

1 **A Science-Based Use of Ensembles of Opportunities for**
2 **Assessment and Scenario study**

3
4 Efisio Solazzo and Stefano Galmarini

5
6
7 European Commission,

8 Joint Research Centre,

9 Institute for Environment and Sustainability, Air and Climate Unit,

10 Ispra (Italy).

11 **KEY NOTES**

12 Multi-model ensembles need inspection prior use

13

14 **ABSTRACT**

15 The multi-model ensemble exercise performed within the HTAP project context [Fiore et al.,
16 2009] is used here as an example of how a *pre-inspection*, diagnosis and selection of an
17 ensemble, can produce more reliable results. The procedure is contrasted with the often-used
18 practice of simply averaging model simulations, assuming different models produce
19 independent results, and using the diversity of simulation as an illusory estimate of model
20 uncertainty. It is further and more importantly demonstrated how conclusions can drastically
21 change when future emission scenarios are analysed using an un-inspected ensemble. The
22 HTAP multi-model ensemble analysis is only taken as an example of a wide spread and
23 common practice in air quality modelling.

24

25 **1. INTRODUCTION**

26 A multi-model (MM) ensemble is defined as a group of simulations of the same case study,
27 produced by formally different models, which are statistically treated in an attempt to
28 improve the quality of the result [Potempski and Galmarini, 2009]. Given the ever increasing
29 collaborations of geophysical modelling communities in joint assessment studies, MM
30 ensembles are becoming very popular and an opportunity to extend and generalize individual
31 deterministic model results [Solazzo et al., 2012 and; 2013; Solazzo and Galmarini, 2014;
32 Galmarini et al., 2004; Vautard et al., 2012; Evans et al., 2013; Bishop and Abramowitz,
33 2013; and many others].

34 In particular in atmospheric sciences, MM ensembles are used extensively in climate and air
35 quality predictions and assessments. While in climate research and applications many of the
36 concepts applied and described here are well known and correctly used, in air quality this is
37 not always the case and several are the examples of direct use of *un-inspected* MM
38 ensembles. We shall describe an *inspected* MM ensemble (opposed to an un-inspected one)
39 as: a set of model results, whose properties and characteristics, have been analysed in an
40 attempt to reduce the presence of redundant information or elements that are not relevant to
41 the determination of an accurate result. An inspected ensemble should is expected to produce

42 a result that is more accurate than the simple average of the multi model results, at least in all
43 the cases when the members of the ensemble are not independent (e.g., Kioutsioukis and
44 Galmarini, 2014).

45 The motivations behind the necessity to inspect a MM ensemble are connected to the way in
46 which MM ensembles are put together and to the nature of the participating models. In fact,
47 the selection of the models whose results are ensembled is not, to the best of our knowledge
48 and at least for air quality applications, regulated by any science based criteria and there is no
49 a-priori specification that defines the characteristics of a model that should or should not take
50 part to an ensemble. The constitution of a MM ensemble is merely based on an opportunity to
51 provide model simulations and to participate to a community activity where anybody is
52 welcome (*ensemble of opportunity*). Regarding the nature of the models producing results for
53 ensemble applications, one should never forget that the best results are those produced by
54 ensembles of independent (and accurate) models [Potemski and Galmarini, 2009;
55 Kioutsioukis and Galmarini, 2014; Weigel et al., 2008; Pirtle et al. 2010; Knutti, 2010; Knutti
56 et al., 2010; Riccio et al., 2012]. Formally, model m_1 is defined independent from m_2 if the
57 joint probability p for a result of m_1 and m_2 can be expressed as $p(m_1, m_2) = p(m_1)p(m_2)$. When
58 many independent models are combined together their bias can be randomly positive or
59 negative increasing the probability of cancelling out and of the sampled uncertainty does not
60 overlap (Knutti et al., 2010; Abramowitz, 2010; Solazzo et al., 2013). Models used in air
61 quality (among others) are not independent, they are often sharing common assumptions,
62 modules, input data, and cannot therefore be considered independent. In most of the cases the
63 models are different (*Phenotypical model difference*, Potemski and Galmarini, [2009]), but
64 are not independent. This leads to the possibility that results obtained from an ensemble,
65 rather than representing a true alternative and independent solution, would just be like in
66 music composition a *variation on the theme*, producing a false sense of variability which
67 could lead to coinciding (diverging) biased results and a false sense of agreement
68 (uncertainty).

