



Error apportionment for atmospheric chemistry-transport models – a new approach to model evaluation

Efisio Solazzo and Stefano Galmarini

European Commission, Joint Research Centre, Institute for Environment and Sustainability, Air and Climate Unit, Ispra, Italy

Correspondence to: Stefano Galmarini (stefano.galmarini@jrc.ec.europa.eu)

Received: 7 January 2016 – Published in Atmos. Chem. Phys. Discuss.: 26 February 2016

Revised: 26 April 2016 – Accepted: 29 April 2016 – Published: 24 May 2016

Abstract. In this study, methods are proposed to diagnose the causes of errors in air quality (AQ) modelling systems. We investigate the deviation between modelled and observed time series of surface ozone through a revised formulation for breaking down the mean square error (MSE) into bias, variance and the minimum achievable MSE (mMSE). The bias measures the accuracy and implies the existence of systematic errors and poor representation of data complexity, the variance measures the precision and provides an estimate of the variability of the modelling results in relation to the observed data, and the mMSE reflects unsystematic errors and provides a measure of the associativity between the modelled and the observed fields through the correlation coefficient. Each of the error components is analysed independently and apportioned to resolved processes based on the corresponding timescale (long scale, synoptic, diurnal, and intra-day) and as a function of model complexity.

The apportionment of the error is applied to the AQMEII (Air Quality Model Evaluation International Initiative) group of models, which embrace the majority of regional AQ modelling systems currently used in Europe and North America.

The proposed technique has proven to be a compact estimator of the operational metrics commonly used for model evaluation (bias, variance, and correlation coefficient), and has the further benefit of apportioning the error to the originating timescale, thus allowing for a clearer diagnosis of the processes that caused the error.

1 Introduction

Due to their use for regulatory applications and to support legislation, air quality (AQ) models must model correctly and be correctly applied, justifying the need for a thorough evaluation. A framework for the operational and scientific evaluation of geophysical models was already envisaged in the early 1980s (Fox, 1981; Wilmott et al., 1985), the former being “a comparison with data exclusively within a particular application context”, and the latter defined as “some understanding of cause-and-effect relationship that relies on testing model components and extensively detailed data collection” (Fox, 1981). Thirty years later, as AQ models became more and more complex and their range of applicability widened, Dennis et al. (2010) further elaborated the concept of model evaluation by proposing a four-level evaluation, according to which different complementary aspects of the models should be tested, namely:

- operational: the level of agreement of model results with observations;
- dynamic: ability of the modelling system to respond to changes (in emissions, or in meteorological events);
- diagnostic: identify and attribute the source of the error to the relevant process;
- probabilistic: confidence and uncertainty levels of the modelled results.

In the framework originally designed by Dennis et al. (2010), the diagnostic component plays a central role. It (i) answers the fundamental issue left open by the operational screening, in other words whether the model provides the right answer for the right reason, (ii) provides feedback to

developers to help make model improvements, and (iii) sets the basis for the probabilistic evaluation (Fig. 1 of Dennis et al., 2010).

Over the years, and despite the increasing relevance of modelling systems for AQ applications, model evaluation continues to rely almost exclusively on operational evaluation, which basically involves gauging the model's performance using distance, variability and associativity metrics. This common practice has little or no impact on model improvement, as it does not target the source of the modelling error and does not discriminate between the reasons for appropriate or inappropriate performance.

Such a requirement is even more pressing these days, with current state-of-the-science AQ modelling systems accounting for an increasing number of coupled physical processes and being described using hundreds of modules, which are the result of decades of targeted and, generally, independent investigations. Furthermore, AQ modelling systems typically depend on external sources for the inputs of meteorology and emissions data, as well as for boundary conditions. These fields are generally produced by other models (which, in turn, depend on external sources for initial and/or boundary conditions) and, after substantial processing, are used by the AQ modelling systems with no guarantee of being unbiased and/or accurate. The bias introduced by these inputs, along with the uncertainty associated with model error, the linearisation of non-linear processes and omitted and unresolved variables and processes, all contribute to the model error. The extensive use of AQ models for AQ assessment and planning is equally important, and requires a good knowledge of the model capabilities and deficiencies that would allow for a more educated use of the modelling systems and their results.

Recently, the AQMEII (Air Quality Model Evaluation International Initiative) activity (Rao et al., 2011) applied the approach proposed by Dennis et al. (2010), by organising model evaluation activities (AQMEII 1, 2 and 3) using operational (Solazzo et al., 2012a, b, 2013a; Im et al., 2015a, b), probabilistic (Solazzo et al., 2013b; Kioutsoukakis and Galmarini, 2014) and diagnostic (Hogrefe et al., 2014; Makar et al., 2015) evaluation frameworks.

The study we present here follows and complements the previous investigations based on the AQMEII models collected in the first and second phases of the activity (AQMEII1 and AQMEII2). The main aim is to introduce a novel method that combines operational and diagnostic evaluations. This method helps apportion the model error to its components, thereby identifying the space/timescale at which it is most relevant and, when possible, to infer which process/es could have generated it. This work is designed to support the analysis of the currently ongoing third phase of the AQMEII activity (Galmarini et al., 2015).

2 Mean square error as a comprehensive metric

For the model evaluation strategy proposed, we start by breaking down the mean square error (MSE) (used here as unique metric to evaluate model performance) into the sum of the variance (and covariance) and the squared bias. The error and its components are then calculated on the spectrally decomposed time series of modelled and observed hourly ozone mixing ratios. The advantage of this evaluation strategy is 2-fold:

- With respect to a conventional operational evaluation, the new method allows for a more detailed assessment of the distance between model results and observations given the breakdown of the error into bias, variance and covariance and their associated interpretations.
- Decomposing the MSE into spectral signals allows for the precise identification of where each portion of the model error predominantly occurs. Given that specific processes are associated with specific scales, the apportionment of the error components to their relevant scales helps to more precisely identify which processes described in the model could be responsible for the error. Information about the nature of the error and the class of process can significantly help modellers and developers to improve model performance.

The data used are produced by the modelling communities participating in AQMEII1 and AQMEII2 over the European (EU) and North American (NA) continental-scale domains for the years 2006 (AQMEII1) and 2010 (AQMEII2).

2.1 Error decomposition

The MSE is the squared difference of the modelled (mod) and observed (obs) values:

$$\text{MSE} = E(\text{mod-obs})^2 = \frac{\sum_{i=1}^{n_t} (\text{mod}_i - \text{obs}_i)^2}{n_t}, \quad (1)$$

where $E(\cdot)$ denotes expectation and n_t is the length of the time series. The bias is

$$\text{bias} = E(\text{mod-obs}) \quad (2)$$

i.e. $\text{bias} = \overline{\text{mod}} - \overline{\text{obs}}$. Thus, the following relationship holds:

$$\text{MSE} = \text{var}(\text{mod-obs}) + \text{bias}^2, \quad (3)$$

which is a well-known property of the MSE, ($\text{var}(\cdot)$ is the variance operator). By using the property of the variance for correlated fields:

$$\text{var}(\text{mod-obs}) = \text{var}(\text{mod}) + \text{var}(\text{obs}) - 2\text{cov}(\text{mod}, \text{obs}), \quad (4)$$

the final formulation for the MSE components reads as follows:

$$\text{MSE} = \text{bias}^2 + \text{var}(\text{mod}) + \text{var}(\text{obs}) - 2\text{cov}(\text{mod}, \text{obs}), \quad (5)$$

where the covariance term (last term on the right-hand side of Eq. 5) accounts for the degree of correlation between the modelled and observed time series. When the covariance term is zero, $\text{var}(\text{obs})$ is referred to as the *incompressible part of the error* and represents the lowest limit that the MSE of the model can achieve. When dealing with model evaluation, the modelled and observed time series are typically highly correlated and therefore, within the limits of the perfect match (correlation coefficient of unity), $\text{cov}(\text{mod}, \text{obs}) = \text{cov}(\text{obs}, \text{obs}) = \text{cov}(\text{mod}, \text{mod}) = \text{var}(\text{mod}) = \text{var}(\text{obs})$ and the MSE can be reduced to only the bias term. That implies that the development of a high-quality model needs to ensure

- the highest possible precision in order to maximise the $\text{cov}(\text{mod}, \text{obs})$ term;
- the highest possible accuracy, in order to minimise the bias.

Elaborating on Eq. (5), Theil (1961) derived the following:

$$\text{MSE} = (\overline{\text{mod}} - \overline{\text{obs}})^2 + (\sigma_{\text{mod}} - \sigma_{\text{obs}})^2 + 2(1-r)\sigma_{\text{mod}}\sigma_{\text{obs}}. \quad (6)$$

In Eq. (6), the variance term is expressed as the difference between the standard deviation of the model and that of the observations, and the covariance term (last term on the right) includes r , the coefficient of correlation between the observed and modelled time series. The ratios of the three terms on the right-hand side of Eq. (6) to the overall MSE are known as *Theil's coefficients* (Pindick and Rubinfeld, 1998). Murphy (1988) provided examples of the scores that can be developed using the components of the MSE.

The bias measures the departure of the modelled from the observed results, and is a measure of systematic error, since it measures the extent to which the average modelled values deviate from the observed ones. The bias is commonly used to express the degree of “trueness”, i.e. “the closeness of agreement between the average value obtained from a large series of measurements and the true value” (Johnson, 2008). The variance shows whether the modelled variability is compatible with that observed. Finally, the covariance term represents the unexplained proportion of the MSE due to the remaining unsystematic errors; i.e. it represents the remaining error after deviations from the mean values have been accounted for. This latter term is a measure of the lack of correlation of the model with comparable observations, and is considered the least “worrisome” portion of the error (Pindick and Rubinfeld, 1998).

Aiming at minimising the MSE, the only controlled variables in Eq. (6) are mod and σ_{mod} , and differentiating with respect to them yields the conditions that minimise the MSE:

$$\begin{cases} \frac{\partial \text{MSE}}{\partial \overline{\text{mod}}} = 2(\overline{\text{mod}} - \overline{\text{obs}}) = 0 \\ \frac{\partial \text{MSE}}{\partial \sigma_{\text{mod}}} = 2(\sigma_{\text{m}} - \sigma_{\text{obs}}) + 2(1-r)\sigma_{\text{obs}} = 0 \end{cases}$$

i.e. the best agreement between modelled and observed values is achieved by

$$\begin{cases} \overline{\text{mod}} = \overline{\text{obs}} \\ \sigma_{\text{m}} = r\sigma_{\text{obs}} \end{cases}, \quad (7)$$

which analytically corresponds to the aforementioned items (a) and (b). By inserting Eq. (7) into Eq. (6), the minimum achievable MSE (mMSE) is

$$\text{mMSE} = \sigma_{\text{obs}}^2(1-r^2), \quad (8)$$

which is the unexplained portion of the error, as it reflects the share of observed variance that is not explained by the model (r^2 is the coefficient of determination). The presence of an unexplained part of the error suggests a modification of the MSE decomposition in Eq. (6) in such a way as to explicitly include mMSE:

$$\text{MSE} = (\overline{\text{mod}} - \overline{\text{obs}})^2 + (\sigma_{\text{mod}} - r\sigma_{\text{obs}})^2 + \text{mMSE}. \quad (9)$$

The decompositions in Eqs. (5), (6) and (9) contain all the relevant operational metrics usually applied to score modelling systems (bias, variance, correlation coefficient), and therefore prove to be a compact estimator of accuracy (bias), precision (variance) and associativity (unexplained portion through the correlation coefficient). Eq. (9) has been explicitly derived in this study to help evaluate AQ models.

