



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53

Two-scale multi-model ensemble: Is a hybrid ensemble of opportunity telling us more?

Stefano Galmarini¹, Ioannis Kioutsioukis¹⁸, Efisio Solazzo¹, Ummugulsum Alyuz⁶, Alessandra Balzarini⁷, Roberto Bellasio², Anna M. K. Benedictow²², Roberto Bianconi², Johannes Bieser⁹, Joergen Brandt¹⁰, Jens H. Christensen¹⁰, Augustin Colette¹¹, Gabriele Curci^{4,5}, Yanko Davila²², Xinyi Dong²⁰, Johannes Flemming¹⁹, Xavier Francis¹², Andrea Fraser¹³, Joshua Fu²⁰, Daven Henze²¹, Christian Hogrefe³, Ulas Im¹⁰, Marta Garcia Vivanco¹⁴, Pedro Jiménez-Guerrero⁸, Jan Eiof Jonson²², Nutthida Kitwiroon¹⁶, Astrid Manders¹⁵, Rohit Mathur³, Laura Palacios-Peña⁸, Guido Pirovano⁷, Luca Pozzoli^{6,1}, Marie Prank¹⁷, Martin Schultz²², Rajeet S. Sokhi¹², Kengo Sudo²⁴, Paolo Tuccella⁵, Toshihiko Takemura²³, Takashi Sekiya²⁴, Alper Unal⁶

1 European Commission, Joint Research Centre (JRC, Ispra (VA), Italy

2 Enviroware srl, Concorezzo, MB, Italy

3 Computational Exposure Division - NERL, ORD, U.S. EPA

4 CETEMPS, University of L'Aquila, Italy

5 Dept. Physical and Chemical Sciences, University of L'Aquila, Italy

6 Eurasia Institute of Earth Sciences, Istanbul Technical University, Turkey

7 Ricerca sul Sistema Energetico (RSE SpA), Milano, Italy

8 University of Murcia, Department of Physics, Physics of the Earth, Facultad de Química, Campus de Espinardo, 30100 Murcia, Spain

9 Institute of Coastal Research, Chemistry Transport Modelling Group, Helmholtz-Zentrum Geesthacht, Germany

10 Aarhus University, Department of Environmental Science, Frederiksborgvej 399, 4000 Roskilde, Denmark

11 INERIS, Institut National de l'Environnement Industriel et des Risques, Parc Alata, 60550 Verneuil-en-Halatte, France

12 Centre for Atmospheric and Instrumentation Research (CAIR), University of Hertfordshire, Hatfield, UK

13 Ricardo Energy & Environment, Gemini Building, Fermi Avenue, Harwell, Oxon, OX11 0QR, UK

14 CIEMAT, Avda. Complutense, 40. 28040. Madrid, Spain

15 Netherlands Organization for Applied Scientific Research (TNO), Utrecht, The Netherlands

16 Environmental Research Group, Kings' College London, London, United Kingdom

17 Finnish Meteorological Institute, Atmospheric Composition Research Unit, Helsinki, Finland

18 University of Patras, Physics Department, Laboratory of Atmospheric Physics, 26504 Rio, Greece

19 European Centre for Medium-Range Weather Forecasts, Reading, UK

20 Department of Civil and Environmental Engineering, The University of Tennessee, Knoxville, TN, 37919, USA

21 Department of Mechanical Engineering, University of Colorado, 1111 Engineering Drive, Boulder, CO, USA.

22 Norwegian Meteorological Institute, Oslo, Norway

23 Research Institute for Applied Mechanics, Kyushu University, Fukuoka, Japan

24 Japan Agency for Marine-Earth Science and Technology, Yokohama, Japan



54

55

56

Abstract

57

58

59

60

61

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

In this study we introduce a *hybrid ensemble* consisting of air quality models operating at both the global and regional scale. The work is motivated by the fact that these different types of models treat specific portions of the atmospheric spectrum with different levels of detail and it is hypothesized that their combination can generate an ensemble that performs better than mono-scale ensembles. A detailed analysis of the hybrid ensemble is carried out in the attempt to investigate this hypothesis and determine the real benefit it produces compared to ensembles constructed from only global scale or only regional scale models. The study utilizes 13 regional and 7 global models participating in the HTAP2/AQMEI3 activity and focuses on surface ozone concentrations over Europe for the year 2010. Observations from 405 monitoring stations are used for the evaluation of the ensemble performance. The analysis first compares the modelled and measured spectra and then assesses the properties of the mono-scale ensembles, particularly their level of redundancy, in order to inform the process of constructing the hybrid ensemble. The main conclusion of this study is that the improvements obtained by the hybrid ensemble relative to the mono-scale ensembles can be attributed to its hybrid nature. Moreover, the optimal set is constructed from an equal number of global and regional models at only 15% of the stations. Finally, the study reaffirms the importance of an in-depth inspection of any ensemble of opportunity in order to extract the maximum amount of information and to have full control over the data used in the construction of the ensemble.



80 **1. Introduction**

81 It has been widely demonstrated (e.g Potempsky and Galmarini, 2009) that when
82 multiple model results are distilled to retain only original and independent
83 contributions (Solazzo et al. 2012) and thereafter statistically combined in what is
84 usually called an ensemble, one obtains results that are systematically superior to
85 the performance of the individual models and therefore can provide more accurate
86 and robust assessments or predictions.

87

88 An additional advantage of using an ensemble treatment resides in the fact that the
89 multiplicity of the results also quantifies the spread of the model solutions, which
90 provides useful information for the subsequent use of the model predictions for
91 planning purposes or more generically decision-making as it is a measure of the
92 variability of the options, scenarios or simply predictions.

93

94 When using ensembles in the realm of air quality modeling and atmospheric
95 dispersion, the general tendency is to combine results of models that belong to the
96 same category. Especially when referring to ensembles of opportunity (e.g.
97 Galmarini et al. 2004; Tebaldi and Knutti al. (2007); Potempsky and Galmarini, 2009,
98 Solazzo et al. 2012; Solazzo and Galmarini, 2015), which combine results from
99 different models applied to the same case study, it is customary to consider as
100 members those obtained from a homogeneous group of models. In particular, the
101 scale at which models operate seems to be a discriminant in all such studies that
102 have been performed to date. Therefore, meso-, regional-, and global-scale model
103 results are grouped in ensembles according to their scale of pertinence. In air quality
104 studies, this has been the case for example in Fiore et al. (2009), Solazzo et al.
105 (2012), Kioutsoukis and Galmarini (2014), and Kioutsoukis et al (2016). Colette et al.
106 (2012) analyzed as part of an analysis of the exposure in Europe, results from an
107 ensemble of opportunity of a total of 6 models, 3 of which where global and 3
108 regional. The focus however was not the analysis of the contribution of neither the
109 hybrid character of the group to the ensemble result nor the role of redundancy and
110 reducibility of the set, but more obtaining a robust assessment of the 2030 air
111 quality in Europe. A potential benefit of the mixed ensemble was spelled out there



112 but never verified in line with the opportunity character of the grouping. Therefore
113 there is no record in the literature of a study of an ensemble of models working at
114 different scales.

115

116 When developing a model, the scale selection is deeply rooted in the approach to
117 atmospheric modeling and it finds a theoretical justification in the alleged scale-
118 separation shown in the energy spectrum of dynamic variables such as horizontal or
119 vertical wind velocities (Van der Hoven, 1957). Although it is now well accepted that
120 the assumed scale separation does not have general validity, (e.g. Galmarini et al.
121 1999, Pielke, 2013) and especially not for scalars (e.g. Galmarini et al., 2000;
122 Michelutti et al., 1999; Jonker et al., 1999; Jonker et al., 2004), it has become a
123 convenient theoretical justification for the development of numerical models at
124 specific scales and to address the challenge that the computational solution of the
125 fundamental equation is imposing. Numerical constraints, in fact, oblige us to
126 identify the portion of the energy spectrum to be explicitly resolved by the model.
127 Larger domains imply larger grid spacing for practical constraints on the number of
128 grid points where the equations are to be solved. Larger domains on the one hand
129 allow us to move the resolved scales up in the atmospheric spectrum but at the
130 same time the coarser resolution leads to the loss of detail in the treatment of sub-
131 grid processes which are represented by parameterizations. Thus, for example, a
132 model that has the entire globe as simulation domain will have to use a horizontal
133 grid spacing of 25 to 100 km and therefore approximate (parameterize) the large
134 number of important processes occurring below those grid sizes. Conversely and
135 under normal conditions, a regional scale model that works with a horizontal grid
136 spacing of approximately 12-15 km will resolve explicitly the dynamics and transport
137 that occurs at scales larger than that distance but will not be able to extend the
138 computational domain to the hemispheric or the global scale. The scale separation
139 hypothesis states that the energy peak of boundary layer processes is isolated from
140 the rest of the spectrum, thus justifying their parameterization in a global model.
141 The same principle holds for a regional scale model. However, in the case of a
142 regional scale model, all the processes with scales falling in between 12-15 km and a
143 global-scale model grid-spacing (25-100 km) are resolved explicitly.



144

145 Although models are developed according to specific scales, nothing prevents us
146 from combining them in an across-scale ensemble. What may appear to be just
147 another attempt to combine model results for the sake of further and diversely
148 populating an ensemble, has in fact a more rigorous motivation. Models working at
149 different scales represent with different degrees of accuracy and precision different
150 portions of the atmospheric spectrum and therefore processes. Our working
151 hypothesis is therefore that by combining global and regional scale models into an
152 ensemble, there is a high probability that they would complement each other across
153 scales and consequently provide an improved ensemble performance compared to
154 single scale ensembles.

155

156 Since in this study we are dealing with chemical transport models (CTM) we should
157 also consider that chemical mechanisms span across a wide range of time scales.
158 This could also constitute an element of diversity among the models working at
159 different space scale although the time resolutions for regional and global scale
160 models are comparable. One could argue that in regional domains in particular,
161 regional models essentially represent in detail the chemistry over a timescale of 10-
162 days which then gets advected out and “reset”. For example, differing
163 representations of organic nitrate lifetimes and how long they sequester NO_x in the
164 system, impacts large scale O₃. Thus the difference in chemical mechanisms related
165 to longer-lived species and multi-day chemistry could also introduce diversity and be
166 another reason for exploring such an “across-scale ensemble”.

167

168 Apparent ancillary elements that could also improve the ensemble results are for
169 example the differences in emission inventories or in general sources of primary
170 information, whose accuracy and precision cannot be guaranteed a priori or
171 evaluated and that could contribute to the development of additional probable
172 solutions.

173

174 As presented in the past, the diversity of modeling approaches is the element that
175 favors a better ensemble product (Kioutsoukis and Galmarini, 2014; Kioutsoukis et



176 al., 2016). In this sense the combination of model results that focus on different
177 scales and that account in a different form for the chemical mechanism has the
178 potential to increase the value of an ensemble to which we will refer from now on as
179 the *hybrid ensemble*.

180

181 The focus in this paper will therefore be on the analysis of the behavior of a hybrid
182 ensemble. The variable considered is the ozone concentration measured and
183 modeled for the year 2010 over the European continent. The analysis takes
184 advantage of the unique opportunity offered by the HTAP2/AQMEI13 activity which
185 brought together global and regional scale models to work on the same case study
186 with a high level of coordination (Galmarini et al., 2017) as far as the input data are
187 concerned.

188

189 In section 2, the observations and model results used in the analysis are presented in
190 detail. In Section 3 the model results are characterized in the phase space to clearly
191 establish whether the two scale groups do indeed account for different portions of
192 the energy spectrum in a distinctly different way. Prior to analyzing the performance
193 of the different ensembles, in Section 4 we evaluate the individual models against
194 the measurements using conventional statistics as well as the newly developed error
195 apportionment analysis presented by Solazzo and Galmarini (2016). Section 5 and 6
196 are dedicated to the analysis of the individual scale ensembles and the hybrid
197 ensemble. Section 7 is dedicated to the comparison hybrid ensemble and single scale
198 ensemble performance. The conclusions are discussed in section 8.

199

200

201 **2. The models used and the case study**

202 The set of models results considered and analyzed in this work are those that
203 contributed to the HTAP2 and AQMEI13 modeling initiatives described in Galmarini et
204 al. (2017).

