



1 ERROR APPORTIONMENT FOR ATMOSPHERIC CHEMISTRY-TRANSPORT
2 MODELS. A NEW APPROACH TO MODEL EVALUATION
3 E. Solazzo, S. Galmarini

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

6 Author for correspondence: S. Galmarini, stefano.galmarini@jrc.ec.europa.eu,
7 Phone: +390332785382

8
9 **Abstract.** In this study, methods are proposed to diagnose the causes of errors in air quality
10 (AQ) modelling systems. We investigate the deviation between modelled and observed time
11 series of surface ozone through a revised formulation for breaking down the mean square
12 error (MSE) into bias, variance, and the minimum achievable MSE (*mMSE*). The bias
13 measures the accuracy and implies the existence of systematic errors and poor
14 representation of data complexity, the variance measures the precision and provides an
15 estimate of the variability of the modelling results in relation to the observed data, and the
16 *mMSE* reflects unsystematic errors and provides a measure of the associativity between the
17 modelled and the observed fields through the correlation coefficient. Each of the error
18 components is analysed independently and apportioned to resolved process based on the
19 corresponding timescale (long scale, synoptic, diurnal, and intra-day) and as a function of
20 model complexity.

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

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

28 **Keywords:** Model evaluation; Time series analysis; Bias-variance decomposition; AQMEII

29 1. INTRODUCTION

30 Due to their use for regulatory applications and to support legislation, air quality (AQ)
31 models must model correctly and be correctly applied, justifying the need for a thorough
32 evaluation. A framework for the operational and scientific evaluation of geophysical models
33 was already envisaged in the early '80s (Fox, 1981; Wilmott et al., 1985), the former being '*a
34 comparison with data exclusively within a particular application context*', and the latter
35 defined as '*some understanding of cause-and-effect relationship that relies on testing model
36 components and extensively detailed data collection*' (Fox, 1981). Thirty years later, as AQ



37 models became more and more complex and their range of applicability widened, Dennis et
38 al. (2010) further elaborated the concept of model evaluation by proposing a four-level
39 evaluation, according to which different complementary aspects of the models should be
40 tested, namely:

- 41 a. Operational: the level of agreement of model results with observations;
- 42 b. Dynamic: ability of the modelling system to respond to changes (in emissions, or in
43 meteorological events);
- 44 c. Diagnostic: identify and attribute the source of the error to the relevant process;
- 45 d. Probabilistic: confidence and uncertainty levels of the modelled results.

46 In the framework originally designed by Dennis et al. (2010), the diagnostic component
47 plays a central role. It *i)* answers the fundamental issue left open by the operational
48 screening, in other words whether the model provides the right answer for the right reason,
49 *ii)* provides feedback to developers to help make model improvements, and *iii)* sets the
50 basis for the probabilistic evaluation (Figure 1 of Dennis et al., 2010).

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

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

71 Recently, the AQMEII (Air Quality Model Evaluation International Initiative) activity (Rao et
72 al., 2011) applied the approach proposed by Dennis et al. (2010), by organising model



73 evaluation activities (AQMEII 1, 2 and 3) using operational (Solazzo et al., 2012a,b; Solazzo
74 et al., 2013a; Im et al., 2015a,b), probabilistic (Solazzo et al., 2013b; Kioutsioukis et al.,
75 2014), and diagnostic (Hogrefe et al., 2014; Makar et al., 2015) evaluation frameworks.

76 The study we present here follows and complements the previous investigations based on
77 the AQMEII models collected in the first and second phases of the activity (AQMEII1 in 2006
78 and AQMEII2 in 2010). The main aim is to introduce a novel method that combines
79 operational and diagnostic evaluations. This method helps apportion the model error to its
80 components, thereby identifying the space/timescale at which it is most relevant and, when
81 possible, to infer which process/es could have generated it. This work is designed to support
82 the analysis of the currently ongoing third phase of the AQMEII activity (Galmarini et al.,
83 2015).

84 2. MEAN SQUARE ERROR AS A COMPREHENSIVE METRIC

85 For the model evaluation strategy proposed, we start by breaking down the Mean Square
86 Error (MSE) (used here as unique metric to evaluate model performance) into the sum of
87 the variance (and covariance) and the squared bias. The error and its components are then
88 calculated on the spectrally decomposed time series of modelled and observed hourly
89 ozone mixing ratios. The advantage of this evaluation strategy is twofold:

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

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

104 2. 1. ERROR DECOMPOSITION

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

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

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



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

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

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

108

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

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

111

112 the final formulation for the MSE components reads:

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

113

114 where the covariance term (last term on the right-hand side of Eq 5) accounts for the
115 degree of correlation between the modelled and observed time series. When the covariance
116 term is zero, $var(obs)$ is referred to as the *incompressible part of the error* and represents
117 the lowest limit that the MSE of the model can achieve. When dealing with model
118 evaluation, the modelled and observed time series are typically highly correlated and
119 therefore, within the limits of the perfect match (correlation coefficient of unity), $cov(mod,$
120 $obs) = cov(obs, obs) = cov(mod, mod) = var(mod) = var(obs)$ and the MSE can be reduced to
121 only the bias term. That implies that the development of a high-quality model needs to
122 ensure:

123 a. the highest possible precision in order to maximise the $cov(mod, obs)$ term, and

124 b. the highest possible accuracy, in order to minimise the bias.

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

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

126

127 In Eq 6, the variance term is expressed as the difference between the standard deviation of
128 the model and that of the observations, and the covariance term (last term on the right)
129 includes r , the coefficient of correlation between the observed and modelled time series.
130 The ratios of the three terms on the right-hand side of Eq 6 to the overall MSE are known as
131 *Theil's coefficients* (Pindick and Rubinfeld, 1998).

132 The bias measures the departure of the modelled from the observed results, and is a
133 measure of systematic error, since it measures the extent to which the average modelled
134 values deviate from the observed ones. The bias is commonly used to express the degree of
135 'trueness', i.e. "the closeness of agreement between the average value obtained from a
136 large series of measurements and the true value" (Johnson, 2008). The variance shows



137 whether the modelled variability is compatible with that observed. Finally, the covariance
138 term represents the unexplained proportion of the MSE due to the remaining unsystematic
139 errors, i.e. it represents the remaining error after deviations from the mean values have
140 been accounted for. This latter term is a measure of the lack of correlation of the model
141 with comparable observations, and is considered the least 'worrisome' portion of the error
142 (Pindick and Rubinfeld, 1998).

143 Elaborating on Eq 6, the conditions that minimise the MSE are:

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

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

145

$$\begin{cases} \overline{mod} = \overline{obs} \\ \sigma_m = r\sigma_{obs} \end{cases} \quad \text{EQ 7}$$

146

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

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

149

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

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

154

155 The decompositions in Eq 5, Eq 6, and Eq 9 contain all the relevant operational metrics
156 usually applied to score modelling systems (bias, variance, correlation coefficient), and
157 therefore prove to be a compact estimator of accuracy (bias), precision (variance) and
158 associativity (unexplained portion through the correlation coefficient). Eq 9 has been
159 explicitly derived in this study to help evaluate AQ models. Murphy (1988) provided
160 examples of the scores that can be developed using the components of the MSE.

161 Ideally, the entire error should be attributable to unsystematic fluctuations. From a model
162 development perspective, the variance and covariance are possibly more revealing of model
163 deficiencies than is the bias term, as they are produced by the AQ model itself, while the
164 bias is also due to external sources (e.g. emissions, boundary conditions). From the



165 application viewpoint, however, it is the overall error that counts, which is mostly made up
166 of the bias.

167 2.2. SPECTRAL DECOMPOSITION OF MODELLED AND OBSERVED TIME SERIES

168 Hourly time series of (modelled and observed) ozone concentrations have been
169 decomposed using an iterative moving average approach known as the Kolmogorov-
170 Zurbenko (kz) low-pass filter (Zurbenko, 1986), whose applications to ozone are vastly
171 documented in the literature (Rao et al., 1997; Wise and Comrie, 2005; Hogrefe et al., 2000
172 and 2014; Galmarini et al., 2013; Kang et al., 2013; Solazzo and Galmarini, 2015). The kz
173 filter depends on two parameters: the length of the moving average window m and the
174 number of iterations k ($kz_{m,k}$). Since the kz is a low-pass filter, the filtered time series
175 consists of the low-frequency fluctuating component, while the difference between two
176 filtered time series provides a band-pass filter. This latter property is used to decompose the
177 ozone concentration time series as:

$$O_3 = LT(O_3) + SY(O_3) + DU(O_3) + ID(O_3) \quad \text{EQ 10}$$

178

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

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

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

- 197 - The DU component accounts for more than half of the total variance, followed by
198 the LT and SY components;
- 199 - The ID component has the smallest influence due to the small amplitude of its
200 fluctuations;



201 - The variance of the spectral component is neither strongly nor systematically
202 associated with the area-type of the monitoring stations (i.e. rural, urban, suburban);

203 - Due to the bias, most of the error is accounted for by the LT component, followed by
204 the DU component. The ID contributes very little to the overall MSE.

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

208 The signal decomposition of Eq 10 is applied to the full-year time series. However, to
209 evaluate the model performance with regard to ozone, the analysis is restricted to the
210 months of May to September, i.e. when the production of ozone due to photochemistry is
211 most relevant.

212 3. DATA AND MODELS USED

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

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

223 There are profound differences between the modelling systems that participated in
224 AQMEII1 and AQMEII2. The two sets of models have been applied to different years (2006
225 for phase 1 and 2010 for phase 2) and are therefore dissimilar with respect to the input data
226 of emissions and boundary conditions for chemistry. The AQ models of the second phase
227 are coupled (online chemistry feedbacks on meteorology), while those of the first phase are
228 not. The effect of using online models for simulating ozone accounts for the impact of
229 aerosols on radiation and therefore on temperature and photolysis rates (Baklanov et al.,
230 2014).

231 The model settings and input data for phase I are described in Solazzo et al. (2012a, b;
232 2013a), Schere et al. (2012), and Pouliot et al. (2012); for phase II, similar information is
233 presented in Im et al. (2015a, b), Brunner et al. (2015), and Pouliot et al. (2015).

234 Table 2 summarises the features of the modelling systems analysed in this study with regard
235 to ozone concentrations in the EU or NA. The modelling contribution to the two phases of
236 AQMEII consists of 12 and 9 models and of 8 and 3 models for EU and NA, respectively.



237 Detailed analysis of the main differences in emissions, boundary conditions, and
238 meteorology between the modelled years of 2006 (AQMEI1) and 2010 (AQMEI2) is
239 presented in Stoekenius et al. (2015). A summary of the performance of the two suites of
240 model runs is provided in Makar et al. (2015), showing that the AQMEI1 models generally
241 performed better than the AQMEI2 models, based on standard operational metrics.
242 However, the use of standard evaluation methods does not allow for the assessment of
243 whether the feedback processes have an effect on the deterioration of model performance,
244 or rather the different sets of emissions and boundary conditions. We try to assess the
245 problem using the error apportionment methods outlined above.

246 4. RESULTS FOR THE SPATIALLY AVERAGED TIME SERIES

247 4.1 MSE OF SPECTRAL COMPONENTS

248 Figure 1 reports the MSE share of the spectral components and cross components for each
249 model, for both phases of AQMEI, spatially averaged over the two continental areas.

250 The LT share of the total MSE is the largest in absolute value for both continents and both
251 simulated years. The LT share ranges between 9.9% (GEM-AQ, AQMEI1, NA) and 86.7%
252 (WRF/Chem, AQMEI1, NA), and averages at ~34% and ~46.5% for the EU and ~50.6% and
253 ~47% for NA (AQMEI1 and AQMEI2, respectively).

254 The second largest share of the total MSE is of the DU component, accounting for ~20% (all
255 cases), followed by the SY component. Depending on the model, the MSE share of the
256 remaining spectral components and cross-components varies significantly. Being the
257 intermediate time scales, the overlap of the DU and SY components is likely to be more
258 significant than the overlap of the LT and ID scales. The contribution of DU_{cc} and SY_{cc} to the
259 total error can be as high as 17% (DU_{cc} for GEM-AQ, AQMEI1, NA) and 16% (SY_{cc} for MM5-
260 CAMx, AQMEI1, EU). Overall, the DU_{cc} terms (interaction of DU with the neighbouring SY
261 and ID scales) are significant in both continents (~10%), while the share of the SY
262 component and cross-components is more significant in the EU.