69

70 MM ensembles derived from simply different models are prone to redundancy and
71 overconfidence. The inspection is therefore primarily finalised at:

- 72 - the identification of the level of diversity (communality) shared by the model results,
- 73 - retaining only those that are contributing with original information
- 74 - removing the redundancy.

75 Techniques exist that allow such screenings that rely on the existence of observations and the
76 comparison of the ensemble variability with the observational variability [Potempski and
77 Galmarini, 2009; Solazzo et al., 2013; Riccio et al., 2012].

78 In this study we aim at demonstrating the importance of using existing good practices in the
79 air quality MM ensemble context. Toward the scope we have selected a case study published
80 in the past which does not exploit the true value of having multiple model results at hand. The
81 case analyzed is the HTAP (Hemispheric Transport of Air Pollution) phase 1 multi-model
82 exercise [Dentener et al. 2010] and in particular the multi-model ensemble activity performed
83 within it and presented by Fiore et al. (2009). The study of Fiore et al. (2009) is used here as
84 mere representative of a wide spread practice in the air quality modelling communities at all
85 scales and it represents just an example on how things could be improved further. The MM
86 ensemble by Fiore et al. (2009) is original in many aspects and, in particular, is used for
87 sensitivity studies with respect to emission reduction options. The inspection of the ensemble
88 can have important consequences also for emission scenarios as shown later, an aspect never
89 considered before in the literature.

90

91 **2. THE CASE STUDY AND MM ENSEMBLE INSPECTION**

92 In 2006 the Task Force on Hemispheric Transport of Air Pollution (<http://www.htap.org/>)
93 organised a comparison exercise of global and hemispheric transport models, focussing on
94 the relationships between regional scale emission perturbations and the response in air
95 quality, ecosystem, and climate related variables. The information was used in an aggregated
96 form to evaluate air pollution abatement strategies and their impact across the Northern
97 Hemisphere. Results of the comparison exercise are summarized in Dentener et al., [2010];
98 Sanderson [2008]; Fry et al. [2012]; Wild et al. [2012]; Jonson et al., [2010]; Anenberg et al.,
99 [2009]; Fiore et al., [2009].

100 We will focus on the MM ensemble analysis by Fiore et al. [2009] (from now FetA09). In
101 FetA09, an average of 21 model results was used to investigate the monthly mean surface
102 ozone concentration in three sub-regions of Europe (Mediterranean, Central Europe with
103 receptors between 0 and 1 km height and Central Europe with receptors between 1 and 2 km
104 height), five North-American sub-regions (North East, South West, South East, Great Lakes,
105 and Mountainous) and one Japanese sub-region (EANET stations). Operational scores (bias,
106 correlation coefficient and standard deviation) were calculated in each sub-region making use

107 of ground-based measurements. The combined spatial and temporal average of the modelled
108 concentration values resulted in smoothed monthly time-series. The analysis of FetA09
109 reveals that the distribution of the results is rather symmetric (Figure 1). Supported by the
110 agreement with observations, the authors considered the MM ensemble mean to be the best
111 possible estimate as it “*generally captures the observed seasonal cycle and is close to the*
112 *observed regional mean*” [FetA09], thus justifying the use of the MM ensemble mean to
113 quantify source-receptor relationships as well as ozone concentration response to changes in
114 the emissions scenarios.

