

1 Dynamic Evaluation of CMAQ Part II: Evaluation of relative response
2 factor metrics for ozone attainment demonstrations

3 Kristen M. Foley*, Patrick Dolwick, Christian Hogrefe, Heather Simon, Brian Timin, Norm
4 Possiel

5
6 United States Environmental Protection Agency, Research Triangle Park, NC 27711, USA

7
8 *Corresponding Author: telephone (919) 541-5367; fax (919) 541-1379; email
9 foley.kristen@epa.gov.

10
11
12 **Abstract**

13 The U.S. Environmental Protection Agency provides guidelines on the use of air quality models
14 for projecting whether an emission reduction strategy will lead to future pollutant levels that are at
15 or below the National Ambient Air Quality Standards (NAAQS). The EPA’s guidance document
16 for ozone attainment demonstrations recommends an attainment test for the 8-hour ozone NAAQS
17 based on using the ratio of output from “future” and “base” model simulations through the
18 calculation of location-specific Relative Response Factors (RRF). The 2007 guidance document
19 as well as other related studies have recommended the use of retrospective evaluation studies in
20 order to evaluate the ability of an air quality model to represent a change in air quality (dynamic
21 evaluation) rather than relying solely on operational evaluation of model predictions under base
22 line conditions. Here simulations from the Community Multiscale Air Quality (CMAQ) modeling
23 system were conducted for 2002 and 2005, a time period characterized by significant emissions
24 reductions associated with the EPA’s Nitrogen Oxides State Implementation Plan Call (NO_x SIP
25 Call) as well as mobile sources. These simulations were used to evaluate the performance of
26 different forms of the RRF metric for projecting 2002 to 2005 against 2005 observed ozone levels.
27 The evaluation study showed that the current form of the RRF calculation is generally well
28 designed for predicting the future 8-hr ozone “design value” metric used for determining
29 attainment. Specifically, the methodology of using air quality model simulations in a relative
30 sense provided better estimates of future ozone design values than using the modeled future year
31 simulation alone. Alternative forms of the RRF metric were found to be very similar to the current
32 methodology in terms of evaluation metrics. However, alternative RRF metrics were sensitive to
33 the number of days used in the calculation of the RRF. Approaches which used more days in the
34 RRF calculation (relative to the 2007 guidance approach) had slightly higher bias and error in
35 predicting 2005 design values compared to approaches using only a subset of high ozone days.

36 **Keywords:** Dynamic Evaluation; NO_x SIP Call; ozone NAAQS; CMAQ version 5.0.1; relative
37 response factor

38
39 **1 Introduction**

40 The U.S. Environmental Protection Agency sets National Ambient Air Quality Standards
41 (NAAQS) for six criteria pollutants considered harmful to public health and the environment.
42 When an area is found to exceed these standards and is subsequently designated as a
43 “nonattainment area” for a particular pollutant, states are required by the Clean Air Act to develop
44 State Implementation Plans (SIP) to describe what emission control strategies will be used to
45 attain the NAAQS by a target date. For ozone, an important aspect of SIP development is the use
46 of air quality models to demonstrate that the planned emission controls will provide the necessary

1 reduction in pollutant levels to achieve attainment of the standard. In 1999 the EPA provided draft
2 guidance on the use of air quality models for SIP development for the 8-hr NAAQS ozone
3 standard (EPA, 1999). The guidance outlined a method for using the ratio of model output from
4 future-case and base-case emission conditions under fixed meteorological conditions to calculate a
5 Relative Response Factors (RRF). The RRF is then used to scale observed ozone mixing ratios
6 representative of the base year to determine if future year ozone levels will be in compliance with
7 the NAAQS. This use of the model in a relative sense was a different methodology than the
8 previous guidance for the 1-hr ozone standard which relied only on absolute modeled future-case
9 concentrations. The RRF approach was updated and finalized in the 2007 guidance document
10 (EPA, 2007).

11
12 One motivation for the use of air quality models in a relative approach is that anchoring the future
13 year prediction on an observed base year value is expected to reduce problems posed by imperfect
14 model performance (EPA, 2007). That is, some types of systematic model errors are expected to
15 “cancel out” using the relative approach. In addition, several studies have shown that the relative
16 approach is much less sensitive to the choice of modeling system or model configuration
17 compared to using the model output directly (Jones et al., 2005; Sistla et al., 2004).

18
19 The use of photochemical grid models to predict changes in air quality motivates the need for
20 additional evaluation, beyond the standard operational evaluation of model predictions under base
21 line conditions. Hogrefe et al. (2008) showed that base model evaluation does not necessarily
22 inform how well the model will estimate the change in air quality due to changes in emissions.
23 Here a dynamic evaluation of the Community Multiscale Air Quality Modeling system (CMAQ)
24 was conducted using a retrospective case study to evaluate the performance of different forms of
25 the RRF metric to accurately predict reductions in ozone levels associated with reductions in NO_x
26 emissions. The case study was based on modeling 2002 and 2005, a time period characterized by
27 significant reductions in ozone precursor emissions. One contributor to these emission reductions
28 was the EPA’s Nitrogen Oxides State Implementation Plan Call (NO_x SIP Call) rule, fully
29 implemented by 2004, and designed to reduce NO_x emissions from power plants in the eastern US
30 (Gilliland et al., 2008). In addition, mobile source NO_x emissions also decreased by roughly 20%
31 during this time period due to motor vehicle fleet turnover to newer lower emitting vehicles (Foley
32 et al., 2014). The combined emissions reductions contributed to a decrease in observed ozone
33 levels up to 20-30% in the eastern US from 2002 to 2005. Additional background on the dynamic
34 evaluation of CMAQv5.0.1 can be found in Part I of this study (Foley et al., 2014). The analysis
35 in Part I shows that this modeling system tends to underestimate changes in ozone between these
36 two years, similar to the finding in other studies (e.g. Gilliland et al., 2008; Zhou et al., 2013).
37 The change in model estimated air quality was found to be most dependent on changes in
38 emissions and somewhat dependent on changes in meteorology from 2002 to 2005. Here we
39 expand on the dynamic evaluation by evaluating model projection metrics used to support
40 regulatory modeled attainment demonstrations for the 8-hr Ozone NAAQS.

41
42 Section 2 provides details on the model simulations conducted for this case study including the
43 development of a model simulation using 2005 emissions under 2002 meteorological conditions
44 which was needed for the calculation of the RRF ratios. Section 3 presents a comparison of
45 model-projected 2005 DVs using model output for base and future scenarios through the
46 calculation of RRF metrics to an approach based solely on model predictions from the “future”

1 year simulation. While EPA guidance and other studies have recommended a relative approach
2 for over a decade, to our knowledge, this is the first time that the difference between these two
3 approaches has been quantified using a retrospective dynamic evaluation study. Five different
4 forms of the RRF metric are compared in section 3 followed by a discussion in Section 4.

6 **2 Model approach for NO_x SIP Call time period**

7 The model simulations used in this analysis are described in detail in sections 2 and 4 in Foley et
8 al. (2014). In summary, four CMAQv5.0.1 simulations were performed over the continental U.S.
9 using a grid with 12km horizontal resolution and 35 vertical layers. Meteorological inputs were
10 based on WRF3.3 with MCIPv4.0 (see Appel et al., 2013 for further details). Emission inputs
11 were developed based on 2002 and 2005 National Emission Inventory data using the Sparse
12 Matrix Operator Kernel Emissions version3.1 emissions processing system (SMOKE; Houyoux et
13 al., 2000), including year specific data for large point sources and year specific mobile emissions
14 derived from MOVESv2010b. Boundary conditions were based on 2005 monthly median values
15 from a GEOS-Chem v9-01-02 simulation (<http://wiki.seas.harvard.edu/geos-chem/>) using v8-02-
16 01 chemistry, GEOS-5 meteorology and ICOADS shipping emissions (Henderson et al., 2014).