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

265 The results of Figure 1 help identify the time-scales and associated processes for which the
266 largest improvement in model accuracy can be achieved. The LT component has the largest
267 share of the error due to the bias (error breakdown is discussed in the next section), but
268 'internal' chemical processes, transport, and deposition also occur at this timescale. Diurnal
269 processes are the second largest source of error, including, among others, chemistry,
270 boundary layer dynamics, radiation forcing, and their interactions. The processes in the SY
271 band bridge meteorological and chemical processes, and discern between the fast-acting
272 diurnal processes and the baseline. As such, although the SY signal is not as strong as that of



273 the DU components (variance of SY is comparable to the variance of ID, see Hogrefe et al.,
274 2014), it accounts for a significant portion of the total error, as discussed next.

275 4.2 THE QUALITY OF THE ERROR: ERROR APPORTIONMENT

276 The error breakdown (Eq 9) of each spectral component complements the analysis
277 presented in the previous section, and is reported in Figure 2. The bias (only included in the LT
278 component) is the average amount by which the modelled time series is displaced with
279 respect to the observed time series, and is the main source of error. The bias can be either
280 due to 'internal' model errors, or inherited from external drivers (emissions, meteorology,
281 boundary conditions). While the former are of interest for model development because
282 they are generated by systematic modelling errors, the bias introduced by external drivers is
283 responsible for the largest share of modelling errors.

284 From the continental average error breakdown of Figure 2 we can conclude that the majority
285 of EU models (in both AQMEII phases) have small bias (continental-wide average), with the
286 important exceptions of CCLM-CMAQ and Muscat models in AQMEII1, and CMAQ in
287 AQMEII2, which introduced large positive biases. The bias for the NA continent is more
288 uniformly distributed across the models (model over-prediction in both AQMEII phases),
289 possibly indicating a common source of (external) bias in the NA models. The error
290 introduced by external fields is reflected by the bias of the baseline component (LT). For the
291 period between May and September, the error in modelled ozone due to the boundary
292 condition is typically small (Solazzo et al., 2012; Im et al., 2015; Giordano et al., 2015;
293 Hogrefe et al., 2014), while the emissions of ozone precursors and VOCs are problematic,
294 especially in the EU (Makar et al., 2015; Brunner et al., 2015). We further notice that the
295 absence of bias in some models may be caused by the presence of compensating bias, i.e.
296 spatially distributed biases of opposite signs. The spatial distribution of the MSE is discussed
297 in the next section. In all cases, the MSE_{best} model is, by definition, the model with lowest
298 MSE and thus the one with the smallest LT bias.

299 The variance share of LT error is generally small (~1 - 2.5 ppb). This is not entirely
300 unexpected, as the LT component has a high signal-to-noise ratio with a well-structured
301 seasonal cycle, peaking in summer. While such a cycle is typically well reproduced by the
302 models, its phase and/or the amplitude are not always well captured (Solazzo et al., 2012;
303 Im et al., 2015), leading to the variance error. In detail, the $mMSE$ error of the LT component
304 outweighs the variance error in most cases (in both the EU and NA), and is due to the
305 unexplained portion of observed variance, thus to the sparseness of the modelled values.
306 The processes responsible for the $mMSE$ error of the LT component (such as deposition,
307 transport, stratospheric mixing and photochemistry) act at timescales of more than 21 days.

308 The DU error (on average 3-4 ppb for AQMEII1 and 2-3 ppb for AQMEII2) makes up the
309 second highest contribution to the total error. The portioning between variance and the
310 $mMSE$ error varies greatly from model to model. However, a comparison of the two AQMEII



311 phases shows that the *mMSE* is predominant for AQMEII2, while the variance error
312 (typically due to model under-prediction of the observed variability) is most relevant in
313 several cases of AQMEII1. Therefore, at the DU scale, the ‘quality’ of the error of the
314 AQMEII2 phase is higher than that of its AQMEII1 counterpart. One possible explanation is
315 the fact that coupled models were used in AQMEII2, while AQMEII1 exclusively used non-
316 coupled models. As already mentioned (end of section 3), Makar et al. (2015) found that
317 AQMEII1 models performed better overall with respect to AQMEII2. An analysis of the LT
318 component showed that the bias in the AQMEII2 models is higher, possibly due to the 2010
319 emission inventory, while an analysis of the DU error found that the variance error in the
320 AQMEII2 models is significantly reduced with respect to the AQMEII1 models, and is almost
321 null. We postulate that the inclusion of feedback effects may have been beneficial, and that
322 the reduced performance of AQMEII2 models is likely due to external bias. The residual
323 *mMSE* error of the DU component (~1-2 ppb on average for both continents) is mostly likely
324 generated by a number of processes, including chemistry, cloudiness, boundary layer
325 transition and vertical mixing.

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

330 Finally, the error of the ID component is less than 1 ppb (on average ~0.2 ppb for AQMEII2)
331 and is generated by both variance (most commonly model over-prediction) and *mMSE*. The
332 fast-acting photochemical processes are, therefore, modelled with satisfactory precision.

333 4.3. SPATIAL DISTRIBUTION OF THE SPECTRAL ERROR COMPONENTS

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

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

344 The Po valley in Italy and the southern part of the EU are the most problematic areas,
345 affected by severe LT errors (Figure 3 and Figure 5). The central and northern parts of the EU are
346 less problematic, especially for AQMEII2. The other components of the error are
347 significantly smaller than the LT error, with some exceptions (especially for the DU



348 component). The length of the segment is in fact normalised to the largest error for each
349 model, to facilitate the interpretation and the relative weight of each error component.

350 Concerning NA (Figure 4 and Figure 6), the DU error has more weight and competes with the LT
351 error in the central and south-eastern parts of the continent. For AQMEII2, the SY error is as
352 significant as the LT error on the East Coast (Wrf/Chem, Figure 6c). The greatest LT error is
353 observed in the coastal areas (east and west) and across the north-eastern border between
354 the US and Canada (due primarily to model underestimation in the east and north, and
355 model overestimation in the west).

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

361 5. MSE DECOMPOSITION AND COMPLEXITY

362 In regression analysis and statistical learning theories, the problem of under- and over-
363 fitting complex systems is at the root of the MSE decomposition into bias and variance. The
364 trade-off between bias and variance is strictly dependent on the complexity of the model.
365 Over-fitting occurs when too many parameters and modules are added to the model: each
366 new module added to describe a process is a new source of variance due to internal
367 parameterisation and linearisation. In other words, over-fitting is associated with the
368 stochasticity inherent to the data/model, and contributes to the increase in variance and
369 consequent decrease in bias. Under-fitting occurs due to an oversimplification of the
370 modelled processes, and is an important source of bias as it is associated with the
371 deterministic property of the modelling activity (Hastie et al., 2009).

372 The problem of the bias-variance trade-off becomes markedly more complicated when
373 dealing with complex models with many degrees of freedom, such as AQ modelling systems.
374 Adding new modules to cope with unexplained physical processes can lead to a reduction in
375 the bias due to that specific process, but also feeds new variance and possibly new bias into
376 the model due to the non-linear interaction of the new module with existing ones, since
377 reducing the bias while preserving the variance is non-trivial. Rao (2005), in the context of
378 dispersion modelling, provided the theoretical variations of the total model uncertainty by
379 exploiting the components of the difference between the modelled and observed variance
380 (Figure 1 of Rao et al., 2005). Rao (2005) used the number of meteorological parameters in
381 the model as a measure of model complexity, and concluded that the optimal model
382 complexity could not be defined a priori, but is a trial-and-error combination of the model,
383 the measurement error and the stochastic uncertainty.

384 In this study we attempt to derive the curves of the MSE components as a function of model
385 complexity. Figure 7 shows an example of the approach used to break down model complexity



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

394 In Figure 8, we report the spatially averaged curves of bias, variance, and covariance according
395 to Eq 6 as a function of model complexity. According to the regression analysis theories
396 outlined above, we would expect the variance to increase according to the complexity
397 ($\frac{d\sigma_m^2}{dcomplexity} > 0$), and the distance between the modelled and observed variance to
398 decrease ($\frac{d(\sigma_m - \sigma_o)^2}{dcomplexity} < 0$), and the opposite for the bias. The curves of variance in Figure 8
399 indeed turn downwards as predicted by the theory, while the curves of bias have a mixed
400 behaviour but are, basically, constant ($\frac{d(mod-obs)^2}{dcomplexity} \approx 0$). More specifically:

- 401 - The $(\sigma_m - \sigma_o)^2$ term decreases steadily but slowly to a timescale of ~ 1 day, after
402 which it drastically drops to significantly lower values. This indicates that *i*) the
403 complexity of the AQ systems increases exponentially at the DU timescales (not
404 entirely surprising, given the day/night behavioural properties of ozone); *ii*) the
405 efforts made to improve the model capabilities on the short-term processes
406 governing the ozone dynamics improve the model precision; *iii*) there is a possible
407 lack of parameterisation and modelling of the processes of transport and chemical
408 transformation over periods longer than 1-2 days.
- 409 - The fact that the bias varies only by small amounts indicates that a fully evolved
410 model, capable of reproducing processes at the shortest timescales (turbulent
411 dispersion, fast chemical reactions, even day/night variability, etc.) is no more
412 accurate than a basic model that only accounts for long-term processes. This might
413 indicate that *i*) the bias at the shorter timescales is introduced entirely by the larger
414 timescales, and/or *ii*) the bias is continuously fed into the model by an external
415 source acting at all scales, as for example the emissions data or boundary conditions.

416 In Figure S4 to Figure S7 we propose the same analysis as that in Figure 8 but replicated for all
417 receptors individually (with no spatial average). In most cases (both continents, both
418 AQMEII phases), the $(\sigma_m - \sigma_o)^2$ term decreases sharply after a timescale of resolved
419 processes of ~ 1 day; the bias term confirms the independency to complexity at all
420 receptors; the covariance is complementary to the variance.

421 5. CONCLUSIONS



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

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

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

- 438 - The bias component of the error is by far the most important source of error, and is
439 mainly associated with long-term processes and/or input fields (likely emissions data
440 or boundary conditions). With regard to the model application, any effort to improve
441 the current capabilities of AQ modelling systems are likely to have little practical
442 impact if this primary issue is not addressed and solved;
- 443 - Most relevant to model development, the variance error (the discrepancy between
444 modelled and observed variance) is mainly associated with the DU component. At
445 timescale of ~1-2 days, the complexity of modelling systems increases substantially
446 and many processes are involved; the fact that the variance error of the DU
447 component for the AQMEII2 runs is reduced with respect to the AQMEII1 runs might
448 indicate the benefits of including feedback in the models. Such a conclusion could
449 not be drawn with simpler operational evaluation strategies;
- 450 - The limited magnitude of the variability of the SY and LT signals produces little
451 variance errors for these two components, and only becomes comparable to the LT
452 or DU error when the bias is negligible or the total MSE is small;
- 453 - The *mMSE* error is predominant in some instances of the analysed models, and is
454 due to the random distribution of modelled values. There are many causes of *mMSE*
455 error, including all 'internal' processes that produce non-systematic errors such as
456 noise, representativeness, the linearisation of non-linear process, and turbulence
457 closure;
- 458 - The analysis of the spatial distribution of the error highlights the diversity in the
459 behaviour of each modelling system. The common spatial structures of the LT error
460 (for example in the central and southern EU) may reveal common sources of error



461 (e.g. emissions data), while the error of the other components (especially DU and SY)
462 are peculiar to each model and need to be assessed individually.
463

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

468

469

470

471 APPENDIX

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

$$\begin{aligned}ID(t) &= \mathbf{x}(t) - kZ_{3,3}(\mathbf{x}(t)) \\DU(t) &= kZ_{3,3}(\mathbf{x}(t)) - kZ_{13,5}(\mathbf{x}(t)) \\SY(t) &= kZ_{13,5}(\mathbf{x}(t)) - kZ_{103,5}(\mathbf{x}(t)) \\LT(t) &= kZ_{103,5}(\mathbf{x}(t)) \\ \mathbf{x}(t) &= ID(t) + DU(t) + SY(t) + LT(t)\end{aligned}\tag{EQ. S.1}$$

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

475 A clear-cut separation of the components of EQ. S.1 cannot be achieved, as the separation is
476 a non-linear function of the parameters m and k (Rao et al., 1997). It follows that the
477 components of EQ. S.1 are not completely orthogonal and that some level of overlapping
478 energy exists (Kang et al., 2013). Galmarini et al. (2013) found that the explained variance by
479 the spectral components account for 75 to 80% of the total variance, the remaining portion
480 being explained by the interactions between the components.