205

206 HTAP2 is the second phase of the modeling activities of the Task Force on
207 Hemispheric Transport of Air Pollutants (TF-HTAP) during which a community of



208 global scale CTMs performed a large number of simulations with the primary goal of
209 investigating the transcontinental exchange of atmospheric pollutants (Dentener et
210 al, 2010; Fiore et al. 2009). AQMEI3 is the third phase of the Air Quality Model
211 Evaluation International Initiative (AQMEI, Rao et al. 2011) which brings together a
212 community of European (EU) and North American (NA) regional scale modelers to
213 work on coordinated case studies over EU and NA. For this third phase, the regional
214 scale air quality modeling activity has been performed within HTAP2 framework. The
215 coordination between HTAP2 and AQMEI3, as detailed in Galmarini et al. (2017),
216 relates to the use of HTAP2 global model results as boundary conditions to the
217 regional scale models and the use of the same anthropogenic emission inventory
218 (Janssens-Maenhout et al., 2015) by both communities. The list of regional and
219 global scale models analyzed in this work is presented in Tables 1 and 2 respectively.
220 The simulations are for the year 2010 and the regional scale models were all initiated
221 and received boundary conditions from the same global chemistry transport model
222 C-IFS (Flemming et.al, 2015). C-IFS is also one of the global models that are part of
223 the global model ensemble. The two sets of models have been extensively evaluated
224 (Solazzo et al. 2017; Solazzo and Galmarini, 2016; Jonson et al., 2018; Galmarini et al.
225 2018).

226

227 The analysis presented here focuses exclusively on ozone over the EU continent for
228 which the largest abundance of models for the two groups is available and for which
229 case we can take stock on the fact that the models' performance has been analyzed
230 with respect to other species elsewhere (Im et al., 2017). In the figures and tables
231 resulting from our analysis, we shall not identify the individual models used since our
232 goal is the identification of possible advantages in using hybrid ensembles rather
233 than evaluating individual model results.

234

235 Hourly modeled concentrations of ozone were extracted by the modeling groups at
236 European routine and non-routine sampling locations presented in Figure 1. Details
237 on the networks used can be found in Solazzo et al. (2012), Im et al. (2015), and
238 Solazzo et al. (2017). Surface data were provided by the European Monitoring and
239 Evaluation Programme (EMEP; <http://www.emep.int/>) and the European Air Quality



240 Database, AirBase (<http://acm.eionet.europa.eu/databases/airbase>). For the
241 purposes of comparing the ensemble performance with observations, only rural
242 stations with data completeness greater than 75% for the entire year and elevation
243 above ground lower than 1000 m have been included in the analysis. The total
244 number of valid time series used is 483.

245

246

247 **3. Spectral analysis of the global and regional model time series of ozone** 248 **concentrations**

249 One year of one-hour resolution ozone data allows us to produce detailed spectra
250 from the two groups of models and the measured concentrations. In Figure 2a the
251 spectrum of the monitoring time series is shown as a function of the frequency and
252 without any smoothing. In Figures 2b and c, smoothed individual power spectra of
253 ozone (plotted against the period in days for easier interpretation) from global (2b)
254 and regional (2c) models are compared with the spectrum of the measured ozone.
255 The time series of the rural monitoring stations have been averaged prior to
256 producing the spectra. In all subsequent results the measured time series should be
257 interpreted as ensemble averages of all available rural monitoring stations.

258

259 Since ozone is a scalar quantity, its spectrum grows monotonically in log-log scale as
260 expected (e.g. Galmarini et al., 2000), showing a distinct peak around a period of 24
261 hours (more visible in the unsmoothed spectrum (Figure 2a)), corresponding to the
262 daily boundary layer evolution and photochemical production of ozone. This peak is
263 captured well by the two groups of model. The global set tends to slightly
264 underestimate the energy associated with this period with only a single model that
265 overestimates it. The regional scale models are evenly distributed around the
266 spectrum of the measured time series. The two groups behave remarkably similarly
267 at scales smaller than the daily peak. The majority of the models overestimate the
268 energy but capture the slope of the measured spectrum. As expected, the spectra of
269 the global models are more scattered but yet very well behaved. A weak second
270 peak is visible between 30 and 50 days, which could be easily attributed to the
271 synoptic variability. Solazzo and Galmarini (2016) demonstrated that it could indeed



272 be connected to meteorology and/or removal by dry deposition. Moving up the
273 period scale, after the daily peak, all regional scale model spectra are below the
274 observed spectra a behavior that continues apart from a few exceptions up until the
275 60-70 day period range. Out of seven global models however, only 3 under- or over-
276 estimate the energy in this period-range while the rest match the observed
277 spectrum. At 70-80 days a new peak appears in the observed time series,
278 corresponding to the seasonal variability. Only one global model captures the
279 observed time series, three models seem to anticipate it at smaller periods and even
280 in the regional scale group there is a variety of behaviors including a monotonic
281 increase of the energy throughout this period range. Beyond the 100-day period the
282 ozone energy spectrum grows monotonically, which the global model group matches
283 the power line very closely whereas the regional scale group shows a more erratic
284 behavior.

285

286 This first test is important to assess the fundamental differences between the two
287 sets of models with respect to the characteristics of the signal, the periodicities
288 present in the latter and the ability to reproduce the power or the variance of the
289 measured signal at the various frequencies (periods). In addition, it can give us an
290 idea of the level of complementarity that exists between the two groups of models
291 in the representation of the measured power spectrum. As clearly evident from
292 Figure 2, both groups of models show an internal coherence in the representation of
293 the power spectra. A remarkable result is the capacity of global models to represent
294 the high frequency part of the ozone spectrum with an accuracy that is comparable
295 with regional models. We can expect a complementarity in the behavior of the two
296 groups in the large-scale energy range, which should be regulating the long-range
297 transport and background values. The global models have a better representation of
298 that portion of the spectrum than the regional one. An element of surprise resides in
299 the fact that the behavior of the two groups is rather similar for ozone as measured
300 by a power spectrum.

301

302

303 **4. Group performance and error apportionment**



304 A characterization of model performance for the individual members of the two
305 groups beyond the information provided in Galmarini et al. (2018), Solazzo et al.
306 (2017), and Jonson et al. (2018) is also appropriate at this stage.

307

308 The Taylor diagrams presented in Figures 3a and b provide an overview of the
309 individual model performance across the year of reference. All model results
310 underwent un-biasing (subtract the annual mean bias from the predicted hourly
311 values, which produces a shift of the annual time series up or down by MB). We
312 notice that the global models show a more scattered behaviour compared to the
313 regional scale models, with performance distributed across a wider range of
314 standard deviation values. Among the global scale models we find a clear outlier
315 (model 5) whereas the rest tend to group in a rather narrow range of standard
316 deviation values and correlations. Among the regional scale models we can also
317 identify an outlier specifically model 9. The RMSE values range from 22.4 to 25.9
318 $\mu\text{g m}^{-3}$ for the global models and 21 to 24.7 $\mu\text{g m}^{-3}$ for the regional models and are
319 thus comparable. Global models overestimate the observed standard deviation while
320 regional scale models with the exception of model 9 are evenly distributed across
321 the observed values. The correlation coefficient is comparable for the two groups of
322 models.

323

324 Figure 4a and b present two classical skill scores for categorical events also applied
325 by Kioutsoukis et al. (2016), namely the probability of detection (POD) and false
326 alarm rate (FAR). The former represents the proportion of occurrences (e.g. events
327 exceeding a threshold value) that were correctly identified, whereas the latter is the
328 proportion of non-occurrences that were incorrectly identified as happening. In
329 other words they measure *true* and *false positives*. In this case the scores are
330 calculated on the basis of the individual model performances at each station. POD
331 and FAR plots are presented as probabilities above a fixed threshold of 100 $\mu\text{g m}^{-3}$
332 (Figure 4a) and as breakdowns for different threshold values (Figure 4b), where the
333 abundance of the observed data per concentration range is also given as histogram.
334 The POD charts show that the global models have a notably higher probability to
335 identify true positives and that this POD is maintained at the various threshold



336 levels. At the same time the global models also have a higher percentage of false
337 positives as can be gleaned from the FAR index in Figure 4a. This analysis is
338 important to establish the capacity of the models to simulate extreme values.

339

340 Using the methodology proposed by Solazzo and Galmarini (2016), in Figure 5 we
341 present the decomposition of the model errors according to specific time scales. In
342 this figure, the individual model errors are shown as decomposed in the diurnal
343 (<6h), inter-diurnal (6h-1d), synoptic (1-10d), and long-term (>10d) time scales and
344 the residual. The decomposition is performed using a Kolmogorov-Zurbenko filter
345 (Rao and Zurbenko, 1997) applied to the Mean Squared Error (MSE) calculated from
346 each model and the observed ozone time series. Such analysis can be very revealing
347 as it identifies the scale and therefore the processes that are mainly responsible for
348 the deviation of the model results from the measurements as well as possible
349 persistence of errors at specific scales.

350

351 The figure reveals that most of the error is contained in the long term and diurnal
352 time scales. For regional-scale models, this is in agreement with the findings of
353 Solazzo and Galmarini (2016) and Solazzo et al. (2017). The same behaviour is also
354 found in the group of global models. What is remarkable is the similarity of the error
355 values at the diurnal time scale across the two groups. This suggests that the
356 difference in spatial resolution between the two sets of models does not seem to
357 influence the error at the scale at which atmospheric boundary layer dynamics and
358 daily emissions of ozone precursors are the dominant processes. Apart from a few
359 exceptions (model 13 and 17 in the regional scale group and model 5 and 1 in the
360 global scale group), all other models have very comparable errors at that scale. A
361 comparable error across the two groups is found at the synoptic scale although this
362 is less surprising because this scale is explicitly resolved by the models in both groups
363 and strongly depends on the quality of the meteorological driver used. Since both
364 global and regional models employ assimilation of meteorological observations, they
365 are able to represent the synoptic scale comparably and are less dependent on
366 parameterizations employed. The long-term components have the largest error and
367 also show the most variability across models. Remarkably, the regional-scale models



368 seem to show smaller long-term error values than the global models although they
369 are highly variable from model to model. The strong dependence of the long-term
370 error on boundary conditions, (specifically lateral boundary conditions and long
371 range transport in the case of a global model, upper air stratospheric intrusions and
372 surface emission of ozone precursors and direct ozone deposition) appears to
373 influence the global scale group concentrations more than the regional scale, though
374 one should consider that almost all regional scale models used boundary conditions
375 from the same global model which nevertheless does not have the smallest long-
376 term error component of the error.

377

378 A useful pre-characterization of an ensemble can be obtained by the construction of
379 the Talagrand diagram (Talagrand et al. 1997). This construction is achieved by
380 binning the range from the minimum to the maximum modelled concentrations with
381 as many bins as the number of ensemble members plus one. The bins are then filled
382 with observed values based on where they fall within the modelled concentration
383 range. For example, if an observed value is lower than the lowest model value, it is
384 assigned to the first bin, if it falls between the lowest and second-lowest model
385 value, it is assigned to the second bin, and so on. If it exceeds the highest model
386 value, it is assigned to the last bin. Figures 6a and 6b show the Talagrand diagrams
387 for the global and regional models respectively. The figures reveal the tendency of
388 the global model ensemble to be over-dispersed as indicated by the accumulation of
389 most of the observed data at the centre of histogram and relatively few observations
390 falling into the more extreme modelled bins. The regional scale model ensemble
391 shows a flat diagram which is an indication of good group performance. A flat
392 Talagrand diagram is an indication of the fact that the group members equally cover
393 (by proportion) all the observed range of values and the group variability does not
394 show an excess or deficiency in the number of predictions in a specific range of
395 observed values.

396

397 The first result obtained for a combined set of model results is shown in Figure 6c,
398 which presents the Talagrand diagram for the combination of the two groups of
399 models. Note that the number of bins (x-axis) has increased corresponding to the



400 new total number of models considered plus 1 (i.e. 7 global models plus 13 regional
401 models plus 1). The diagram for the combined group of models qualitatively
402 constitutes an improvement compared to those of the individual group ensembles.
403 The combination of the bell shaped diagram of the global set with the relatively flat
404 shape of the regional set produces an extension of the range correctly modelled by
405 the new combined set of models (flat region between bins 5 and 18) and an under
406 prediction between bins 1 and 5 and 19 and 21, which now account for lower and
407 higher values respectively compared to the extreme bins of the global and regional
408 sets.

409

410 **5. Ensemble analysis per scale group**

411 Prior to analyzing the performance of the hybrid multi-model ensemble (mme_GR),
412 let us concentrate on the individual ensembles (mme_R and mme_G) of the two
413 groups for the sake of having a term of comparison beyond the measured
414 concentrations against which to compare the mme_GR one. In this study, we would
415 also like to build upon the research performed in other multi-model ensembles over
416 the years and rather than calculating only the classical model average or median
417 ensemble (mme) we shall also calculate three ensembles based on the findings from
418 Potempski and Galmarini (2009), Riccio et al. (2012), Solazzo et al. (2012); Solazzo et
419 al. (2013); Galmarini et al. (2013), and Kioutsoukis and Galmarini (2014). We shall
420 therefore refer to mmeS (Solazzo et al., 2012) as the ensemble made by the optimal
421 subset of models that produce the minimum RMSE; kzFO (Galmarini et al., 2013) as
422 the ensemble produced by filtering measurements and all model results using the
423 Kolmogorov-Zurbenko decomposition presented earlier and recombining the four
424 components that best compare with the observed components into a new model
425 set; and the optimally weighted combination mmeW (Potempski and Galmarini,
426 2009, Kioutsoukis and Galmarini, 2014, Kioutsoukis et al., 2016).