17
18 Two base year simulations were conducted for June 1 through September 30th 2002 and 2005
19 (Sim02e02m, Sim05e05m). In addition, two “cross” simulations were used to simulate air quality
20 under 2005 emissions with 2002 meteorology (Sim05e02m) and 2002 emissions with 2005
21 meteorology (Sim02e05m). Modeling for attainment demonstrations is based on a base year and a
22 future year where the meteorology is held constant across both simulations. Thus the focus of this
23 analysis was on the Sim02e02m and Sim05e02m simulations. The other two simulations are used
24 briefly to illustrate the change in ozone levels due to changes in meteorology from 2002 to 2005.

25
26 Many methods exist for creating future year emissions for attainment demonstrations and different
27 methods have been developed for different emission sources. Here the NEI-based 2005 emission
28 estimates are modified, as described below, to create the emission inputs for Sim05e02m which is
29 used as the ‘future’ simulation in the RRF calculations. The focus of the current study is to isolate
30 the errors and biases in predicted ozone levels that are associated with either the form of the RRF
31 metric or the air quality modeling system itself. The amount of error in the predicted future case
32 ozone levels that can be attributed to errors in projecting future emission levels will depend on the
33 methodology used. An analysis of this type of modeling uncertainty is beyond the scope of the
34 current work.

35
36 For the simulation with 2005 emissions and 2002 meteorology (Sim05e02m), the emissions from
37 electrical generating units (EGUs) with available continuous emission monitoring systems (CEM)
38 data in 2005 were adjusted to account for the impact of different meteorology in 2002.

39 Summertime 2005 NO_x emissions are generally lower than 2002 emissions due to emission
40 reductions. Temporal fluctuations are different due to differences in electricity demand which is
41 heavily influenced by year-specific meteorology. 2005 emission levels with 2002 meteorological
42 patterns (EMIS05e02m) were estimated by scaling the hourly 2002 CEM emissions (CEM2002)
43 based on the ratio of summer total CEM emissions (STCEM) for a particular EGU unit in 2005
44 versus 2002: $EMIS05e02m = CEM2002 \times (STCEM2005/STCEM2002)$. An analogous
45 calculation is made to estimate 2002 emissions with 2005 meteorological patterns for Sim02e05m.
46 Emissions from electrical generating point sources with no CEMS data used 2005 annual total

1 emissions scaled with 2004-2006 annual-to-month ratios and 2002 day-to-month ratios in order to
2 match the daily temporal fluctuations of the point sources with CEMS data in 2002. Additional
3 details on the calculation of the scaling ratios is provided in Supplemental material S2.

4
5 Mobile emissions for Sim05e02m were based on a MOVES simulation using 2005 VMT and
6 emission factors and 2002 meteorological inputs. Emissions from nonroad (e.g. construction),
7 industrial point and commercial marine sectors are based on 2005 emission levels but shifted to
8 match the day-of-the week in 2002. Emissions from fertilizer application, biogenic sources, NO_x
9 from lightning, fires and dust are based on 2002 estimates since these sources are associated with
10 meteorological conditions. All other sectors have identical emissions inputs for both Sim02e02m
11 and Sim05e02m.

12
13 Section 3 outlines how the base year (Sim02e02m) and future year (Sim05e02m) simulations were
14 used to calculate different RRF metrics in order to estimate 2005 ozone “design values”. As part
15 of determining which areas are attaining the ozone standard, EPA calculates ozone design values
16 (DVs) as the 3-year average of the annual 4th highest maximum daily 8-hour average (MDA8)
17 ozone concentration. Determining attainment based on multiple years of observations reduces the
18 impact of interannual meteorological variability on pollutant concentrations. For attainment
19 demonstrations in which states show how they plan to bring future ozone levels in line with the
20 standard, EPA recommends establishing site-specific base DVs to be the average of three
21 consecutive 3-year DVs, centered about the modeling year (which is effectively a 5-year weighted
22 average DV) to further reduce the efforts of meteorological variability on design values.
23 However, in this retrospective study the base and ‘future’ years are very close together in order to
24 isolate the impact of a specific emission control event (i.e. the NO_x SIP Call). As a result, the
25 standard calculation of the 2002 DVs would have been based on summers both before and after
26 the controls were implemented in 2003 and 2004. To avoid this issue, the following analysis is
27 based on 3-year DVs using 2000-2002 as the base year average (2002 DV) and 2004-2006 as the
28 future year average (2005 DV).

29
30 Ozone observations from 644 monitoring stations from the EPA’s Air Quality System (AQS;
31 <http://www.epa.gov/ttn/airs/airsaqs>) across the U.S. were used to evaluate model predicted 2005
32 design values. These monitoring sites were selected from the entire AQS network based on
33 observed 2002 DVs. Only sites with 2002 DVs greater than 75ppb were included in this study.
34 Sites with lower base year ozone are of less interest because these areas would not be expected to
35 have attainment issues for this case study.

36 37 **3 Comparison of 5 RRF metrics**

38
39 In this section the form of the 8-hr ozone NAAQS attainment test is reviewed followed by a
40 comparison of alternate RRF metrics. The RRF metrics are also compared to using the model in
41 an absolute sense rather than in a relative sense.

42
43 The 8-hr ozone model attainment test for monitoring site *i* provides an estimate of the future DV
44 for the year in which attainment is required (DVF_i) based on the relative response factor (RRF_i)
45 and the base design value for the site (DVB_i):

$$46 \quad DVF_i = RRF_i \times DVB_i, \quad (1)$$

1 where RRF_i is a ratio of modeled future ozone to modeled base ozone. In the model attainment
2 test suggested by the 2007 guidance (EPA, 2007), the modeled base ozone value (the denominator
3 of the RRF) is an average of MDA8 ozone on all days in the base simulation greater than or equal
4 to 85ppb. If there are less than ten days ≥ 85 ppb in the base year, then the average of the ten
5 highest days in the simulation are taken as long as each of these values is ≥ 70 ppb. If there are at
6 least five days but less than ten days with MDA8 ozone ≥ 70 ppb then the average of all days
7 ≥ 70 ppb are used. In the case when there are less than five days with MDA8 ozone ≥ 70 ppb the
8 RRF is not calculated. The modeled future ozone (the numerator of the RRF) is the average of the
9 MDA8 ozone in the future year simulation over the same calendar days used in the base ozone
10 calculation. This is referred to as the TH85 (“Threshold 85”) approach for the remainder of the
11 paper.

12
13 Before exploring alternate methods for calculating the RRF at a given site, the NO_x SIP Call case
14 study was used to evaluate two underlying assumptions of the ozone attainment test. First is the
15 idea that the change in the average of high ozone values is a good predictor of the change in the 4th
16 highest ozone value across several summers. The motivation for using a threshold approach is
17 that days with high base-case ozone levels are expected to be the days most representative of
18 nonattainment conditions, which is what emission controls are designed to limit. Specifically,
19 high base-case ozone levels have been found to be more responsive to emission reductions
20 compared to days with lower ozone (Hogrefe et al., 2008). To investigate this with the current
21 dataset, “daily RRF” values were calculated for each summer day in 2002 (June – September).
22 The daily RRF at a given grid cell containing an AQS monitor is the ratio of the ozone from the
23 2005 simulation using 2002 meteorology (Sim05e02m) to the ozone value from the 2002

24 simulation (Sim02e02m) at that grid cell. That is, for location i , day j , $dailyRRF_{ij} = \frac{MDA8_{ij}^{05e02m}}{MDA8_{ij}^{02e02m}}$.

25 Figure 1 shows daily RRF values as a function of 2002 modeled MDA8 ozone values. The top
26 left figure shows that the daily RRF values across all sites tend to steadily decrease with
27 increasing base level ozone values and then level off around 80ppb. The daily RRFs for several
28 urban sites are also shown with the individual days that would be used in the calculation of the
29 TH85 RRF shown in red.