481

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

$$\begin{aligned}MSE(O3) &= MSE(LT + SY + DU + ID) \\ &= \sum MSE(spec\ comp) + \sum MSE(cross\ comp)\end{aligned}\tag{EQ. S.2}$$

484

485 Where *spec comp* are the diagonal terms, and *LT*, *SY*, *DU*, *ID* and *cross comp* are the off-
486 diagonal terms deriving from the squared nature of the MSE: LT_oSY_m , SY_oLT_m , SY_oDU_m ,



487 $DU_oSY_m, DU_oID_m, ID_oDU_m, LT_mSY_m, LT_oSY_o, DU_mSY_m, DU_mID_m, DU_oSY_o, DU_oID_o$ (o and m
488 represent observed and modelled fields, respectively). For simplicity, the cross-components
489 are assumed to be symmetric, so the o and m subscripts are dropped. This simplification has
490 little impact on the MSE breakdown since, as shown by Galmarini et al. (2013), the diagonal
491 terms alone account for over 80% of the total variance.

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

$$MSE(O_3-LT(O_3)) = MSE(SY)+MSE(DU)+MSE(ID)+2MSE(IDDU)+2MSE(IDSY)+2MSE(DUSY) \quad \text{EQ. S.3}$$

494

495 By removing EQ. S.3 from EQ. S.2, the contribution of LT and its cross-component is isolated:

$$\text{EQ. S.2- EQ. S.3} = MSE(LT) + MSE(LTID) + MSE(LTSY) + MSE(LTDU) \quad \text{EQ. S.4}$$

496

497 We can further elaborate on EQ. S.4 to isolate the contribution of each cross-component.

498 For instance, the case of SYLT:

499

$$MSE(SY-ID-DU)-MSE(SY)-MSE(LT) = [MSE(SY)+MSE(LT)+ 2MSE(SYLT)] - MSE(SY) - MSE(LT) = 2MSE(SYLT) \quad \text{EQ. S.5}$$

500

501 The procedure in EQ. S.5 has been applied to derive the contribution of all cross-
502 components.

503

504 ACKNOWLEDGEMENTS

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

507

508

509

510



511 REFERENCES

- 512 Baklanov, A., and et al., 2014. Online coupled regional meteorology chemistry models in Europe: current status
513 and prospects. *Atmospheric Chemistry and Physics* 14, 317-398.
- 514 Brunner, D., Jorba, O., Savage, N., Eder, B., Makar, P., Giordano, L., Badia, A., Balzarini, A., Baro, R., Bianconi,
515 R., Chemel, C., Forkel, R., Jimenez-Guerrero, P., Hirtl, M., Hodzic, A., Honzak, L., Im, U., Knote, C., Kuenen,
516 J.J.P., Makar, P.A., Manders-Groot, A., Neal, L., Perez, J.L., Pirovano, G., San Jose, R., Savage, N., Schroder,
517 W., Sokhi, R.S., Syrakov, D., Torian, A., Werhahn, K., Wolke, R., van Meijgaard, E., Yahya, K., Zabkar, R.,
518 Zhang, Y., Zhang, J., Hogrefe, C., Galmarini, S., 2015. Evaluation of the meteorological performance of
519 coupled chemistry-meteorology models in phase 2 of the air quality model evaluation international
520 initiative. *Atmos. Environ*
- 521 Dennis, R., Fox, T., Fuentes, M., Gilliland, A., Hanna, S., Hogrefe, C., Irwin, J., Rao, S.T., Scheffe, R., Schere, K.,
522 Steyn, D., Venkatram, A., 2010. A framework for evaluating regional-scale numerical photochemical
523 modeling systems. *Environ. Fluid Mech. (Dordr.)* 10, 471-489. [http://dx.doi.org/10.1007/s10652-009-](http://dx.doi.org/10.1007/s10652-009-9163e2)
524 [9163e2](http://dx.doi.org/10.1007/s10652-009-9163e2).
- 525 Fox, D.G., 1981. Judging air quality model performance. *Bulletin of the American Meteorological Society* 62,
526 No.5, 599-609.
- 527 Galmarini, S. Solazzo, E., Im, U., Kioutsioukis, I., 2015. AQMEII 1, 2 and 3: Direct and Indirect Benefits of
528 Community Model Evaluation Exercises. 34th International Technical Meeting on Air Pollution Modelling
529 and its Application, Montpellier (France) 4-8 May 2015.
- 530 Galmarini, S., Kioutsioukis, I., Solazzo, E., 2013. E pluribus unum: ensemble air quality predictions. *Atmos.*
531 *Chem. Phys.* 13, 7153-7182.
- 532 Giordano, L., Brunner, D., Flemming, J., Hogrefe, C., Im, U., Bianconi, R., and et al., 2015. Assessment of the
533 MACC reanalysis and its influence as chemical boundary conditions for regional air quality modelling in
534 AQMEII-2. *Atmospheric Environment* 115, 371-388.
- 535 Hastie, T., Tibshirani, R., Friedman, J., 2009. *The elements of statistical learning (2nd edition)*. Springer-Verlag.
536 763 pages.
- 537 Hogrefe, C., Rao, S.T., Zurbenko, I.G., Porter, P.S., 2000. Interpreting the information in ozone observations and
538 model predictions relevant to regulatory policies in the Eastern United States. *Bull. Am. Meteorol. Soc.*
539 81, 2083e2106. [http:// dx.doi.org/10.1175/1520-0477\(2000\)0812.3.CO;2](http://dx.doi.org/10.1175/1520-0477(2000)0812.3.CO;2).
- 540 Hogrefe, C., Roselle, S., Mathur, R., Rao, S.T., Galmarini, S., 2014. Space-time analysis of the Air Quality Model
541 Evaluation International Initiative (AQMEII) phase 1 air quality simulation. *J. Air Waste Manag. Assoc.* 64,
542 388-405.
- 543 Im, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., Baro, R., Bellasio, R., Brunner, D.,
544 Chemel, C., Curci, G., Denier van der Gon, H., Flemming, J., Forkel, R., Giordano, L., Jimenez-Guerrero, P.,
545 Hirtl, M., Hodzic, A., Honzak, L., Jorba, O., Knote, C., et al., 2015a Evaluation of operational online-coupled
546 regional air quality models over Europe and North America in the context of AQMEII phase 2. Part II:
547 particulate matter. *Atmos. Environ.* 115, 421-441
- 548 Im, U., Bianconi, R., Solazzo, E., Kioutsioukis, I., Badia, A., Balzarini, A., Baro, R., Bellasio, R., Brunner, D.,
549 Chemel, C., Curci, G., Flemming, J., Forkel, R., Giordano, L., Jimenez-Guerrero, P., Hirtl, M., Hodzic, A.,
550 Honzak, L., Jorba, O., Knote, C., Kuenen, J. J.P., et al., 2015b. Evaluation of operational on-line-coupled



- 551 regional air quality models over Europe and North America in the context of AQMEII phase 2. Part I:
552 ozone. *Atmos. Environ.* 115, 404-420
- 553 Johnson, R. 2008 Assessment of Bias with Emphasis on Method Comparison. *Clin Biochem Rev Vol 29 Suppl (i)*
554 S37–S42.
- 555 Kang, D., Hogrefe, C., Foley, K.L., Napelenok, S.L., Mathur, R., Rao, S.T., 2013. Application of the Kolmogorov-
556 Zurbenko filter and the decoupled direct 3D method for the dynamic evaluation of a regional air quality
557 model. *Atmos. Environ.* 80, 58-69.
- 558 Kioutsioukis, I., Galmarini, S., 2014. De praeceptis ferendis: good practice in multi-model ensembles.
559 *Atmospheric Chemistry and Physics* 14, 11791–11815.
- 560 Makar, P.A., Gong, W., Hogrefe, C., and et al., 2015. Feedbacks between air pollution and weather, part 2:
561 effects on chemistry. *Atmospheric Environment* 115, 499-526
- 562 Murphy, A.H., 1988. Skill scores based on the mean square error and their relationship to the correlation
563 coefficient. *Monthly Weather Review* 116, 2417-2424
- 564 Pindyck, R.S., Rubinfeld, D.L., 1998. *Econometric Models and Economic Forecast*, Irwin/McGraw-Hill,
565 Singapore, 388 pg
- 566 Pouliot, G., Denier van der Gon, H., Kuenen, J., Makar, P., Zhang, J., Moran, M., 2015. Analysis of the emission
567 inventories and model-ready emission datasets of Europe and North America for phase 2 of the AQMEII
568 project. *Atmos. Environ.* 115, 345-360.
- 569 Pouliot, G., Pierce, T., Denier van der Gon, H., Schaap, M., Moran, M., and Nopmongcol, U., 2012. Comparing
570 Emissions Inventories and Model-Ready Emissions Datasets between Europe and North America for the
571 AQMEII Project. *Atmos. Environ.* 53, 4–14.
- 572 Rao, K.S., 2005. Uncertainty analysis in atmospheric dispersion modelling. *Pure and Applied Geophysics* 162,
573 1893-1917.
- 574 Rao, S.T., Galmarini, S., Puckett, K., 2011. Air quality model evaluation international initiative (AQMEII). *Bull.*
575 *Am. Meteorol. Soc.* 92, 23-30. <http://dx.doi.org/10.1175/2010BAMS3069.1>.
- 576 Rao, S.T., Zurbenko, I.G., Neagu, R., Porter, P.S., Ku, J.Y., Henry, R.F., 1997. Space and time scales in ambient
577 ozone data. *Bull. Am. Meteorol. Soc.* 78, 2153e2166. [http://dx.doi.org/10.1175/1520-0477\(1997\)0782.0.CO;2](http://dx.doi.org/10.1175/1520-0477(1997)0782.0.CO;2).
- 579 Schere, K., Flemming, J., Vautard, R., Chemel, C., Colette, A., Hogrefe, C., Bessagnet, B., Meleux, F., Mathur, R.,
580 Roselle, S., Hu, R.-M., Sokhi, R. S., Rao, S. T., and Galmarini, S.: Trace gas/aerosol concentrations and their
581 impacts on continental-scale AQMEII modelling sub-regions, *Atmos. Environ.*, 53, 38–50, 2012.
- 582 Solazzo, E., Bianconi, R., Vautard, R., Appel, K.W., Moran, M.D., Hogrefe, C., Bessagnet, B., Brandt, J.,
583 Christensen, J.H., Chemel, C., Coll, I., van der Gon, H.D., Ferreira, J., Forkel, R., Francis, X.V., Grell, G.,
584 Grossi, P., Hansen, A.B., Jericevic, A., Kraljevic, L., Miranda, A.I., Nopmongcol, U., Pirovano, G., Prank, M.,
585 Riccio, A., Sartelet, K.N., Schaap, M., Silver, J.D., Sokhi, R.S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G.,
586 Zhang, J., Rao, S.T., Galmarini, S., 2012a. Model evaluation and ensemble modelling and for surface-level
587 ozone in Europe and North America. *Atmos. Environ.* 53, 60-74.
- 588 Solazzo, E., Bianconi, R., Pirovano, G., Matthias, V., Vautard, R., Moran, M.D., Appel, K.W., Bessagnet, B.,
589 Brandt, J., Christensen, J.H., Chemel, C., Coll, I., Ferreira, J., Forkel, R., Francis, X.V., Grell, G., Grossi, P.,



- 590 Hansen, A.B., Hogrefe, C., Miranda, A.I., Nopmongco, U., Prank, M., Sartelet, K.N., Schaap, M., Silver, J.D.,
591 Sokhi, R.S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S.T., Galmarini, S., 2012b.
592 Operational model evaluation for particulate matter in Europe and North America. Atmos. Environ. 53,
593 75-92.
- 594 Solazzo, E., Bianconi, R., Pirovano, G., Moran, M., Vautard, R., Hogrefe, C., Appel, K.W., Matthias, V., Grossi, P.,
595 Bessagnet, B., Brandt, J., Chemel, C., Christensen, J.H., Forkel, R., Francis, X.V., Hansen, A., McKeen, S.,
596 Nopmongcol, U., Prank, M., Sartelet, K.N., Segers, A., Silver, J.D., Yarwood, G., Werhahn, J., Zhang, J., Rao,
597 S.T., Galmarini, S., 2013a. Evaluating the capabilities of regional scale air quality models to capture the
598 vertical distribution of pollutants. Geophys. Model Dev. 6, 791-818.
- 599 Solazzo, E., Riccio, A., Kioutsioukis, I., Galmarini, S., 2013b. *Pauci ex tanto numero*: reduce redundancy in multi-
600 model ensemble. Atmos. Chem. Phys. 13, 8315-8333.
- 601 Solazzo, E., Galmarini, S., 2015. Comparing apples with apples: Using spatially distributed time series of
602 monitoring data for model evaluation. Atmos. Environ. 112, 234-245
- 603 Stoeckenius, T.E., Hogrefe, C., Zagunis, J., Sturtz, T.M., Wells, B., Sakulyanontvittaya, T., 2015. A comparison
604 between 2010 and 2006 air quality and meteorological conditions, and emissions and boundary
605 conditions used in simulations of the AQMEII2 North American domain. Atmospheric Environment, 115,
606 389-403.
- 607 Theil, H., 1961. Economic forecast and policy. North-Holland, Amsterdam
- 608 Willmott, C.J., and et al., 1985. Statistics for the evaluation and comparison of models. Journal of Geophysical
609 research 90, No. C5, 8995-9005.
- 610 Wise, E.K., Comrie, A.C., 2005. Extending the KZ filter: application to ozone, particulate matter, and
611 meteorological trends. J. Air Waste Manag. Assoc. 55 (8), 1208e1216.
- 612 Zurbenko, I.G., 1986. The Spectral Analysis of Time Series. North-Holland, Amsterdam, 236 pp.
- 613
- 614
- 615
- 616
- 617
- 618
- 619
- 620
- 621
- 622
- 623
- 624



625 FIGURES

626 **Figure 1.** Share (in %) of the total MSE in the main spectral components and the cross components (see Appendix for
627 detail) for *a)* AQMEI1 and *b)* AQMEI2. Top panel: EU; lower panel: NA.