427

428 Figures 7a and b show the effect of the various ensemble treatments for the two
429 groups of models separately and presented as Taylor diagram. The correlation has
430 increased and narrowed between 0.90 and 0.95 for both groups. As expected, the
431 best ensemble treatment of the two individual groups is mmeW which in the case of



432 the global models is comparable to mmeS and in the case of the regional scale
433 models is farther apart from mmeS. The fact that the optimal partition of the error in
434 terms of accuracy and diversity in an equal weighted sub-ensemble (mmeS) and the
435 analytical optimization of the error in a weighted full-ensemble (mmeW) are
436 comparable for the global models implies that this group better replicates the
437 behavior of an independent and identically distributed (i.i.d.) ensemble around the
438 true state set (on average). The range of improvement of the RMSE is comparable
439 for the two groups of models.

440

441 Of the entire set of ensemble treatments proposed, mmeS is the only one that works
442 with an identified subset of elements. The elements chosen in this context are those
443 that minimize a specific metric (e.g. RMSE). The combination of all possible
444 permutations of a pre-defined subset and for all possible subsets allows us to
445 identify the subgroup of models that performs best (Solazzo et al. 2012). This group
446 is the one that best reduces the redundancies and optimizes the complementarity of
447 the model results (Kioutsioukis and Galmarini, 2014). Other methods have been
448 devised to determine the optimal number of models (Bretherton et al., 1999; Riccio
449 et al. 2012) that are equally effective as the one used here, though they do not allow
450 identifying the members of the subset. Beyond the use of the mmeS for the current
451 analysis, given the diversity in the number of models comprising the two ensembles
452 we have calculated the effective numbers of models (Bretherton et al., 1999) for the
453 regional and global sets in the attempt to verify whether the effective numbers were
454 close for the two sets. Figure 8a shows the values obtained for the global set and the
455 regional set. At over two third of the stations, the mmeS used 3-4 global models and
456 3-5 regional models. In other words, roughly half of the global models (3-4 out of 7)
457 produce the best result when constructing the mmeS globally while in the case of
458 the regional scale models less than half (3-5 out of 13) of all models are required.
459 Figure 8b provides the frequency of contribution of the individual models to the
460 mmeS thus confirming the dominance of 3 global and 4 regional models determined
461 with the N_{eff} analysis. What is presented in Figure 8 is the analysis for the aggregated
462 set of model results at all available monitoring points. We also would like to
463 determine the spatial variability of this result, i.e. to answer the question whether



464 N_{eff} is uniform throughout the domain or whether there are sub regions that require
465 more or less models to construct mmeS.

466

467 In order to have a more objective assessment of the result presented in Figure 8 we
468 introduce a metric which samples only the diversity of the model results (see section
469 6). Following Pennel and Reichler (2011) and Solazzo et al. (2013) we introduce the
470 metric d_m defined for M models at location i as:

471

$$472 \quad d_{m,i} = e_{m,i}^* - R_{m,mme} mme_i^* \quad (1)$$

473 where

$$474 \quad mme_i = \frac{1}{M} \sum_{m=1}^M e_{m,i} \quad (2)$$

475

$$476 \quad e_{m,i} = \frac{mod_{m,i} - obs_i}{\sigma_{obs}} \quad (3)$$

477

478 and the * version of $e_{m,i}$ and mme_i is obtained by normalizing them with σ_e and
479 σ_{mme_i} respectively. $R_{m,mme}$ is the correlation between the individual and average
480 model results. Therefore only the uncorrelated portion of the individual result is
481 retained in d as measure of the diversity whereas the correlated portion is filtered
482 out. Applying this metric, the model results have been decomposed by means of the
483 Kolmogorov-Zurbenko filter described earlier and N_{eff} has been calculated across the
484 domain for the most relevant components LT, SY, and DU. Figure 9 presents the
485 results for the two groups of models. For the long-term component, N_{eff} results
486 shown in Figure 8a are largely confirmed with an overall spatial homogeneity of N_{eff} .
487 The global model set appears to require a larger number of models than the average
488 in critical areas like Northern Italy where the resolution would be insufficient to
489 capture the inhomogeneity of the concentration field due to the complex terrain in
490 that region (similarly in the western part of the domain). At the synoptic scale, the
491 regional scale models require slightly more models on average than the numbers
492 presented in Figure 8 and in some portions of the domain almost all available models
493 are required. The number of required models increases even further at the diurnal
494 scale. In the case of the global set, the average N_{eff} is the same across these two



495 scales and more models are required in the Po valley (Italy) at the synoptic scale and
496 western Poland at the diurnal scale.

497

498

499 **6. Building the hybrid ensemble**

500 Given the fact that there is redundancy in the two groups of models and a disparity
501 exists in the overall and effective number of models in the two groups, a strategy has
502 to be devised so that no pre-determined weight is assigned to one of the two groups
503 thus masking the potential outcome of this study or creating false results. This goal is
504 accomplished by applying the following strategy.

505

506 We want to compare three equally populated ensembles of just global, just regional,
507 and mixed global and regional models. We will therefore reduce the ensemble of
508 regional-scale models and include extra models in the ensemble of global models
509 beyond the effective number calculated in Figures 8 and 9 so that the joint ensemble
510 will not be too small. In order to accomplish this, we select the global models
511 contributing most to the global ensemble beyond those identified by N_{eff} . We begin
512 by assuming that six is a reasonably abundant ensemble (as also indicated by the
513 effective number of regional scale models) and as such the single-scale ensembles
514 will be based on six members. Taking advantage of the various techniques developed
515 to build an ensemble presented earlier we define the following sets:

- 516 - (mme_GR) hybrid ensemble of rank 6 (ensemble of 6 members) composed
517 of the best three global models and the best three regional models
- 518 - (mme_G) global ensemble of best six global models
- 519 - (mme_R) regional ensemble of best six regional models
- 520 - (mmeS_GR) optimally generated hybrid ensemble of rank 6 from the pool of
521 the best six global models and the best six regional models
- 522 - (mmeS_G) optimal global ensemble of rank 6
- 523 - (mmeS_R) optimal regional ensemble of rank 6
- 524 - (mmeW_GR) weighted hybrid ensemble composed from the best three
525 global models and the best three regional models
- 526 - (mmeW_G) weighted global ensemble of best six global models



527 - (mmeW_R) weighted regional ensemble of best six regional models

528

529 **7. Comparing the single scale multi model ensembles with the hybrid one**

530 The comparison of the ensemble performances will be restricted to the months of
531 June -August when the photochemical production of ozone is at its maximum and
532 the number of exceedances is expected to peak throughout the continent. The
533 results of the comparison of the mme, mmeS and mmeW for the regional (_R),
534 global (_G) and hybrid cases (_GR) are shown in Figures 10a,b, and c and 11 a, b and
535 c. The elements common to the three figures are:

536

- 537 • The hybrid ensemble of rank 6 composed of the three best global models and
538 the three best regional models (mme_GR) when compared to mme_G (best
539 six global models) and mme_R (best six regional models) does not show
540 improved performance, rather its skill is inferior to both mme_G and mme_R.
- 541 • For the other two kinds of ensemble treatments (mmeS and mmeW), the
542 combination of global and regional models produces a systematic
543 improvement compared to just the global or regional ensembles in terms of
544 correlation coefficients, standard deviations and RMSE.
- 545 • The partition of global and regional models in mmeS (Figure 11) shows that
546 the contribution of regional models is more frequent. Specifically, at two
547 thirds of the stations, the optimum hybrid ensemble of rank 6 consists of one
548 or two global models and five or four regional models, respectively. At only
549 15% of the stations, mmeS consists of an equal number of global and regional
550 models. The maximum number of global models in the mmeS_GR ensemble
551 is four, achieved at roughly 1% of the stations. Conversely, at around 10% of
552 the stations the hybrid ensemble utilized only regional models.
- 553 • POD and FAR (Figures 12 a and b) show a net improvement over the
554 mmeW_G results when the hybrid ensemble is considered, with a minimum
555 in false positives and a maximum in true positives that closely match the
556 mmeW_R results.

557



558 The real improvement of the hybrid ensemble with respect to the single scale model
559 ensembles becomes evident when analyzing Figure 13. The plots in the figure are the
560 collective representation of three of the most important characteristics of an
561 ensemble as proposed by Kioutsioukis and Galmarini (2014), i.e. diversity, accuracy
562 and error. On the x and y axes respectively “*diversity*” and “*accuracy*” are presented.
563 The former represents the average square deviation of the single models from the
564 mean of the models, whereas the latter is the square of the average deviation of the
565 individual model results from the observed value. As presented by Krogh and
566 Vedelsby (1995), the difference of the *diversity* and *accuracy* defines the quadratic
567 deviation of the ensemble average from the observed value. From the definition it
568 follows that in order for the ensemble result to be closer to the observed value one
569 has to find the right trade off between *accuracy* and *diversity* (A-D). A mere increase
570 in *diversity* does not guarantee a minimization of the ensemble error since it will also
571 produce a reduction in the *accuracy*. What one hopes to obtain is the right
572 combination of models that provides the maximum *accuracy* and maximum
573 *diversity*. In the plot of Figure 13, the optimal condition is achieved when the model
574 results concentrate in the upper left quadrant of the plot toward the
575 ($x=100/(\text{Number of Models}), y=1$) point. In the plot, the accuracy parameter is
576 presented as deviation from the best model performance. The dots represent the
577 estimate of the two parameters at every location where measurements are
578 available. The colour scale is based on the RMSE. The two upper panels (13a and
579 13b) give the A-D mapping for the mme_R and mme_G ensembles; the lower two
580 panels give the map for the hybrid ensembles, i.e. mme_GR (13c) and mmeS_GR
581 (13d). The difference in nature of the two ensembles is clear from the two panels.
582 Ensemble mme_G is more diverse and accurate than mme_R (y values: 69 in G and
583 66 for R, x: 0.75 in G, 0.66 in R). The combination of the two produces a decrease in
584 the two parameters (13c). However, if the models are selected as in mmeS_GR, both
585 accuracy and diversity increase. The real advantage of the combination is visible in a
586 slight increase of the diversity as compared to mme_GR and a marked improvement
587 of the accuracy from 0.71 to 0.81. The error decreases from a median value of 17.9
588 to 15.6 and from an Inter Quartile Range of 5.1 to 3.8.



589 In Figures 14 a, b and c the spectra of the ensembles are presented. For the just
590 global and just regional-scale ensembles, the spectra of mme, mmeS, mmeW and
591 kFO are shown in Figures 14 a and b and the ensembles are based on the entire set
592 of available members per group. Figure 14 c shows the spectra of the four
593 ensembles, mme_R6, mme_G6, mme_GR6 and mmeS_GR6 for which the largest six
594 contributors from the regional models, the six global, and three regional plus three
595 global models were used. From the picture we see that regardless of the treatment,
596 the ensemble data captures the ozone power spectrum with no notable deviation
597 from the measured spectrum. It is important to note that an ensemble treatment is
598 a purely statistical treatment that does not consider any physics constraints. The
599 deficiencies that were originally present in the individual model spectra are still
600 present in the ensemble results, particularly the large power deficit in the range
601 from 0.8 days to 100 days. The mme_GR spectrum appears to produce a slight
602 improvement toward filling this energy gap, but the change is very small.

603

604 **8. Discussion and conclusions. How much is the improvement attributable to the** 605 **hybrid character of the ensemble?**

606 The analysis presented above gives us clear indications that the combination of the
607 two sets of models analysed produces an improvement in the ensemble
608 performance. In particular, the hybrid ensemble appears to be superior to any
609 single-scale ensemble in the optimum setting. For example, given six global, six
610 regional and three global and three regional ensembles, the optimization always
611 favours the hybrid ensemble. This was repeated for all examined cases: the annual
612 hourly records, the JJA hourly records and the annual daily maximum records.

- 613 - The improvement is in the range 1-5% compared to single scale optimum
614 ensembles
- 615 - POD/FAR show a remarkable improvement, with a steep increase in the
616 largest POD values, though comparable to the other for the hybrid ensemble
617 and comparatively smallest values of FAR across the concentration ranges.