Ideally, the entire error should be attributable to unsystematic fluctuations. From a model development perspective, the variance and covariance are possibly more revealing of model deficiencies than is the bias term, as they are produced by the AQ model itself, while the bias is also due to external sources (e.g. emissions, boundary conditions). From the application viewpoint, however, it is the overall error that counts, which is mostly made up of the bias.

2.2 Spectral decomposition of modelled and observed time series

Hourly time series of (modelled and observed) ozone concentrations have been decomposed using an iterative moving average approach known as the Kolmogorov–Zurbenko (kz) low-pass filter (Zurbenko, 1986), whose applications to ozone are vastly documented in the literature (Rao et al., 1997; Wise and Comrie, 2005; Hogrefe et al., 2000, 2014; Galmarini et al., 2013; Kang et al., 2013; Solazzo and Galmarini, 2015). The kz filter depends on two parameters: the length of the moving average window m and the number of

iterations $k(kz_{m,k})$. Since the kz is a low-pass filter, the filtered time series consists of the low-frequency fluctuating component, while the difference between two filtered time series provides a band-pass filter. This latter property is used to decompose the ozone concentration time series as

$$O_3 = LT(O_3) + SY(O_3) + DU(O_3) + ID(O_3), \quad (10)$$

where LT is the long-term component (periods longer than 21 days), SY is the synoptic component (weather processes that last between 2.5 and 21 days), DU is the diurnal component (day/night alternation period between 0.5 and 2.5 days) and ID is the intra-day component accounting for fast-acting processes (less than 12 h). The decomposition presented in Eq. (10) is such that the original time series is perfectly returned by the summation of the components (see Appendix A for details). Dealing with 1 year of data, any filter longer than the LT component would not be meaningful. The periods of the components correspond to well-defined peaks in the power spectrum of ozone, e.g. as detailed in Rao et al. (1997) and Hogrefe et al. (2000).

The LT component is the baseline and incorporates the bias of the original (un-decomposed time series). The other components (SY , DU and ID) are zero-mean fluctuations around the LT time series and are therefore unbiased. The band-pass nature of the SY , DU and ID components is such that they only account for the processes occurring in the time window the filter allows the signal to “pass”. For instance, the DU component is insensitive to processes outside the range of 0.5 to 2.5 days.

Further properties of the spectrally decomposed ozone time series of AQMEII derived by Galmarini et al. (2013), Hogrefe et al. (2014) and Solazzo and Galmarini (2015) are as follows:

- The DU component accounts for more than half of the total variance, followed by the LT and SY components.
- The ID component has the smallest influence due to the small amplitude of its fluctuations.
- The variance of the spectral component is neither strongly nor systematically associated with the area-type of the monitoring stations (i.e. rural, urban, sub-urban).
- Due to the bias, most of the error is accounted for by the LT component, followed by the DU component. The ID contributes very little to the overall MSE .

Further important technicalities of the spectral decomposition, including a method to estimate the contribution of the spectral cross-components (the overlapping regions of the power spectrum) to the total error, are reported in the Appendix A.

The signal decomposition of Eq. (10) is applied to the full-year time series. However, to evaluate the model performance with regard to ozone, the analysis is restricted to

the months of May to September, i.e. when the production of ozone due to photochemistry is most relevant.

3 Data and models used

The observational data set derived from the surface AQ monitoring networks operating in the EU and NA constitutes the same data set used in the first and second phases of AQMEII to support model evaluation. Only stations with over 75 % valid records for the whole periods and located at altitudes below 1000 m have been used for this analysis. Details of the modelled regions and number of receptor stations are reported in Table 1.

Since the main scope of this study is to introduce the error apportionment methodology (rather than to strictly evaluate the models), the analysis is presented for continental areas for convenience and easier display of the results. However, given the size of the domains and the heterogeneity of climatic and emission conditions, dedicated analyses for three sub-regions in both continents are proposed in the Supplement (Figs. S1 to S3).

There are profound differences between the modelling systems that participated in AQMEII1 and AQMEII2. The two sets of models have been applied to different years (2006 for phase 1 and 2010 for phase 2) and are therefore dissimilar with respect to the input data of emissions and boundary conditions for chemistry. The AQ models of the second phase are coupled (online chemistry feedbacks on meteorology), while those of the first phase are not. The effect of using online models for simulating ozone accounts for the impact of aerosols on radiation and therefore on temperature and photolysis rates (Baklanov et al., 2014).

The model settings and input data for phase I are described in Solazzo et al. (2012a, b, 2013a), Schere et al. (2012) and Pouliot et al. (2012); for phase II, similar information is presented in Im et al. (2015a, b), Brunner et al. (2015) and Pouliot et al. (2015).

Table 2 summarises the features of the modelling systems analysed in this study with regard to ozone concentrations in the EU or NA. The modelling contribution to the two phases of AQMEII consists of 12 and 9 models and of 8 and 3 models for EU and NA, respectively. Solazzo et al. (2012a, 2013b) showed the existence of a subset of models, whose ensemble mean, MSE_{best} , optimises the accuracy (minimum error over all possible ensemble mean combinations). The set and number of models composing MSE_{best} varies by pollutant and, for the same pollutant, by the examined period (year, season, etc.). In this study MSE_{best} is identified for the continent-wide-averaged time series of ozone concentration and for the un-decomposed ozone time series (i.e. not spectrally decomposed) during the period May–September. There are circumstances where a single model outperforms any combination of models. In such cases the MSE_{best} is identified with the best single model.

Table 1. Features of the modelled domains

	Europe		North America	
	phase 1	phase 2	phase 1	phase 2
Simulated year	2006	2010	2006	2010
Extension	$(-10,39)^{\circ}$ W; $(30,65)^{\circ}$ N		$(-125,-55)^{\circ}$ W; $(26,51)^{\circ}$ N	
Number of receptors (min validity = 75 %; max altitude = 1000 m)	1339	1360	672	652

Detailed analysis of the main differences in emissions, boundary conditions, and meteorology between the modelled years of 2006 (AQMEII1) and 2010 (AQMEII2) is presented in Stoeckenius et al. (2015). A summary of the performance of the two suites of model runs is provided in Makar et al. (2015), showing that the AQMEII1 models generally performed better than the AQMEII2 models, based on standard operational metrics. However, the use of standard evaluation methods does not allow for the assessment of whether the feedback processes have an effect on the deterioration of model performance, or rather the different sets of emissions and boundary conditions. We try to assess the problem using the error apportionment method outlined above.

4 Results for the spatially averaged time series

4.1 MSE of spectral components

Figure 1 reports the MSE share of the spectral components and cross-components for each model, for both phases of AQMEII, derived from the ozone time series spatially averaged over each continental area. The spatial average is carried out prior to the spectral decomposition and error apportionment.

The LT share of the total MSE is the largest in absolute value for both continents and both simulated years. The LT share ranges between 9.9 % (GEM-AQ, AQMEII1, NA) and 86.7 % (WRF/Chem, AQMEII1, NA), and averages at ~ 34 and ~ 46.5 % for the EU and ~ 50.6 and ~ 47 % for NA (AQMEII1 and AQMEII2, respectively).

The second largest share of the total MSE is of the DU component, accounting for ~ 20 % (all cases), followed by the SY component. Depending on the model, the MSE share of the remaining spectral components and cross-components varies significantly. Being the intermediate timescales, the overlap of the DU and SY components is likely to be more significant than the overlap of the LT and ID scales. The contribution of DU_{cc} and SY_{cc} to the total error can be as high as 17 % (DU_{cc} for GEM-AQ, AQMEII1, NA) and 16 % (SY_{cc} for MM5-CAMx, AQMEII1, EU). Overall, the DU_{cc} terms (interaction of DU with the neighbouring SY and ID scales) are significant in both continents (~ 10 %), while the share of

the SY component and cross-components is more significant in the EU.

The ID component has a little impact on the total MSE (negligible in some instances), exceeding the 3 % share only for the two EU instances of the L.-Euros model.

The results of Fig. 1 help identify the timescales and associated processes for which the largest improvement in model accuracy can be achieved. The LT component has the largest share of the error due to the bias (error breakdown is discussed in the next section), but “internal” chemical processes, transport and deposition also occur at this timescale. Diurnal processes are the second largest source of error, including, among others, chemistry, boundary layer dynamics, radiation forcing and their interactions. The processes in the SY band bridge meteorological and chemical processes, and discern between the fast-acting diurnal processes and the baseline. As such, although the SY signal is not as strong as that of the DU components (variance of SY is comparable to the variance of ID; see Hogrefe et al., 2014), it accounts for a significant portion of the total error, as discussed next.

4.2 The quality of the error: error apportionment

The error breakdown (Eq. 9) of each spectral component complements the analysis presented in the previous section, and is reported in Fig. 2 (please note that results in Fig. 2 are reported in ppb^2 for reason of clarity). The bias (only included in the LT component) is the average amount by which the modelled time series is displaced with respect to the observed time series, and is the main source of error. The bias can be either due to “internal” model errors, or inherited from external drivers (emissions, meteorology, boundary conditions). Based on the experience matured within AQMEII, while the internal model errors are of interest for model development because they are generated by systematic modelling errors, the bias introduced by external drivers is responsible for the largest share of modelling errors.

From the continental average error breakdown of Fig. 2 we can conclude that the majority of EU models (in both AQMEII phases) have small bias (continental-wide average), with the important exceptions of CCLM-CMAQ and Muscat models in AQMEII1, and CMAQ in AQMEII2, which introduced large positive biases. The bias for the NA continent is

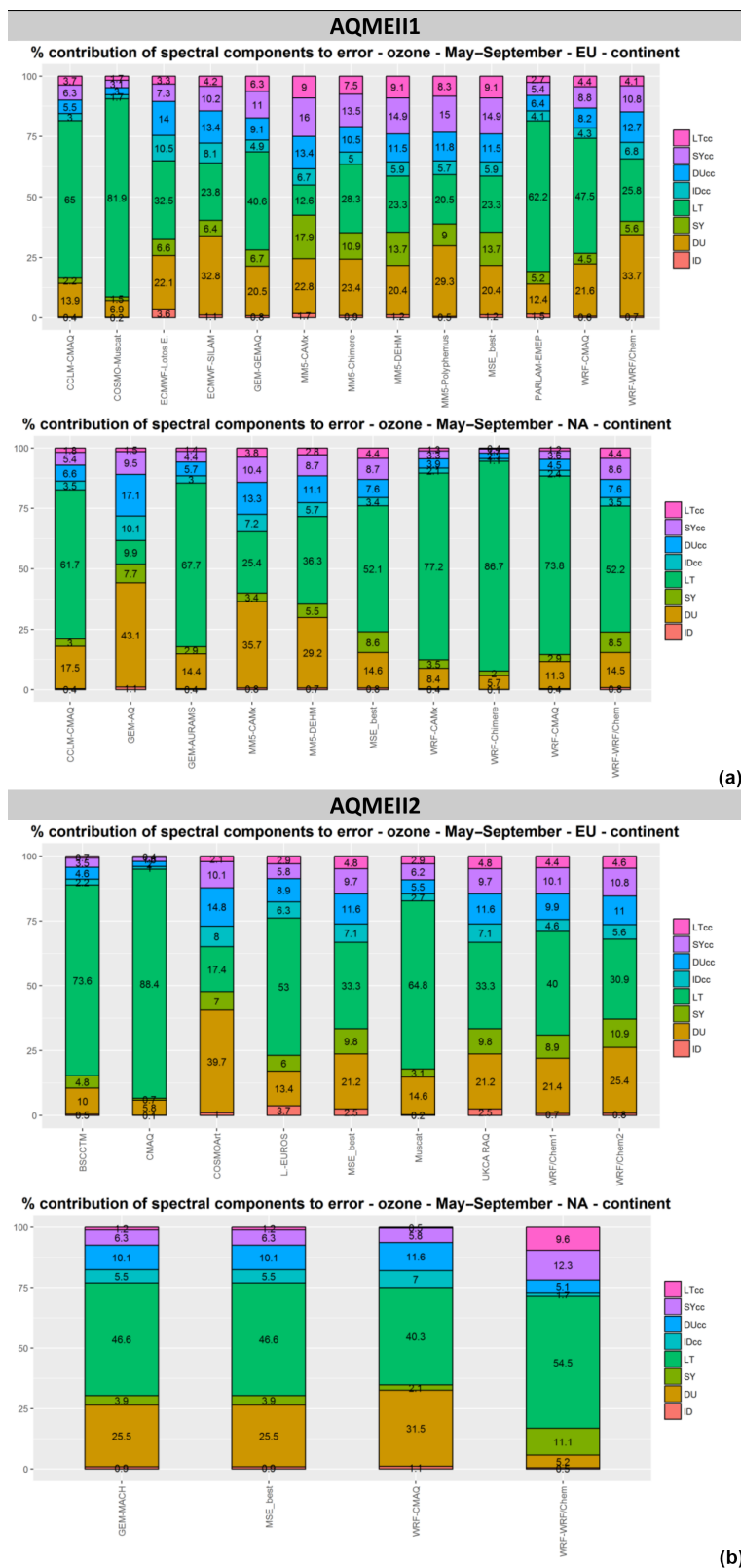


Figure 1. Share (in %) of the total MSE in the main spectral components and the cross-components (see Appendix for detail) for (a) AQMEII1 and (b) AQMEII2. Top panel: EU; lower panel: NA.