628 **Figure 2.** MSE (ppb^2) breakdown in bias, variance and mMSE of the spectral components ID, DU, SY, LT, based on Eq 9. The
629 sign within the share of bias and variance indicates model overestimation (+) or underestimation (-) of mean concentration
630 (bias) and variance. *a)* AQMEI1 and *b)* AQMEI2. Top panel: EU; lower panel: NA.

631 **Figure 3** Spatial distribution of the MSE in the spectral components for the EU models of AQMEI1. The segments are
632 centred at the rural receptors' position (clockwise from north: MSE of ID, DU, SY, and LT). Their length is proportional to
633 the MSE magnitude, coded according to the colour scale. For each model, the colour scale extends from zero up to the 75th
634 percentile, and the last value of the scale is the maximum MSE. The colour of the MSE values above the 75th percentile
635 represents the maximum value. The thick dashed LT segment indicates model underestimation (low model bias).

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

637 **Figure 5.** As in **Figure 3**, but for the EU models of AQMEI2.

638 **Figure 6** As in **Figure 3**, but for the NA models of AQMEI2.

639 **Figure 7** Example of the model complexity as time-resolved scale of the transport and dispersion processes: the minimum
640 complexity (far right) is a poor time-resolving time series obtained as $kz(250,5)$. The complexity increases towards the left,
641 with the scale of resolved processes becoming finer up to the maximum complexity (far left), which represents the full time
642 series.

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

646
647

648
649
650
651

652

653
654
655

656

657

658

659

660

661

662



663 TABLES

664 **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)W; (30,65)N		(-125,-55)W; (26,51)N	
Number of receptors (min validity=75%; max altitude = 1000 m)	1339	1360	672	652

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685



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

688 a)

Model			Grid(km)	Emissions	Chemical BC
Code	Met	AQ			
EUROPE					
DK1	MM5	DEHM	50	Global emission databases, EMEP	Satellite measurements
FR3	MM5	Polyphemus	24	Standard [§]	Standard
HR1	PARLAM-PS	EMEP	50	EMEP model	From ECMWF and forecasts
UK2	WRF	CMAQ	18	Standard [§]	Standard
US4	WRF	WRF/Chem	22.5	Standard [§]	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 setup (no boundary conditions)
NL1	ECMWF	Lotos-EUROS	25	Standard [§]	Standard
DE1	COSMO	Muscat	24	Standard [§]	Standard
US3	MM5	CAMx	15	MEGAN, Standard	Standard
DE3	COSMO-CLM	CMAQ	24	Standard [§]	Standard
NORTH AMERICA					
CA1	GEM	AURAMS	45	Standard*	Climatology
PL1	GEM	GEM-AQ	25	Standard over AQMEII region; Global EDGAR/GEIA over the rest of the global domain	Global variable grid setup (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 [§]	Standard
ES3	WRF	WRF/Chem	23	Standard	Standard

689
 690

[§] Standard anthropogenic emissions and biogenic emissions derived from meteorology (temperature and solar radiation) and land use distribution implemented in the meteorological driver.



691 *Standard anthropogenic inventory but independent emission processing, exclusion of wildfires, and different versions of BEIS(v3.09)
 692 used.
 693 Refer to Solazzo et al. (2012a-b) and references therein for details.
 694
 695

b)

Model			Grid	Emissions	Chemical BC
Code	Met	AQ			
EUROPE					
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° x 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					
ES1	WRF	WRF/Chem	36 km	Standard	Standard
US6	WRF	CMAQ	12km	Standard	Standard
CA2f	GEM	MACH	15 km	Standard	Standard

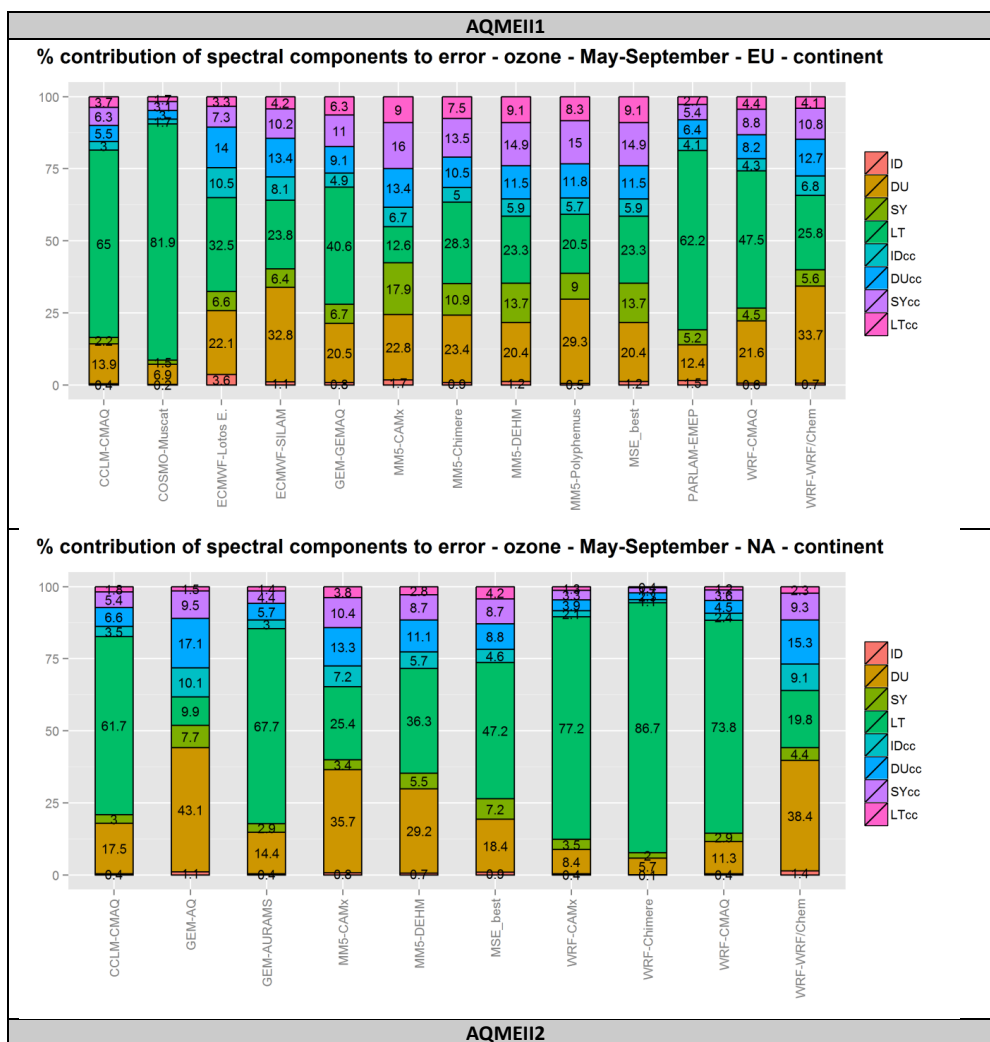
696 Standard Boundary conditions: 3-D daily chemical boundary conditions were provided by the ECMWF IFS-MOZART model run in the
 697 context of the MACC-II project (Monitoring Atmospheric Composition and Climate - Interim Implementation) at 3-hourly and 1.125 spatial
 698 resolution. Refer to Im et al. (2015a-b) for details.

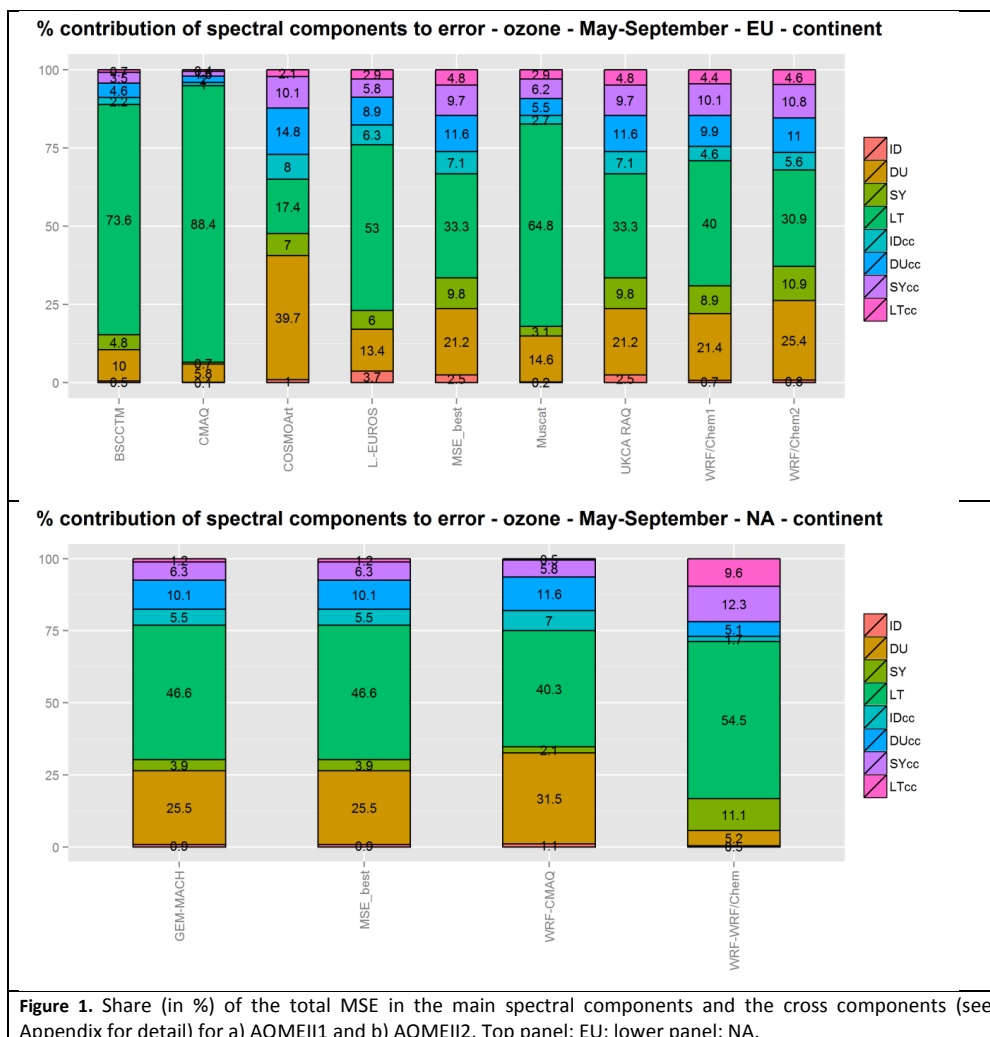
699 Standard Emissions: based on the TNO-MACC-II (Netherlands Organization for Applied Scientific Research, Monitoring Atmospheric
 700 Composition and Climate - Interim Implementation) framework for Europe and by the US EPA (Environmental Protection Agency) and
 701 Environment Canada for North America. The 2008 National Emissions Inventory (<http://www.epa.gov/ttn/chief/net/2008inventory.html>)
 702 and the 2008 Emissions Modeling Platform (<http://www.epa.gov/ttn/chief/emch/index.html#2008>) with year-specific updates for 2006
 703 and 2010 were used for the US portion of the modelling domain. Canadian emissions were derived from the Canadian National Pollutant
 704 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.
 705