30
31 To further evaluate the use of a threshold approach, an “observed” RRF was calculated at each
32 AQS site based on the ratio of the average of observed MDA8 ozone in June-September 2005 to
33 2002. At site i , the denominator of the ratio is the average of the N_i days in 2002 with observed
34 MDA8 ozone ≥ 85 ppb and the numerator is the average of the top N_i MDA8 observed ozone days
35 in 2005. The value of N_i is the same in the calculation of the numerator and denominator at a
36 given site but can change across AQS sites. For example, if there were 15 observed MDA8 ozone
37 values greater than 85ppb at a particular AQS site in 2002, then the denominator is the average of
38 these 15 MDA8 ozone values and the numerator is the average of the top 15 observations at that
39 site in 2005. An observation based estimate of the 2005 DV was then calculated by multiplying
40 the RRF times the 2002 DV. Figure 2(a) shows the observed RRF approach does indeed evaluate
41 well against the actual observed 2005 DVs. Note that unlike the model-based RRF, this
42 calculation captures changes in both meteorology and emissions across these years. This result is
43 included to provide context for the evaluation metrics used to evaluate the model-based results.
44 Using a model-based RRF approach would not be expected to improve upon the observation-
45 based metrics (e.g. MB = -1.9ppb, RMSE=5.6ppb, $R^2 = 0.68$).

1
2 The second issue evaluated was whether it is more appropriate to use the model output in a
3 relative sense rather than comparing the absolute future year model prediction (e.g. output directly
4 from Sim05e02m) to the level of the 8-hr ozone NAAQS. The second and third panels in Figure 2
5 show the observed 2005 DV compared to (b) the 4th highest MDA8 ozone value from Sim05e02m
6 and (c) the predicted 2005 DV using the model-based TH85 RRF approach. Using the TH85 RRF
7 approach results in smaller positive bias (3.8ppb vs. 4.1ppb) and much higher correlation ($R^2 =$
8 0.67 vs. 0.37) compared to the Sim05e02m prediction. The implication for attainment testing is
9 that the relative approach anchored to 2002 observed design values is better able to predict
10 exceedances of both the 75ppb (threshold for current 8-hr O₃ NAAQS) and 84ppb (threshold for
11 NAAQS during 1997-2008) as reflected by the accuracy scores (Acc75 and Acc84). Note that the
12 1997 NAAQS is 0.08ppm with DVs rounded to two significant digits, i.e., any DV below
13 0.085ppm (85ppb) is considered in attainment. The 2008 standard is 0.075ppm with DVs
14 truncated to three significant digits, i.e., any DV below 0.076ppm (76ppb) is considered in
15 attainment. These results indicate that for the current O₃ NAAQS, the methodology of using high
16 base-case ozone days and using the air quality model simulations in a relative sense is well suited
17 for determining attainment and can mitigate some biases in the modeling system being applied.
18

19 Several alternate RRF calculations were compared in order to see if an alternate approach could be
20 adopted to improve the accuracy of model-based DV predictions. To test whether it makes sense
21 to tie the RRF calculation to the level of the ozone NAAQS, two additional threshold approaches
22 were compared to the TH85 RRF. TH76 is based on averaging model days in 2002 with MDA8
23 ozone ≥ 76 ppb. If there are fewer than ten days ≥ 76 ppb in 2002 then all days ≥ 60 ppb are used. If
24 there are fewer than five days ≥ 60 ppb then this metric is not calculated. TH71 is constructed in an
25 analogous manner based on days greater than or equal to 71ppb. To avoid the use of a threshold
26 all together, two “number of days” tests were also compared. The Top10 RRF averages the ten
27 highest MDA8 ozone days in 2002 that are above 60ppb. If there are less than ten days ≥ 60 ppb in
28 2002 then all days ≥ 60 ppb are used. If there are less than five days ≥ 60 ppb then this metric is not
29 calculated at that site. The Top20 RRF is constructed in the same fashion except that at the first
30 step the twenty highest MDA8 ozone days are used. The numerator of the RRF for all of these
31 metrics is calculated by averaging the MDA8 ozone in the future year simulation over the same
32 calendar days that were used in the base ozone calculation. Table 1 provides an overview of the
33 five metrics that were compared.
34

35 Table 2 provides evaluation statistics for the five methods compared to using the output from the
36 future model simulation directly (i.e., Sim05e02m). The 644 AQS sites are subset to sites where
37 all six methods could be calculated. Twenty-five of these sites did not have five days with MDA8
38 ozone ≥ 70 ppb needed for the calculation of the TH85 metric, reducing the final dataset to n=619.
39 (Note that Tables S1 and S2 in the Supplemental material also show these statistics segregated by
40 NOx SIP call states and non-NOx SIP Call states.)
41

42 The evaluation statistics of the five RRF type methods are almost identical with the TH85 and
43 Top10 metrics having slightly lower bias and RMSE compared to the other 3 RRF metrics. All of
44 the methods have a positive bias (3.8-4.2ppb) which means the model did not predict as large a
45 decrease in the ozone DVs from 2002 to 2005 as was observed. Part of this bias can be attributed
46 to the fact that the model based predictions do not account for changes in meteorology across these

1 summers while the observed changes in DV, though averaged over three years, do not completely
2 eliminate the effects of such changes in meteorology. Analysis of the second cross run,
3 Sim02e05m, showed that the modeled change in meteorology across these years did lead to a
4 decrease in ozone levels in many parts of the domain (see Supplemental Figure S1). For example,
5 using the Sim05e05m output for the numerator of the model-based RRF calculations (i.e. creating
6 a model RRF that captured both meteorological and emissions changes) reduced the mean bias in
7 the DV predictions to 2.4ppb as shown in Figure 2(d). The remaining underestimation of the
8 observed change is a result of the model predictions for high summer time MDA8 ozone being too
9 low in 2002 and too high in 2005 (the dynamic evaluation of these two years is discussed in detail
10 in Part I, Foley et al. 2014). In other scenarios when the bias in the base and future year
11 simulations are more comparable (e.g. the MDA8 ozone predictions are too high in both years)
12 using the relative approach would be expected to provide an even greater improvement in
13 prediction performance compared to using a single model simulation in an absolute sense.

14
15 Table 2 also provides two metrics for evaluating the 2005 predictions in a categorical sense. The
16 modeled 2005 DVs were compared to the ozone NAAQS level to determine whether a site
17 exceeded the standard or was in attainment. The model prediction at a given site was compared to
18 the observed 2005 DV, producing four outcomes: true exceedance, false exceedance, false
19 attainment, true attainment. For example, for the 84ppb NAAQS if *Mod* is the model predicted
20 2005 DV and *Obs* is the observed 2005 DV, the four outcomes are: $Mod \geq 85\text{ppb}$ and $Obs \geq 85\text{ppb}$
21 (true exceedance); $Mod \geq 85\text{ppb}$ and $Obs < 85\text{ppb}$ (false exceedance); $Mod < 85\text{ppb}$ and Obs
22 $\geq 85\text{ppb}$ (false attainment); $Mod < 85\text{ppb}$ and $Obs < 85\text{ppb}$ (true attainment). Spatial plots of these
23 categorical outcomes for the future model output and the TH85 RRF approach are shown in Figure
24 3 and Supplemental Figure S3 for both the 84ppb and 75ppb standards for the entire modeling
25 domain and the Eastern U.S., respectively. Note that the locations where false attainments or false
26 exceedances occurred are specific to model simulations used in this case study and would be
27 expected to change for a different model setup or a different set of simulation years.

