than at a specific time. GTG2.3, however, is valid at specific times, such as 2100 UTC. To address this difference, AIRMETs are converted into a series of six 1-h forecasts for each issuance time and are valid at specific times, much like the GTG2.3 algorithm. For each of the six intermediate forecasts, the original spatial bounds of the forecast areas (i.e., polygons) from the AIRMET issuance are used. This conversion of the four AIRMET issuances each day provides a new set of AIRMET forecasts that are valid for each hour of the day rather than every six hours. Because AIRMETs are only issued four times per day, a limited number of forecasts are available for intercomparison with GTG2.3, which include the 1500 UTC issuance 6-h lead time and the 2100 UTC issue 3-h lead time. Owing to the limited number of available times for the intercomparison, results will be presented for the 3- and 6-h lead times combined.

It is important to restate the limitations of this comparison. AIRMETs are not intended to provide forecasts for a given valid time, instead they are intended to provide forecasts for valid periods. Despite this, users must make decisions at certain instances in time and AIRMETs are often treated as though they represent snapshots in time. GTG2.3E is intended to provide forecasts at valid times. The mechanics of the intercomparison have been designed to provide as fair a comparison as possible to AIRMETs.

GTG2.3E outperforms AIRMETs at upper levels (Fig. 24). The ROC AUC value for GTG2.3E is 0.783 for the National region. AIRMETs have the overall statistics of MOG PODy equal to 0.66 and MOG PODn equal to 0.57. Aside from the AUC values, the AIRMETs can be compared to GTG2.3E dichotomously through the choice of a threshold value that differentiates forecasts of turbulence from forecasts of no turbulence. One obvious threshold for GTG2.3E is 0.475, which is the threshold that will be used on ADDS to depict MOG turbulence. A second threshold of interest is the GTG2.3E value