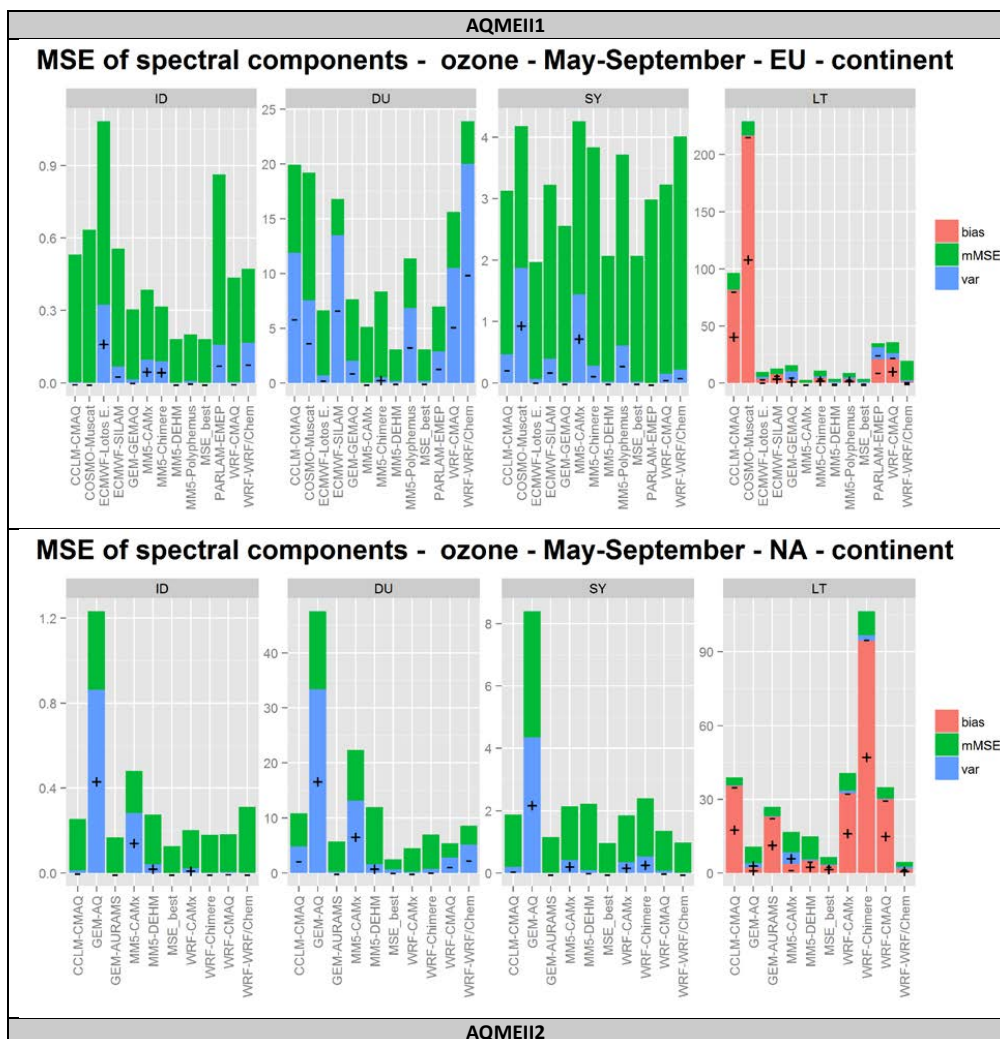


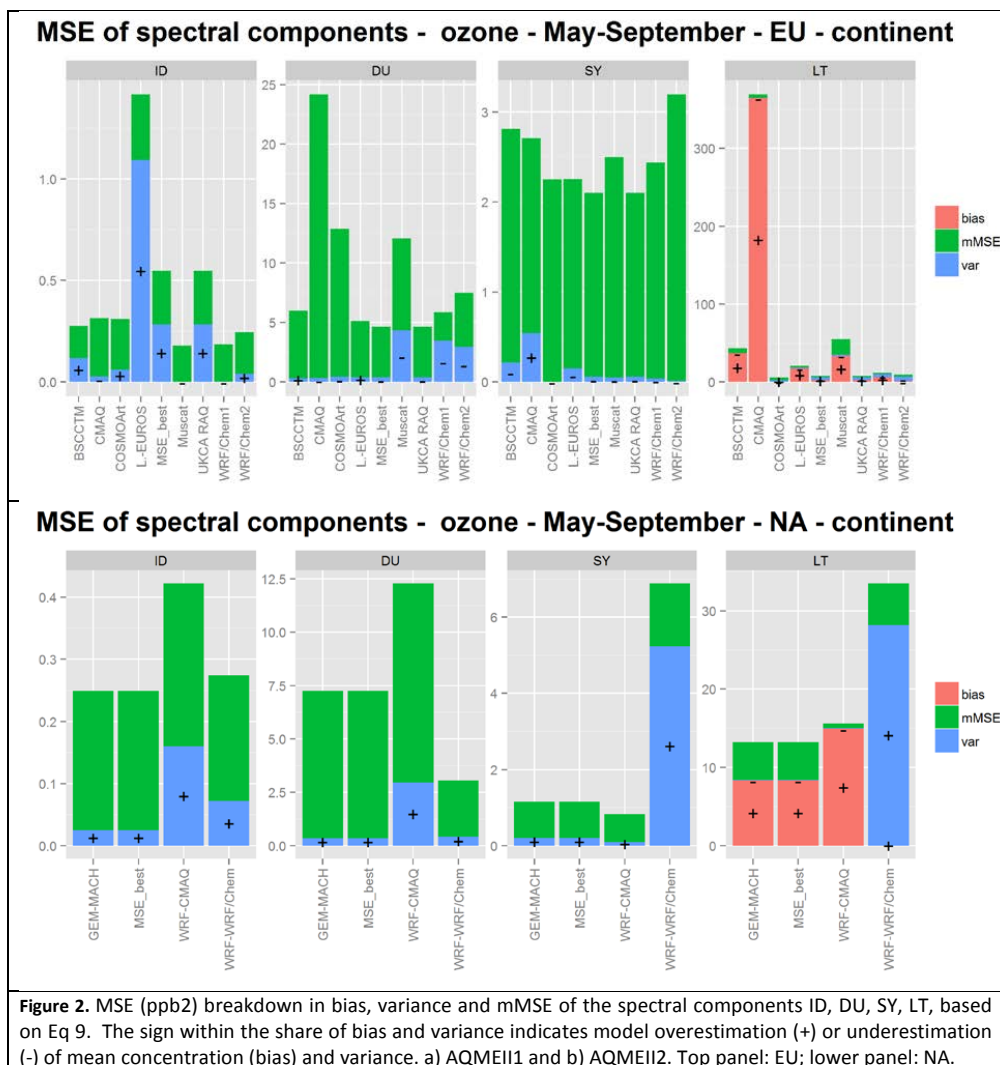
ERROR APPORTIONMENT FOR ATMOSPHERIC CHEMISTRY-TRANSPORT MODELS. A NEW APPROACH TO
 MODEL EVALUATION, BY E. SOLAZZO, S. GALMARINI