618

619 Some important considerations need to be made at this point. It is difficult to find
620 quantitative evidence for the fact that the hybrid ensemble improvement can be



621 unequivocally attributed to the multi-scale nature of the ensemble. We have no
622 evidence, nor guarantee, that the same kind of improvement could be reached by
623 adding more regional-scale models to the regional-scale ensemble, or more global
624 models to the global-scale ensemble. However, what is clear is that the regional-
625 scale ensemble is characterised by a higher level of redundancy in the members than
626 the global ensemble, since less than half of the members produced the optimal
627 ensemble, and that the use of the three best members from the regional-scale
628 ensemble and three best global-scale models produces an improvement in the
629 ensemble performance. This last argument suggests that the addition of more model
630 results of the same “nature” would just contribute to further increase the level of
631 redundancy, while on the other hand, the improvement obtained could indeed be
632 attributed to the different “nature” of the global-scale models compared to the
633 regional-scale models.

634

635 Therefore, considering:

- 636 • the large number of regional scale models and the spectrum of diversity in
637 their nature (only a small number of the same models were used by multiple
638 groups and there was an abundance of models developed independently
639 from one another);
- 640 • the relatively smaller number of global model results compared to the
641 regional models and also their different nature;
- 642 • the fact that the two groups of models used the same emission inventories
643 and all the regional scale models used boundary conditions from the same
644 global model;

645 one could attribute the improvement of the *mmeS_GR* ensemble performance to
646 the difference in nature of the two groups and a complementary contribution of the
647 two toward an improved result.

648

649

650 **Acknowledgments**

651 The group from University of L’Aquila kindly thanks the EuroMediterranean Centre
652 on Climate Change (CMCC) for the computational resources. P.T. is beneficiary of an



653 AXA Research Fund postdoctoral grant. We acknowledge the EC FP7 financial
654 support for the TRANSPHORM project (grant agreement 243406). CIEMAT has been
655 financed by the Spanish Ministry of Agriculture and Fishing, Food and Environment.
656 DKH and YD recognize support from NASA HAQAST. The UMU group acknowledges
657 the Project REPAIR-CGL2014-59677-R of Spanish Ministry of the Economy and
658 Competitiveness and the FEDER European program for support to conduct this
659 research. The views expressed in this article are those of the authors and do not
660 necessarily represent the views or policies of the U.S. Environmental Protection
661 Agency. The MetNo work has been partially funded by EMEP under UNECE.
662 Computer time for EMEP model runs was supported by the Research Council of
663 Norway through the NOTUR project EMEP (NN2890K) for CPU, and NorStore project
664 European Monitoring and Evaluation Programme (NS9005K) for storage of data. RSE
665 contribution to this work has been financed by the research fund for the Italian
666 Electrical System under the contract agreement between RSE S.p.A. and the Ministry
667 of Economic Development – General Directorate for Nuclear Energy, Renewable
668 Energy and Energy Efficiency in compliance with the decree of 8 March 2006.
669
670



671

672 **References**

673

674 Bretherton, C. S., Widmann, M., Dymnikov, V. P., Wallace, J. M., and Bladè, I.: The
675 effective number of spatial degrees of freedom of a time-varying field, *J.*
676 *Climate*, 12, 1990–2009, 1999.

677 Colette A., C. Granier, O. Hodnebrog, H. Jakobs, A. Maurizi, A. Nyiri, S. Rao, M.
678 Amann, B. Bessagnet, A. D'Angiola, M. Gauss, C. Heyes, Z. Klimont, F. Meleux,
679 M. Memmesheimer, A. Mieville, L. Rouil, F. Russo, S. Schucht, D. Simpson, F.
680 Stordal, F. Tampieri and M. Vrac, Future air quality in Europe: a multi-model
681 assessment of projected exposure to ozone, *Atmos. Chem. Phys.* 12(2012a),
682 pp. 10613-10630.

683 Fiore A. M., Dentener F. J., Wild O., Cuvelier C., Schultz M. G., Hess P., Textor C.,
684 Schulz M., Doherty R. M., Horowitz L. W., MacKenzie I. A., Sanderson M. G.,
685 Shindell D. T., Stevenson D. S., Szopa S., Van Dingenen R., Zeng G., Atherton
686 C., Bergmann D., Bey I., Carmichael G., Collins W. J., Duncan B. N., Faluvegi
687 G., Folberth G., Gauss M., Gong S., Hauglustaine D., Holloway T., Isaksen I.
688 S. A., Jacob D. J., Jonson J. E., Kaminski J. W., Keating T. J., Lupu A., Marmer
689 E., Montanaro V., Park R. J., Pitari G., Pringle K. J., Pyle J. A., Schroeder S.,
690 Vivanco M. G., Wind P., Wojcik G. and Wu S., Zuber A., (2009), Multimodel
691 estimates of intercontinental source-receptor relationships for ozone
692 pollution, *J. Geophys. Res.*, 114, D04301, doi:[10.1029/2008JD010816](https://doi.org/10.1029/2008JD010816).

693 Flemming, J., Huijnen, V., Arteta, J., Bechtold, P., Beljaars, A., Blechschmidt, A.-M.,
694 Diamantakis, M., Engelen, R. J., Gaudel, A., Inness, A., Jones, L., Josse, B.,
695 Katragkou, E., Marecal, V., Peuch, V.-H., Richter, A., Schultz, M. G., Stein, O.,
696 and Tsikerdekis, A.: Tropospheric chemistry in the Integrated Forecasting
697 System of ECMWF, *Geosci. Model Dev.*, 8, 975-1003,
698 <https://doi.org/10.5194/gmd-8-975-2015>, 2015.

699 Galmarini, S., and P. Thunis, 2000: Estimating the contribution of Leonard and cross
700 terms to the subfilter scale from atmospheric data. *J. Atmos. Sci.*, **57**, 1785–
701 1796.

702 Galmarini, S., F. Michelutti, and P. Thunis, 1999: Evaluation of Leonard and cross
703 terms from atmospheric data. *13th Symp. Boundary Layer Turbulence*, Dallas,
704 TX, Amer. Meteor. Soc., 115-118

705 Galmarini S., R Bianconi, W Klug, T Mikkelsen, R Addis, S Andronopoulos, P Astrup, A
706 Baklanov, J Bartniki, JC Bartzis, R Bellasio, F Bompay, R Buckley, M Bouzom, H
707 Champion, R D'Amours, E Davakis, H Eleveld, GT Geertsema, H Glaab, M Kollax,
708 M Ilvonen, A Manning, U Pechinger, Christer Persson, E Polreich, S Potemski, M
709 Prodanova, J Saltbones, H Slaper, MA Sofiev, D Syrakov, JH Sørensen, L Van der
710 Auwera, I Valkama, R Zelazny, 2004: Ensemble dispersion forecasting—Part I:
711 concept, approach and indicators, *Atmospheric Environment* 38 (28), 4607-
712 4617

713 Galmarini, S., Kioutsioukis, I., and Solazzo, E.: E pluribus unum: ensemble air quality
714 predictions, *Atmos. Chem. Phys.*, 13, 7153– 7182, doi:[10.5194/acp-13-7153-](https://doi.org/10.5194/acp-13-7153-2013)
715 2013, 2013.

716 Galmarini, S., Koffi, B., Solazzo, E., Keating, T., Hogrefe, C., Schulz, M., Benedictow,
717 A., Griesfeller, J. J., Janssens-Maenhout, G., Carmichael, G., Fu, J., and



- 718 Dentener, F. 2017: Technical note: Coordination and harmonization of the
719 multi-scale, multi-model activities HTAP2, AQMEII3, and MICS-Asia3:
720 simulations, emission inventories, boundary conditions, and
721 model output formats, *Atmos. Chem. Phys.*, 17, 1543-1555
- 722 HEMISPHERIC TRANSPORT OF AIR POLLUTION 2010 PART A: OZONE AND
723 PARTICULATE MATTER, Edtrs F. Dentener, T. Keating, and H. Akimoto, AIR
724 POLLUTION STUDIES No. 17, ECONOMIC COMMISSION FOR EUROPE
- 725 Henze, D. K., A. Hakami and J. H. Seinfeld (2007), Development of the adjoint of
726 GEOS-Chem, *Atmos. Chem. Phys.*, 7, 2413-2433
- 727 Im U., Roberto Bianconi, Efisio Solazzo, Ioannis Kioutsioukis, Alba Badia, Alessandra
728 Balzarini, Rocío Baró, Roberto Bellasio, Dominik Brunner, Charles Chemel,
729 Gabriele Curci, Johannes Flemming, Renate Forkel, Lea Giordano, Pedro
730 Jiménez-Guerrero, Marcus Hirtl, Alma Hodzic, Luka Honzak, Oriol Jorba,
731 Christoph Knote, Jeroen J.P. Kuenen, Paul A. Makar, Astrid Manders-Groot,
732 Lucy Neal, Juan L. Pérez, Guido Pirovano, George Pouliot, Roberto San Jose,
733 Nicholas Savage, Wolfram Schroder, Ranjeet S. Sokhi, Dimiter Syrakov, Alfreida
734 Torian, Paolo Tuccella, Johannes Werhahn, Ralf Wolke, Khairunnisa Yahya,
735 Rahela Zabkar, Yang Zhang, Junhua Zhang, Christian Hogrefe, Stefano
736 Galmarini, Evaluation of operational on-line-coupled regional air quality
737 models over Europe and North America in the context of AQMEII phase 2. Part
738 I: Ozone, In *Atmospheric Environment*, Volume 115, 2015, Pages 404-420, ISSN
739 1352-2310
- 740 Im, U., Brandt, J., Geels, C., Hansen, K. M., Christensen, J. H., Andersen, M. S.,
741 Solazzo, E., Kioutsioukis, I., Alyuz, U., Balzarini, A., Baro, R., Bellasio, R.,
742 Bianconi, R., Bieser, J., Colette, A., Curci, G., Farrow, A., Flemming, J., Fraser, A.,
743 Jimenez-Guerrero, P., Kitwiroon, N., Liang, C.-K., Pirovano, G., Pozzoli, L., Prank,
744 M., Rose, R., Sokhi, R., Tuccella, P., Unal, A., Vivanco, M. G., West, J., Yarwood,
745 G., Hogrefe, C., and Galmarini, S.: Assessment and economic valuation of air
746 pollution impacts on human health over Europe and the United States as
747 calculated by a multi-model ensemble in the frame work of AQMEII3, *Atmos.*
748 *Chem. Phys. Discuss.*, <https://doi.org/10.5194/acp-2017-751>, in review, 2017.
- 749 Janssens-Maenhout, G., Crippa, M., Guizzardi, D., Dentener, F., Muntean, M.,
750 Pouliot, G., Keating, T., Zhang, Q., Kurokawa, J., Wankmüller, R., Denier van der
751 Gon, H., Kuenen, J. J. P., Klimont, Z., Frost, G., Darras, S., Koffi, B., and Li, M.:
752 HTAP_v2.2: a mosaic of regional and global emission grid maps for 2008 and
753 2010 to study hemispheric transport of air pollution, *Atmos. Chem. Phys.*, 15,
754 11411-11432, 2015
- 755 Jonker, H. J., J. Vilà-Guerau de Arellano, and P. G. Duynkerke. (2004) Characteristic
756 Length Scales of Reactive Species in a Convective Boundary Layer. *Journal of*
757 *the Atmospheric Sciences* 61:1, 41-56.
- 758 Jonker, H. J., J. W. Cuijpers, and P. G. Duynkerke, 1999: Mesoscale fluctuations in
759 scalars generated by boundary layer convection. *J. Atmos. Sci.*, 56, 801-808.
- 760 Kioutsioukis I., and S. Galmarini. De praeceptis ferendis: good practice in multi-model
761 ensembles, *Atmos. Chem. Phys.*, 14, 11791-11815, 2014
- 762 Kioutsioukis, I., Im, U., Solazzo, E., Bianconi, R., Badia, A., Balzarini, A., Baró, R.,
763 Bellasio, R., Brunner, D., Chemel, C., Curci, G., van der Gon, H. D., Flemming, J.,
764 Forkel, R., Giordano, L., Jiménez-Guerrero, P., Hirtl, M., Jorba, O.,