28
29 The accuracy of the model predictions is the fraction of sites where the model produced a correct
30 prediction (true exceedance or true attainment). For the 75ppb standard, the accuracy of all of the
31 RRF methods is higher than the future model output (0.78-0.79. vs. 0.72). The accuracy scores for
32 the 84ppb standard are even higher with the TH85 and Top10 metrics having the highest score
33 (0.82) of all of the methods. Figure 3 shows the TH85 approach decreases both false attainment
34 predictions (17 vs. 28) and false exceedance predictions (97 vs 114) compared to the future model.
35 Because the TH85 RRF approach defaults to the top ten days when there are ten or fewer modeled
36 days $\geq 85\text{ppb}$, the Top10 RRF approach and the TH85 approach were identical at 77% of the AQS
37 sites. At the remaining sites, the Top10 RRF had slightly lower bias (MB of 5.3ppb vs. 5.4ppb)
38 and error (RMSE of 7.9ppb vs. 8.0ppb) than the TH85 approach. Moreover, the Top10 RRF and
39 the TH85 RRF approaches predict the same attainment designations (i.e. exceedance or
40 attainment) at all but a single AQS site for both the 84ppb and the 75ppb standards.

41
42 This analysis showed that all of the tested RRF metrics were similar. However methods that used
43 a larger of number of days in the RRF calculation (see Figure 4) all led to higher MB and RMSE
44 compared to the Top10 RRF and TH85 RRF approaches. The Top10 metric was found to be a
45 comparable metric to the current TH85 approach and offers the advantage that it is independent of
46 the threshold used in the NAAQS standard, which could be revised again in the future. In

1 addition, the Top10 metric simplifies the RRF calculation and ensures that only the top percentile
2 of ozone days are included in the model averages. For example, some locations with very high
3 ozone had more than twenty or even thirty days above 84ppb in 2002, i.e. 15-25% of the entire
4 simulation period.

6 **4 Summary**

7 The dynamic evaluation case study based on the 2002-2005 time period characterized by
8 emissions reductions associated with the NOx SIP Call and mobile source emissions offered a
9 valuable opportunity to assess the methodology used in the ozone attainment test, as well as the
10 suitability of the CMAQ modeling system for use in this type of regulatory modeling application.
11 The RRF approach was found to be generally well designed to capture changes in the observed
12 DVs across these years, and offered several advantages over using output from the air quality
13 model directly, in terms of lower bias, higher correlation and higher accuracy in predicting ozone
14 exceedances. There were no large differences found in the evaluation of the five different RRF
15 metrics. Basing the metrics on the highest modeled base-year ozone days yielded a slightly lower
16 bias and higher correlation compared to the metrics that averaged over a larger percentage of the
17 simulation days.

18
19 The performance of the RRF based approach was quantified here based on aggregate statistics for
20 monitoring sites across the contiguous US domain. For a given SIP there may be local
21 circumstances where large uncertainties in model output and observations would necessitate
22 analysis beyond the RRF methodology, referred to by the EPA guidance as a “weight of evidence”
23 analysis. For example, Vizuite et al. 2010 and 2011 provide detailed discussion and evaluation of
24 RRF calculations for the Houston, TX area. Such area-specific analysis will continue to be an
25 important part of the SIP process regardless of the specific RRF methodology being applied.

26
27 This study also demonstrated how inter-annual meteorological variability can impact the
28 attainment test and be a confounding factor in the dynamic evaluation of the air quality modeling
29 system being used. In this case, the change in meteorological conditions from 2002 to 2005
30 caused ozone to decrease in many locations resulting in an increase in the positive bias in the
31 model predicted DVs. The use of 3-year DVs helps to reduce the impacts of year-to-year changes
32 in meteorology on observed pollutant levels but likely does not eliminate these effects. Future
33 work using longer-term records of observations and model simulations is needed to assess whether
34 using model output from multiple years in calculating RRFs may improve the ability of the
35 modeling system to reproduce observed changes in DVs.

36
37 **Disclaimer:** Although this work was reviewed by EPA and approved for publication, it may not
38 necessarily reflect official Agency policy.

39 **Acknowledgements**

40 The authors would like to recognize the many contributions of others in their support in processing
41 the large suite of emission, meteorology and boundary condition inputs needed for the
42 CMAQv5.0.1 simulations, in assisting with observed datasets and in providing many thoughtful
43 suggestions on how to improve the final analysis: Wyatt Appel, Rob Gilliam, Jim Godowitch, Rob
44 Pinder, George Pouliot, Shawn Roselle (EPA/ORD/NERL); Kirk Baker, Alison Eyth, Sharon
45 Philips, Benjamin Wells, Alexis Zubrow (EPA/OAR/OAQPS); Allan Beidler, Lucille Bender,
46

1 Ryan Cleary (Computer Sciences Corporation); Barron Henderson (now at Univ. of FL); and
2 Farhan Akhtar (now at U.S. State Dept.).
3
4

5 **References**

- 6
7 Appel, K.W., Pouliot, G.A., Simon, H., Sarwar, G., Pye, H.O.T., Napelenok, S.L., Akhtar, F.,
8 Roselle, S.J. (2013) Evaluation of dust and trace metal estimates from the Community
9 Multiscale Air Quality (CMAQ) model version 5.0. *Geoscientific Model Development* 6, 883-
10 899.
- 11 Foley, K.M., Hogrefe, C., Pouliot, G., Possiel, N., Roselle, S., Simon, H., Timin, B. (2014)
12 Dynamic evaluation of CMAQv5.0.1 Part 1: Quantifying changes in ozone associated with the
13 EPA's NO_x SIP Call Rule, *in review*.
- 14 Gilliland, A.B., Hogrefe, C., Pinder, R.W., Godowitch, J.M., Foley, K.M., Rao, S.T. (2008)
15 Dynamic evaluation of regional air quality models: Assessing changes in O₃ stemming from
16 changes in emissions and meteorology. *Atmospheric Environment* 42, 5110-5123.
- 17 Henderson, B.H., Akhtar, F., Pye, H.O.T., Napelenok, S.L., Hutzell, W.T. (2014) A database and
18 tool for boundary conditions for regional air quality modeling: description and evaluation.
19 *Geoscientific Model Development*, 7, 339-360.
- 20 Hogrefe, C.; Civerolo, K.L.; Hao, W., Ku, J.Y.; Zalewsky, E.E.; Sistla, G. (2008) Rethinking the
21 assessment of photochemical modeling systems in air quality planning applications. *J. Air &*
22 *Waste Manage. Assoc.*, 58, 1086-1099.
- 23 Houyoux, M. R., Vukovich, J. M., Coats Jr., C. J., Wheeler, N. J. M., Kasibhatla, P. (2000)
24 Emission inventory development and processing for the seasonal model for regional air quality,
25 *Journal of Geophysical Research: Atmospheres*, 105, 9079–9090.
- 26 Jones, J.M.; Hogrefe, C.; Henry, R.F.; Ku, J.-Y.; Sistla, G. (2005) An assessment of the sensitivity
27 and reliability of the relative reduction factor (RRF) approach in the development of 8-hr
28 Ozone Attainment Plans; *J. Air & Waste Manage. Assoc.*, 55, 13-19.
- 29 Sistla, G.; Hogrefe, C.; Hao, W.; Ku, J.-Y.; Zalewsky, E.; Henry, R.F.; Civerolo, K. (2004) An
30 operational assessment of the application of the relative reduction factors (RRF) in
31 demonstration of attainment of the 8-hr ozone National Ambient Air Quality Standard
32 (NAAQS); *J. Air & Waste Manage. Assoc.*, 54, 950-959.
- 33 United States Environmental Protection Agency, 1999. Draft report on the use of models and
34 other analyses in attainment demonstrations for the 8-hour ozone NAAQS. EPA-44/R-99-
35 0001, May 1999, *United States Environmental Protection Agency*, Research Triangle Park, NC
36 27711.
- 37 United States Environmental Protection Agency, 2007. Guidance on the use of models and other
38 analyses for demonstrating attainment of air quality goals for ozone, PM_{2.5}, and regional haze.
39 EPA-454/B-07-002, April 2007, *U.S. Environmental Protection Agency*, Research Triangle
40 Park, NC, 27711.
- 41 Vizuite, W., Biton, L., Jeffries, H., Couzo, E. (2010) Evaluation of relative response factor
42 methodology for demonstrating attainment of ozone in Houston, Texas. *J. Air & Waste*
43 *Manage. Assoc.*, 60, 838-848.
- 44 Vizuite, W., Jeffries, H.E., Tesche, T.W., Olague, E.P., Couzo, E. (2011) Issues with ozone
45 attainment methodology in Houston, TX. *J. Air & Waste Manage. Assoc.*, 61, 238-253.