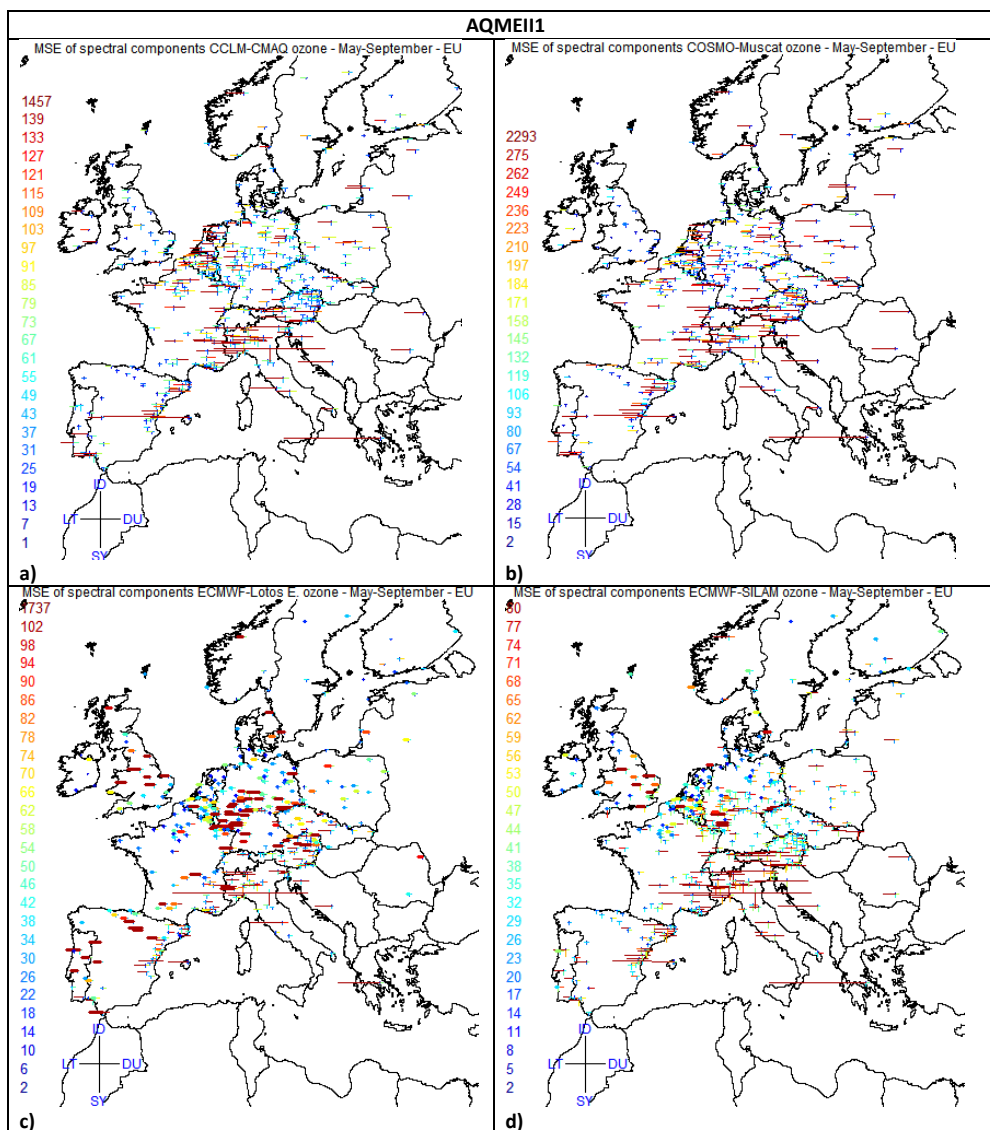
FIGURES

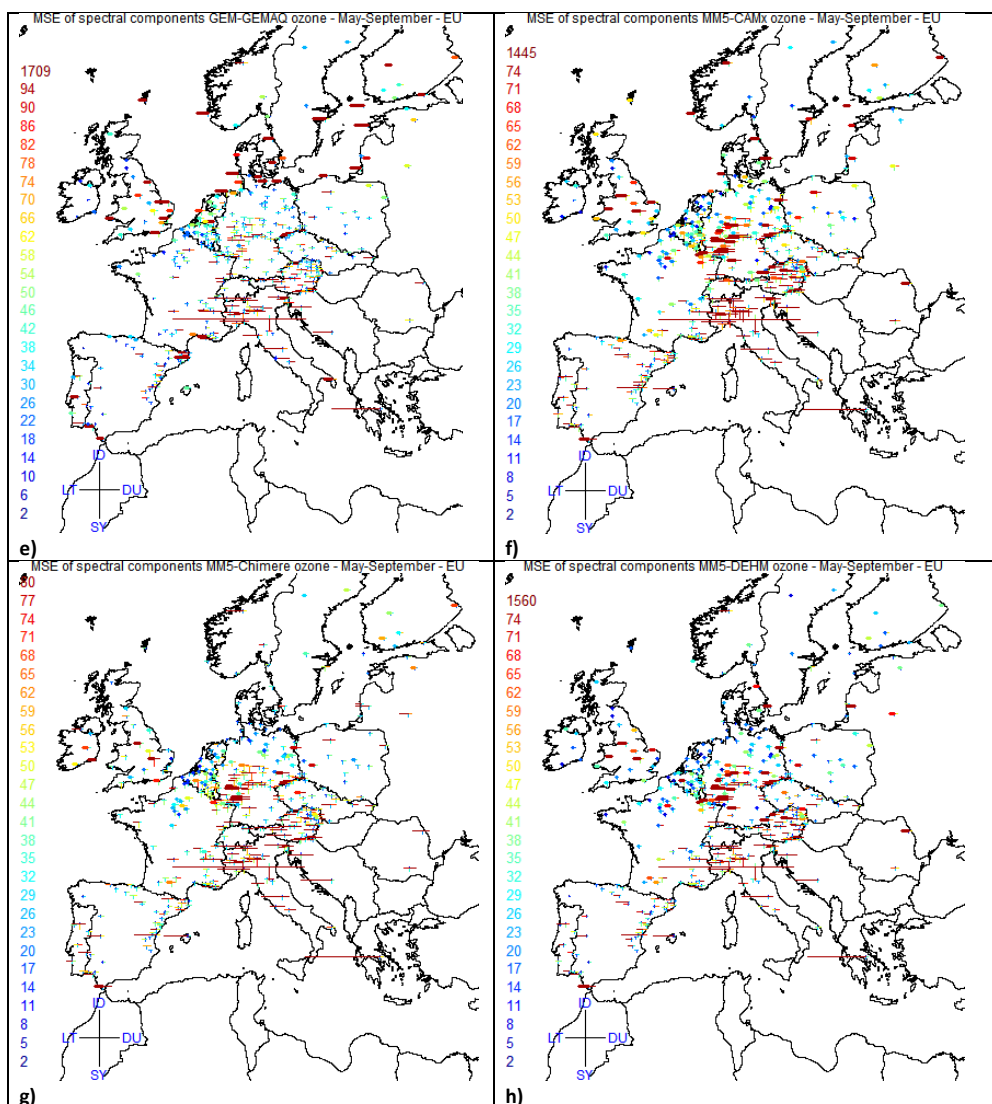












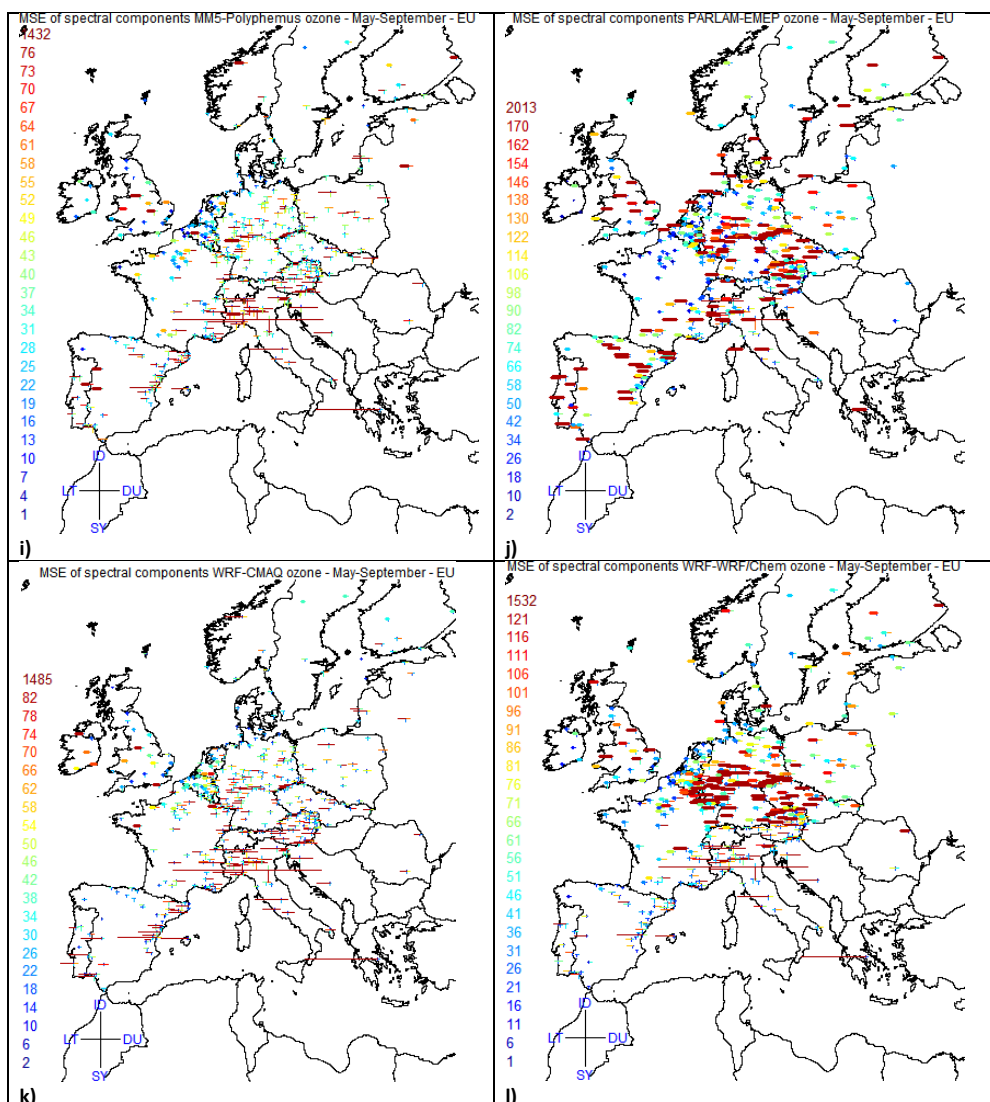
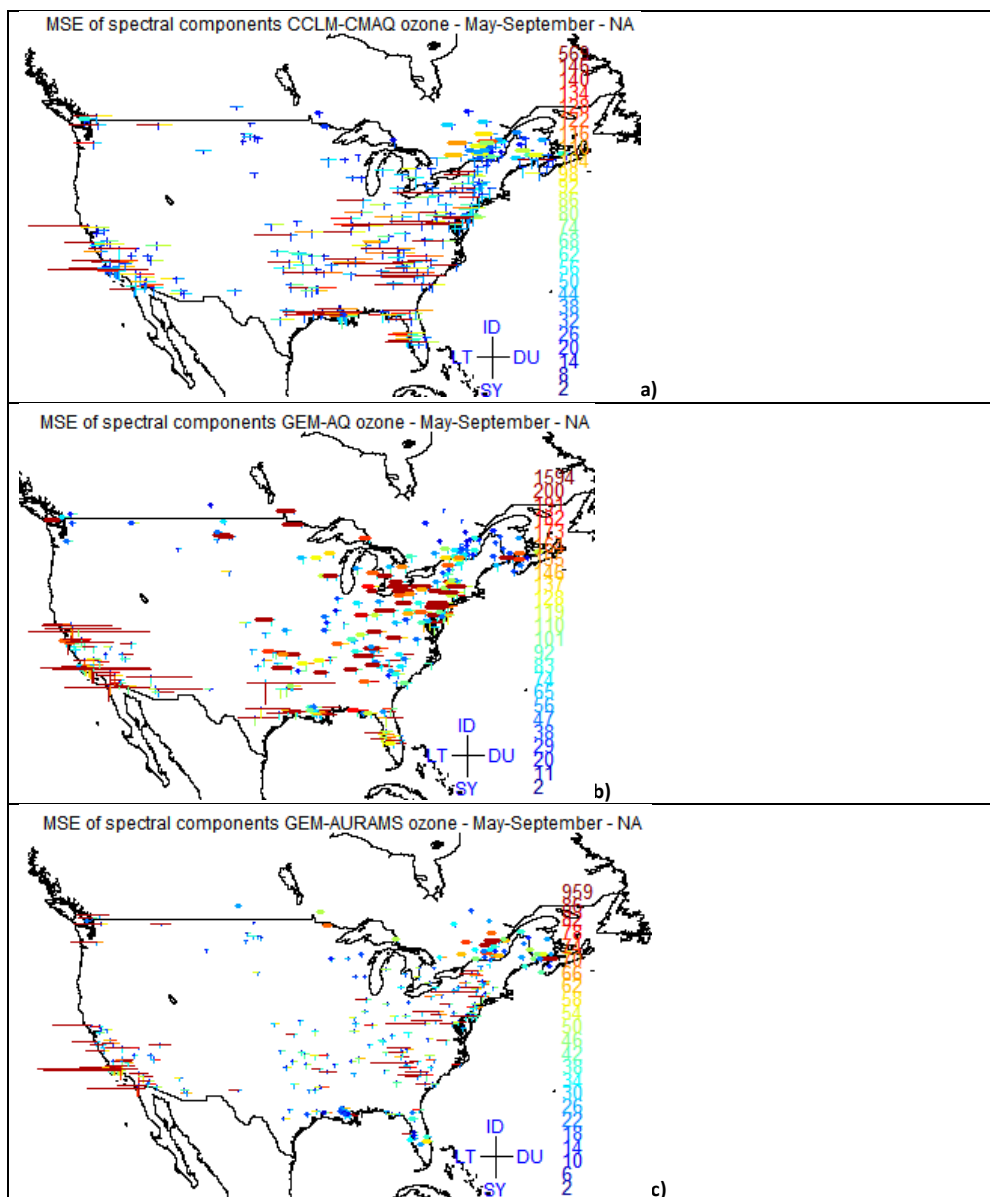
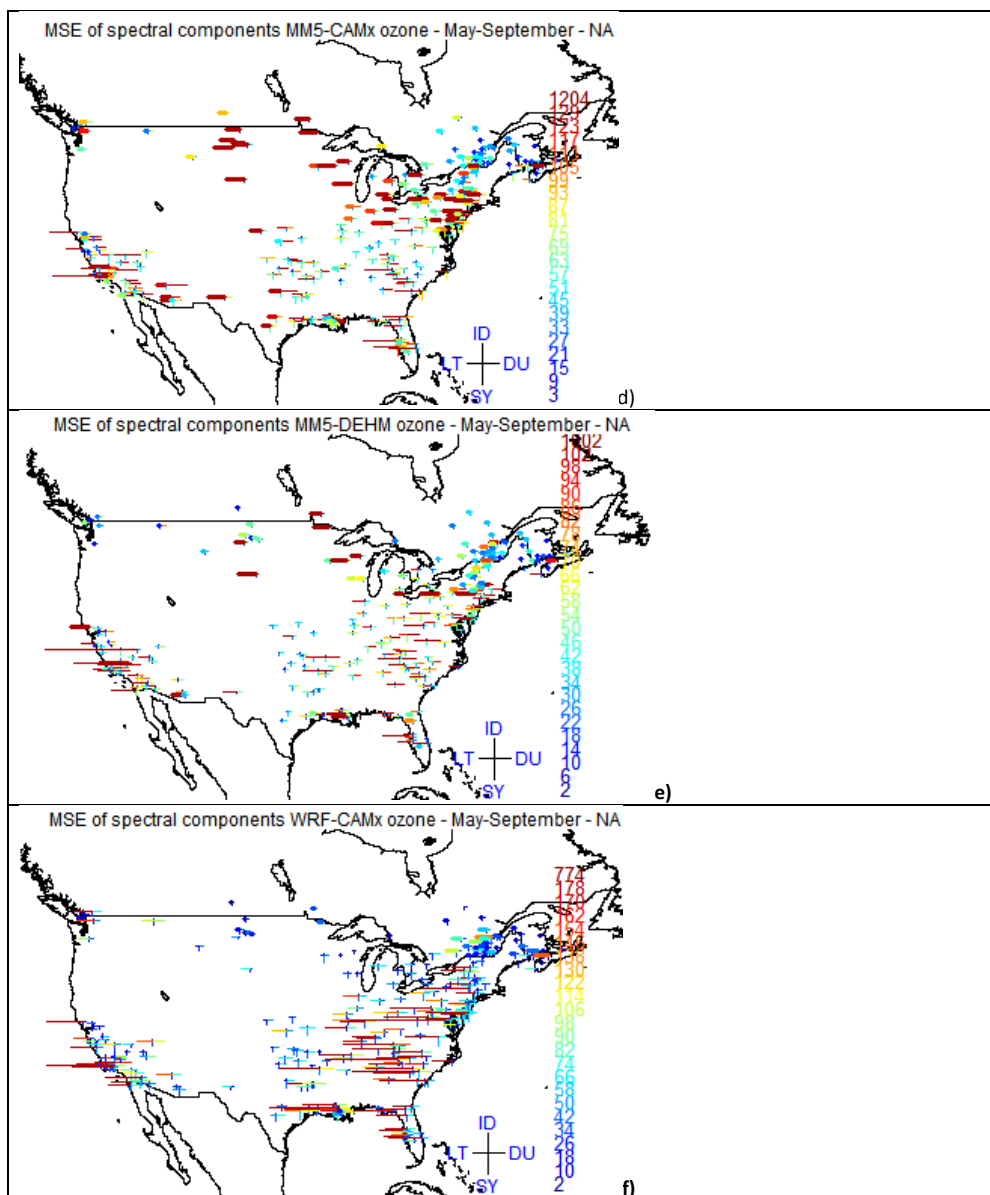
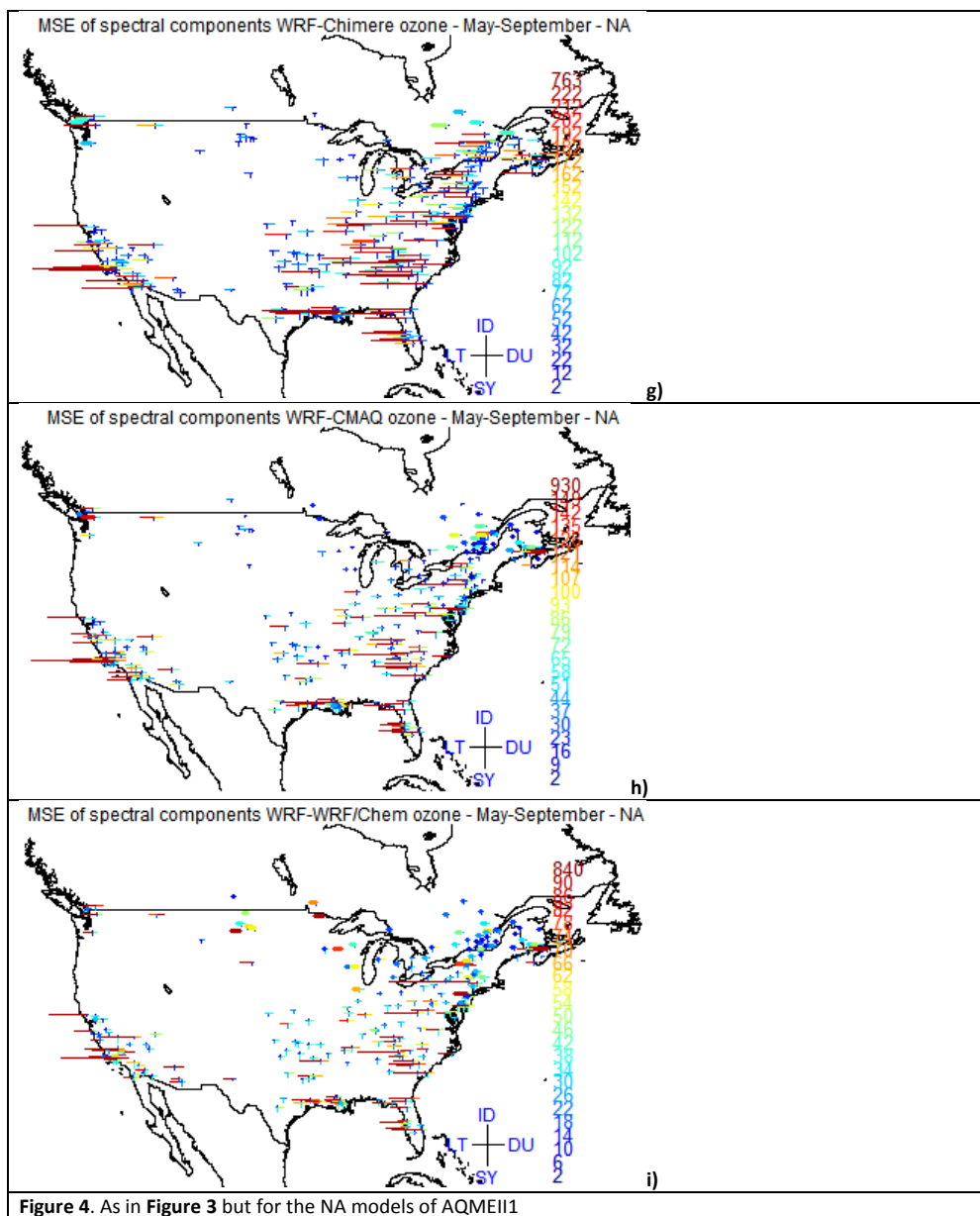
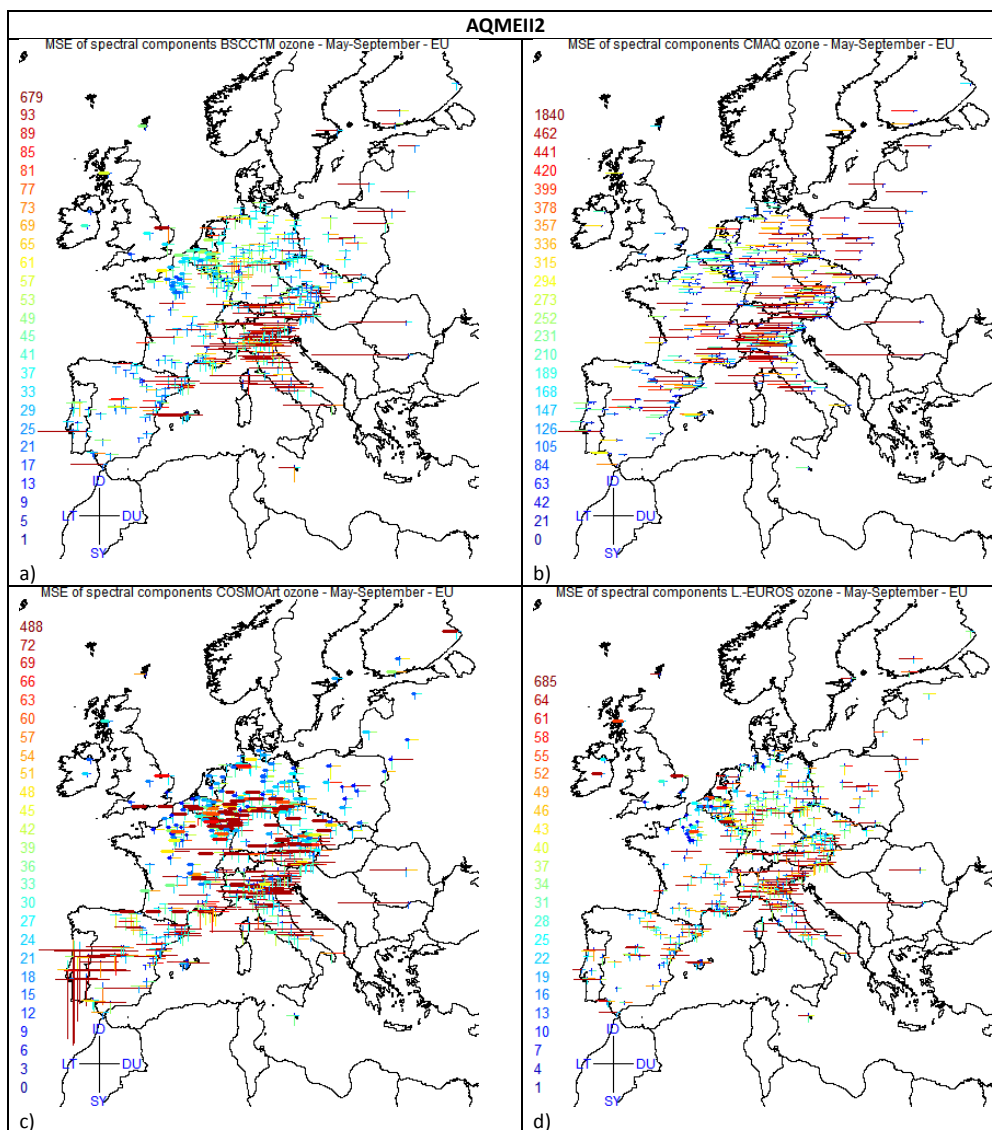