- 765 MandersGroot, A., Neal, L., Pérez, J. L., Pirovano, G., San Jose, R., Savage, N.,
766 Schroder, W., Sokhi, R. S., Syrakov, D., Tuccella, P., Werhahn, J., Wolke, R.,
767 Hogrefe, C., and Galmarini, S.: Insights into the deterministic skill of air quality
768 ensembles from the analysis of AQMEII data, *Atmos. Chem. Phys.*, 16, 15629–
769 15652, doi:10.5194/acp-16-15629-2016, 2016.
- 770 Mathur, R., Xing, J., Gilliam, R., Sarwar, G., Hogrefe, C., Pleim, J., Pouliot, G., Roselle,
771 S., Spero, T. L., Wong, D. C., and Young, J.: Extending the Community Multiscale
772 Air Quality (CMAQ) modeling system to hemispheric scales: overview of
773 process considerations and initial applications, *Atmos. Chem. Phys.*, 17, 12449–
774 12474, <https://doi.org/10.5194/acp-17-12449-2017>, 2017.
- 775 Pennel, C. and Reichler, T.: On the effective numbers of climate models, *J. Climate*,
776 24, 2358–2367, 2011.
- 777 Pielke R. A. Sr, *Mesoscale Meteorological Modeling*
778 *Volume 98 di International Geophysics*, ISBN 0123852382 and 9780123852380,
779 Academic Press, 2013, pp760
- 780 Potempski, S. and Galmarini, S.: Est modus in rebus: analytical properties of multi-
781 model ensembles, *Atmos. Chem. Phys.*, 9, 9471–9489, doi:10.5194/acp-9-
782 9471-2009, 2009.
- 783 Rao, S. T., Zurbenko, I. G., Neagu, R., Porter, P. S., Ku, J. Y., and Henry, R. F.: Space
784 and time scales in ambient ozone data, *B. Am. Meteorol. Soc.*, 78, 2153,
785 doi:10.1175/1520-0477(1997)0782.0.CO;2, 1997
- 786 Rao, S. T., Galmarini, S., and Puckett, K.: Air quality model evaluation international
787 initiative (AQMEII): Advancing the state of the science in regional
788 photochemical modelling and its applications, *B. Am. Meteorol. Soc.*, 92, 23–
789 30, 2011
- 790 Riccio, A., Ciaramella, A., Giunta, G., Galmarini, S., Solazzo, E., and Potempski, S.: On
791 the systematic reduction of data complexity in multi-model ensemble
792 atmospheric dispersion modelling, *J. Geophys. Res.*, 117, D05314,
793 doi:10.1029/2011JD016503, 2012.
- 794 Simpson, D., Benedictow, A., Berge, H., Bergström, R., Emberson, L., Fagerli, H.,
795 Flechard, C., Hayman, G., Gauss, M., Jonson, J., Jenkin, M., Nyíri, A., Richter, C.,
796 Semeena, V., Tsyro, S., Tuovinen, J.-P., Valdebenito, A. and Wind, P. (2012).
797 The EMEP MSC-W chemical transport model technical description, *Atmos.*
798 *Chem. Phys.* 12: 7825–7865
- 799 Solazzo, E. and Galmarini, S.: A science-based use of ensembles of opportunities for
800 assessment and scenario studies, *Atmos. Chem. Phys.*, 15, 2535-2544,
801 <https://doi.org/10.5194/acp-15-2535-2015>, 2015.
- 802 Solazzo, E., Bianconi, R., Hogrefe, C., Curci, G., Tuccella, P., Alyuz, U., Balzarini, A.,
803 Baró, R., Bellasio, R., Bieser, J., Brandt, J., Christensen, J. H., Colette, A., Francis,
804 X., Fraser, A., Vivanco, M. G., Jiménez-Guerrero, P., Im, U., Manders, A.,
805 Nopmongcol, U., Kitwiroon, N., Pirovano, G., Pozzoli, L., Prank, M., Sokhi, R. S.,
806 Unal, A., Yarwood, G., and Galmarini, S.: Evaluation and error apportionment
807 of an ensemble of atmospheric chemistry transport modeling systems:
808 multivariable temporal and spatial breakdown, *Atmos. Chem. Phys.*, 17, 3001-
809 3054, 2017.
- 810 Solazzo, E., Bianconi, R., Vautard, R., Appel, K. W., Moran, M. D., Hogrefe, C.,
811 Bessagnet, B., Brandt, J., Christensen, J. H., Chemel, C., Coll, I., Denier van der



- 812 Gon, H., Ferreira, J., Forkel, R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B.,
813 Jericvić, A., Kraljević, L., Miranda, A. I., Nopmongkol, U., Pirovano, G.,
814 Prank, M., Riccio, A., Sartelet, K. N., Schaap, M., Silver, J. D., Sokhi, R. S., Vira, J.,
815 Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S., and Galmarini, S.:
816 Model evaluation and ensemble modelling of surface-level ozone in Europe
817 and North America in the context of AQMEII, *Atmos. Environ.*, 53, 60–74,
818 2012a.
- 819 Solazzo, E., Bianconi, R., Vautard, R., Appel, K. W., Moran, M. D., Hogrefe, C.,
820 Bessagnet, B., Brandt, J., Christensen, J. H., Chemel, C., Coll, I., Denier van der
821 Gon, H., Ferreira, J., Forkel, R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B.,
822 Jericvić, A., Kraljević, L., Miranda, A. I., Nopmongkol, U., Pirovano, G.,
823 Prank, M., Riccio, A., Sartelet, K. N., Schaap, M., Silver, J. D., Sokhi, R. S., Vira, J.,
824 Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S., and Galmarini, S.:
825 Model evaluation and ensemble modelling of surface-level ozone in Europe
826 and North America in the context of AQMEII, *Atmos. Environ.*, 53, 60–74,
827 2012a.
- 828 Solazzo, E., Bianconi, R., Pirovano, G., Matthias, V., Vautard, R., Moran, M. D., Wyat
829 Appel, K., Bessagnet, B., Brandt, J., Christensen, J. H., Chemel, C., Coll, I.,
830 Ferreira, J., Forkel, R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B., Miranda,
831 A. I., Nopmongkol, U., Prank, M., Sartelet, K. N., Schaap, M., Silver, J. D., Sokhi,
832 R. S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S. T., and
833 Galmarini, S.: Operational model evaluation for particulate matter in Europe
834 and North America in the context of AQMEII, *Atmos. Environ.*, 53, 75–92,
835 2012b.
- 836 Solazzo, E., Riccio, A., Kioutsioukis, I., and Galmarini, S.: Pauci ex tanto numero:
837 reduce redundancy in multi-model ensembles, *Atmos. Chem. Phys.*, 13, 8315–
838 8333, doi:10.5194/acp-13-8315-2013, 2013.
- 839 Sudo, K., M. Takahashi, J. Kurokawa, and H. Akimoto, Chaser: A global chemical
840 model of the troposphere, 1. Model description, *J. Geophys. Res.*, 107(D17),
841 4339, doi:10.1029/2001JD001113, 2002.
- 842 Talagrand, O., R. Vautard, B. Strauss: Evaluation of probabilistic prediction systems,
843 Workshop proceedings "Workshop on predictability", 20-22 October 1997,
844 ECMWF, Reading, UK, 1999
- 845 Tebaldi C. and R. Knutti (2007), The use of the multimodel ensemble in probabilistic
846 climate projections. *Philosophical Transactions of the Royal Society (special*
847 *issue on Probabilistic Climate Change Projections)*, Vol. 365, pp. 2053-2075.
- 848 UNITED NATIONS New York and Geneva, 2010, pp304
- 849 van der Hoven, I. V., 1957: Power spectrum of horizontal wind speed in the
850 frequency range from 0.0007 to 900 cycles per hour. *J. Meteor.*, **14**, 160–164.
- 851 Watanabe, S., Hajima, T., Sudo, K., Nagashima, T., Takemura, T., Okajima, H.,
852 Nozawa, T., Kawase, H., Abe, M., Yokohata, T., Ise, T., Sato, H., Kato, E., Takata,
853 K., Emori, S., and Kawamiya, M.: MIROC-ESM 2010: model description and
854 basic results of CMIP5-20c3m experiments, *Geosci. Model Dev.*, 4, 845-872,
855 <https://doi.org/10.5194/gmd-4-845-2011>, 2011.
- 856 Xing, J., Mathur, R., Pleim, J., Hogrefe, C., Gan, C.-M., Wong, D. C., and Wei, C.: Can a
857 coupled meteorology–chemistry model reproduce the historical trend in



858 aerosol direct radiative effects over the Northern Hemisphere?, Atmos. Chem.
859 Phys., 15, 9997-10018, <https://doi.org/10.5194/acp-15-9997-2015>, 2015.

860 **Figure Captions**

861

862 Figure 1: Spatial distribution of the 405 rural monitoring stations where ozone model
863 results were produced and observations were available

864

865 Figure 2: (a) Power spectrum of observed ozone obtained from the average one year
866 time series across all measuring locations. Running averaged spectrum of
867 observations (thick red line), global models (b) and regional models (c)

868

869 Figure 3: Taylor diagram of Global models (a) and regional models (b)

870

871 Figure 4: Cumulated (a) Probability of detection (POD) and False alarm ration (FAR)
872 for Global and regional models; (b) POD and FAR for ozone concentration ranges

873

874 Figure 5: Distribution of the Mean Square Error (MSE) across the models of the two
875 communities and the scales in which the signal has been decoposed (LT, long term;
876 SY synoptic; DU diurnal; ID inter diurnal; see text for definition)

877

878 Figure 6: Talagrand diagrams of Global (a) and Regional (b) models and Global +
879 Regional set of model results (c)

880

881 Figure 7: Taylor diagram of the four ensemble treatments considered in the text
882 obtained from the global (a) and regional (b) models

883

884 Figure 8: Effective number (N_{eff}) of models calculated according to Bretherton et al.
885 (1999) for the two groups of models (a); and frequency of contribution of each
886 model to the relative ensemble (b)

887

888 Figure 9: Number of effective models for the two groups obtained at all monitoring
889 locations considered thus giving the spatial structure of the ensemble size and for



890 three of the four components in which the modelled time series have been
891 decomposed, namely: LT, SY and DU.

892

893 Figure 10: Comparison of the performance of three ensemble treatments (mme,
894 mmeS and mmeW) for three groupings of models (regional *R*, global *_G*, and mixed
895 global and regional *_RG*)

896

897 Figure 11: Contribution of Global models to mmeS

898

899 Figure 12: POD and FAR for the best performing ensemble treatment (mmeW) and
900 for three ensemble grouping (regional *R*, global *_G*, and mixed global and regional
901 *_RG*)

902

903 Figure 13: Representation of the accuracy (y-axis) vs diversity (x-axis) and RMSE for
904 and ensemble of the most present 6 global and regional models respectively and an
905 hybrid ensemble of three most frequently present global and 3 regional models.

906

907 Figure 14: Spectra behaviour of the ensemble treatments: (a) full global ensemble;
908 (b) full regional ensemble; (c) mme of 6 most frequently present global and regional
909 models and the hybrid ensemble

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925
926
TABLE 1. PARTICIPATING REGIONAL MODELLING SYSTEMS AND KEY FEATURES. THE DARK SHADED CELLS CONTAIN INFORMATION ON MODELS THAT WORKED OVER THE NA DOMAIN THE OTHERS ON THE EU ONE
927
928

Operated by	Modelling system	Horizontal grid	Vertical grid	Global meteo data provider	Gaseous chemistry module
Finnish Meteorological Institute (working with 2 versions)	ECMWF-SILAM_H, SILAM_M	0.25 x 0.25 deg (LatLon)	12 uneven layers up to 13km. First layer ~30m	ECMWF (nudging within the PBL)	CBM-IV
Netherlands Organization for Applied Scientific Research	ECMWF-L-EUROS	0.5 x 0.25 deg (latlon)	Surface layer (~25m depth), mixing layer, 2 reservoir layers up to 3.5km.	Direct interpolation from ECMWF	CBM-IV
University of L'Aquila	WRF-WRF/Chem1	23 km	33 levels up to 50hPa. 12 layers below 1km. First layer ~12m	ECMWF (nudging above the PBL)	RACM-ESRL
University of Murcia	WRF-WRF/Chem2	23 x 23 km ²	33 levels, from ~24m to 50hPa	ECMWF (nudging above the PBL)	RADM2
Ricerca Sistema Energetico	WRF-CAMx	23 x 23 km ²	14 layers up to 8km. First layer ~25m.	ECMWF (nudging within the PBL)	CB05
University of Aarhus	WRF-DEHM	50 x 50 km ²	29 layers up to 100hPa	ECMWF (no nudging within the PBL)	Brandt et al. (2012)
Istanbul Technical University	WRF-CMAQ1	30 x 30 km ²	24 layers up to 10hPa	NCEP (nudging within PBL)	CB05
Kings College	WRF-CMAQ4	15 x 15 km ²	23 layers up to 100hPa, 7 layer below 1km. First layer ~14m	NCEP (Nudging within the PBL)	CB05
Ricardo E&E	WRF-CMAQ2	30 x 30 km ²	23 VL up to 100hPa, 7 layers < 1km. 1 st L @ ~15m	NCEP (nudging above the PBL)	CB05-TUCL
Helmholtz-Zentrum Geesthacht	CCLM-CMAQ	24 x 24 km ²	30 VL from ~40m to 50hPa	NCEP (spectral nudging above f. troposphere)	CB05-TUCL
University of Hertfordshire	WRF-CMAQ3	18 x 18 km ²	35 VL from ~20m to ~16km	ECMWF (nudging above PBL)	CB05-TUCL
INERIS/CIEMAT	ECMWF-Chimere_H Chimere_M	0.25 x 0.25 deg	9 VL up to 500hPa. 1 st L @ ~20m	Direct interpolation from ECMWF	MELCHIOR2

929
930
931
932
933
934
935
936
937
938
939
940
941



942

943

TABLE 2. PARTICIPATING GLOBAL MODELLING SYSTEMS AND KEY FEATURES.