1 Zhou, W., Cohan, D.S., Napelenok, S.L. (2013) Reconciling NOx emissions reductions and ozone
2 trends in the U.S., 2002-2006, *Atmospheric Environment*, 70, 236-244.
3
4

1 List of Figures

2
3 Figure 1. Daily RRF ratios versus 2002 modeled MDA8 ozone values at all AQS sites (top left)
4 and at 5 select AQS sites in very urban areas. The color scale of the density scatter plot in the top
5 right figure represents the percent of the data that fall within a particular pixel in the plot, e.g., the
6 red pixels in the plot indicate areas where more than 1.5% of the data (>10 data points) would be
7 clustered in a regular scatterplot. The green curve shows a spline fit to the daily RRF values. The
8 red points in the scatter plot indicate what days were used in the calculation of the TH85 RRF
9 metric for each site.

10
11 Figure 2. Comparison of 2005 observed versus predicted DV at n=619 AQS sites. Predicted DVs
12 are based on using the (a) observed TH85 RRF approach, (b) the 4th highest MDA8 ozone value
13 from the Sim05e02m simulation (c) the modeled TH85 approach and (d) the modeled TH85
14 approach modified by using 2005 meteorology in the “future” model simulation (i.e. output from
15 Sim05e05m rather than Sim05e02m is used in the numerator). Evaluation statistics include mean
16 bias (MB), root mean square error (RMSE), R2, and accuracy scores for predicting attainment of
17 the 75ppb and 84ppb NAAQS (Acc75, Acc84).

18
19 Figure 3. Observed 2002 and 2005 design values (top row) across the entire model domain. The
20 remaining plots show model predicted 2005 DVs in a categorical sense based on output from
21 Sim05e02m (left column) and the TH85 RRF based approach (right column) for both the 84ppb
22 standard (middle row) and the 75ppb standard (bottom row). Four categories are depicted: true
23 exceedance (green circle); false exceedance (purple triangle); false attainment (red inverted
24 triangle); true attainment (yellow square).

25
26 Figure 4. The number of days used in the calculation of each of the five RRF metrics across the
27 619 AQS sites used in Table 2.

28

1 List of Tables

2
3 Table 1. Summary of five RRF metrics compared in section 3. TH85 was the method described
4 in the 2007 guidance document.

5
6 Table 2. Evaluation of different approaches for predicting the observed 2005 DVs at 619 AQS
7 sites. The evaluation statistics include mean bias (MB), root mean square error (RMSE), R2, and
8 accuracy for predicting attainment for the 75ppb and 84ppb NAAQS.

9

1 Table 1. Summary of five RRF metrics compared in section 3. TH85 was the method described
2 in the 2007 guidance document.

3

RRF metric	Type of metric	Days averaged	Minimum ozone level used in the average	Minimum # days to be averaged
TH85	Threshold	Days \geq 85	70	5
TH76	Threshold	Days \geq 76	60	5
TH71	Threshold	Days \geq 71	60	5
Top10	Number of days	Highest 10 days	60	5
Top20	Number of days	Highest 20 days	60	5

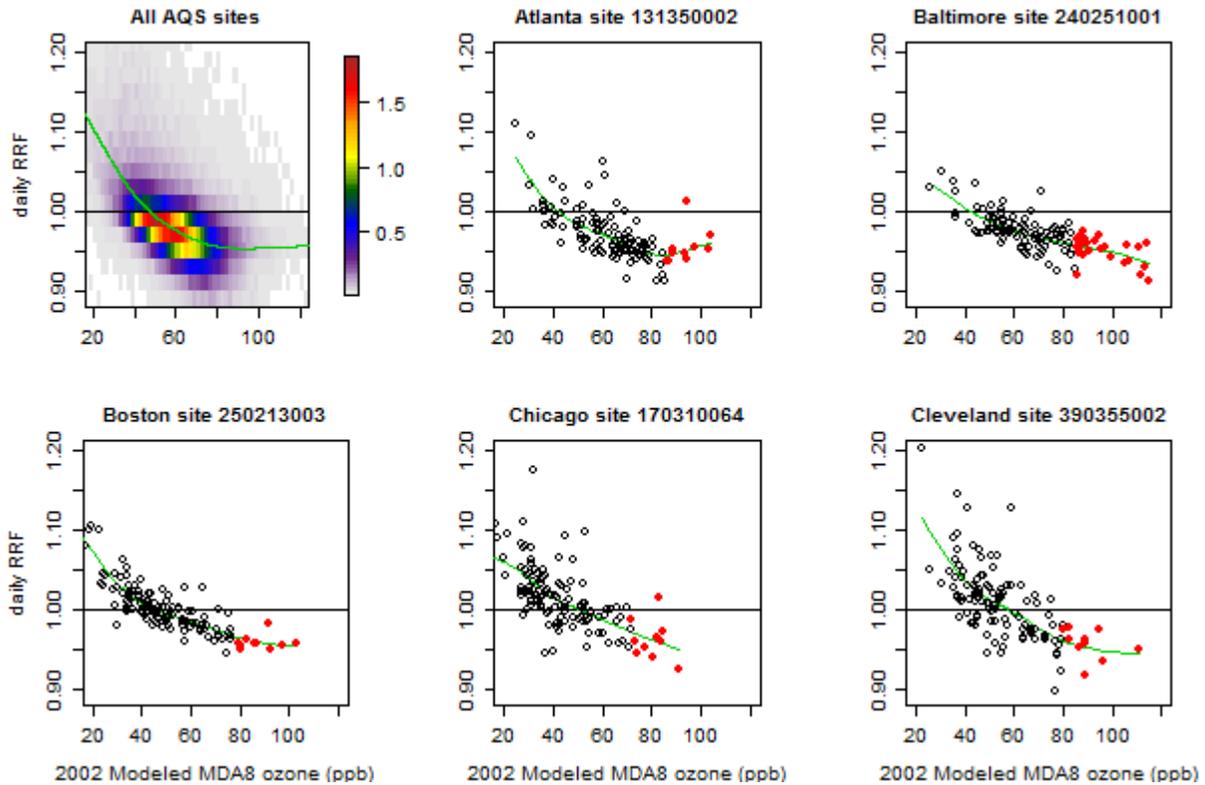
4

1 Table 2. Evaluation of different approaches for predicting the observed 2005 DVs at 619 AQS
 2 sites. The evaluation statistics include mean bias (MB), root mean square error (RMSE), R2, and
 3 accuracy for predicting attainment for the 75ppb and 84ppb NAAQS.
 4

	MB (ppb)	RMSE (ppb)	R²	Accuracy for 75ppb	Accuracy for 84ppb
Future Model	4.1	8.5	.37	.72	.77
TH85 RRF	3.8	6.6	.67	.79	.82
TH76 RRF	4.2	6.9	.66	.78	.80
TH71 RRF	4.2	6.9	.67	.78	.80
Top10 RRF	3.8	6.6	.67	.78	.82
Top20 RRF	4.1	6.8	.66	.78	.81