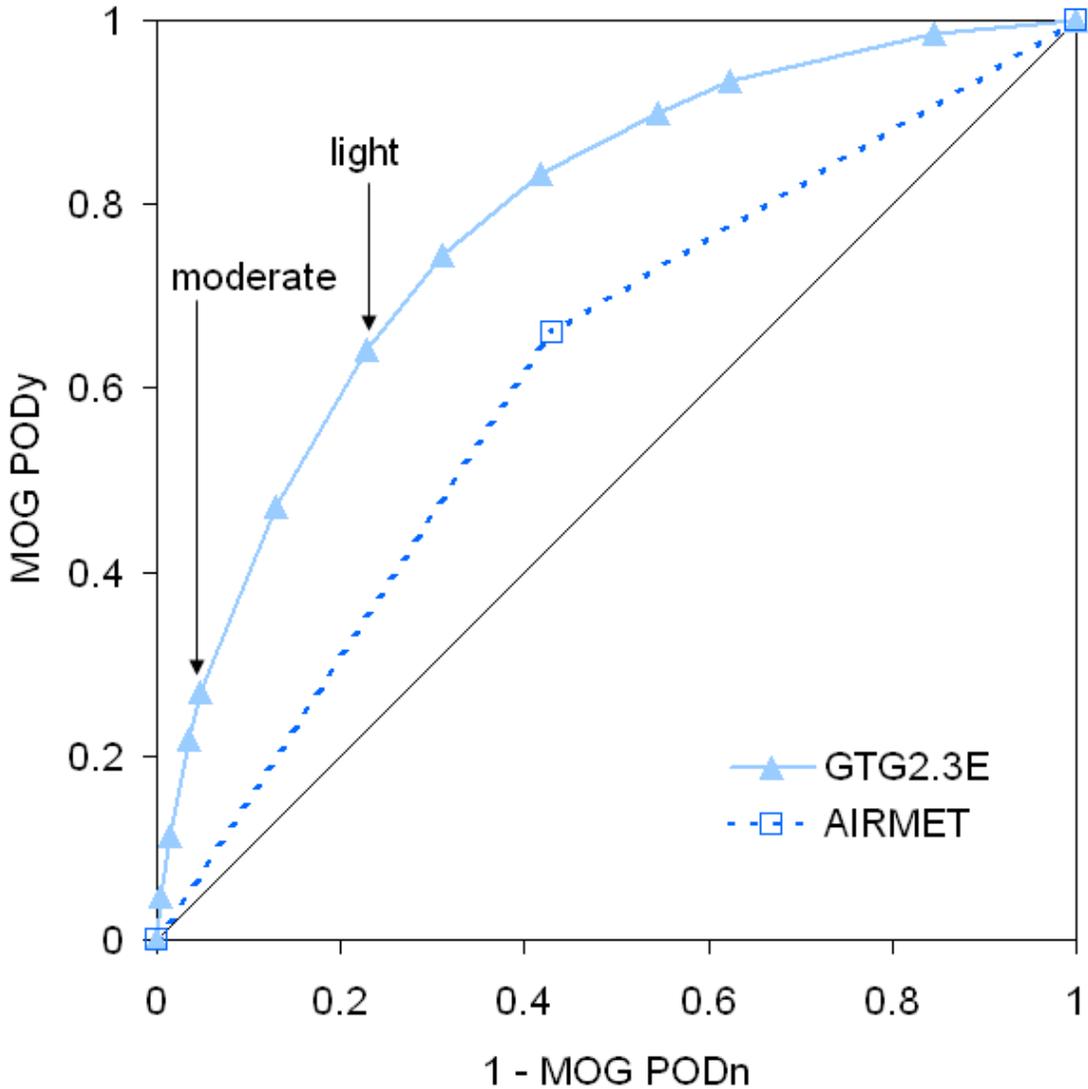


Fig. 24: ROC diagram for GTG2.3E and AIRMETs at upper levels. GTG2.3 points corresponding to the light and moderate thresholds are highlighted.

that is nearest to the location of the AIRMET point on the ROC diagram. Using this approach, the appropriate threshold is 0.250, which gives a MOG PODy value that is larger than the AIRMET value, while also increasing the MOG PODn value. A third comparison uses the threshold that allows a similar MOG PODy to AIRMETs. The GTG2.3E threshold of 0.300 provides a comparable MOG PODy value with a significantly increased MOG PODn value compared to AIRMETs. For each of these thresholds, it is instructive to consider the volume of airspace where moderate or greater turbulence is forecast by GTG2.3E relative to that of the AIRMETs (Fig. 25). The 0.475 threshold for GTG2.3E results in the smallest forecast volumes. Correspondingly, the MOG PODy value is the smallest of the group at 0.27 at this threshold. The 0.250 threshold has an associated median forecast volume that is nearly identical to the median

AIRMET forecast volume with a MOG POD_y value of 0.75 and MOG POD_n of 0.31. The threshold of 0.3 for GTG2.3E, which results in a nearly identical MOG POD_y value with AIRMETs, has a median forecast volume that is close to 18% of the possible volume compared to the median forecast volume of 27% for AIRMETs. These results are summarized in Table 13.

Table 13: MOG POD_y, MOG POD_n, and median % Volume values for GTG2.3E and AIRMETs at upper levels.

	Threshold	MOG POD _y	MOG POD _n	% Volume _{median}
AIRMETs	-	0.66	0.57	27.8
GTG2.3E	0.475	0.27	0.95	3.00
	0.300	0.64	0.77	18.2
	0.250	0.75	0.31	26.5

These results suggest several key differences between GTG2.3E and AIRMETs. First and foremost, for the 0.475 threshold from GTG2.3E, which is used to indicate regions of moderate or greater turbulence, the MOG POD_y for GTG2.3E was significantly smaller than the value for AIRMETs and the MOG POD_n was appreciably larger compared to the AIRMET value. The larger MOG POD_y value for AIRMETs is associated with much larger forecast volumes than those from GTG2.3E. The lack of specificity of AIRMETs is further reflected in its smaller MOG POD_n value. Additionally, GTG2.3E is able to achieve nearly identical MOG POD_y values as AIRMETs with median forecast volumes that are approximately 33% smaller than the AIRMET volumes, while at the same time improving upon correct differentiation between MOG and less-than-MOG PIREPs as compared to AIRMETs.