Table 2. Modelling systems participating in the first (a) and second (b) phases of AQMEII for Europe and North America.

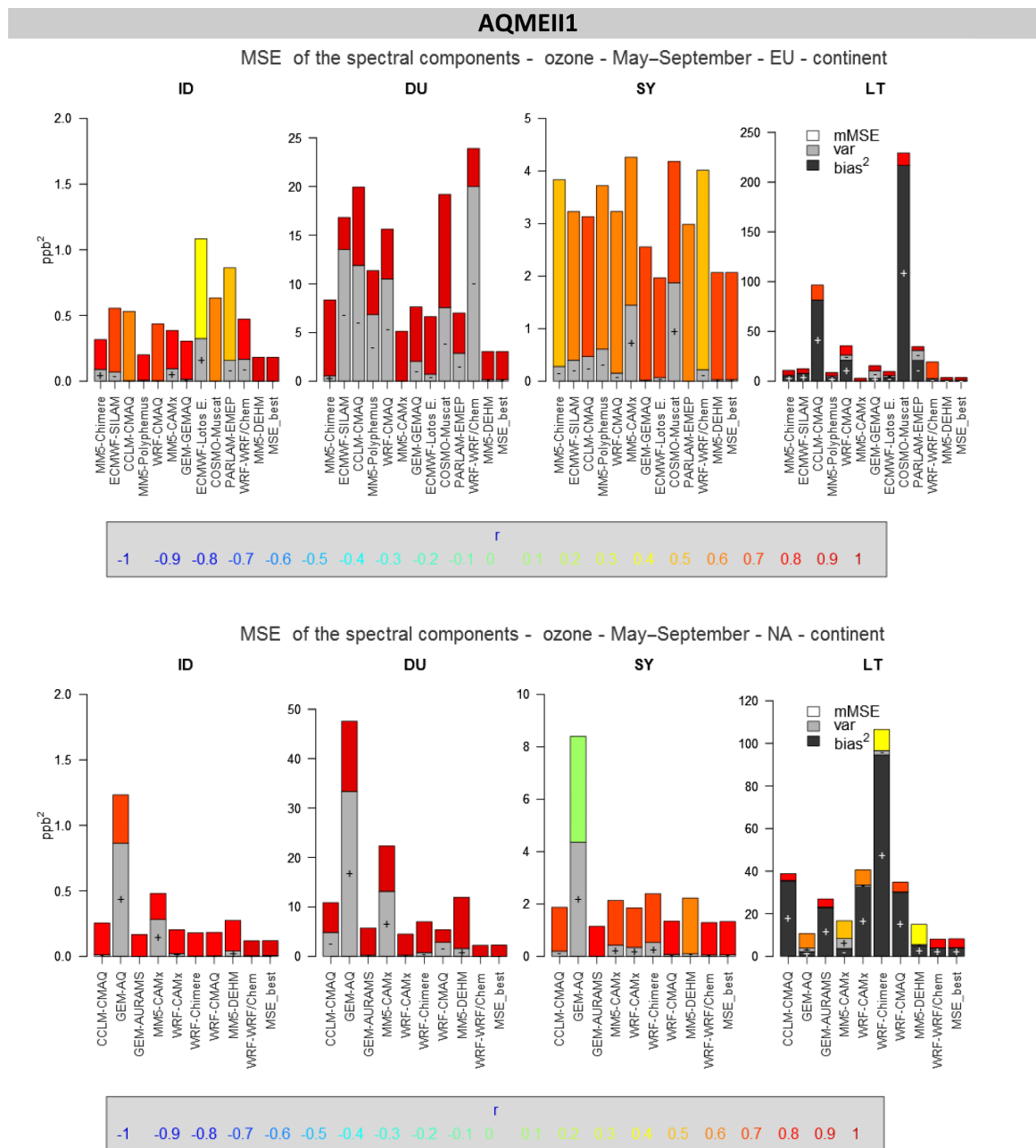
(a) Model			Grid (km)	Emissions	Chemical BC
Code	Met	AQ			
EUROPE – AQMEII 1					
DK1	MM5	DEHM	50	Global emission databases, EMEP	Satellite measurements
FR3	MM5	Polyphemus	24	Standard ^a	Standard
HR1	PARLAM-PS	EMEP	50	EMEP model	From ECMWF and forecasts
UK2	WRF	CMAQ	18	Standard ^a	Standard
US4	WRF	WRF/Chem	22.5	Standard ^a	Standard
FI1	ECMWF	SILAM	24	Standard anthropogenic; In-house biogenic	Standard
FR4	MM5	Chimere	25	MEGAN, Standard	Standard
PL1	GEM	GEM-AQ	25	Standard over AQMEII region; Global EDGAR/GEIA over the rest of the global domain	Global variable grid set-up (no boundary conditions)
NL1	ECMWF	Lotos-EUROS	25	Standard ^a	Standard
DE1	COSMO	Muscat	24	Standard ^a	Standard
US3	MM5	CAMx	15	MEGAN, Standard	Standard
DE3	COSMO-CLM	CMAQ	24	Standard ^a	Standard
NORTH AMERICA – AQMEII 1					
CA1	GEM	AURAMS	45	Standard ^b	Climatology
PL1	GEM	GEM-AQ	25	Standard over AQMEII region; Global EDGAR/GEIA over the rest of the global domain	Global variable grid set-up (no boundary conditions)
PT1	MM5	CAMx	24	Standard	LMDZ-INCA
US1	WRF	CAMQ	12	Standard	Standard
US3	WRF	CAMx	12	Standard	Standard
FR4b	WRF	CHIMERE			
DK1	MM5	DEHM	50	Global emission databases, EMEP	Satellite measurements
DE3	COSMO-CLM	CMAQ	24	Standard ^a	Standard
ES3	WRF	WRF/Chem	23	Standard	Standard
(b) Model			Grid	Emissions	Chemical BC
Code	Met	AQ			
EUROPE – AQMEII 2					
AT1	WRF	WRF/Chem	23 km	Standard	Standard
CH1	COSMO	Cosmo-ART	0.22°	Standard	Standard
ES2a	NMMB	BSCCTM	0.20°	Standard	Standard
ES3	WRF	WRF/Chem	23 km	Standard	Standard
NL2	RACMO	LOTOS-EUROS	0.5° × 0.25°	Standard	Standard
UK5	WRF	CMAQ	18 km	Standard	Standard
UK4	MetUM	UKCA RAQ	0.22°	Standard	Standard
DE3	COSMO	Muscat	0.25°	Standard	Standard
NORTH AMERICA – AQMEII 2					
ES1	WRF	WRF/Chem	36 km	Standard	Standard
US6	WRF	CMAQ	12 km	Standard	Standard
CA2f	GEM	MACH	15 km	Standard	Standard

Footnotes for (a): ^a standard anthropogenic emissions and biogenic emissions derived from meteorology (temperature and solar radiation) and land use distribution implemented in the meteorological driver.

^b Standard anthropogenic inventory but independent emission processing, exclusion of wildfires, and different versions of BEIS(v3.09) used.

Refer to Solazzo et al. (2012a, b) and references therein for details.

Footnotes for (b): standard boundary conditions: 3-D daily chemical boundary conditions were provided by the ECMWF IFS-MOZART model run in the context of the MACC-II project (Monitoring Atmospheric Composition and Climate – Interim Implementation) at 3-hourly and 1.125 spatial resolution. Refer to Im et al. (2015a, b) for details. Standard emissions: based on the TNO-MACC-II (Netherlands Organization for Applied Scientific Research, Monitoring Atmospheric Composition and Climate - Interim Implementation) framework for Europe and by the US EPA (Environmental Protection Agency) and Environment Canada for North America. The 2008 National Emissions Inventory (<https://www.epa.gov/air-emissions-inventories>) and the 2008 Emissions Modeling Platform (<https://www.epa.gov/air-emissions-modeling/20072008-version-5-air-emissions-modeling-platforms>) with year-specific updates for 2006 and 2010 were used for the US portion of the modelling domain. Canadian emissions were derived from the Canadian National Pollutant Release Inventory (<http://www.ec.gc.ca/inrp-npri/>) and Air Pollutant Emissions Inventory (<http://www.ec.gc.ca/inrp-npri/donnees-data/ap/index.cfm?lang=En>) values for the year 2006. Refer to Im et al. (2015a, b) for details.



(a)

Figure 2.

more uniformly distributed across the models (model over-prediction in both AQMEII phases), possibly indicating a common source of (external) bias in the NA models. The bias introduced by external fields is reflected by the bias of the baseline component (LT). For the period between May and September, the error in modelled ozone due to the boundary condition is typically small (Solazzo et al., 2012; Im et al., 2015; Giordano et al., 2015; Hogrefe et al., 2014), while the emissions of ozone precursors and VOCs (volatile organic

compounds) are problematic, especially in the EU (Makar et al., 2015; Brunner et al., 2015). We further notice that the absence of bias in some models may be caused by the presence of compensating bias, i.e. spatially distributed biases of opposite signs. The spatial distribution of the MSE is discussed in the next section. In all cases, the MSE_{best} model is, by definition, the model with the lowest MSE and thus the one with the smallest LT bias.