5
 6

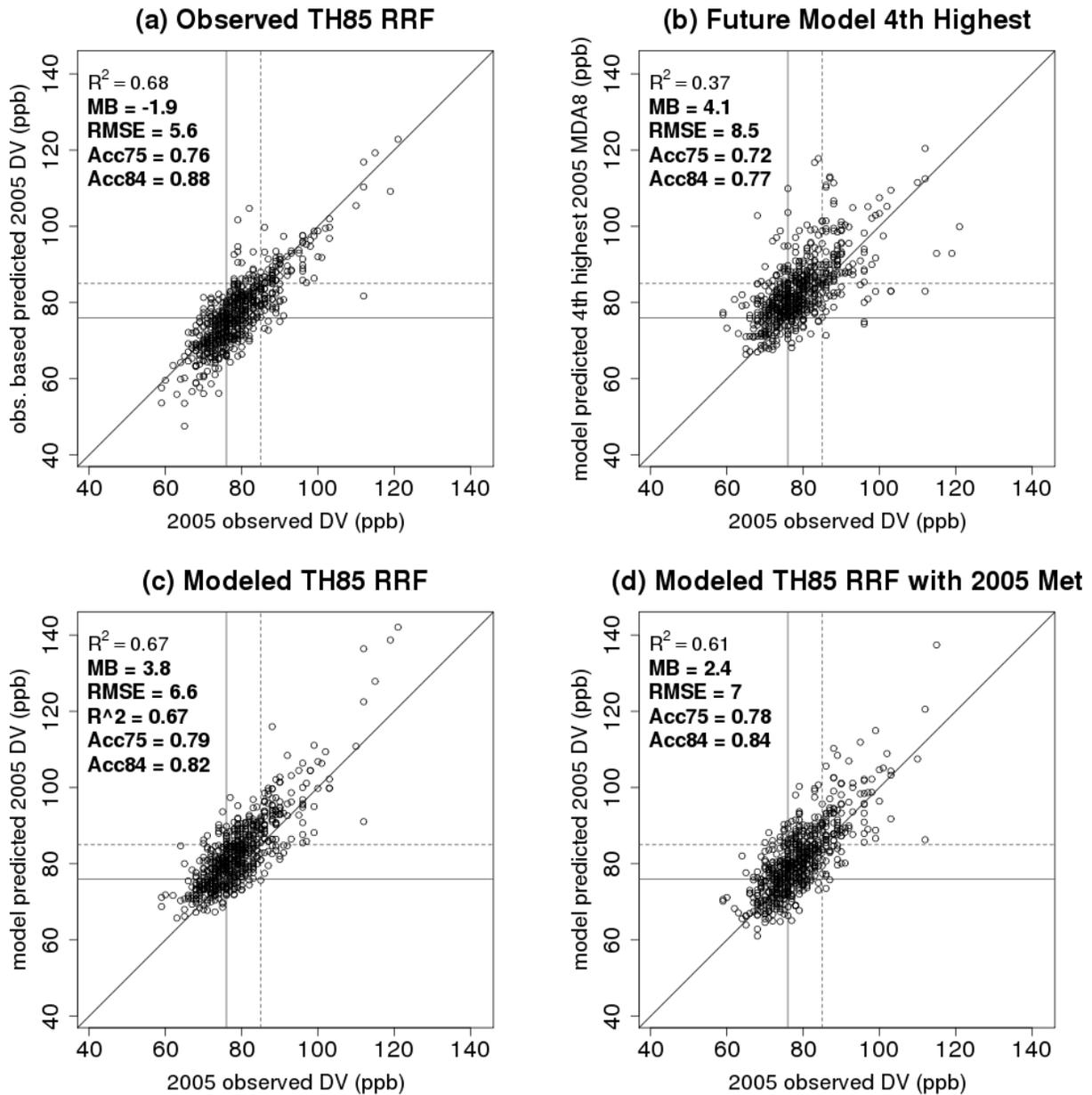
1



2
3
4
5
6
7
8
9
10
11
12
13

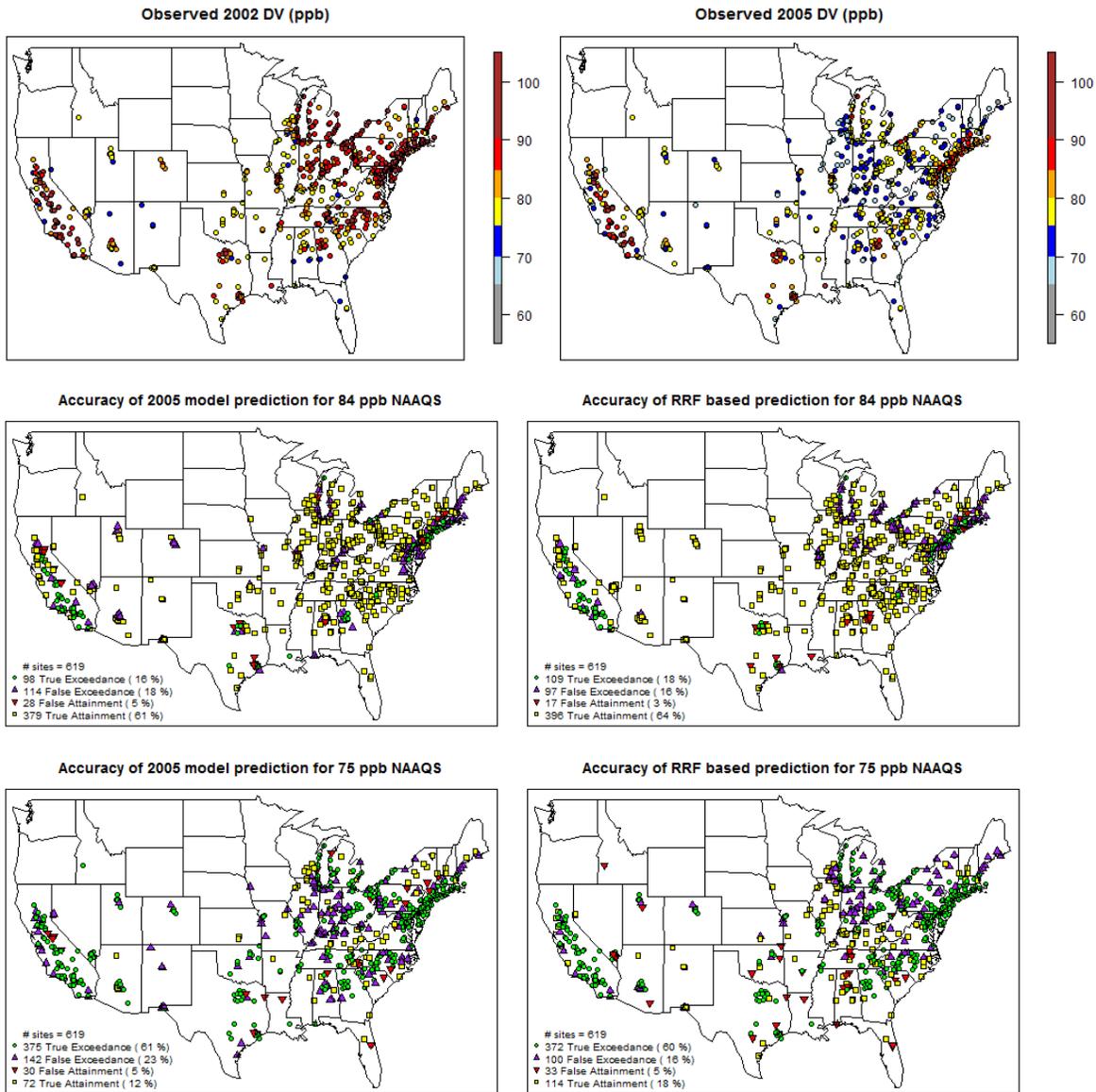
Figure 1. Daily RRF ratios versus 2002 modeled MDA8 ozone values at all AQS sites (top left) and at 5 select AQS sites in very urban areas. The color scale of the density scatter plot in the top right figure represents the percent of the data that fall within a particular pixel in the plot, e.g., the red pixels in the plot indicate areas where more than 1.5% of the data (>10 data points) would be clustered in a regular scatterplot. The green curve shows a spline fit to the daily RRF values. The red points in the scatter plot indicate what days were used in the calculation of the TH85 RRF metric for each site.