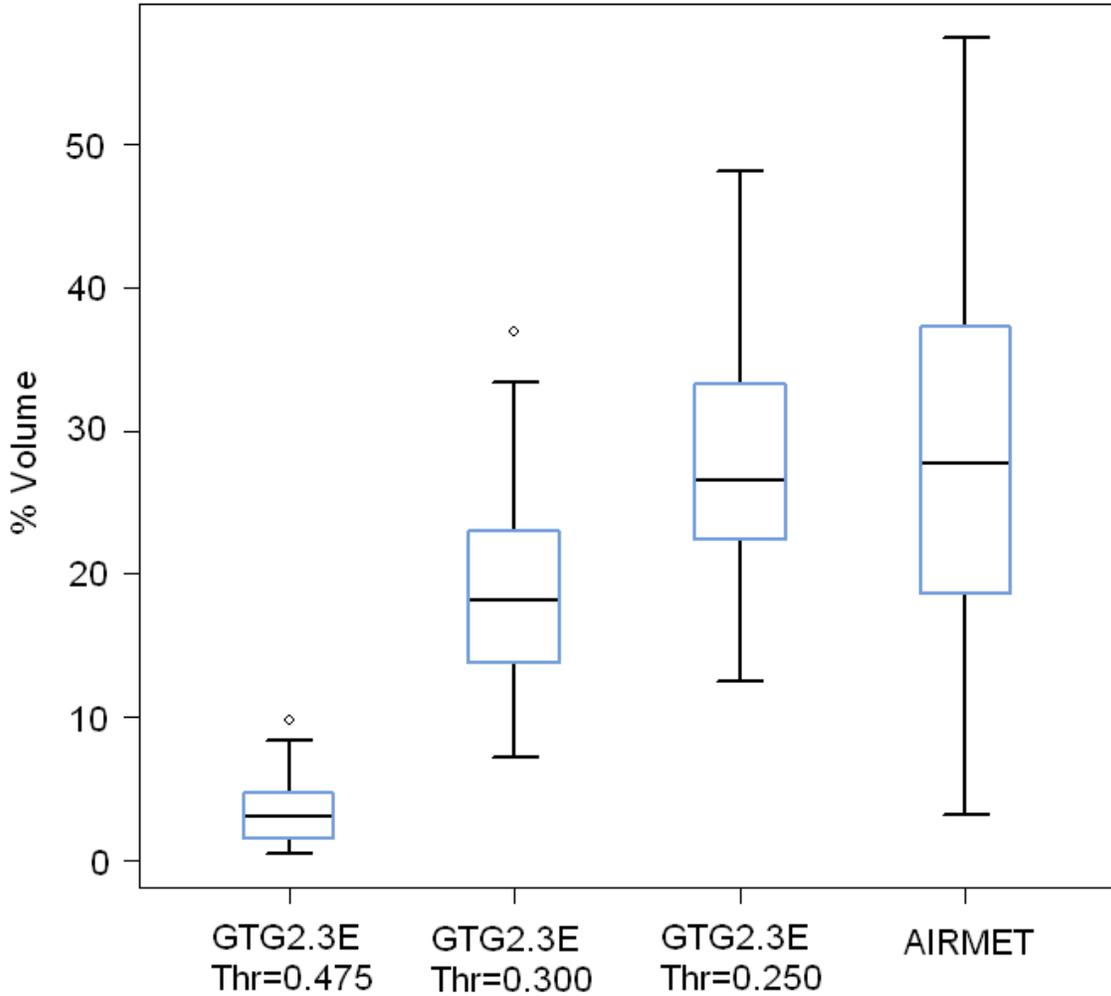


Fig. 25: Boxplot of % Volume values for GTG2.3E with thresholds 0.475, 0.300, 0.250, and AIRMETs at upper levels.

In the midlevels, the ROC diagram is very similar to the upper-level ROC presented previously (Fig. 26). GTG2.3E again outperforms AIRMETs with AUC values for the two forecast systems of 0.738 and 0.597, respectively. The scores for AIRMETs place its performance between the 0.375 and 0.475 thresholds for GTG2.3E. GTG2.3E MOG PODy values for the 0.375 and 0.475 thresholds are 0.67 and 0.45, respectively compared to 0.58 for AIRMETs. The MOG PODn values for GTG2.3E for both the 0.375 and 0.475 thresholds are higher than the AIRMET value of 0.62. At the 0.375 threshold, median forecast volumes from GTG2.3E are slightly larger than those of AIRMETs while volumes using the 0.475 threshold for GTG2.3E are approximately 50% smaller than the AIRMET volumes. These results are summarized in Table 14.