115 The scope of the analysis by FetA09 was not to prove the robustness of the MM ensemble
116 mean, and provides an example of the widespread practice of averaging all available
117 members, assuming that the average of many model results is always a better result than that
118 of one model. That would be true if the models were independent but there is no a-priori
119 proof of that. Some questions arise: how robust are the results if the members are not
120 independent models? How different the result would be should some model not taking part to
121 the activity or more outliers (like the one present in the Figure 1) would be present? How
122 generalised is the result since the selection of the ensemble members is based on the
123 voluntary participation to a joint activity and the MM ensemble does not contain all possible
124 results? Is there any duplication of information? Is all the information contained in a MM
125 ensemble relevant and necessary? Since the construction of a MM ensemble is not governed
126 by scientific selection criteria, so it happens that the subsequent ensemble result strictly
127 depends on *aleatory* factors and one can presume that it lacks generality as it is supported by
128 assumptions known to be valid for independent members only.

129 The screening methodology we propose and that we apply as an example to the FetA09 set,
130 is a good way to exploit an abundance of model results in the best way, to transform the
131 aleatory gathering of information into a more robust result that is based on general selection
132 criteria. The large ensemble of model results becomes an opportunity to *cherry-pick* those
133 models whose combination produce the most accurate MM ensemble and use only those to
134 drive conclusions. The analysis will help identifying the size of the non-redundant ensemble
135 and the subsets of members to produce skilled results.

136

137 **2.1 INSPECTING A MULTI MODEL ENSEMBLE**

138 In this section the MM ensemble of FetA09 is inspected. We will concentrate on the ozone
139 simulations over the same regions presented in FetA09 and we will make use of exactly the

140 same model data and observations used in by FetA09 as the main point of the investigation here is to
141 use the same available information of Fetal09 to show that results are different when, an inspected MM
142 ensemble is adopted. The inspection is based on the following steps:

- 143 - determine to what extent the variability (standard deviation about the ensemble mean
144 as in Fortin et al., 2014) present in the observation is reproduced by the ensemble
- 145 - determine the minimum number of models necessary to represent the observed
146 variability
- 147 - identification of the models forming the reduced MM ensemble used for subsequent
148 analysis.

149

150

151 **2.1.1 THE “ACCOUNTED FOR” VARIABILITY: EIGEN-ANALYSIS AND RANKED** 152 **HISTOGRAM TECHNIQUE**

153 The goal of this first analysis is to determine to what extent the observational variability is
154 reproduced by the ensemble. An optimal situation is the one in which the variability of
155 observations coincides with that produced by the ensemble of models, in other words the
156 ensemble of the results all together covers the same range of variation of the measurements.
157 Any deviation from this condition, namely a smaller or a larger variability of the MM
158 ensemble with respect to the observed one would show, on one side, the incapacity of the
159 ensemble to span the observed reality, or on the other, the addition of irrelevant information
160 to the simulation of the observed situation. Therefore considering that a MM ensemble is
161 assembled on an opportunity basis rather than results characteristics, this first step is of
162 primary importance to estimate to what extent the gathered set is appropriate for the case
163 study.

164 A technique to assess the variability and to estimate the redundancy of the MM ensemble
165 with respect to that of the observations, was suggested by Annan and Hargreaves [2010] and
166 applied in several MM ensemble modelling contexts (see, e.g. Solazzo et al., [2013]; Solazzo
167 and Galmarini [2014]). It consists of projecting the observation anomalies (the element-wise
168 difference between the observations and their mean) onto the principal components (PCs) of
169 the covariance matrix of the deviation of the ensemble of models from the MM mean (the
170 element-wise difference between each model realisation and the MM ensemble mean).
171 Principal component analysis [Jolliffe, 2002] is probably the most well-known and wide-

172 spread dimension-reduction technique. It is based on eigen-analysis to select uncorrelated
173 directions associated with the largest variances.