1
2



3
4 Figure 2. Comparison of 2005 observed versus predicted DV at n=619 AQS sites. Predicted DVs
5 are based on using the (a) observed TH85 RRF approach, (b) the 4th highest MDA8 ozone value
6 from the Sim05e02m simulation (c) the modeled TH85 approach and (d) the modeled TH85
7 approach modified by using 2005 meteorology in the “future” model simulation (i.e. output from
8 Sim05e05m rather than Sim05e02m is used in the numerator). Evaluation statistics include mean
9 bias (MB), root mean square error (RMSE), R2, and accuracy scores for predicting attainment of
10 the 75ppb and 84ppb NAAQS (Acc75, Acc84).

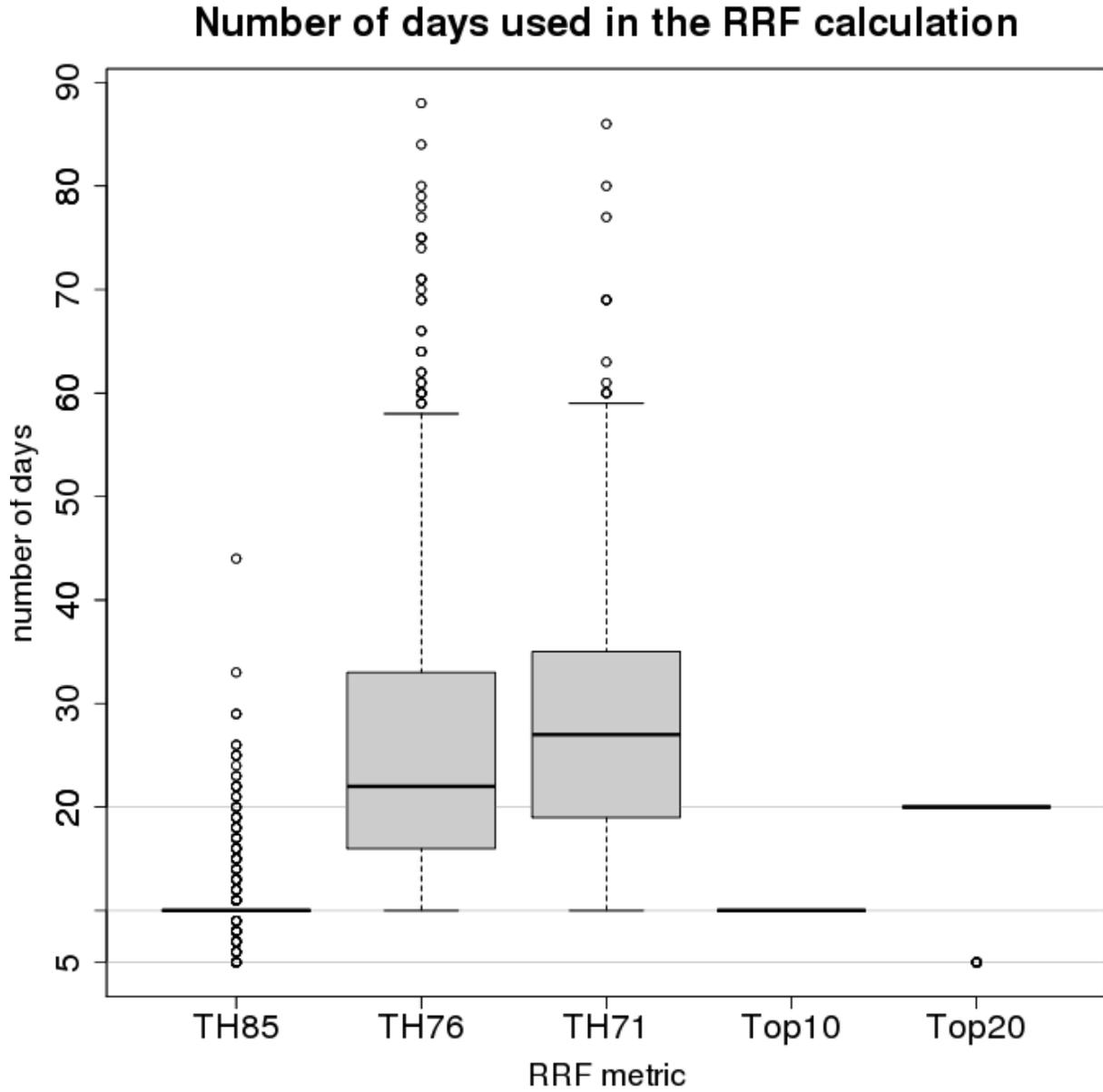
11
12



1
 2 Figure 3. Observed 2002 and 2005 design values (top row) across the entire model domain. The
 3 remaining plots show model predicted 2005 DVs in a categorical sense based on output from
 4 Sim05e02m (left column) and the TH85 RRF based approach (right column) for both the 84ppb
 5 standard (middle row) and the 75ppb standard (bottom row). Four categories are depicted: true
 6 exceedance (green circle); false exceedance (purple triangle); false attainment (red inverted
 7 triangle); true attainment (yellow square).

8
 9

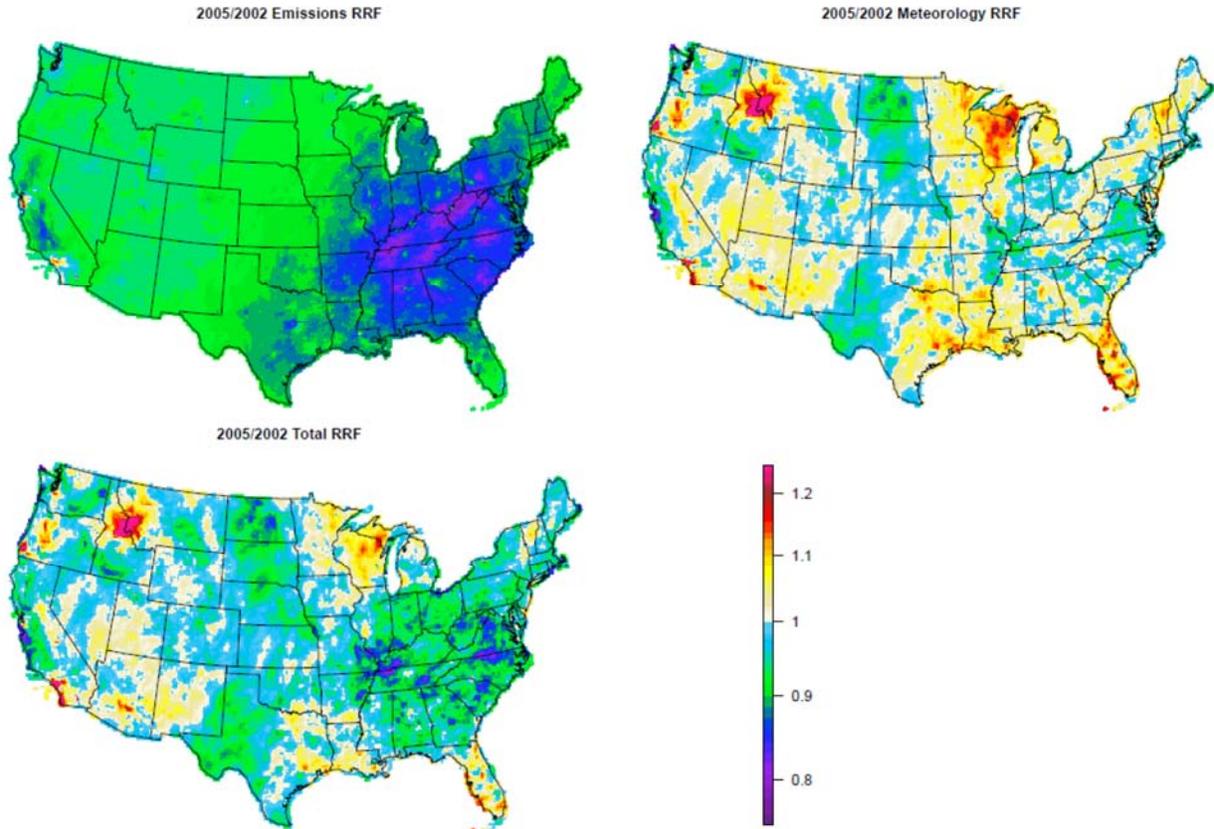
1
2



3
4
5
6
7

Figure 4. The number of days used in the calculation of each of the five RRF metrics across the 619 AQS sites used in Table 2.