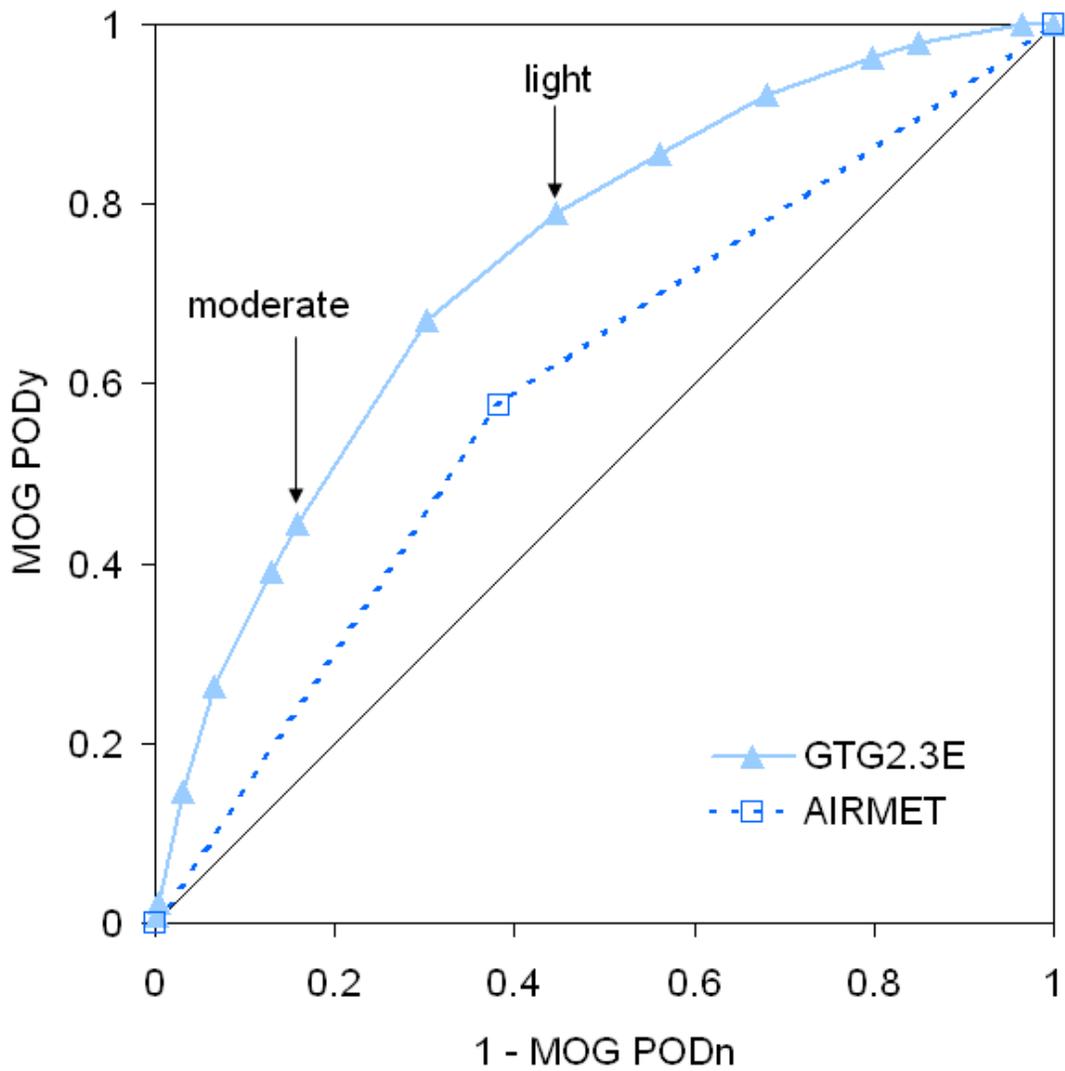


Fig. 26: As in Fig. 24 except for midlevels.

Table 14: MOG PODy, MOG PODn, and median % Volume for GTG2.3E and AIRMETs in midlevels.

	Threshold	MOG PODy	MOG PODn	% Volume _{median}
AIRMETs	-	0.58	0.62	10.0
GTG2.3E	0.475	0.45	0.84	4.15
	0.375	0.67	0.70	11.1

Regionally, at upper levels, GTG2.3E performs best in the Central region followed by the East and West regions while performing best in the West at midlevels and worst in

the East (Table 15). There is no consistent regional pattern of performance between mid- and upper levels for GTG2.3E. AUC values for GTG2.3E are larger in all AWC forecast regions than the corresponding AIRMET values.

Table 15: Regional AUC values

		<i>Upper Level</i>		<i>Midlevel</i>	
		<i>GTG2.3E</i>	<i>AIRMETs</i>	<i>GTG2.3E</i>	<i>AIRMETs</i>
Region	West	0.739	0.598	0.764	0.625
	Central	0.811	0.616	0.729	0.575
	East	0.799	0.632	0.690	0.589

Vertically, GTG2.3E using the MOG threshold of 0.475 and AIRMETs exhibit different behavior. AIRMET MOG PODy values increase very slightly with height with all values exceeding 0.6 except for the 15,000 to 20,000-ft layer (Fig. 27). GTG2.3E MOG PODy values decrease with height, with a maximum value of 0.47 in the 15,000 to 20,000-ft layer and a minimum from 30,000 to 35,000 ft of 0.21. The results again suggest a strong link to the forecast volume of each product as a function of height. Forecast volumes per 5,000 ft vertical layer were unavailable for this study. As a proxy, consider the integral forecast volumes for the mid- and upper-level layers shown in Fig. 28. The large forecast volumes for AIRMETs relative to GTG2.3E are the trade-off associated with the large MOG PODy values achieved by AIRMETs in both mid- and upper levels. For the 0.475 threshold, GTG2.3E % Volume values are quite small for the midlevels with median forecast volumes of less than 4% of the possible airspace. The GTG2.3E % Volume values are even smaller at upper levels. While it is possible that individual 5,000 ft layers within the aggregate mid- and upper layers exhibit differing behavior, it is likely that the volumes within each of these layers decreases with height in a similar manner as the MOG PODy curve. The larger forecast volumes for AIRMETs also contribute to much smaller MOG PODn scores for these forecasts (Fig. 29). The smaller forecast volumes for GTG2.3E correspond to larger MOG PODn scores for GTG2.3E.