174 When applied to the HTAP 21-member ensemble analysed by FetA09, this method shows
175 that the first (largest) eigenvalue already explains more than 90% of the observational
176 variability in most regions, the only exception being Japan with 60%. In other words, most of
177 the ensemble members have a significant projection onto the first eigen-vector defining the
178 major component, thus explaining the same portion of variance. If too many models are
179 projected on the same eigenvector, it means that there are too many models producing
180 repeating or 'overlapping' solutions (thus, the MM ensemble is redundant and
181 overconfident). A well-behaved MM ensemble (not necessarily the theoretical case of
182 independent models) should be made of a number of models whose eigenvalues contribute to
183 the explanation of as many different components as the observational variability and the ratio
184 model-to-observed variance should be close to unity. In the case of the HTAP MM ensemble,
185 when all eigen-values are taken into account (and all of the associated eigen-vectors), the
186 MM ensemble variance is 4.7, 6.0, 8.7 times the variance of the observation anomalies for the
187 EU Mediterranean, Central 0-1 km, and Central 1-2 km, regions respectively. Concerning
188 the US Mountains, Great Lakes, SE, NE, SW regions, the full MM ensemble mean accounts
189 for 25.4, 9.1, 20.6, 10.7, 5.6 times the observed variability, respectively, and finally 4.7 times
190 for the Japanese sub-region. According to the definition of Annan and Hargreaves [2010] the
191 ensemble is therefore *wide*, i.e. its variability is larger than the observed one. Dealing with a
192 wide ensemble implies that there is a substantial amount of redundant variability, i.e.
193 variability already accounted for by other models. Not all information contained in the
194 ensemble is needed in principle and needs to be reduced.

195 An alternative method to diagnose the variability spanned by an ensemble of models to the
196 eigenvalues used is the Talagrand or Ranked Histogram (RH) [Talagrand et al., 1998], which
197 provides an evaluation of the consistency of the ensemble with an observed quantity. In a RH
198 the observations are ranked into a number of bins equal to the number of models making up
199 the ensemble plus one for the extremes. The ensemble members are sorted to define ranges or
200 "bins" of the modeled variable such that the probability of occurrence of the observation
201 within each bin is, ideally, equal. The bins are determined by ranking the ensemble member
202 from lowest to highest. The interval between each pair of ranked values forms a bin. To a N -
203 member ensemble correspond $N+1$ bins [Hamill, 2001]. The underlying assumption is that
204 each ensemble member in principle introduces an independent degree of variability. An

205 indication of an ill-constructed ensemble is the ratio between the number of elements and the
 206 number of data available per model. If there are N models with time series each of size n_t
 207 (elements of the time series), the implication of $N > n_t$ is that there will be at least $N - n_t$ empty
 208 bins in the RH, indicating redundancy of the ensemble and that the ensemble is inappropriate
 209 for the case analyzed. This same result could be visualized by looking at the load factors
 210 resulting from the decomposition in PCs: many projections would be null, as the number of
 211 eigen-vector is larger than the number of data to project. The HTAP MM ensemble used in
 212 this example, $N = 21$ and $n_t = 12$. The RH for the nine sub-regions is reported in Fig. 2. Six
 213 (NA NE) to nine (NA SW) bins out of 22 are populated, (i.e. contain non-zero values), due to
 214 insufficient data and excess of redundant information. The use of the RH reveals another
 215 important problem with the FetA09 MM ensemble. Good ensemble practice would require n_t
 216 $\gg N$. The plots clearly show that there are many empty bins (so degrees of freedom in the
 217 process that are not part of the reality as no observations are present in that range). The
 218 uneven distribution of the histograms shows that much emphasis (overconfidence) is given to
 219 some aspects of the process description, while others are neglected, that is another way of
 220 representing the redundancy obtained with PC analysis presented earlier

221

222 **2.1.2 EFFECTIVE NUMBER OF MODELS**

223

224 Having assessed that the ensemble is redundant it is important to determine the minimum
 225 number of models from those available in the MM ensemble that would suffice to describe
 226 the observational variability. A method developed by Bretherton et al. [1999], and firstly
 227 applied to air quality models by Solazzo et al. [2013], quantifies the effective number of
 228 models sufficient to reproduce the variability of the observation as:

$$229 \quad N_{eff} = \frac{(\sum_{k=1}^N \lambda_k)^2}{\sum_{k=1}^N \lambda_k^2} \quad \text{Eq (1)}$$