Figure 3. Spatial distribution of the MSE in the spectral components for the EU models of AQMEII1. The segments are centred at the rural receptors' position (clockwise from north: MSE of ID, DU, SY, and LT). Their length is proportional to the MSE magnitude, coded according to the colour scale. For each 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).









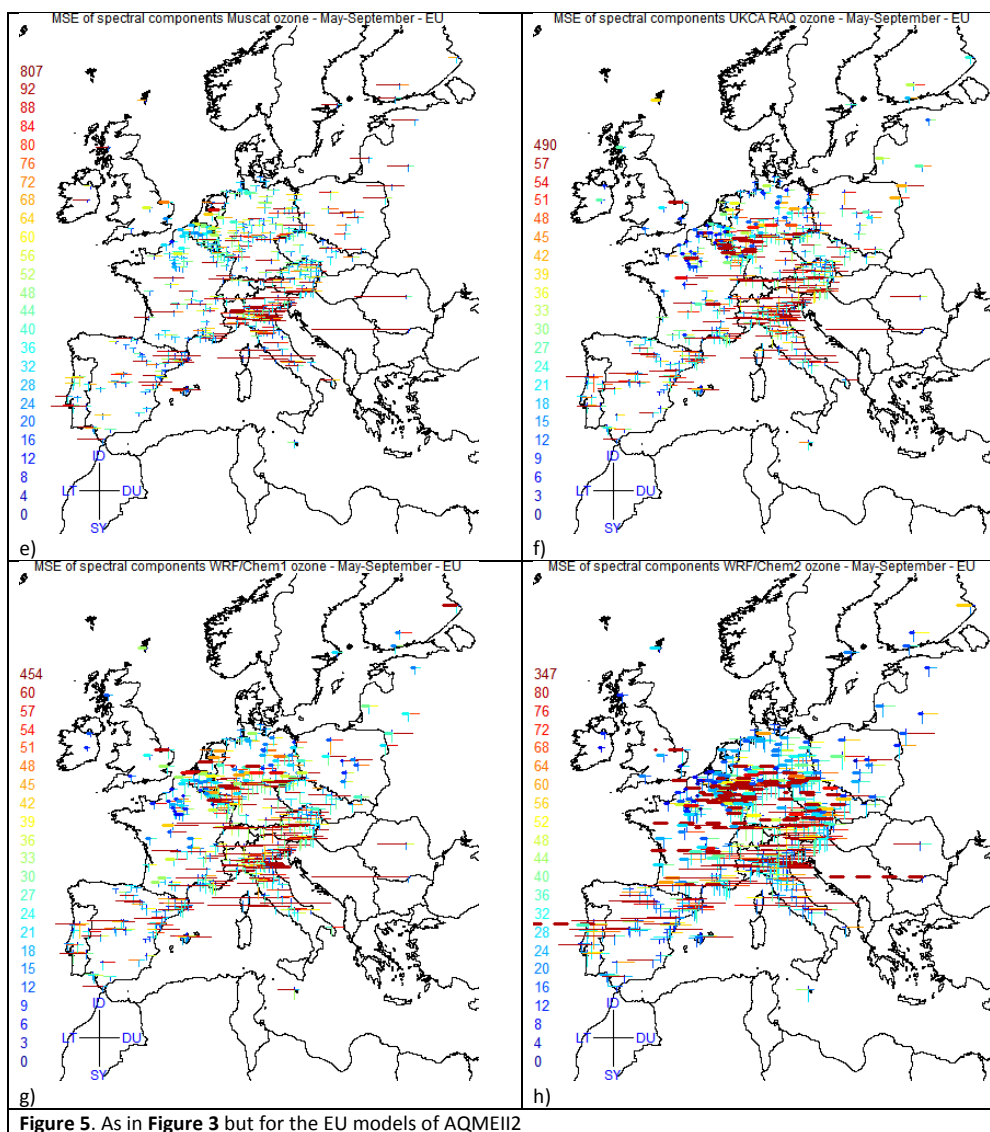
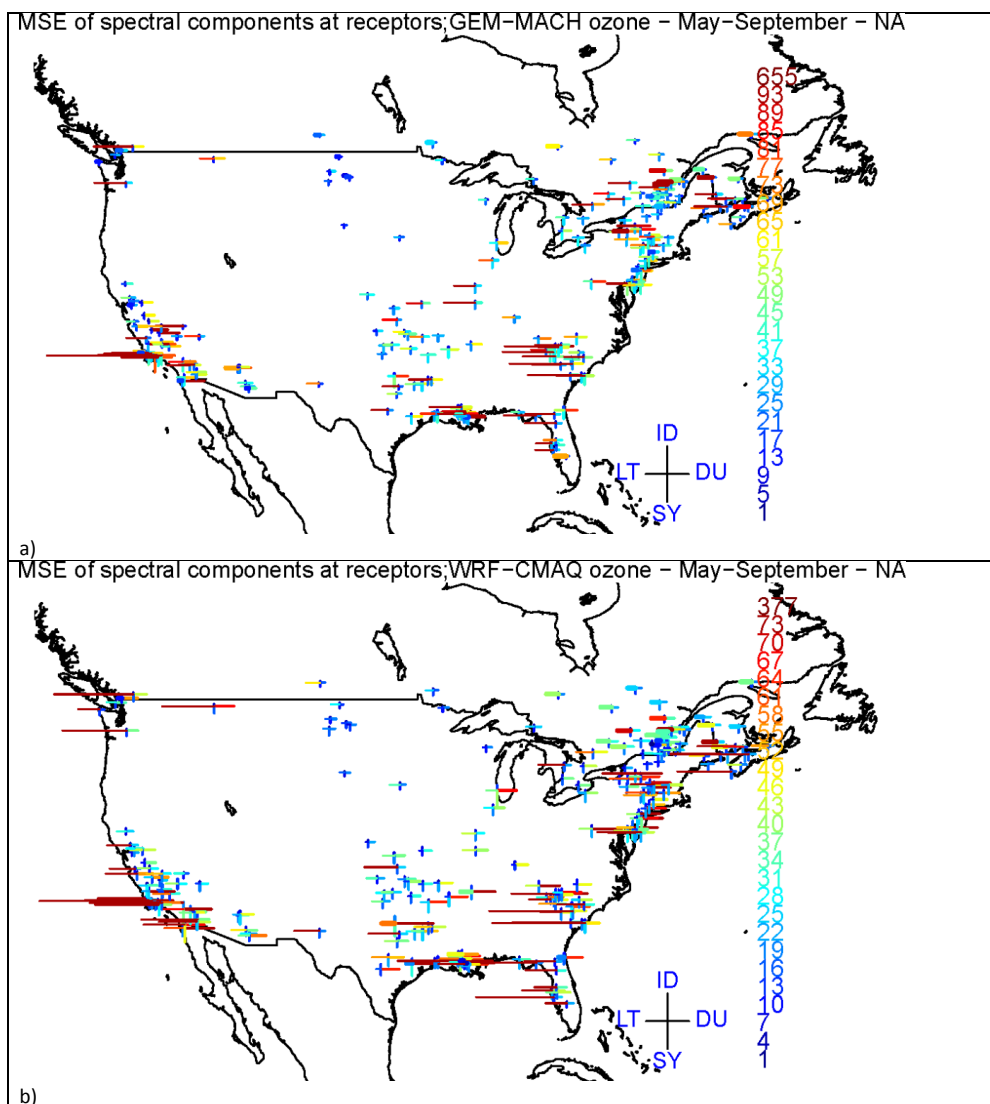
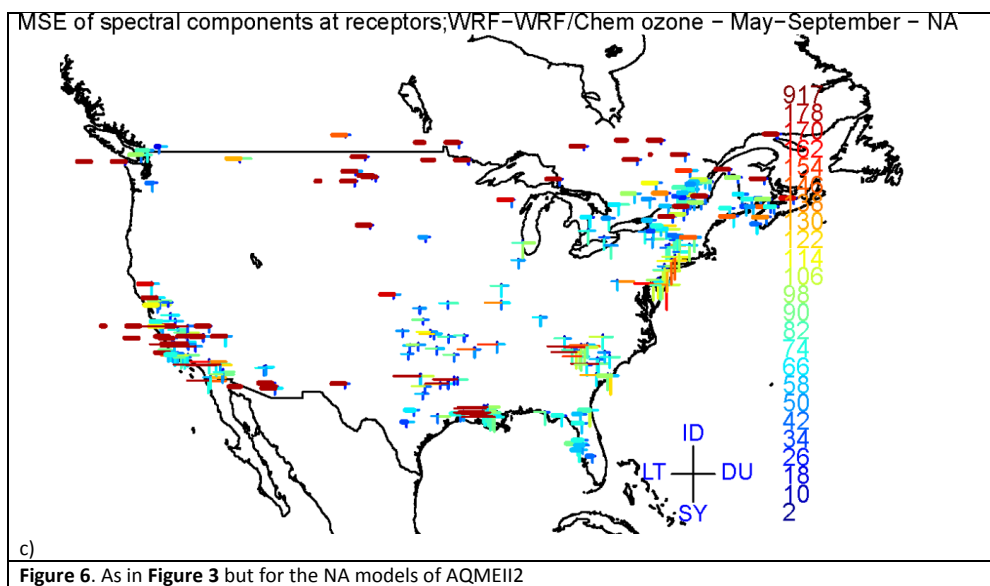