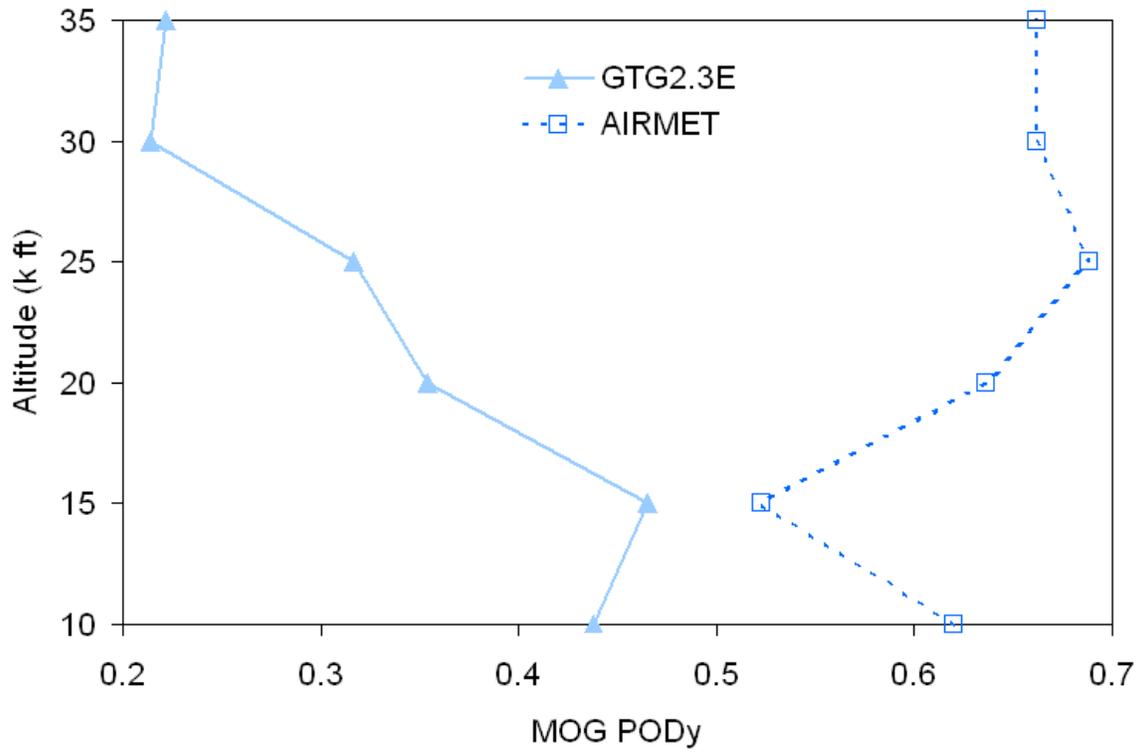


Fig. 27: Height series of MOG PODy for GTG2.3E and AIRMETs. GTG2.3E threshold is 0.475.