230

231 with λ eigenvalue of the $corr(d_i, d_j)$ matrix, which contains the linear correlation coefficient
 232 between any pair d_i, d_j ($i, j = 1, \dots, N$). d is a metric defined accordingly to Pennel and Reichler
 233 [2011]:

$$234 \quad d_m = e_m - R \text{ MME} \quad \text{Eq (2)}$$

235

236 where the index m identifies the model, MME is the multi model error (the average of all
237 individual model's errors) and R is the Pearson correlation coefficient between e_m , the error
238 of model m and the MME. The removal of MME in Eq. (2) makes model errors more
239 dissimilar from one another and uncovers "hidden" trends that are outweighed by overarching
240 commonalities. Indeed the scope of the metric d_m is to determine similarities between models
241 beyond the dominating ones induced by shared inputs and/or common parameterisations to
242 the extent that the former are accounted for in the average. The relationship (1) should be
243 interpreted as: only if all eigenvalues were equal to unity, Eq. (1) would take a value of N_{eff}
244 $=N$, which corresponds to the situation where all directions are equally important and all
245 models add independent contributions to the explanation of the observational variability. On
246 the other hand, if all error fields were similar, only one eigenvalue would be non-zero and N_{eff}
247 $= 1$. Equation (1) provides an analytical estimate of the dimensions of the subspace of models
248 necessary to produce the information of the whole ensemble.

249

250 For the HTAP MM ensemble of FetA09, Eq. (1) gives N_{eff} ranging between ~ 2 and 4 for the
251 regions analysed by FetA09 compared to the original 21 models. Thus, approximately three
252 quarter of the available members participate to the ensemble with already 'accounted for'
253 information. This is a revealing result that indicates paradigmatically the relevance of a pre-
254 inspection of an ensemble. What seemed like a largely populated ensemble turns out to be
255 incapable of capturing several degrees of freedom of observations and 2 to 4 members of 21
256 are sufficient to describe the observational variability. One may ask: if so, why is the average
257 of the 21 models fitting so well with the observations as presented in FetA09? The answers
258 could be: pure chance, since finally the model results participated out of good will, and
259 happened to be there in the right mixture. Just consider what would have happened to the
260 mean of the models should one of the two most evident outliers in Figure 1 decide to
261 withdraw from the exercise. Alternatively an explanation could be the massive smoothing
262 due to the monthly averaging along with the high level of tuning of the models around
263 specific solutions that are normally distributed around the average observed data.

264

265 **2.1.3 REDUCING ENSEMBLES**

266 As demonstrated in the previous sections, the HTAP MM ensemble is redundant and in
267 particular 2 to 4 members are sufficient to represent the observational variability while the
268 rest do not add any new information. Similarly, the extra elements are likely to deteriorate

269 any evaluation metrics applied to the ensemble. At this point we know that the number of
270 models that are necessary and sufficient is smaller than 21 but we do not know which
271 combination of members for every grouping produces the optimal ensemble.

272 Given N members, there are $G=N!/[(r!(N-r)!]$ possible groups of r elements. A straight
273 forward way to identify the optimal ensemble (optimal sub set) and maximize the accuracy of
274 the ensemble is to analyse all the G combinations of subsets of models and identify the one
275 that minimize the Root Mean Square Error (RMSE). The latter is a measure of the accuracy
276 (the even distribution of model results from the observed value), and high accuracy also
277 improves precision (a reduced spread/scatter of the model results around the observed value).
278 In principle measurement errors should be also taken into account in the procedure for
279 reducing the ensemble, but in case where they are significantly smaller than the model ones,
280 RMSE is sufficient measure.