944

Operated by	Modelling system	Horizontal grid	Vertical grid	Global meteo data provider	Gaseous chemistry module	References
NAGOYA, JAMSTEC, NIES	CHASER_re1	128x64 cells, Approximately 2.8x2.8deg	32 VL up to 40 km	ECMWF (nudging above PBL)	Sudo et al. (2002)	Sudo et al. (2002), Watanabe et al. (2011)
NAGOYA, JAMSTEC, NIES	CHASER_t106	320x160 cells, Approximately 1.1x1.1deg	32 VL up to 40 km	ECMWF (nudging above PBL)	Sudo et al. (2002)	Sudo et al. (2002), Watanabe et al. (2011)
ECMWF	C-IFS	Ca. 80 km	60 VL from surface to 0.1 hPa – lowest level 15 m	IFS	CB05	Flemming et al. 2015 http://emep.int/mscw/mscw_publications.html
MetNo	EMEP_rv4.8	0.5 x 0.5 deg Lat x Lon	20 uneven layers up to 100hpa. First layer ~90m	ECMWF IFS dedicated model run	EMEP	Simpson et al. 2012
Univ. Tennessee	H-CMAQ	108 km x 108 km	44 layers up to 50hPa	WRF	CB05	Xing et al. (2015)
Univ.Col. Boulder	GEOSCHEM-ADJOINT	2° lat x 2.5° lon	47 levels up to 0.066 mb	GEOS-5	GEOS-Chem	Henze et al. (2007)
US-EPA	Hemispheric CMAQ	108kmx108km	44 lev to 50hPa	WRF nudged with NCEP/NCAR	CB05TUCL	Mathur et al. (2017)

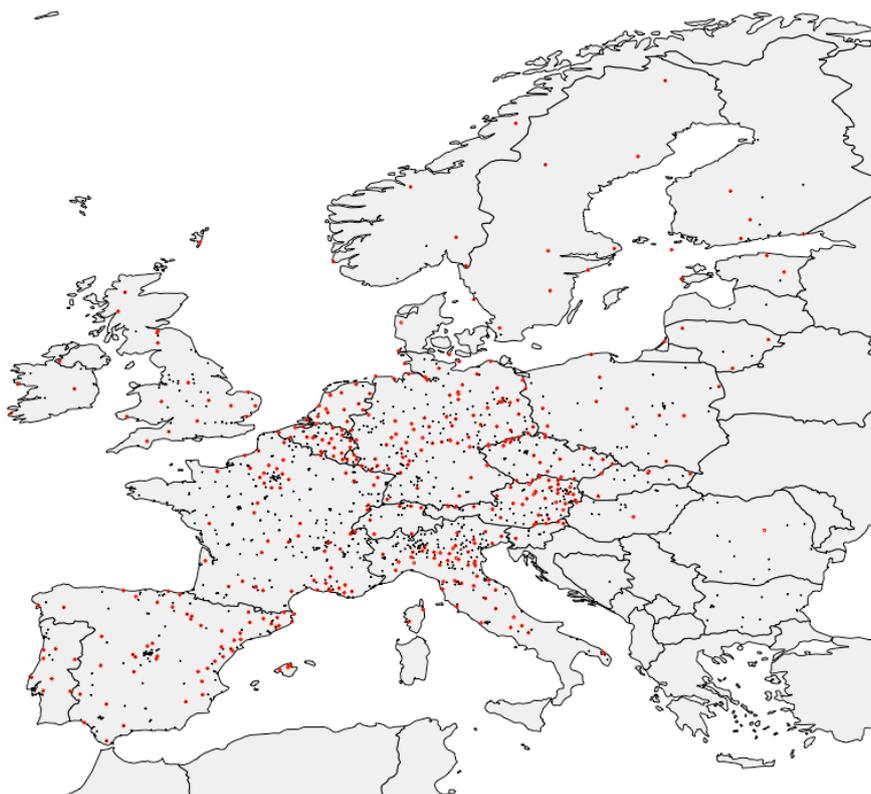
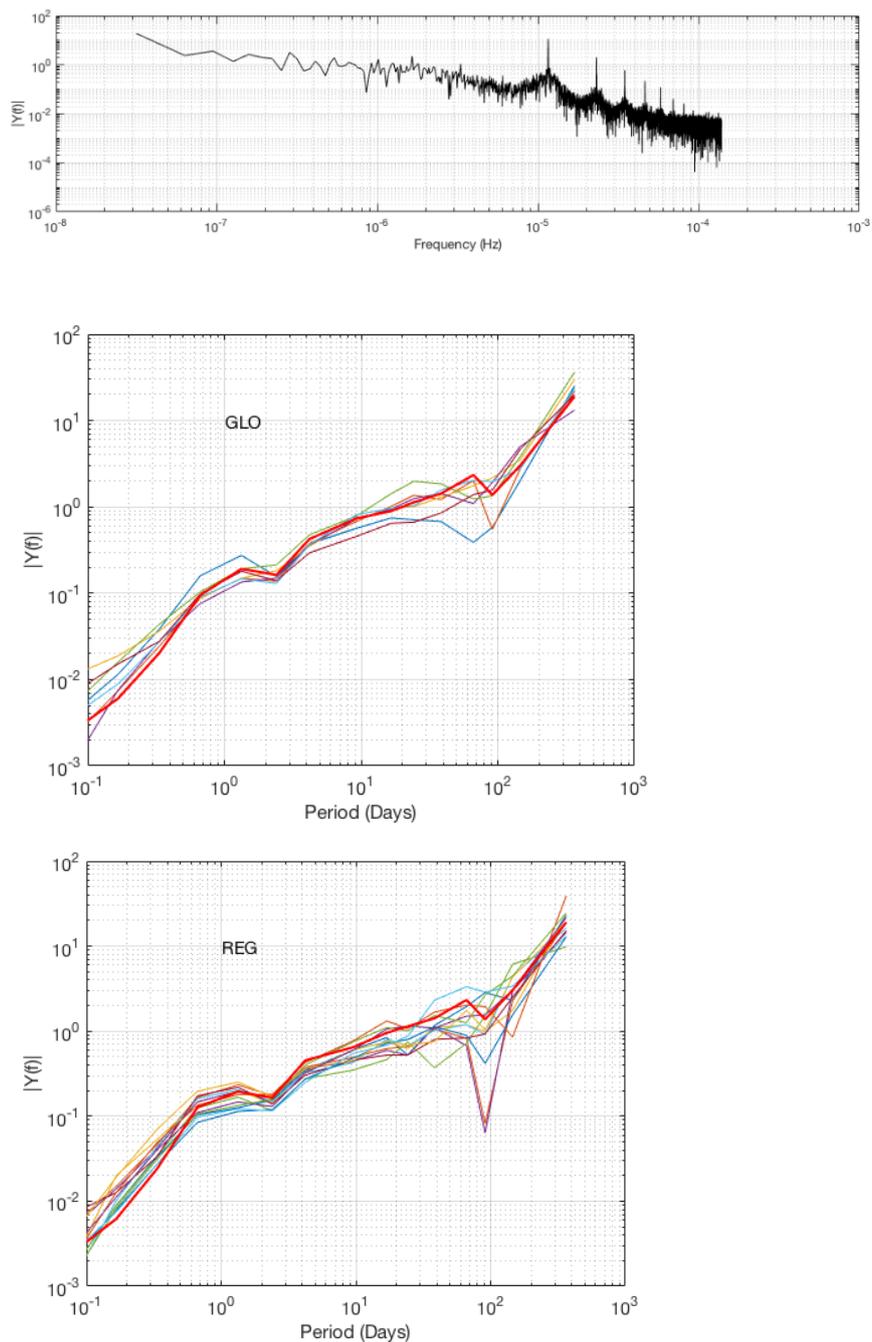