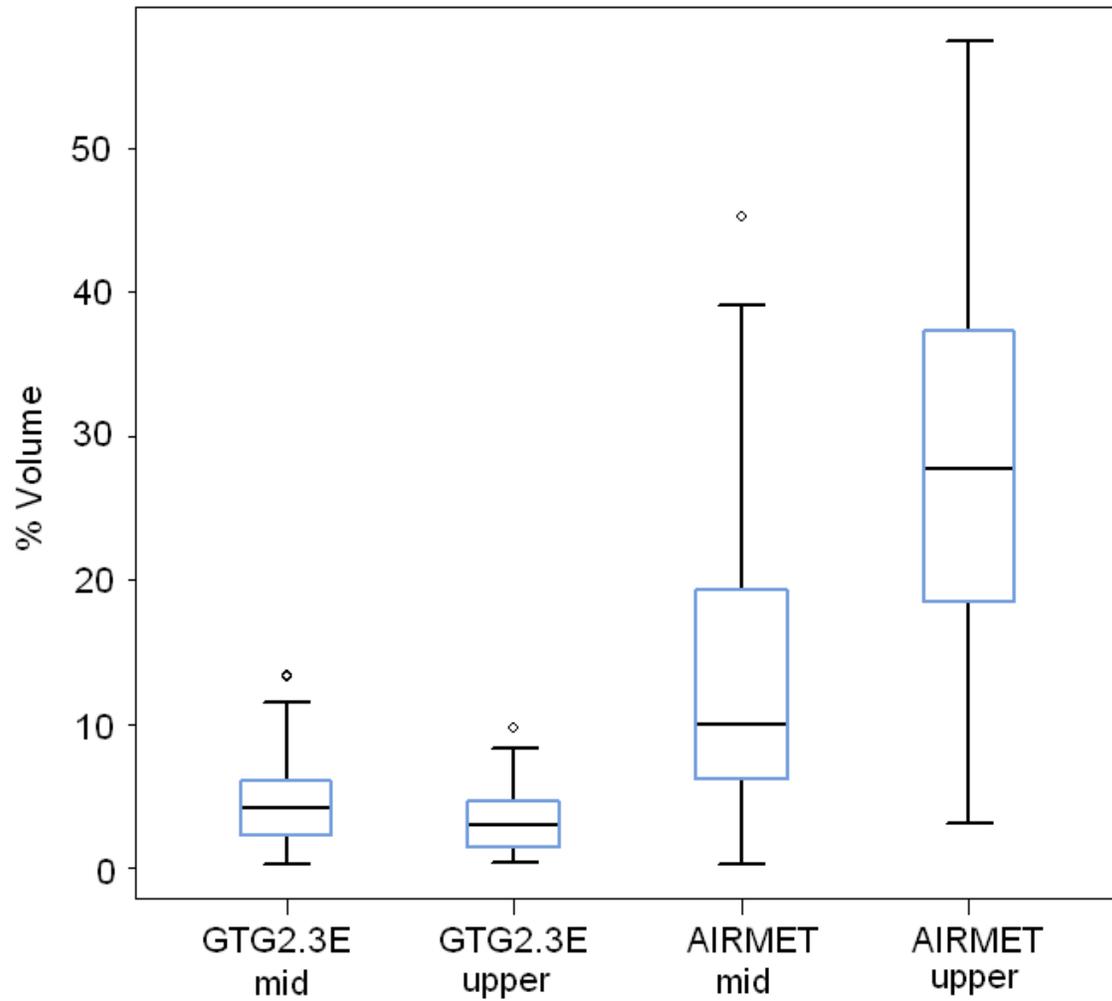


Fig. 28: Boxplot of % Volume values for GTG2.3E and AIRMETs in the mid- and upper levels. GTG2.3E threshold is 0.475.

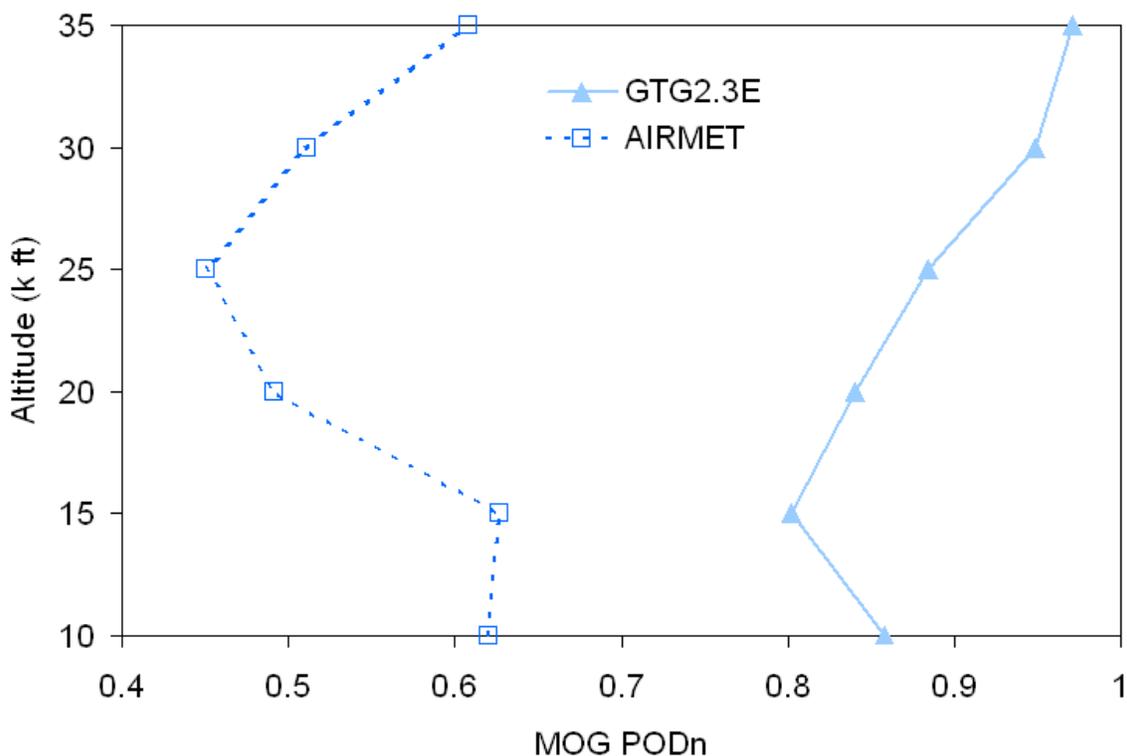


Fig. 29: Height series of MOG PODn for GTG2.3E and AIRMETs. GTG2.3E threshold is 0.475.

6.5 GTG2.3E vs. SIGMETs Comparison

An overview of the capabilities of GTG2.3E and SIGMETs to forecast severe turbulence is considered next. SIGMETs are issued for areas of severe turbulence, and are often once severe conditions have already been encountered. They are valid for up to four hours, but may be canceled at any time. These characteristics of the SIGMETs make it impossible to make a direct comparison of the two types of forecasts. Therefore, any differences that may be identified between GTG2.3E and SIGMETs can only provide general guidance regarding forecast quality rather than define specific differences in the performance of the two types of forecasts.

The following approach was taken for this intercomparison. Forecast/observation pairs were generated from both GTG2.3E and SIGMETs for forecasts valid times 1500, 1800, and 0000 UTC. Data valid at 2100 UTC were also used for SIGMETs. Only 3-h lead time forecasts were used for GTG2.3E. The temporal window around each valid time was ± 60 min. A 0.8 forecast threshold was used to define severe turbulence regions. Additional, lower-, thresholds of 0.625 and 0.475 are also presented below.