281 In Fig (3) we report the curves of minimum, mean, and maximum RMSE for the nine sub-
282 regions used by FetA09 as a function of the number of members of ensembles ($r=2,\dots,21$).
283 The figure confirms the results on the number of models necessary to maximize the ensemble
284 performance and tells us that which combination of the 2 to 4 models out of 21 produces such
285 improvement. The scores of the reduced ensemble are reported in Table 2 and are compared
286 against the ones produced by the full ensemble mean. In all cases the mean of the reduced
287 ensemble improves the accuracy (from 31% for NA NW to 71% for NA Mountain and NA
288 Lakes) and precision (most notably for NA SE and NA NE). As it can be seen in several
289 regions the use of the full MM ensemble of opportunity produces a clear deterioration in the
290 ensemble statistics. In Table 2 we report also the ranking of the models contributing to
291 minimize the error in the sub-regions. As from the table it is often the case that the error is
292 minimized by mix-ranked (good performing and bad performing) of members. In fact, if the
293 two best models have a high chance of being also highly correlated then they would share
294 some portion of information, thus resulting redundant. Therefore when considering the
295 ensemble mean of these two models, very little decrease in error would be found compared to
296 the individual models. Mathematically, the theorems by Elashoff et al. [1967] and Cover
297 [1974] have proven two important results on the selection of member and evaluation of
298 individual scores: the best two models are seldom the combination of two models that
299 maximises the score of an ensemble average, and furthermore, that the best single model may
300 not appear in the ensemble maximising the feature score. As a result, the simple method of
301 making ranked combinations of models with the best individual features may prove

302 unsuccessful, as also demonstrated by e.g. Solazzo et al. [2013], Hannan and Hargreaves
303 [2011], Kioutsioukis and Galmarini [2014], Knutti et al., [2010], and others. This confirms
304 the importance of the inspection of the available results prior to their use and of having at
305 disposal a large pool of models from which optimal subsets can be extracted.

306

307 **3. IMPACT ON THE RESULTS OF EMISSION SENSITIVITY ANALYSIS OF AN INSPECTED** 308 **VS UNINSPECTED ENSEMBLE**

309 An important part of FetA09 relates to the sensitivity study on emission reduction. As part of
310 the HTAP program the consequences of an emission reduction of 20% anthropogenic NO_x in
311 specific part of the globe where investigated using the MM ensemble available. Since we
312 have demonstrated that the MM ensemble used in FetA09 is redundant and having identified
313 the optimal number of elements and the most accurate set of models, one may wonder how
314 the predicted consequences of the emission reduction on ozone concentration would change if
315 we used the reduced ensemble.

316 We focused the analysis on the North-American region only. In FetA09 the use the mean of
317 the full ensemble produced an average response in ozone concentration of -0.76 ppb in the
318 NA region as a consequence of the reduction of NO_x emission by 20%. We shall note that the
319 NA region is subjected to the emission reduction and therefore the investigation includes the
320 whole of the US and part of Mexico (Figure 1 of FetA09), and thus it has a spatial extension
321 that includes the five NA sub-regions described in section 2 for the evaluation. Furthermore,
322 of the 21 models participating to the evaluation part of the exercise, only 14 models results
323 were made available for the simulation with reduced emission scenarios. Therefore, for the
324 sake of consistency, we repeated the redundancy inspection for the 14-member ensemble and
325 calculated the most accurate set through the minimization of RMSE described section 2.1.3.
326 The size of the newly calculated subsets ranges between three for the Lakes, North-East,
327 South-West, South- East of USA, and four for the Mountainous region. The newly calculated
328 set obtained from the original 14 member ensemble produced an ozone concentration
329 reduction of 2.32 ppb on average across all regions. That is 300% more than that found by
330 FetA09. The largest variation is obtained for the South-East region of USA, with an ozone
331 concentration decrease of 5.30 ppb that is a 5-fold than what obtained by FetA09. Such an
332 analysis demonstrates how conclusions could change if the ensemble is not inspected a priori
333 and reduced if necessary.