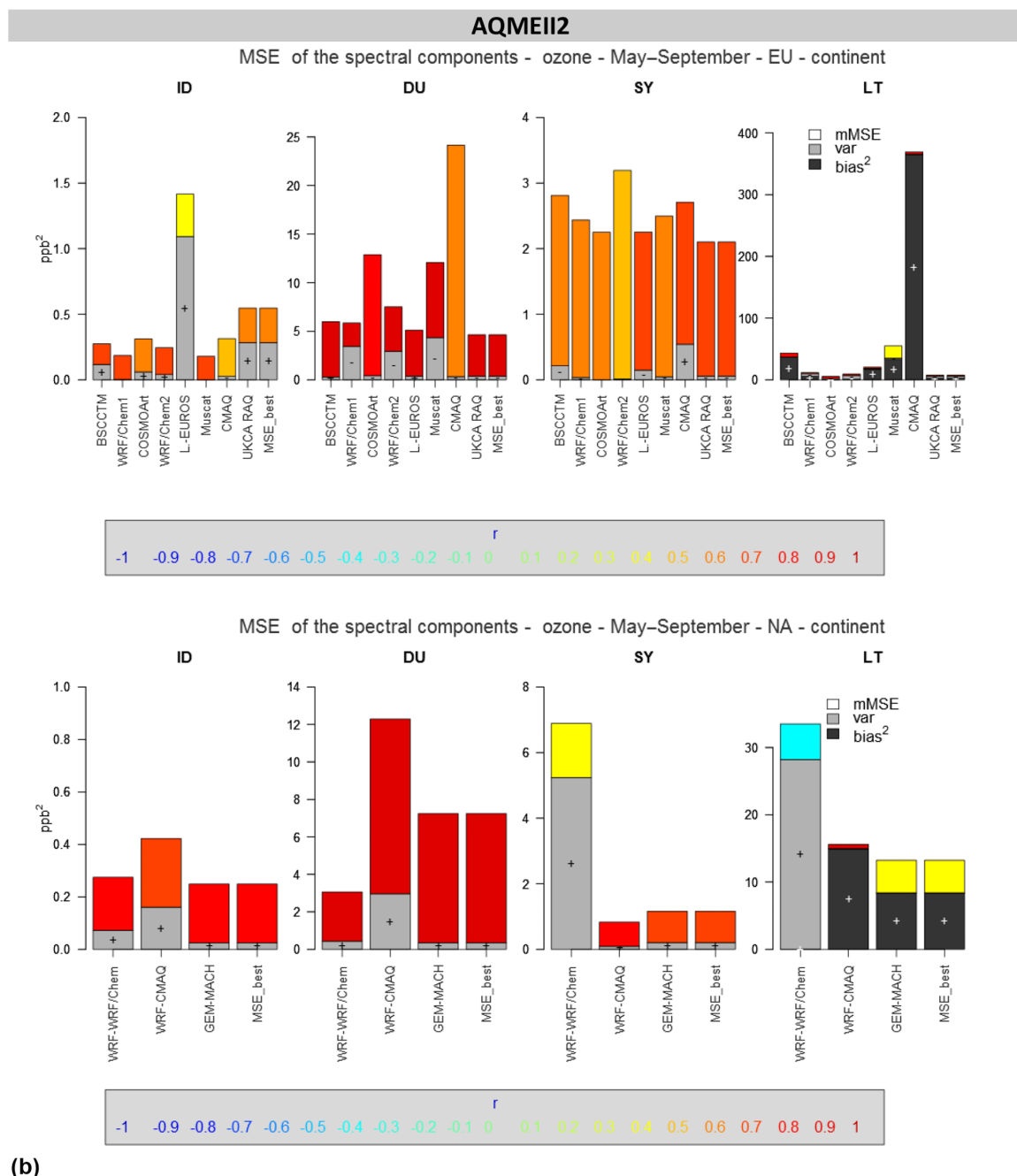


Figure 2. MSE (ppb^2) breakdown in bias squared, variance and mMSE of the spectral components ID, DU, SY, LT, based on Eq. (9). The bias is entirely accounted for by the LT component. The sign within the share of bias and variance indicates model overestimation (+) or underestimation (–) of mean concentration (bias) and variance. The colour of the mMSE share of the error is coded based on the values of r , the correlation coefficient, according to the colour scale at the bottom of each plot. (a) AQMEII1 and (b) AQMEII2. Top panel: EU; lower panel: NA.

The variance share of LT error is generally small (~ 1 – 2.5 ppb). This is not entirely unexpected, as the LT component has a high signal-to-noise ratio with a well-structured seasonal cycle, peaking in summer. While such a cycle is typically well reproduced by the models, its phase and/or the

amplitude are not always well captured (Solazzo et al., 2012; Im et al., 2015), leading to the variance error. The variance error also originates from the different spatial support (incommensurability) of point measurements vs. gridded model outputs. The latter have typically larger spatial support, while

receptors are more likely to detect local-scale effects that enhance the observed variance.

The mMSE error of the LT component outweighs the variance error in most cases (in both the EU and NA), and is due to the unexplained portion of observed variance. The processes responsible for the mMSE error of the LT component (such as deposition, transport, stratospheric mixing and photochemistry) act at timescales of more than 21 days.

The DU error (on average 3–4 ppb for AQMEII1 and 2–3 ppb for AQMEII2) makes up the second highest contribution to the total error. The partitioning between variance and the mMSE error varies greatly from model to model. However, a comparison of the two AQMEII phases shows that the mMSE is predominant for AQMEII2, while the variance error (typically due to model under-prediction of the observed variability) is most relevant in several cases of AQMEII1. Therefore, at the DU scale, the “quality” of the error of the AQMEII2 phase is higher than that of its AQMEII1 counterpart. One possible explanation is the fact that coupled models were used in AQMEII2, while AQMEII1 exclusively used non-coupled models. As already mentioned (end of Sect. 3), Makar et al. (2015) found that AQMEII1 models performed better overall with respect to AQMEII2. An analysis of the LT component showed that the bias in the AQMEII2 models is higher, possibly due to the 2010 emission inventory, while an analysis of the DU error found that the variance error in the AQMEII2 models is significantly reduced with respect to the AQMEII1 models, and is almost null. We postulate that the inclusion of feedback effects may have been beneficial, and that the reduced performance of AQMEII2 models is likely due to external bias. The residual mMSE error of the DU component (~ 1 –2 ppb on average for both continents) is mostly likely generated by a number of processes, including chemistry, cloudiness, boundary layer transition and vertical mixing. From Fig. 2, the values of the correlation coefficient for the DU component are very high (exceeding 0.8 in the majority of the cases). Such a high performance can be misleadingly optimistic though, because it mostly reflects the 24 h and annual forcing embedded in both the observations and model values. Further analysis on the amplitude and phase of the error can be more informative.

The SY error (almost entirely due to mMSE in AQMEII2) is comparable across all models applied to the same continental domain (except for GEM-AQ and WRF/Chem, NA), indicating that a possible common source of error may be due to missing processes in the models related to the interaction between chemistry and transport.

Finally, the error of the ID component is less than 1 ppb (on average ~ 0.2 ppb for AQMEII2) and is generated by both variance (most commonly model over-prediction) and mMSE. The fast-acting photochemical processes are, therefore, modelled with satisfactory precision, although the small errors in the ID component can be quite large relative to the total amount of ID variability. Furthermore, the spatial averaging of the time series prior to the spectral decomposition

might have suppressed the variance of the ID component, which is more uncorrelated in space than the other components and would therefore tend to average out.

4.3 Spatial distribution of the spectral error components

Maps of MSE by spectral components are reported in Figs. 3 to 6. As anticipated by the error analysis, the LT is the most problematic source of error for both continents, although the variety in the models’ behaviour does not allow for generalisation.

Some of the cases presented in Fig. 2, where the bias was null (MM5-CAMx, MM5-DEHM for AQMEII1 and CosmoArt for AQMEII2, both in EU), show bias compensation, typically due to model underestimation in the central part of the EU (Germany, eastern France) and model overestimation in the rest of the continent. The case of the CosmoArt model (Fig. 5c) clearly shows the effect of the spatial averaging in masking the error that is only cancelled when a continental average is calculated. The model is in fact affected by severe bias and component errors.

The Po Valley in Italy and the southern part of the EU are the most problematic areas, affected by severe LT errors (Figs. 3 and 5). The central and northern parts of the EU are less problematic, especially for AQMEII2. The other components of the error are significantly smaller than the LT error, with some exceptions (especially for the DU component). The length of the segment is in fact normalised to the largest error for each model, to facilitate the interpretation and the relative weight of each error component.

Concerning NA (Figs. 4 and 6), the DU error has more weight and competes with the LT error in the central and south-eastern parts of the continent. For AQMEII2, the SY error is as significant as the LT error on the east coast (Wrf/Chem, Fig. 6c). The greatest LT error is observed in the coastal areas (east and west) and across the north-eastern border between the USA and Canada (due primarily to model underestimation in the east and north, and model overestimation in the west).

The analysis presented provides a detailed breakdown of the error in terms of error components, spectral decomposition and spatial distribution, thereby avoiding the pitfalls of extreme averaging and providing a comprehensive analysis of where the error occurs and the associated timescales and processes, and whether the error is internally generated or stems from the model’s input data.

5 MSE decomposition and complexity

In regression analysis and statistical learning theories, the problem of under- and over-fitting complex systems is at the root of the MSE decomposition into bias and variance. The trade-off between bias and variance is strictly dependent on

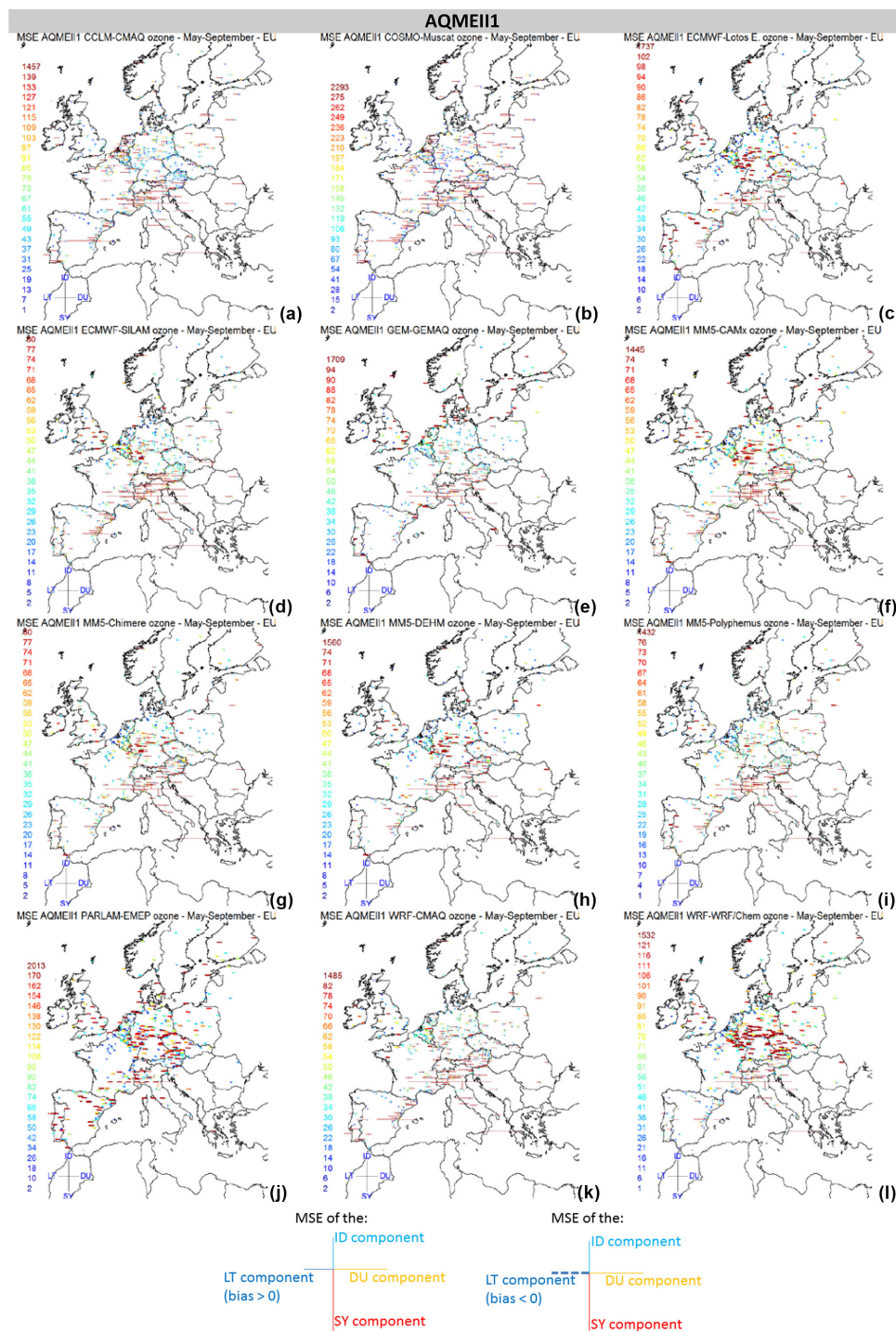


Figure 3. Spatial distribution of the MSE in the spectral components for the EU models of AQMEII1. The segments are centred at the rural receptors' position (clockwise from north: MSE of ID, DU, SY, and LT). Their length is proportional to the MSE magnitude, coded according to the colour scale. For each model, the colour scale extends from zero up to the 75th percentile, and the last value of the scale is the maximum MSE. The colour of the MSE values above the 75th percentile represents the maximum value. The tick-dashed LT segment indicates model underestimation (low model bias), while thin continuous segment indicates model overestimation (high model bias). The example in the last panel indicates how the maps reports the error of the spectral components at each receptor (the colours are arbitrary). The example on the left represents the error at a receptor where the LT component is biased high, while the example on the right refers to a case where the bias is negative. The other components do not change.