The results of the analysis are presented in Table 16. At upper levels, GTG2.3E did not capture any of the 107 severe intensity PIREPs using the Severe category threshold of

0.8. SIGMETs captured a larger fraction of observed severe PIREPs than GTG2.3E with approximately 20% falling inside a SIGMET (27 out of 122). If the forecast threshold is relaxed to 0.625 or 0.475 (the threshold for MOG turbulence) the situation improves slightly. GTG2.3E captures only half of the severe PIREPs when a threshold of 0.475 is used. At a threshold of 0.475, the mean percent volume forecast from GTG2.3E is approximately 2.5 times larger than the mean volume forecast by SIGMETs.

For the midlevel layer from 10,000 to 20,000 ft, a similar pattern of performance is evident. GTG2.3E captured only 7 of 155 severe intensity PIREPs while SIGMETs correctly captured 48 out of 170 reports. If the threshold for forecasts of severe turbulence for GTG2.3E is decreased toward those used to indicate MOG turbulence, detection rates improve at the expense of significantly increased volumes.

As with earlier analyses, the higher detection rate of SIGMETs is likely associated with an increase in forecast volume as compared to the volumes produced by the GTG2.3E algorithm at a threshold of 0.8. Mean forecast volumes for areas of severe turbulence from SIGMETs are more than one hundred times larger than the corresponding volumes from GTG2.3E in both mid- and upper levels. No rigorous statements can be made however about the true amount of overforecasting by either the SIGMETs or GTG2.3E due to the incomplete, non-systematic sampling of the troposphere and lower stratosphere by commercial aircraft.

Table 16: Distribution of forecast/observation pairs and mean % Volume for GTG2.3E with thresholds 0.475, 0.625, and 0.8 and SIGMETs when severe intensity PIREPs were reported for both mid- and upper levels.

<i>Layer</i>	<i>Forecast</i>	<i>Threshold</i>	<i>YY</i>	<i>NY</i>	<i>Mean(% Volume)</i>
Upper Level	GTG2.3E	0.800	0	107	0.001
	GTG2.3E	0.625	6	101	0.244
	GTG2.3E	0.475	52	55	3.19
	SIGMETs	-	27	95	1.2
Midlevel	GTG2.3E	0.800	7	148	0.016
	GTG2.3E	0.625	49	106	0.586
	GTG2.3E	0.475	85	70	4.40
	SIGMETs	-	48	122	1.0

7. CONCLUSIONS

This evaluation was performed to assess the quality of forecasts produced by GTG2.3. Forecasts were assessed for the period 1 November 2005 to 31 January 2006. In addition, GTG2.3 was compared to several operational forecast systems which included GTG

version 1.0, AIRMETs, and SIGMETs. All comparisons were performed using PIREPs as the observational dataset. The limitations of PIREPs, which remain the best operational dataset for many weather phenomena that are hazardous to aviation, preclude a complete evaluation of the forecasting systems. It is important to keep these limitations in mind when interpreting the results and conclusions of this study. The results are only applicable to situations where PIREPs were received. Therefore, no statements can be made about the true amount of over- or underforecasting that may have occurred for any of the forecasts considered here.

GTG2.3 versions E and P performed nearly identically during the evaluation period. The introduction of the *in situ* EDR data into the GTG2.3 algorithm did not appear to decrease the skill of the algorithm in any way. The *in situ* EDR measurements, which have several important attributes not possessed by PIREPs, are currently a relatively small dataset. The limited number of EDR values may have constrained any noticeable effect that the data may have contributed to GTG2.3E. In addition, the evaluation of the algorithms was based on PIREPs, and did not include any *in situ* EDR observations.

Overall, GTG2.3E performs well at forecasting moderate or greater turbulence in both mid- and upper-levels. The new version of the algorithm improves upon the operational version with comparable discrimination of Yes and No PIREPs with significantly reduced volumes. GTG2.3E improved upon operational GTG in all AWC forecast regions. GTG2.3E, like its predecessor, performs only modestly at predicting specific categories of turbulence intensity. It is plausible that the increased availability of *in situ* EDR measurements in the future will allow later versions of GTG to perform better in this regard. In addition, GTG2.3E also improves upon the AIRMET forecasts by reducing the forecast volumes and improving the detection of Yes and No PIREPs. Vertically, GTG2.3E appears to have much less ability to correctly capture moderate or greater PIREPs than AIRMETs. However, to achieve this result, the AIRMET volumes are consistently larger than those produced by the algorithm. GTG2.3E was also compared to SIGMETs to determine how well it performed at forecasting severe turbulence. This comparison was hampered by the small number of severe reports of turbulence received during the evaluation period. Additionally, pilots tend to avoid areas where severe turbulence is either forecast or where a severe PIREP was received previously. Both of these factors limit what can be concluded about forecast quality for either type of forecast. Neither forecast performed particularly well at capturing severe turbulence reports. GTG2.3E forecast volumes for severe turbulence were significantly smaller than those of SIGMETs. However, the true amount of over-, or underforecasting cannot be known owing to the non-systematic nature of PIREPs.

Additional verification information and analyses are available through the RTVS web site (<http://www-ad.fsl.noaa.gov/fvb/rtps/turb/>).