Figure 1



Figure 2



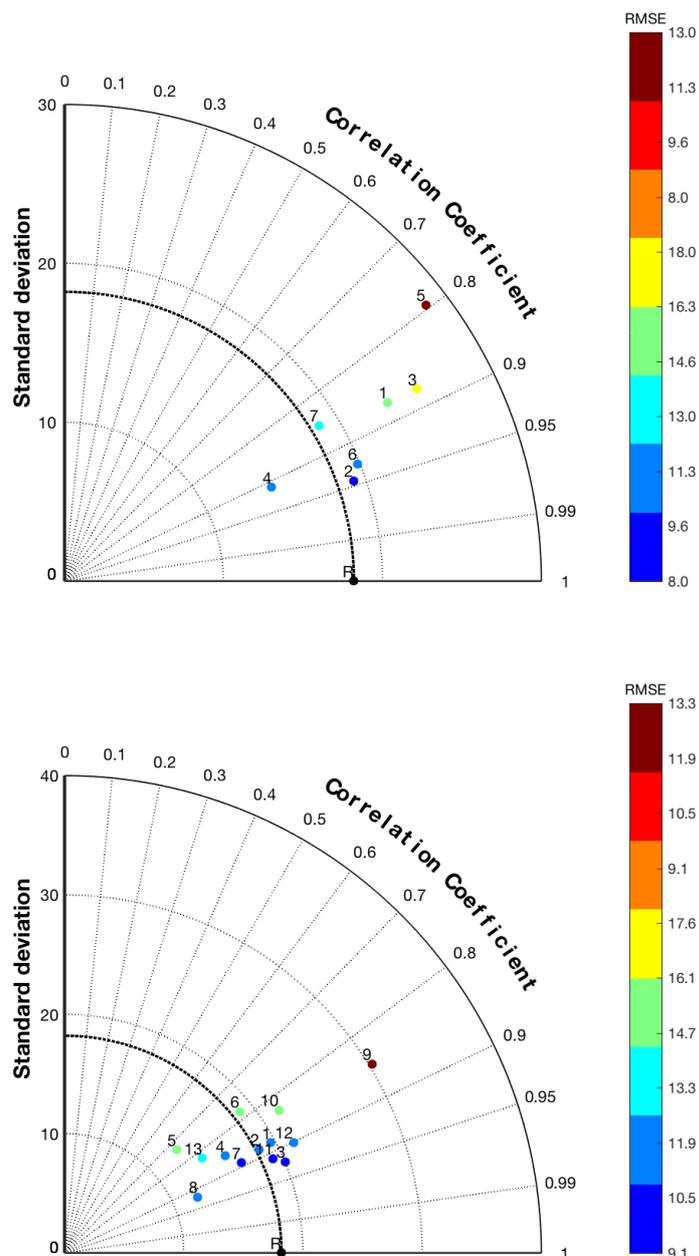


Figure 3 a and b

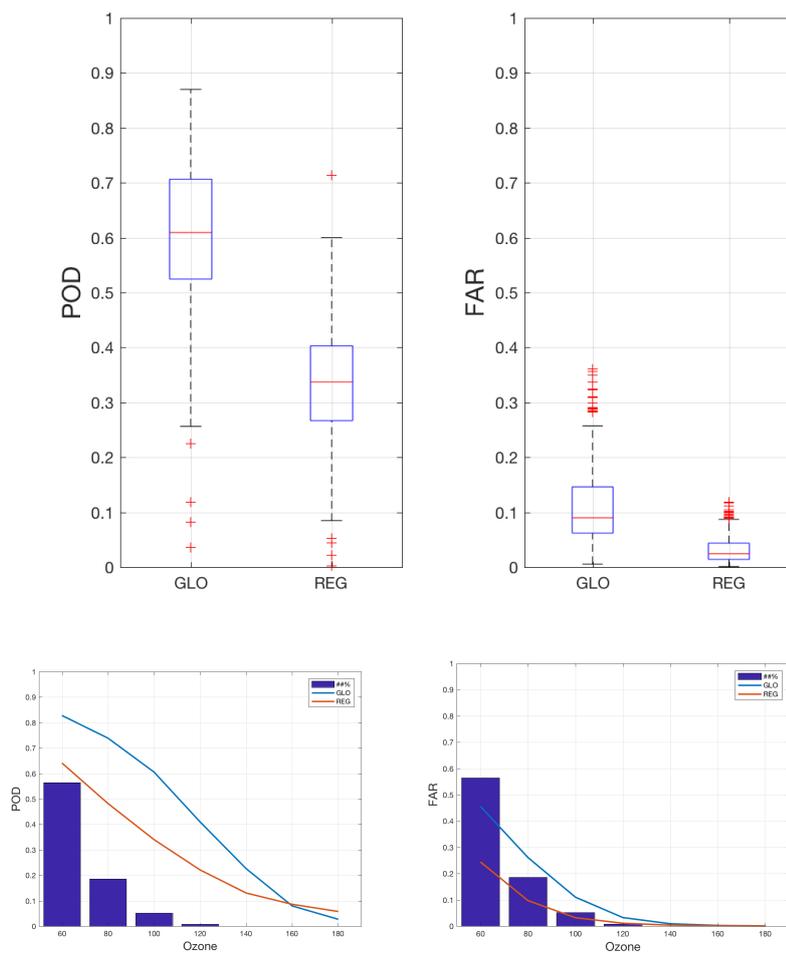


Figure 4 a and b

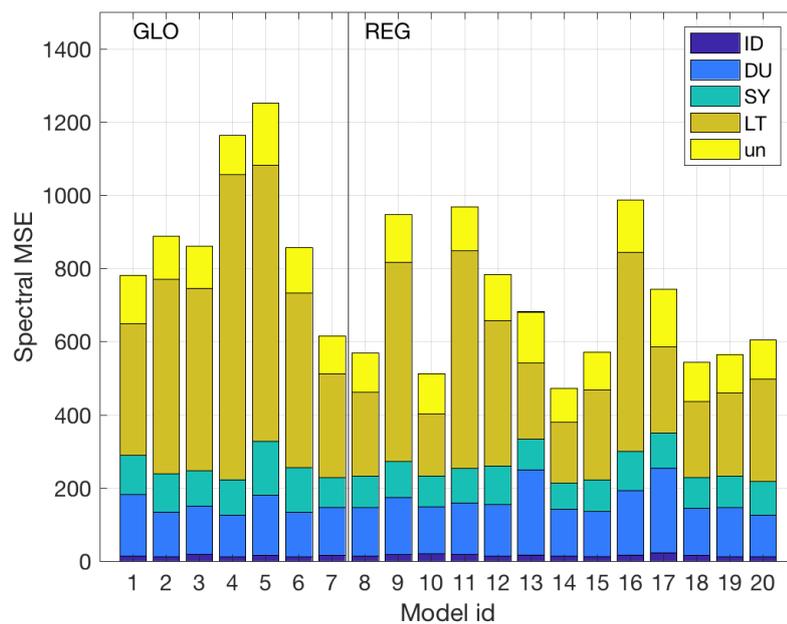


Figure 5

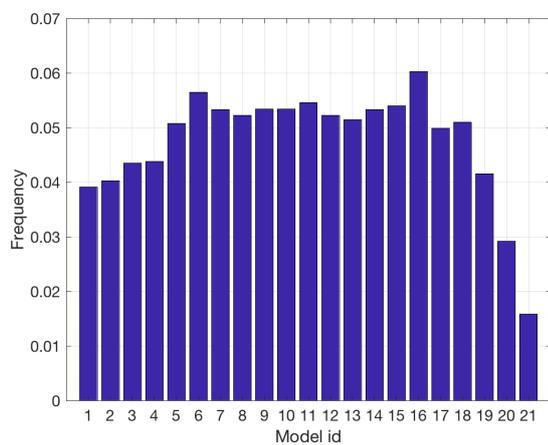
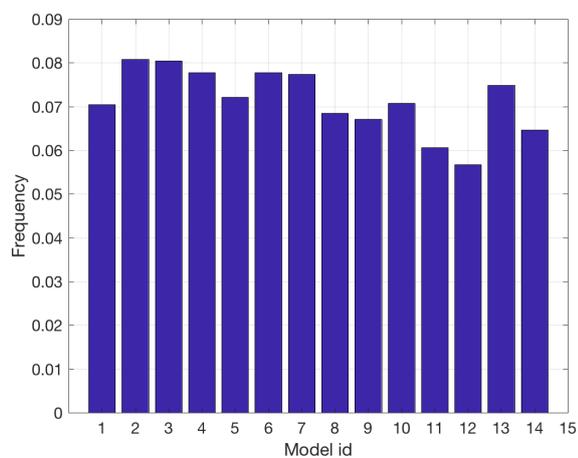
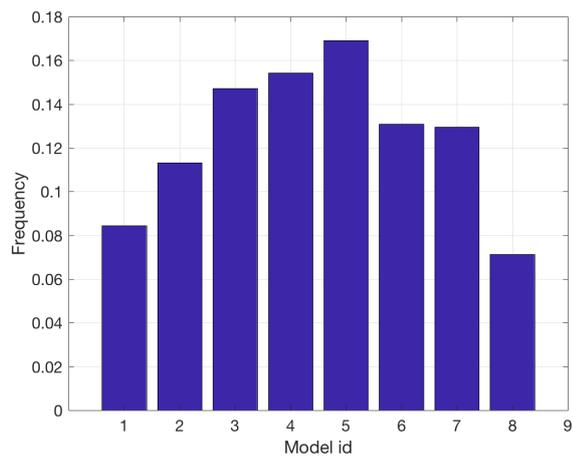


Figure 6 a,b and c

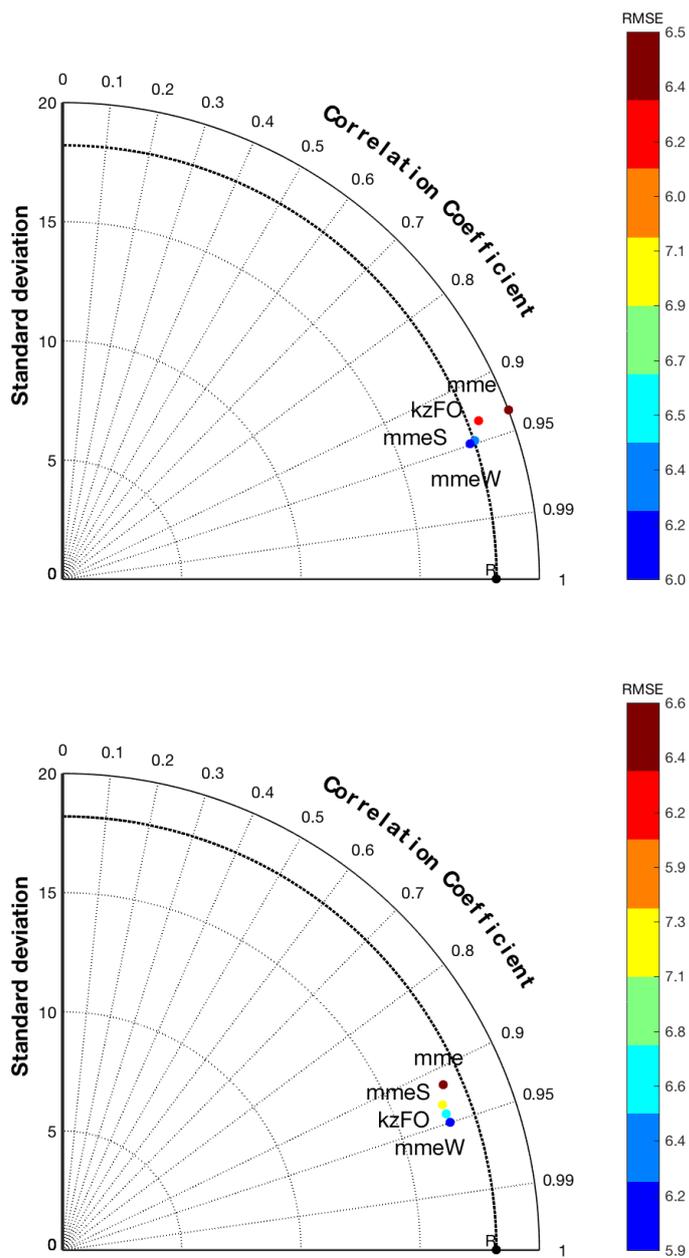


Figure 7 a and b

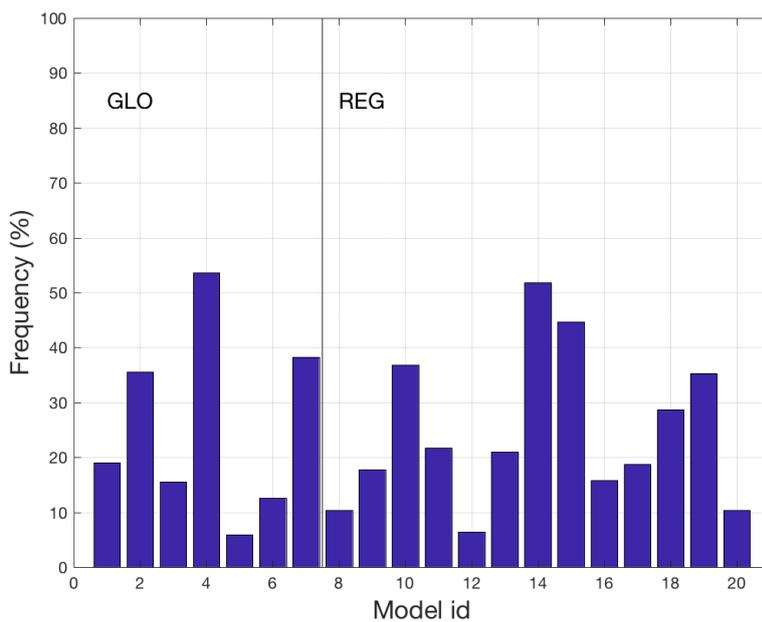
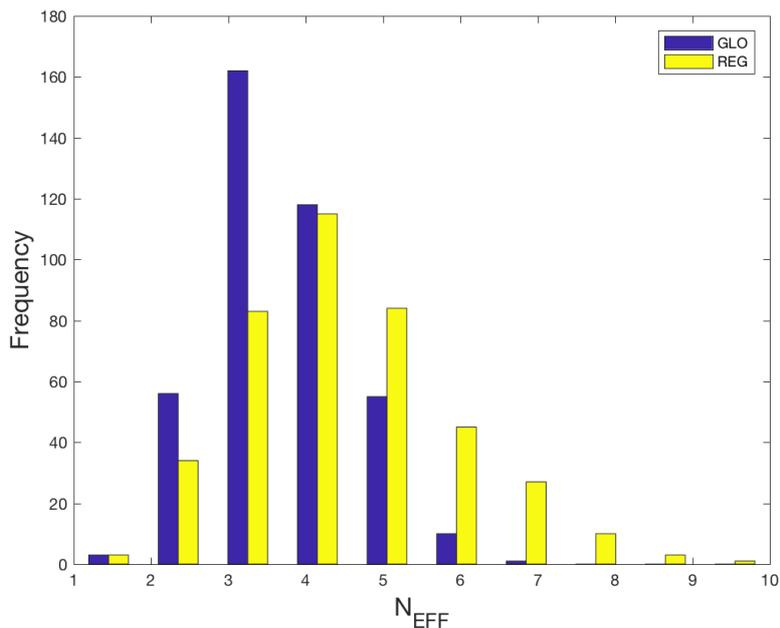