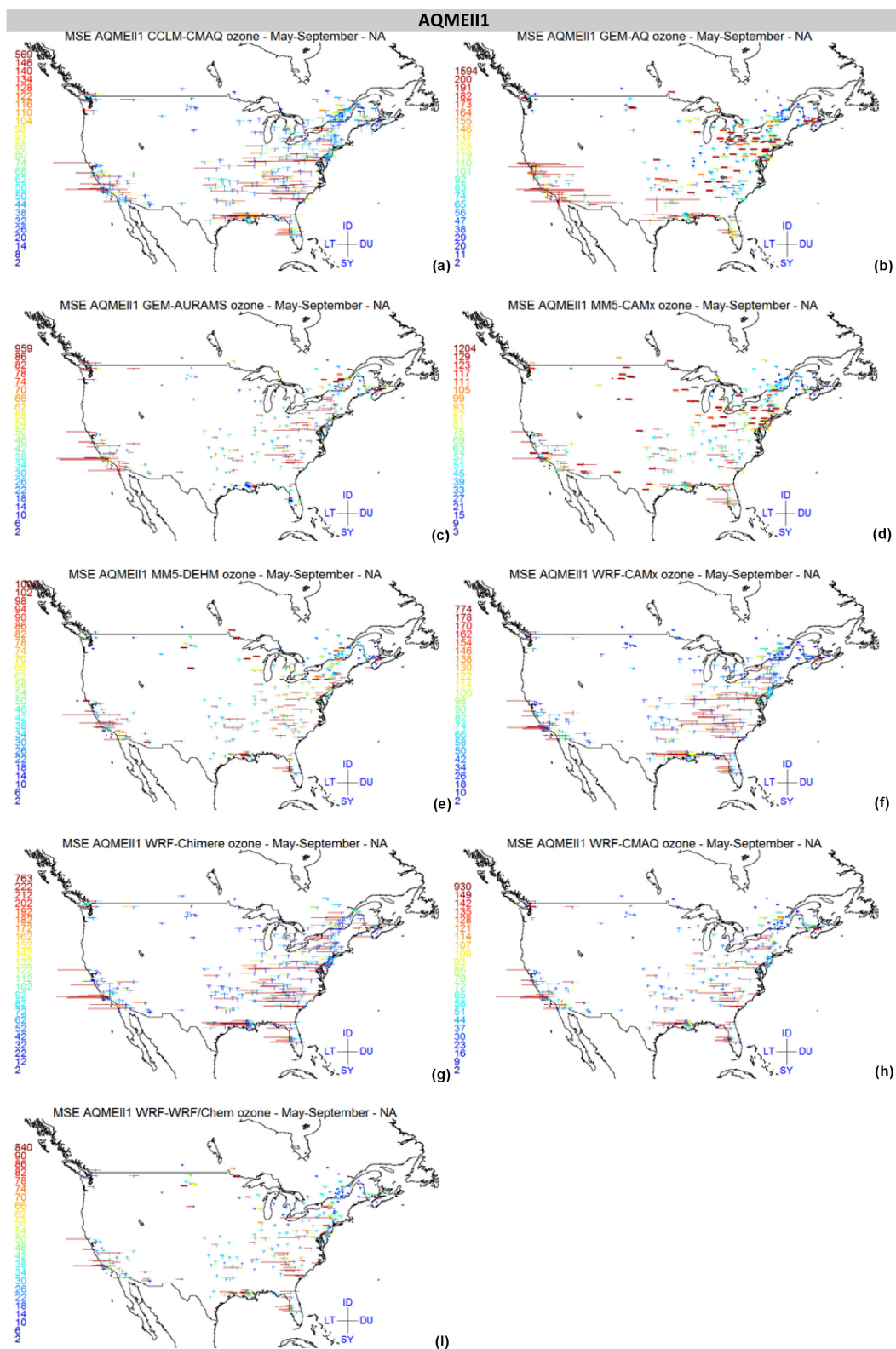


Figure 4. As in Fig. 3, but for the NA models of AQMEI1.

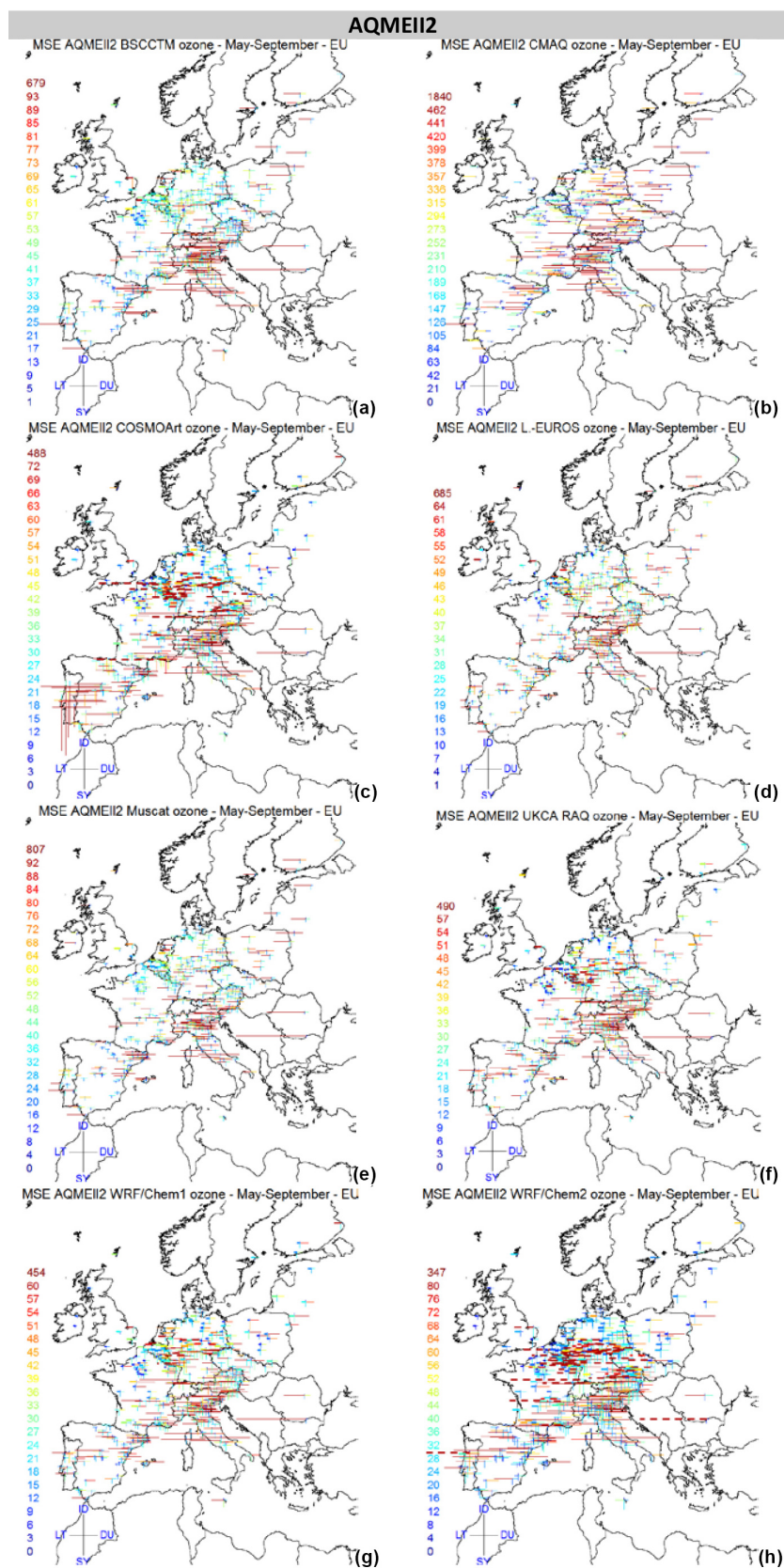


Figure 5. As in Fig. 3, but for the EU models of AQMEII2.

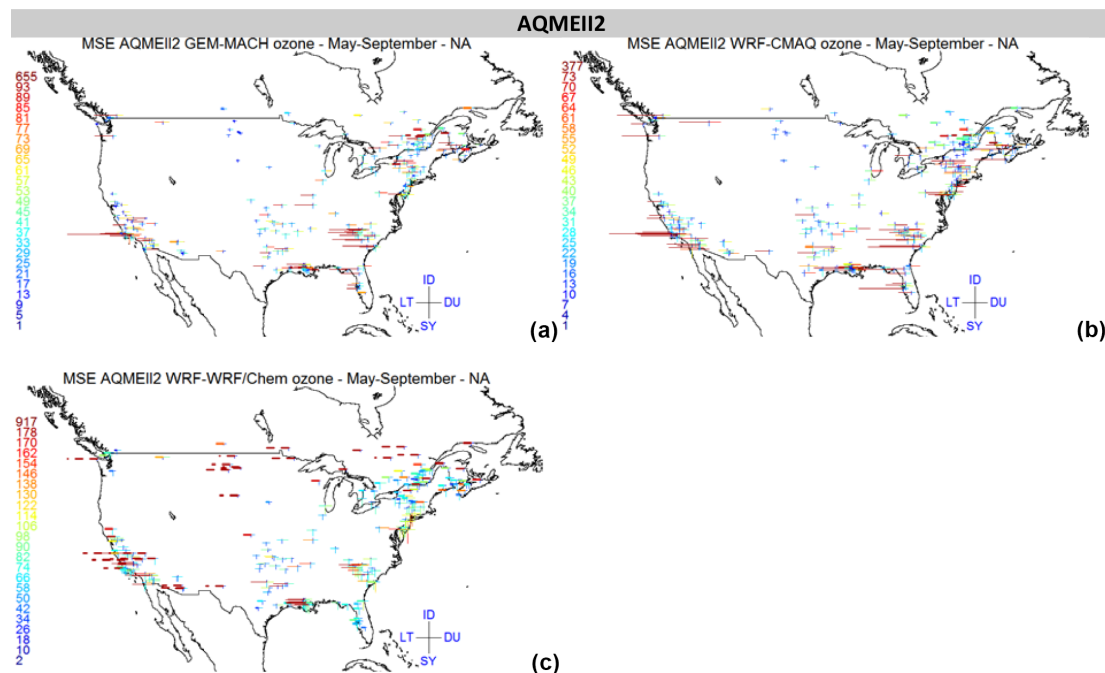


Figure 6. As in Fig. 3, but for the NA models of AQMEII2.

the complexity of the model. Over-fitting occurs when too many parameters and modules are added to the model: each new module added to describe a process is a new source of variance due to internal parameterisation and linearisation. In other words, over-fitting is associated with the stochasticity inherent to the data/model, and contributes to the increase in variance and consequent decrease in bias. Under-fitting occurs due to an oversimplification of the modelled processes, and is an important source of bias as it is associated with the deterministic property of the modelling activity (Hastie et al., 2009).