1 **Supplemental Material**
2 **Figure S1.** Top10 RRF values for the entire US domain. The top left figure is based on the ratio
3 of Sim05e02m output to Sim02e02m output, representing the RRFs used in the attainment
4 evaluation in this paper. The top right figure shows the ratio of Sim02e05m output to
5 Sim02e02m. This figure represents the model-predicted change in ozone due to changes in
6 meteorology across these two years under 2002 emission levels. The final figure on the bottom
7 row is a “total” RRF based on the ratio of Sim05e05m to Sim02e02m output and shows the model
8 predicted change in ozone levels from 2002 to 2005 due to changes in both emissions and
9 meteorology.
10



11
12
13
14
15

1 **S2. Creation of hourly emissions for point sources without CEMS data (ptipm)**

2 Pseudo-CEMS data were created for point sources that only have annual total emissions available
3 in the NEI. For Sim02e02m, state-specific month-to-annual ratios were created by calculating
4 three-year averages of 2001-2003 CEMS data for each month and dividing these monthly averages
5 by the three year annual average for the state. The annual total ptipm emissions for each unit were
6 then allocated to month totals using these state-specific monthly ratios. To allocate the monthly
7 emissions to each day, state-specific day-to-month ratios were calculated using daily 2005 CEMS
8 data divided by the monthly average for 2005. These state-specific day-to-month factors were
9 then multiplied by the monthly total emissions for a given unit to calculate the total emissions on
10 each day for that unit. The resulting daily emissions were input into the SMOKE processing
11 system and hourly-to-daily allocation was performed using diurnal profiles. An analogous
12 calculation was made to estimate hourly 2005 emissions for the Sim05e05m simulation based on
13 2004-2006 CEMS data. This is the standard method used in regulatory applications for creating
14 simulations based on future emissions levels under current, base line meteorological conditions.
15 Emissions inputs for Sim02e05m used 2002 ptipm unit annual total emissions scaled with 2001-
16 2003 annual-to-month ratios to preserve the NO_x SIP call seasonal distribution and 2005 day-to-
17 month ratios to preserve the meteorological patterns of the meteorological year. Emissions inputs
18 for Sim05e02m used 2005 ptipm unit annual total emissions scaled with 2004-2006 annual-to-
19 month ratios and 2002 day-to-month ratios.

20
21
22

1 Table S1. Evaluation of different approaches for predicting the observed 2005 DVs at 388 AQS
 2 sites within NOx SIP Call States. The evaluation statistics include mean bias (MB), root mean
 3 square error (RMSE), R², and accuracy for predicting attainment for the 75ppb and 84ppb
 4 NAAQS.

	MB (ppb)	RMSE (ppb)	R²	Accuracy for 75ppb	Accuracy for 84ppb
Future Model	4.8	8.1	.35	.70	.80
TH85 RRF	4.6	6.4	.60	.78	.79
TH76 RRF	5.0	6.8	.59	.77	.77
TH71 RRF	5.1	6.8	.61	.77	.78
Top10 RRF	4.6	6.4	.60	.78	.79
Top20 RRF	4.9	6.6	.60	.77	.78

5
 6 Table S2. Evaluation of different approaches for predicting the observed 2005 DVs at 231 AQS
 7 sites that are not within NOx SIP Call States. The evaluation statistics include mean bias (MB),
 8 root mean square error (RMSE), R², and accuracy for predicting attainment for the 75ppb and
 9 84ppb NAAQS.

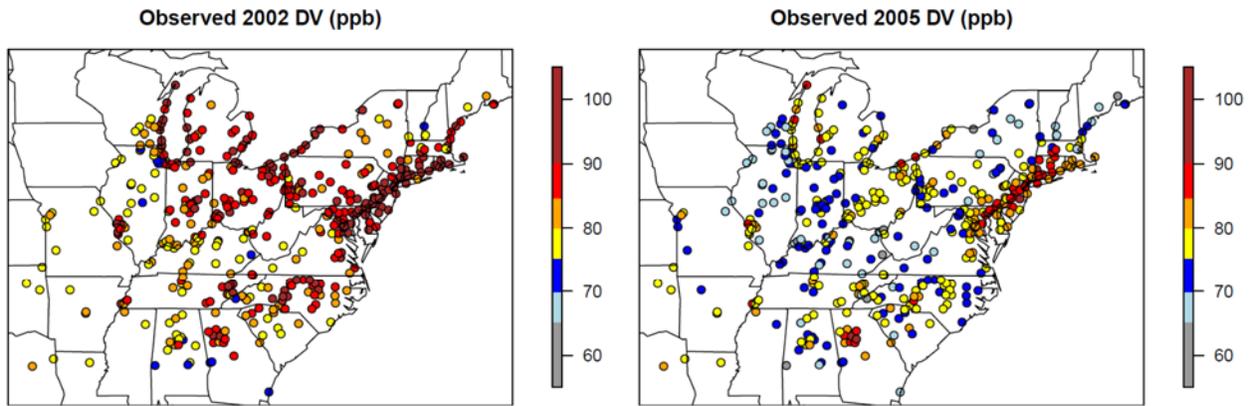
	MB (ppb)	RMSE (ppb)	R²	Accuracy for 75ppb	Accuracy for 84ppb
Future Model	3.0	9.3	.37	.77	.71
TH85 RRF	2.5	7.0	.71	.79	.86
TH76 RRF	2.7	7.1	.71	.79	.85
TH71 RRF	2.7	7.1	.71	.79	.85
Top10 RRF	2.5	6.9	.71	.79	.86
Top20 RRF	2.6	7.0	.71	.79	.85

10

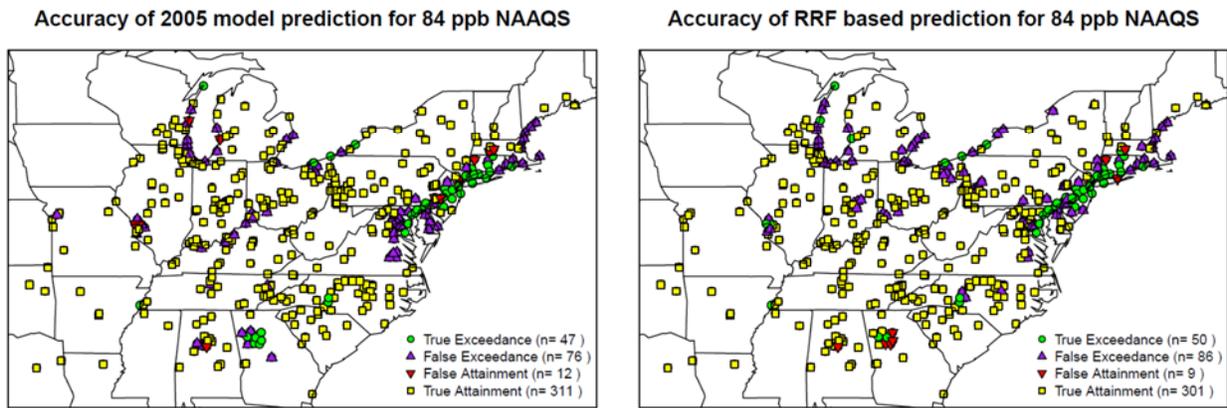
1
2
3
4
5
6
7
8

Figure S3. Observed 2002 and 2005 design values (top row) in the eastern US. The remaining plots show model predicted 2005 DVs in a categorical sense based on output from Sim05e02m (left column) and the TH85 RRF based approach (right column) for both the 84ppb standard (middle row) and the 75ppb standard (bottom row). Four categories are depicted: true exceedance (green circle); false exceedance (purple triangle); false attainment (red inverted triangle); true attainment (yellow square).

9



10



11
12
13