Figure 8 a and b

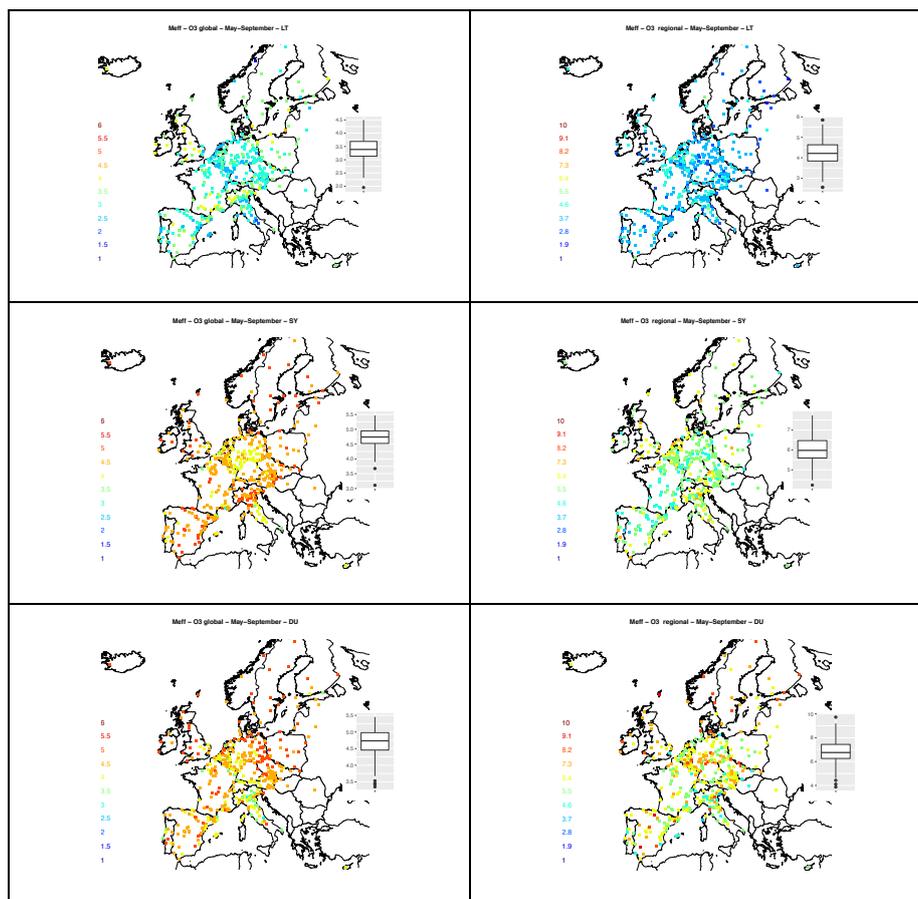


Figure 9

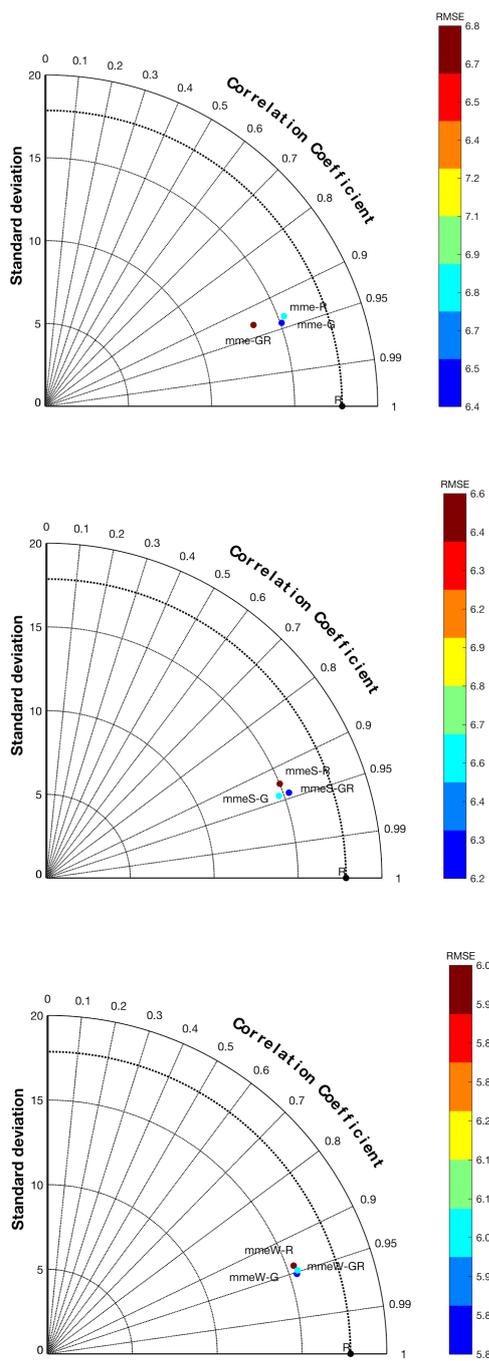


Figure 10 a b and c

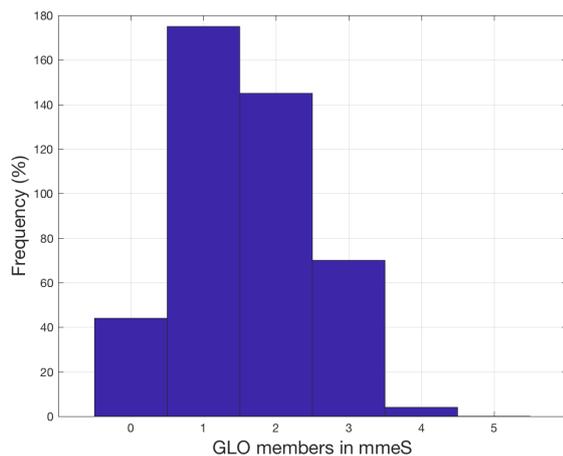


Figure 11

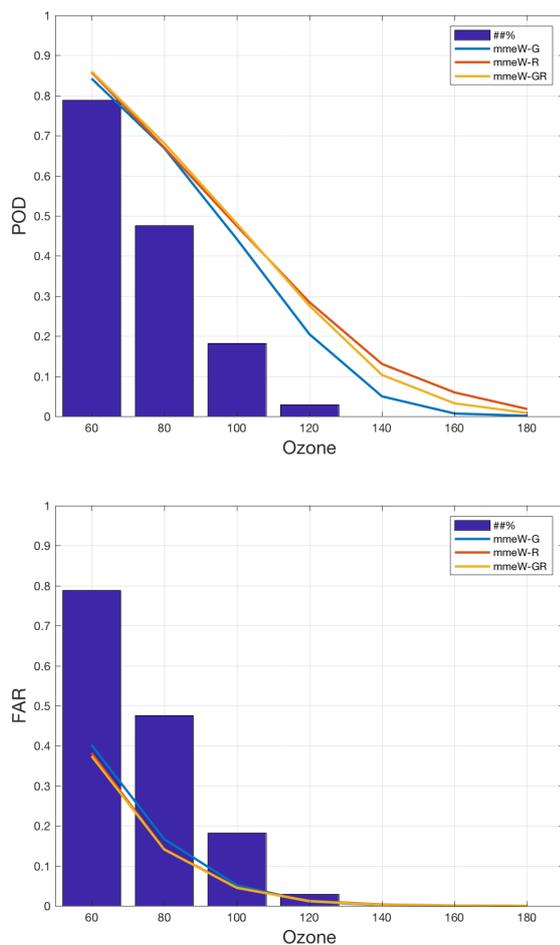


Figure 12 a b

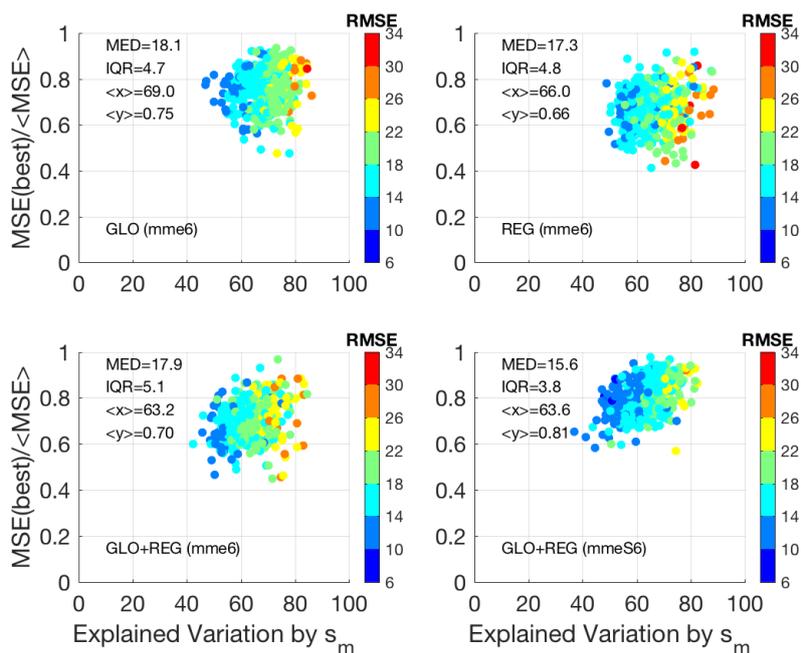
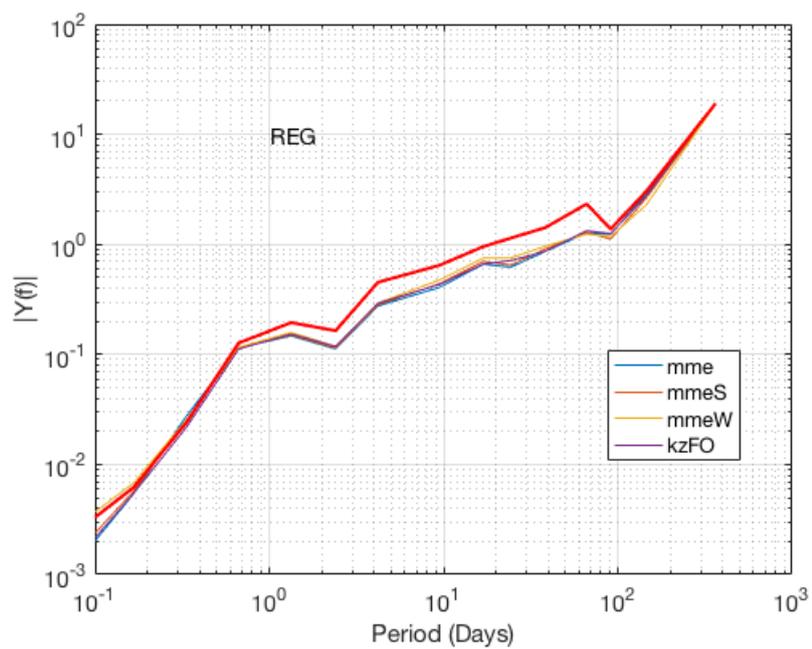
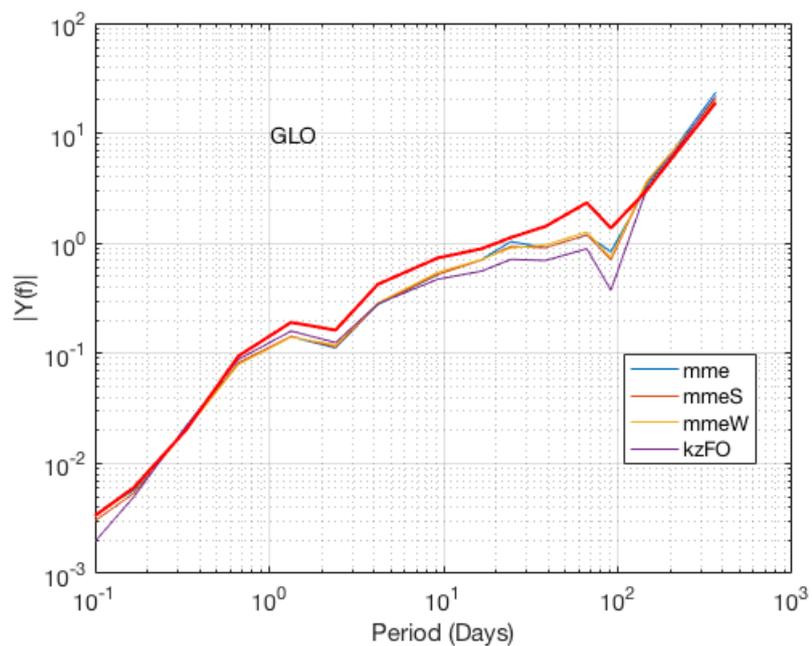


Figure 13



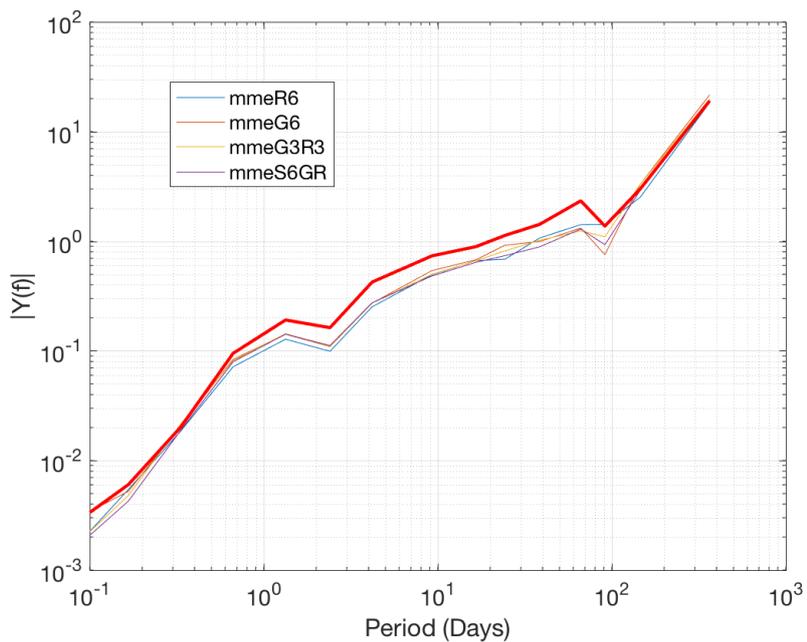


Figure 14 a,b and c