The problem of the bias-variance trade-off becomes markedly more complicated when dealing with complex models with many degrees of freedom, such as AQ modelling systems. Adding new modules to cope with unexplained physical processes can lead to a reduction in the bias due to that specific process, but also feeds new variance and possibly new bias into the model due to the non-linear interaction of the new module with existing ones, since reducing the bias while preserving the variance is non-trivial.

Rao (2005), in the context of dispersion modelling, provided the theoretical variations of the total model uncertainty by exploiting the components of the difference between the modelled and observed variance (Fig. 1 of Rao et al., 2005). Rao (2005) used the number of meteorological parameters in the model as a measure of model complexity, and concluded that the optimal model complexity could not be defined a priori, but is a trial-and-error combination of the model, the measurement error and the stochastic uncertainty.

In this study we attempt to derive the curves of the MSE components (bias, variance and covariance) as a function of model complexity, providing a first-time attempt to analyse the error of a regional AQ model as function of its complexity. The aim is to find the timescale dominated by the error (and what type of error) and, if it exists, the time window where the error decreases. The information obtained is of immediate usefulness for model development, as it provides a clear temporal cut-off that distinguishes the dynamics of the error.

Figure 7 shows an example of the approach used to break down model complexity, which basically relies on the resolved timescale of the model. The complexity of the model is assumed to increase when the resolved timescale is shortened: the shorter the timescale, the more complex the model. The timescale of the resolved processes is thus used as a measure of the complexity, and is obtained by recursively applying the kz filter to the ozone time series. The minimum complexity is assumed to be represented by a model that cannot resolve any temporal scale below ~ 1 month (far right of Fig. 7), while the maximum complexity corresponds to the hourly time series, i.e. the standard model's output (far left of Fig. 7).

In Fig. 8, we report the spatially averaged curves of bias, variance and covariance according to Eq. (6) as a function of model complexity. According to the regression analysis theories outlined above, we would expect the variance to increase according to the complexity ($\frac{d\sigma_m^2}{d\text{complexity}} > 0$), and the distance between the modelled and observed variance to de-

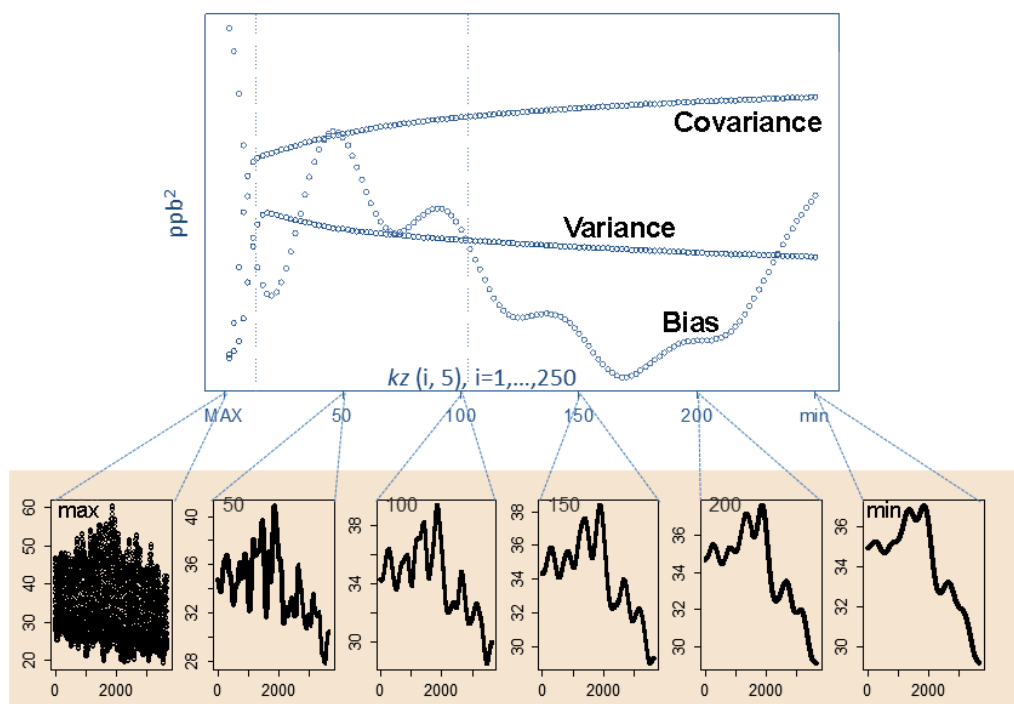


Figure 7. Example of the model complexity as time-resolved scale of the transport and dispersion processes: the minimum complexity (far right) is a poor time-resolving time series obtained as $kz(250,5)$ (~ 1 month). The complexity increases towards the left, with the scale of resolved processes becoming finer up to the maximum complexity (far left), which represents the full time series. The upper panel shows an example of how the curves of the error for covariance, variance and bias vary according to complexity.

crease $\left(\frac{d(\sigma_m - \sigma_o)^2}{d\text{complexity}} < 0\right)$, and the opposite for the bias. The curves of variance in Fig. 8 indeed turn downwards as predicted by the theory, while the curves of bias have a mixed behaviour but are, basically, constant $\left(\frac{d(\text{mod} - \text{obs})^2}{d\text{complexity}} \approx 0\right)$.

More specifically:

- The $(\sigma_m - \sigma_o)^2$ term decreases steadily but slowly to a timescale of ~ 1 day, after which it drastically drops to significantly lower values. This indicates that (i) the complexity of the AQ systems increases exponentially at the DU timescales (not entirely surprising, given the day/night behavioural properties of ozone), (ii) the efforts made to improve the model capabilities on the short-term processes governing the ozone dynamics improve the model precision and (iii) there is a possible lack of parameterisation and modelling of the processes of transport and chemical transformation over periods longer than 1–2 days.
- The fact that the bias varies only by small amounts indicates that a fully evolved model, capable of reproducing processes at the shortest timescales (turbulent dispersion, fast chemical reactions, even day/night variability, etc.) is no more accurate than a basic model that only accounts for long-term processes. This might indicate that (i) the bias at the shorter timescales is introduced

entirely by the larger timescales, and/or (ii) the bias is continuously fed into the model by an external source acting at all scales, as for example the emissions data or boundary conditions.

Summarising, in most cases (both continents, both AQMEII phases), the $(\sigma_m - \sigma_o)^2$ term decreases sharply after a timescale of resolved processes of ~ 1 day; the bias term is surprisingly independent on complexity; the covariance is complementary to the variance. Thus, the bias is the error term to which efforts should be expanded first, and current studies are carried out to diagnose more precisely its origin within AQ modelling systems.

6 Conclusions

This study presents a novel approach to model evaluation, and aims to combine standard operational statistics with the time allocation of the component error. The methodology we propose tackles the issue of diagnostic evaluation from the angle of the spectral decomposition and error breakdown of model/data signals, introducing a compact operator for the quantification of bias, variance and the correlation coefficient.

When the analytical decomposition of the error into bias, variance and mMSE is applied to the decomposition of the

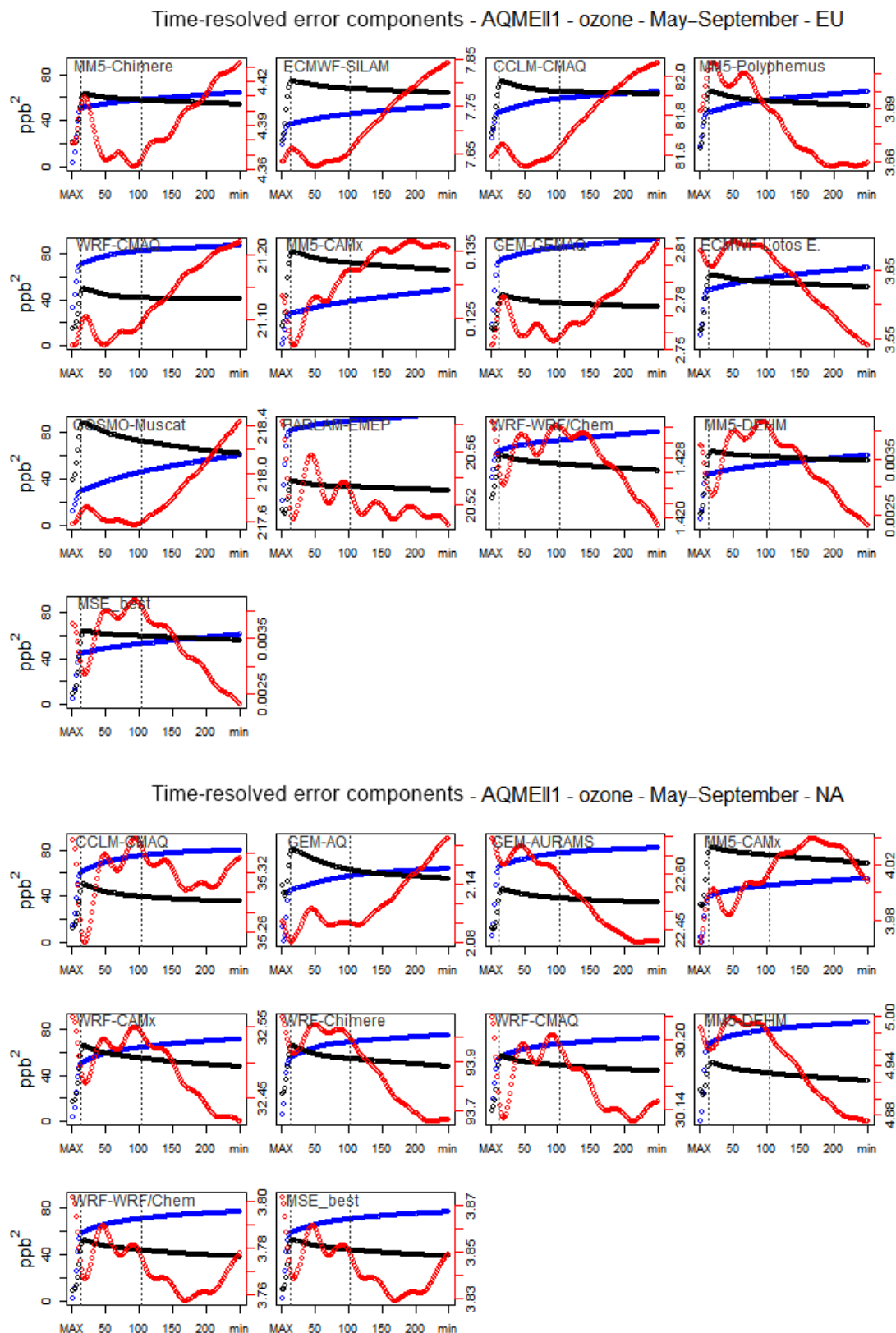


Figure 8.