ACKNOWLEDGMENTS

This research is in response to requirements and funding by the Federal Aviation Administration (FAA). The views expressed are those of the authors and do not necessarily represent the official policy and position of the U.S. Government.

We would like to thank the members of the Forecast Verification Section at Earth Systems Research Laboratory for their work on integrating the GTG processing into RTVS. We also thank Jamie Wolff for making the GTG data available to the QA PDT for this report.

8. REFERENCES

- Benjamin, S.G., J.M. Brown, K.J. Brundage, B.E. Schwartz, T.G. Smirnova, and T.L. Smith, 1998: The operational RUC-2. *Preprints, 16th Conference on Weather Analysis and Forecasting*, Phoenix, AZ, American Meteorological Society (Boston), 249-252.
- Brown, B.G. and J.L. Mahoney, 1998: Verification of turbulence algorithms. Report, available from B.G. Brown, NCAR, P.O. Box 3000, Boulder, CO 80307-3000, 9pp.
- Brown, B.G., and G.S. Young, 2000: Verification of icing and turbulence forecasts: Why some verification statistics can't be computed using PIREPs. *Preprints, 9th Conference on Aviation, Range, and Aerospace Meteorology*, 11-15 Sept., Orlando, FL, Amer. Met. Soc., 393-398.
- Cornman, L., G. Meymaris, and M. Limber, 2004: An update on the FAA Aviation Weather Research Program's in situ turbulence measurement and reporting system. *Preprints, 11th Conference on Aviation, Range, and Aerospace Meteorology*, 4-8 Oct., Hyannis, MA, Amer. Met. Soc.
- Kay, M.P., J., J.K. Henderson, and J.L. Mahoney, 2006: GTG version 1.0 proposed threshold changes for Aviation Digital Data Service displays. Report available from the author (mike.kay@noaa.gov), 9pp.
- Mahoney, J.L., J.K. Henderson, B.G. Brown, J.E. Hart, A. Loughe, C. Fischer, and B. Sigren, 2002: The Real-Time Verification System (RTVS) and its application to aviation weather forecasts. *Preprints, 10th Conference on Aviation, Range, and Aerospace Meteorology*, 13-16 May, Portland, OR, Amer. Met. Soc., 323-326.
- Mason, I., 1982, A model for assessment of weather forecasts. *Australian Meteorological Magazine*, 30, 291-303.
- Murphy, A.H. and R.L. Winkler, 1987: A general framework for forecast verification. *Monthly Weather Review*, **115**, 1330-1338.

National Weather Service, 2003: Airman's Meteorological Advisories (AIRMET) Product Description Document. Available from <http://products.weather.gov/PDD/AIRMETv2.pdf>.

National Weather Service, 2003: Significant Meteorological Information (SIGMET) Product Description Document. Available from <http://products.weather.gov/PDD/SIGMET.pdf>.

Schwartz, B.E., 1996: The quantitative use of PIREPS in developing aviation weather guidance products. *Weather and Forecasting*, **11**, 372-384.

Sharman, R., C. Tebaldi, G. Wiener, and J. Wolff, 2006: An integrated approach to mid- and upper-level turbulence forecasting. *Weather and Forecasting*, **21**, 268-287.

Sharman, R., J. Wolff, G. Wiener, and C. Tebaldi, 2004: Technical Description Document for the Graphical Turbulence Guidance Product 2 (GTG2). Report, submitted to the Federal Aviation Administration Aviation Weather Research Program; available from the author (sharman@ucar.edu).

Sharman, R., J. Wolff, G. Wiener, and C. Tebaldi, 2002: Technical Description Document for the Integrated Turbulence Forecast Algorithm (ITFA). Report, submitted to the Federal Aviation Administration Aviation Weather Research Program; available from the author (sharman@ucar.edu).

Takacs, A., L. Holland, M. Chapman, B.G. Brown, J.L. Mahoney, and C. Fischer, 2004: Graphical Turbulence Guidance 2 (GTG2): Quality Assessment Report. Submitted to Aviation Weather Technology Transfer (AWTT) Technical Review Panel.

APPENDIX: CATEGORICAL DISTRIBUTIONS OF FORECASTS AND OBSERVATIONS FOR GTG

The following tables provide the categorical forecast and observation information for operational GTG upper level (20,000 to 40,000 ft) forecasts.

Joint distribution of GTG forecasts and observed PIREPs.

		<i>Observed</i>				Total
		None	Light	Moderate	Severe	
<i>Forecast</i>	None	16122	980	1922	40	19064
	Light	18796	2520	5569	156	27041
	Moderate	4755	1263	3371	101	9490
	Severe	181	118	409	12	720
	total	39854	4881	11271	309	56315

Conditional probability of GTG forecasts given observed PIREPs (p(fl|x)).

		<i>Observed</i>			
		None	Light	Moderate	Severe
<i>Forecast</i>	None	0.405	0.201	0.171	0.129
	Light	0.472	0.516	0.494	0.505
	Moderate	0.119	0.259	0.299	0.327
	Severe	0.005	0.024	0.036	0.039
	p(x)	0.708	0.087	0.200	0.005