Figure 5. As in Figure 3 but for the EU models of AQMEII2





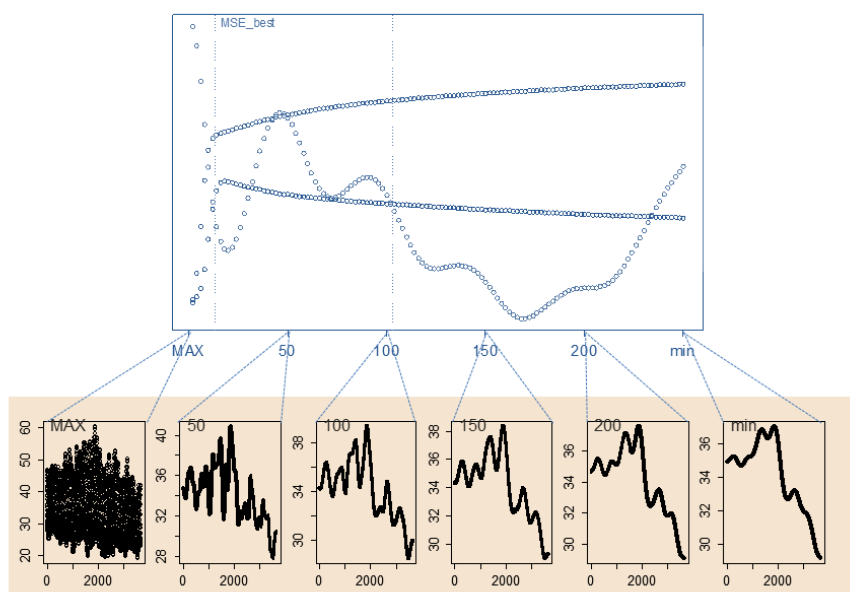
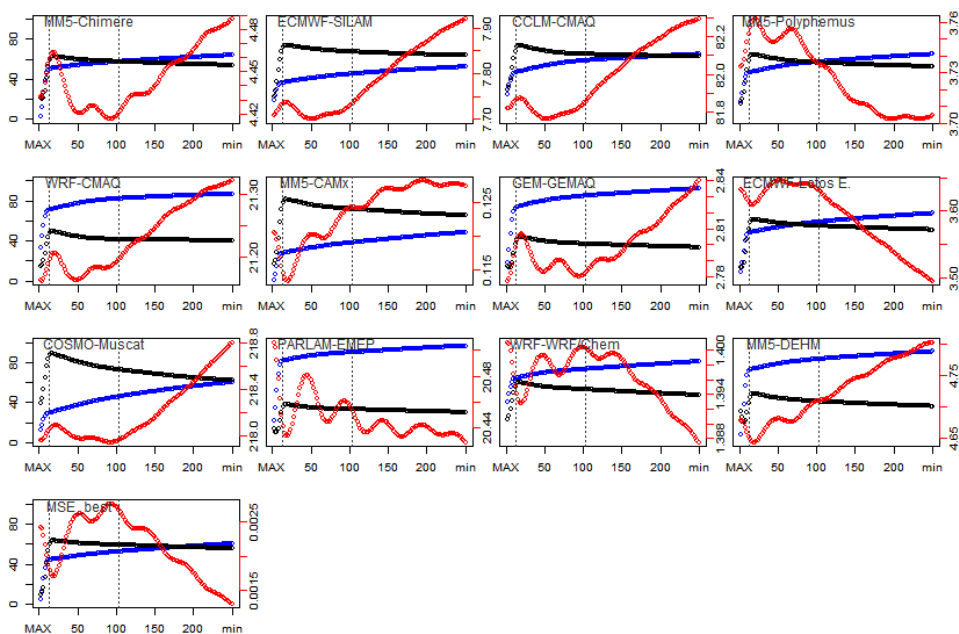


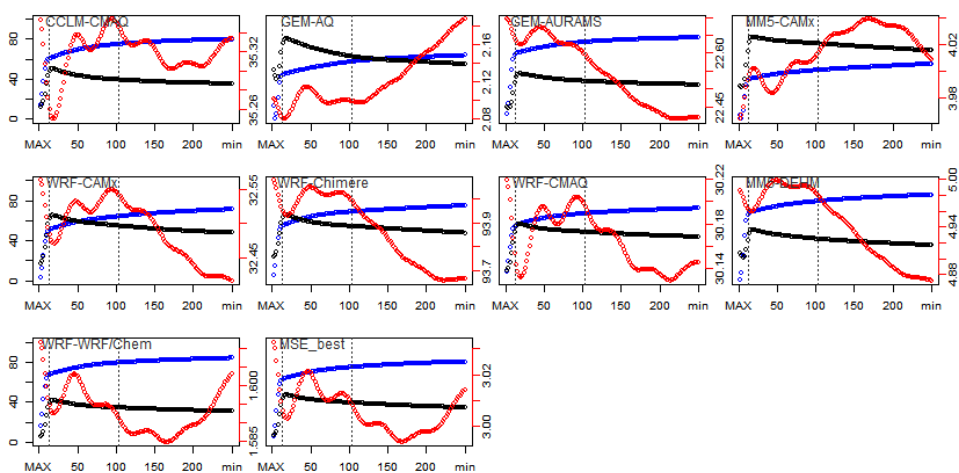
FIGURE 7 Example of the model complexity as time-resolved scale of the transport and dispersion processes: the minimum complexity (far right) is a poor time-resolving time series obtained as $kz(250,5)$. The complexity increases towards the left, with the scale of resolved processes becoming finer up to the maximum complexity (far left), which represents the full time series.



TIME-RESOLVED ERROR COMPONENTS (PPB2) - AQMEI1 - ozone - May-September - EU



TIME-RESOLVED ERROR COMPONENTS (PPB2) - AQMEI1 - ozone - May-September - NA



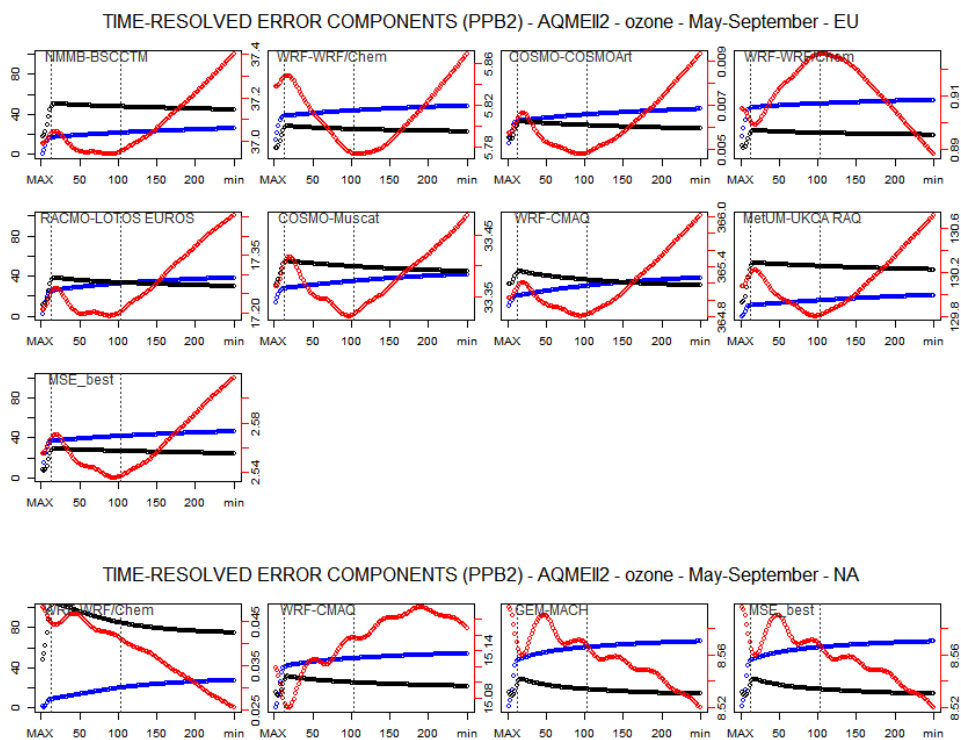


FIGURE 8 Evolution of error components (red: bias; Blue: variance; Black: covariance) as a function of model complexity. Complexity increases from left (min.) to right (max.) and is calculated as the temporal scale of the resolved process using the k_z filter on the modelled signal: $k_z(i,5)$, $i=2,\dots,250$.