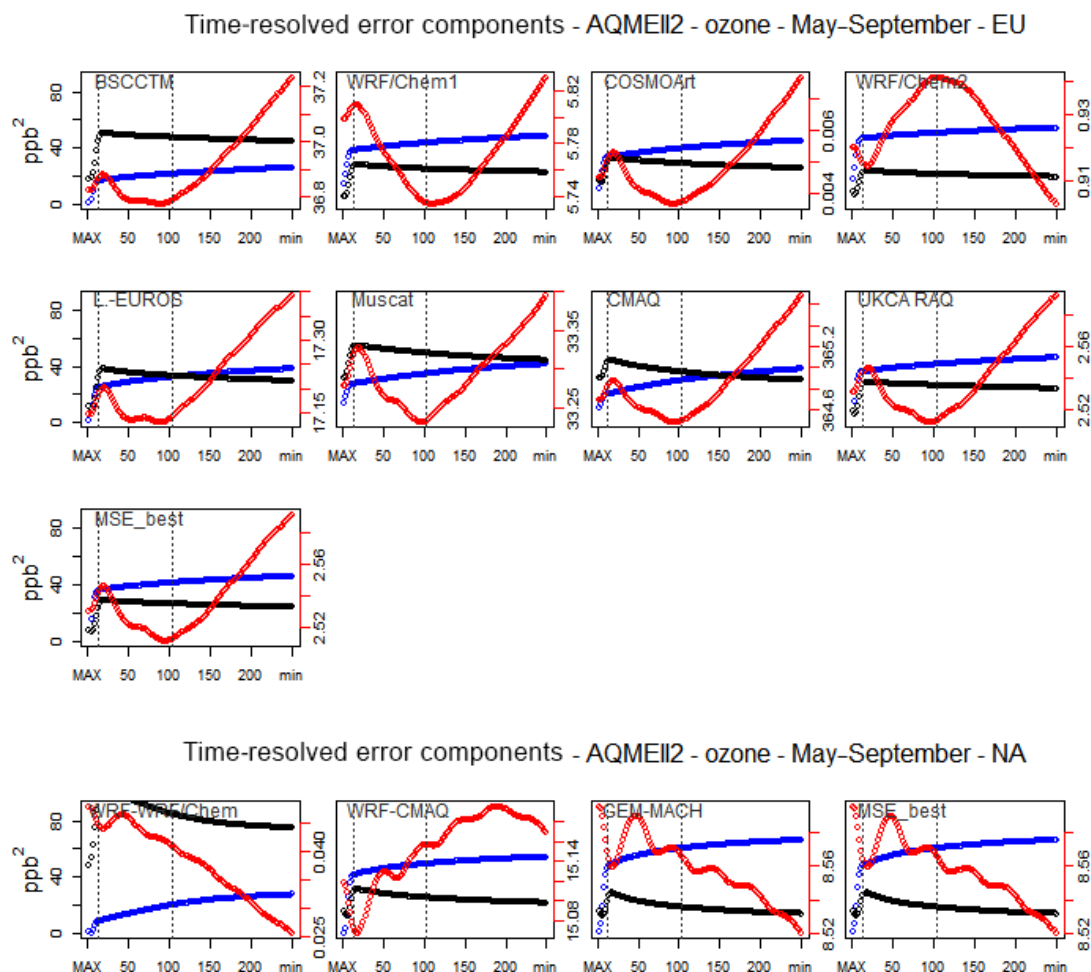


Figure 8. Evolution of error components (red: bias; blue: variance; black: covariance) as a function of model complexity. Complexity increases from right (min) to left (MAX) and is calculated as the temporal scale of the resolved process using the kz filter on the modelled signal: $kz(i,5)$, $i = 2, \dots, 250$.

signals into long-term, synoptic, inter-diurnal and diurnal components, information can be gathered that helps reduce the spectrum of possible sources of errors and pinpoint the processes that are most active at a particular scale that need to be improved. The procedure is denoted here as *error apportionment* and provides an improved and more powerful capacity to identify the nature of the error and associate it with a specific part of the spectrum of the model/measurement signal. The AQMEII set of models and measurements have been used in the evaluation procedure.

After analysing the ozone concentrations gathered in the two phases of AQMEII, which cover a number of modelling systems in two different years and geographical areas, we conclude that

- The bias component of the error is by far the most important source of error, and is mainly associated with long-term processes and/or input fields (likely emissions data or boundary conditions). With regard to the

model application, any effort to improve the current capabilities of AQ modelling systems are likely to have little practical impact if this primary issue is not addressed and solved.

- Most relevant to model development, the variance error (the discrepancy between modelled and observed variance) is mainly associated with the DU component. At a timescale of ~ 1 –2 days, the complexity of modelling systems increases substantially and many processes are involved; the fact that the variance error of the DU component for the AQMEII2 runs is reduced with respect to the AQMEII1 runs might indicate the benefits of including feedback in the models. Such a conclusion could not be drawn with simpler operational evaluation strategies.
- The limited magnitude of the variability of the SY and LT signals produces little variance errors for these two components, and only becomes comparable to the LT or

DU error when the bias is negligible or the total MSE is small.

- The mMSE error is predominant in some instances of the analysed models, and is due to the random distribution of modelled values. There are many causes of mMSE error, including all “internal” processes that produce non-systematic errors such as noise, representativeness, the linearisation of non-linear process and turbulence closure.
- The analysis of the spatial distribution of the error highlights the diversity in the behaviour of each modelling system. The common spatial structures of the LT error (for example in the central and southern EU) may reveal common sources of error (e.g. emissions data), while the error of the other components (especially DU and SY) are peculiar to each model and need to be assessed individually.

Analyses of the modelling results for the third phase of AQMEII are currently building on the methodology outlined in this study, with specific attention being given to the diagnostic of the error of the LT component in relation to external forcing (emissions and boundary conditions) and of the DU component with respect to the variance error.

Appendix A

As in Hogrefe et al. (2000) and Galmarini et al. (2013), the time windows (m) and the smoothing parameter (k) have been selected as follows:

$$\begin{aligned} \text{ID}(t) &= \mathbf{x}(t) - k\mathbf{z}_{3,3}(\mathbf{x}(t)) \\ \text{DU}(t) &= k\mathbf{z}_{3,3}(\mathbf{x}(t)) - k\mathbf{z}_{13,5}(\mathbf{x}(t)) \\ \text{SY}(t) &= k\mathbf{z}_{13,5}(\mathbf{x}(t)) - k\mathbf{z}_{103,5}(\mathbf{x}(t)) \\ \text{LT}(t) &= k\mathbf{z}_{103,5}(\mathbf{x}(t)) \\ \mathbf{x}(t) &= \text{ID}(t) + \text{DU}(t) + \text{SY}(t) + \text{LT}(t), \end{aligned} \quad (\text{A1})$$

where $\mathbf{x}(t)$ is the time series vector.

A clear-cut separation of the components of Eq. (A1) cannot be achieved, as the separation is a non-linear function of the parameters m and k (Rao et al., 1997). It follows that the components of Eq. (A1) are not completely orthogonal and that some level of overlapping energy exists (Kang et al., 2013). Galmarini et al. (2013) found that the explained variance by the spectral components account for 75 to 80 % of the total variance, the remaining portion being explained by the interactions between the components.

Assuming a spectral decomposition, which is valid for the modelling and the observational time series, the MSE formulation outlined in Galmarini et al. (2013) holds:

$$\begin{aligned} \text{MSE}(\text{O}_3) &= \text{MSE}(\text{LT} + \text{SY} + \text{DU} + \text{ID}) = \\ &= \sum \text{MSE}(\text{spec comp}) + \sum \text{MSE}(\text{cc}), \end{aligned} \quad (\text{A2})$$

where spec comp are the diagonal terms, and LT, SY, DU, ID and cc identifies the cross-components, i.e. the off-diagonal terms deriving from the squared nature of the MSE: LT_oSY_m , SY_oLT_m , SY_oDU_m , DU_oSY_m , DU_oID_m , ID_oDU_m , LT_mSY_o , LT_oSY_o , DU_mSY_m , DU_mID_m , DU_oSY_o , DU_oID_o (o and m represent observed and modelled fields, respectively). For simplicity, the cross-components are assumed to be symmetric, so the o and m subscripts are dropped.

To isolate the contribution to MSE of a single spectral component, we proceed as follows. We subtract a component (e.g. LT) from the whole time series:

$$\begin{aligned} \text{MSE}(\text{O}_3 - \text{LT}(\text{O}_3)) &= \text{MSE}(\text{SY}) + \text{MSE}(\text{DU}) \\ &+ \text{MSE}(\text{ID}) + 2\text{MSE}(\text{IDDU}) + 2\text{MSE}(\text{IDSY}) \\ &+ 2\text{MSE}(\text{DUSY}). \end{aligned} \quad (\text{A3})$$

By removing Eq. (A3) from Eq. (A2), the contribution of LT and its cross-component is isolated:

$$\begin{aligned} \text{Eq. A2} - \text{Eq. A3} &= \text{MSE}(\text{LT}) + \text{MSE}(\text{LTID}) \\ &+ \text{MSE}(\text{LTSY}) + \text{MSE}(\text{LTDU}). \end{aligned} \quad (\text{A4})$$

We can further elaborate on Eq. (A4) to isolate the contribution of each cross-component. For instance, the case of SYLT:

$$\begin{aligned} \text{MSE}(\text{SY-ID-DU}) - \text{MSE}(\text{SY}) - \text{MSE}(\text{LT}) &= \\ [\text{MSE}(\text{SY}) + \text{MSE}(\text{LT}) + 2\text{MSE}(\text{SYLT})] - \\ \text{MSE}(\text{SY}) - \text{MSE}(\text{LT}) &= 2\text{MSE}(\text{SYLT}). \end{aligned} \quad (\text{A5})$$

The procedure in Eq. (A5) has been applied to derive the contribution of all cross-components.

The Supplement related to this article is available online at doi:10.5194/acp-16-6263-2016-supplement.

Acknowledgements. We would like to thank the community of modellers and data providers of the first and second phases of AQMEII.

Edited by: C. Hogrefe

References

- Baklanov, A., Schlünzen, K., Suppan, P., Baldasano, J., Brunner, D., Aksoyoglu, S., Carmichael, G., Douros, J., Flemming, J., Forkel, R., Galmarini, S., Gauss, M., Grell, G., Hirtl, M., Joffe, S., Jorba, O., Kaas, E., Kaasik, M., Kallos, G., Kong, X., Kørsholm, U., Kurganskiy, A., Kushta, J., Lohmann, U., Mahura, A., Manders-Groot, A., Maurizi, A., Moussiopoulos, N., Rao, S. T., Savage, N., Seigneur, C., Sokhi, R. S., Solazzo, E., Solomos, S., Sørensen, B., Tsegas, G., Vignati, E., Vogel, B., and Zhang, Y.: Online coupled regional meteorology chemistry models in Europe: current status and prospects, *Atmos. Chem. Phys.*, 14, 317–398, doi:10.5194/acp-14-317-2014, 2014.
- Brunner, D., Jorba, O., Savage, N., Eder, B., Makar, P., Giordano, L., Badia, A., Balzarini, A., Baro, R., Bianconi, R., Chemel, C., Forkel, R., Jimenez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Im, U., Knote, C., Kuenen, J. J. P., Makar, P. A., Manders-Groot, A., Neal, L., Perez, J. L., Pirovano, G., San Jose, R., Savage, N., Schroder, W., Sokhi, R. S., Syrakov, D., Torian, A., Werhahn, K., Wolke, R., van Meijgaard, E., Yahya, K., Zabkar, R., Zhang, Y., Zhang, J., Hogrefe, C., and Galmarini, S.: Evaluation of the meteorological performance of coupled chemistry meteorology models in phase 2 of the air quality model evaluation international initiative, *Atmos. Environ.*, 115, 470–498 doi:10.1016/j.atmosenv.2014.12.032, 2015.
- Dennis, R., Fox, T., Fuentes, M., Gilliland, A., Hanna, S., Hogrefe, C., Irwin, J., Rao, S. T., Scheffe, R., Schere, K., Steyn, D., and Venkatram, A.: A framework for evaluating regional-scale numerical photochemical modeling systems, *Environ. Fluid Mech.*, 10, 471–489, doi:10.1007/s10652-009-9163-2, 2010.
- Fox, D. G.: Judging air quality model performance, *B. Am. Meteorol. Soc.*, 62, 599–609, 1981.
- Galmarini, S., Kioutsioukis, I., and Solazzo, E.: *E pluribus unum**: ensemble air quality predictions, *Atmos. Chem. Phys.*, 13, 7153–7182, doi:10.5194/acp-13-7153-2013, 2013.
- Galmarini, S., Solazzo, E., Im, U., and Kioutsioukis, I.: AQMEII 1, 2 and 3: Direct and Indirect Benefits of Community Model Evaluation Exercises, 34th International Technical Meeting on Air Pollution Modelling and its Application, Montpellier, France, 4–8 May 2015.
- Giordano, L., Brunner, D., Flemming, J., Hogrefe, C., Im, U., Bianconi, R., Badia, A., Balzarini, A., Baró, R., Chemel, C., Curci, G., Forkel, R., Jiménez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Jorba, O., Knote, C., Kuenen, J. J. P., Makar, P. A., Manders-Groot, A., Neal, L., Pérez, J. L., Pirovano, G., Pouliot, G., San José, R., Savage, N., Schröder, W., Sokhi, R. S., Syrakov, D., Torian, A., Tuccella, P., Werhahn, J., Wolke, R., Yahya, K., Žabkar, R., Zhang, Y., and Galmarini, S.: Assessment of the MACC re-analysis and its influence as chemical boundary conditions for regional air quality modelling in AQMEII-2, *Atmos. Environ.*, 115, 371–388, 2015.
- Hastie, T., Tibshirani, R., and Friedman, J.: The elements of statistical learning, 2nd Edn., Springer-Verlag, 763 pp., 2009.
- Hogrefe, C., Rao, S. T., Zurbenko, I. G., and Porter, P. S.: Interpreting the information in ozone observations and model predictions relevant to regulatory policies in the Eastern United States, *B. Am. Meteorol. Soc.*, 81, 2083e2106, doi:10.1175/1520-0477(2000)0812.3.CO;2, 2000.
- Hogrefe, C., Roselle, S., Mathur, R., Rao, S. T., and Galmarini, S.: Space-time analysis of the Air Quality Model Evaluation International Initiative (AQMEII) phase 1 air quality simulation, *J. Air Waste Manage.*, 64, 388–405, 2014.
- Im, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., Baro, R., Bellasio, R., Brunner, D., Chemel, C., Curci, G., Denier van der Gon, H., Flemming, J., Forkel, R., Giordano, L., Jimenez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Jorba, O., Knote, C., Makar, P. A., Manders-Groot, A., Neal, L., Pérez, J. L., Pirovano, G., Pouliot, G., San Jose, R., Savage, N., Schroder, W., Sokhi, R. S., Syrakov, D., Torian, A., Tuccella, P., Wang, K., Werhahn, J., Wolke, R., Zabkar, R., Zhang, Y., Zhang, J., Hogrefe, C., and Galmarini, S.: Evaluation of operational onlinecoupled regional air quality models over Europe and North America in the context of AQMEII phase 2. Part II: particulate matter, *Atmos. Environ.*, 115, 421–441, 2015a.
- Im, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., Baro, R., Bellasio, R., Brunner, D., Chemel, C., Curci, G., Flemming, J., Forkel, R., Giordano, L., Jimenez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Jorba, O., Knote, C., Kuenen, J. J. P., Makar, P. A., Manders-Groot, A., Neal, L., Pérez, J. L., Pirovano, G., Pouliot, G., San Jose, R., Savage, N., Schroder, W., Sokhi, R. S., Syrakov, D., Torian, A., Tuccella, P., Werhahn, J., Wolke, R., Yahya, K., Zabkar, R., Zhang, Y., Zhang, J., Hogrefe, C., and Galmarini, S.: Evaluation of operational online-coupled regional air quality models over Europe and North America in the context of AQMEII phase 2. Part I: ozone, *Atmos. Environ.*, 115, 404–420, 2015b.
- Johnson, R.: Assessment of Bias with Emphasis on Method Comparison, *Clin. Biochem.*, 29, S37–S42, 2008.
- Kang, D., Hogrefe, C., Foley, K. L., Napelenok, S. L., Mathur, R., and Rao, S. T.: Application of the Kolmogorov-Zurbenko filter and the decoupled direct 3D method for the dynamic evaluation of a regional air quality model, *Atmos. Environ.*, 80, 58–69, 2013.
- Kioutsioukis, I. and Galmarini, S.: *De praeceptis ferendis*: good practice in multi-model ensembles, *Atmos. Chem. Phys.*, 14, 11791–11815, doi:10.5194/acp-14-11791-2014, 2014.
- Makar, P. A., Gong, W., Hogrefe, C., Zhang, Y., Curci, G., Žabkar, R., Milbrandt, J., Im, U., Balzarini, A., Baró, R., Bianconi, R., Cheung, P., Forkel, R., Gravel, S., Hirtl, M., Honzak, L., Hou, A., Jiménez-Guerrero, P., Langer, M., Moran, M. D., Pabla, B., Pérez, J. L., Pirovano, G., San José, R., Tuccella, P., Werhahn, J., Zhang, J., and Galmarini, S.: Feedbacks between air pollution and weather, part 2: effects on chemistry, *Atmos. Environ.*, 115, 499–526, 2015.

- Murphy, A. H.: Skill scores based on the mean square error and their relationship to the correlation coefficient, *Mon. Weather Rev.*, 116, 2417–2424, 1988.
- Pindyck, R. S. and Rubinfeld, D. L.: *Econometric Models and Economic Forecast*, Irwin/McGraw-Hill, Singapore, 388 pp., 1998.
- Pouliot, G., Pierce, T., Denier van der Gon, H., Schaap, M., Moran, M., and Nopmongcol, U.: Comparing Emissions Inventories and Model-Ready Emissions Datasets between Europe and North America for the AQMEII Project, *Atmos. Environ.*, 53, 4–14, 2012.
- Pouliot, G., Denier van der Gon, H., Kuenen, J., Makar, P., Zhang, J., and Moran, M.: Analysis of the emission inventories and model-ready emission datasets of Europe and North America for phase 2 of the AQMEII project, *Atmos. Environ.*, 115, 345–360, 2015.
- Rao, K. S.: Uncertainty analysis in atmospheric dispersion modelling, *Pure Appl. Geophys.*, 162, 1893–1917, 2005.
- Rao, S. T., Zurbenko, I. G., Neagu, R., Porter, P. S., Ku, J. Y., and Henry, R. F.: Space and time scales in ambient ozone data, *B. Am. Meteorol. Soc.*, 78, 2153e2166, doi:10.1175/1520-0477(1997)078<2153:SATSIA>2.0.CO;2, 1997.
- Rao, S. T., Galmarini, S., and Puckett, K.: Air quality model evaluation international initiative (AQMEII), *B. Am. Meteorol. Soc.*, 92, 23–30, doi:10.1175/2010BAMS3069.1, 2011.
- Schere, K., Flemming, J., Vautard, R., Chemel, C., Colette, A., Hogrefe, C., Bessagnet, B., Meleux, F., Mathur, R., Roselle, S., Hu, R.-M., Sokhi, R. S., Rao, S. T., and Galmarini, S.: Trace gas/aerosol concentrations and their impacts on continental-scale AQMEII modelling sub-regions, *Atmos. Environ.*, 53, 38–50, 2012.
- Solazzo, E. and Galmarini, S.: Comparing apples with apples: Using spatially distributed time series of monitoring data for model evaluation, *Atmos. Environ.*, 112, 234–245, 2015.
- Solazzo, E., Bianconi, R., Vautard, R., Appel, K. W., Moran, M. D., Hogrefe, C., Bessagnet, B., Brandt, J., Christensen, J. H., Chemel, C., Coll, I., van der Gon, H. D., Ferreira, J., Forkel, R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B., Jericevic, A., Kraljevic, L., Miranda, A. I., Nopmongcol, U., Pirovano, G., Prank, M., Riccio, A., Sartelet, K. N., Schaap, M., Silver, J. D., Sokhi, R. S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S. T., and Galmarini, S.: Model evaluation and ensemble modelling and for surface-level ozone in Europe and North America, *Atmos. Environ.*, 53, 60–74, 2012a.
- Solazzo, E., Bianconi, R., Pirovano, G., Matthias, V., Vautard, R., Moran, M. D., Appel, K. W., Bessagnet, B., Brandt, J., Christensen, J. H., Chemel, C., Coll, I., Ferreira, J., Forkel, R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B., Hogrefe, C., Miranda, A. I., Nopmongcol, U., Prank, M., Sartelet, K. N., Schaap, M., Silver, J. D., Sokhi, R. S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S. T., and Galmarini, S.: Operational model evaluation for particulate matter in Europe and North America, *Atmos. Environ.*, 53, 75–92, 2012b.
- Solazzo, E., Bianconi, R., Pirovano, G., Moran, M. D., Vautard, R., Hogrefe, C., Appel, K. W., Matthias, V., Grossi, P., Bessagnet, B., Brandt, J., Chemel, C., Christensen, J. H., Forkel, R., Francis, X. V., Hansen, A. B., McKeen, S., Nopmongcol, U., Prank, M., Sartelet, K. N., Segers, A., Silver, J. D., Yarwood, G., Werhahn, J., Zhang, J., Rao, S. T., and Galmarini, S.: Evaluating the capability of regional-scale air quality models to capture the vertical distribution of pollutants, *Geosci. Model Dev.*, 6, 791–818, doi:10.5194/gmd-6-791-2013, 2013a.
- Solazzo, E., Riccio, A., Kioutsioukis, I., and Galmarini, S.: Pauci ex tanto numero: reduce redundancy in multi-model ensembles, *Atmos. Chem. Phys.*, 13, 8315–8333, doi:10.5194/acp-13-8315-2013, 2013b.
- Stoeckenius, T. E., Hogrefe, C., Zagunis, J., Sturtz, T. M., Wells, B., and Sakulyanontvittaya, T.: A comparison between 2010 and 2006 air quality and meteorological conditions, and emissions and boundary conditions used in simulations of the AQMEII2 North American domain, *Atmos. Environ.*, 115, 389–403, 2015.
- Theil, H.: *Economic forecast and policy*, North-Holland, Amsterdam, 1961.
- Willmott, C. J., Ackleson, S. G., Davis, R. E., Feddema, J. J., Klink, K. M., Legates, D. R., O'Donnell, J., and Rowe, C. M.: Statistics for the evaluation and comparison of models, *J. Geophys. Res.*, 90, 8995–9005, 1985.
- Wise, E. K. and Comrie, A. C.: Extending the KZ filter: application to ozone, particulate matter, and meteorological trends, *J. Air Waste Manage.*, 55, 1208e1216, doi:10.1080/10473289.2005.10464718, 2005.
- Zurbenko, I. G.: *The Spectral Analysis of Time Series*, North-Holland, Amsterdam, 236 pp., 1986.