334 In the exploration of scenario or sensitivity to ideal conditions like that presented in HTAP,
335 one may be tempted to construct an ensemble that only groups the best performing models
336 results in the evaluation against measurements and using only those in the sensitivity or
337 scenario case study grouping them in an ensemble. This would be wrong in principle or in
338 other words would not produce the best ensemble by definition as demonstrated by the
339 already cited theorems of Elashoff et al. [1967] and Cover [1974].

340

341 4. CONCLUSIONS

342 Multi-model ensemble is becoming very popular in geophysical studies. In this paper we
343 have been contrasting the results from an *ensemble of opportunity* where casually assembled
344 model *phenotypical different* are the driving elements, with the results obtained when the
345 same pool of model is screened to eliminate redundancy and the optimal combination is used.

346 The case of HTAP phase 1 is taken here as an example of a practice that is wide spread,
347 especially in the realm of air quality, atmospheric dispersion at all scales. A very limited
348 amount of studies apply correctly the technique. The HTAP case has been selected for two
349 main reasons:

- 350 - The very large number of models that participated to the initiative and that were
351 available for the ensemble analysis;
- 352 - the ensemble results were also used as basis to assess the consequences of an emission
353 reduction strategy on ozone in several regions of the world.

354 The HTAP ensemble has been assessed against available measurements and the following
355 conclusion were obtained:

- 356 - In spite of the large number of participating models, the scarcity of time steps
357 produces an important level of redundancy as from the simple analysis of a ranked
358 histogram.
- 359 - At smaller subset of model perform much better when compared to measurements and
360 it is statistically more significant.
- 361 - In the case of HTAP [FetA09] the objective of the study was to determine, through a
362 MM ensemble, the impact of emission changes produced in one continent on another.
363 The analysis conducted on the impact over the same continent where the emissions

364 are produced, reveals that the conclusions remain the same as those produced by
365 FetA09 but the values found are between 3 to 5 times higher when using a non-
366 redundant ensemble.

367 These are problems that are common to many multi model studies and for which a minimum
368 set of good practice rules should be taken into account (Kioutsioukis and Galmarini, 2014).
369 Among these, we point out that in order to have any reasonable statistics the number of
370 measurement should be much greater than the number of ensemble members. Otherwise rank
371 histogram is simply not a proper tool for the analysis.

372 On a more general level, it is clear that the use of un-inspected ensembles of opportunities is
373 a miss-practice that could lead to under-exploitation of the latter and in some case even
374 wrong conclusions. Quantitative practices guarantee for the best possible diagnosis of the
375 ensemble potential and its full exploitation. The availability of monitoring information is
376 essential for the performance of the analysis presented here and it could be argued that the
377 optimal ensemble identification is prone to the time and spatial representativity of the
378 observations. This is true but as much as it is for the evaluation of any individual model result
379 that depends on the space and time distribution of observation and the phenomenology
380 represented.

381 The hemispheric transport case analyzed here brings to the attention also the issue of the
382 space and timescale at which a set of model verified in a certain area could be used. The
383 verification of the effect of the selection of an optimal set out of an ensemble based on data
384 pertaining to a specific region and time frame, produces over another region, remains an
385 important element of research. In other words, whether an optimal set selected for region A
386 using observation in region A can be used for a region B and in a scenario or sensitivity
387 analysis mode. Scale dependence of the atmospheric processes involved could become an
388 issue in this case and will have to be verified. On the other end we consider the use of the
389 optimal set for scenario and sensitivity study in the area where the observation used for its
390 selection have been collected much more appropriate than the use of a full ensemble of
391 opportunity. The selection of the optimal set through observations on a base case scenario is
392 equivalent to the evolution of a single deterministic model and its application for speculative
393 scenario analysis or forecast applications.

394 The representativity of the ensemble compared to observation and the minimization of the
395 redundancy remain an important issue. In the light of that we speculate here, the use of multi-

396 scale multi-model ensembles, constructed with the combinations of models covering different
397 portions of the atmospheric power spectrum, could greatly improve the representativity and
398 provide coverage of the problem in a much more detailed form. The combination of global
399 and regional scale results, for example, in one ensemble is a possibility that will be explored
400 in the framework of the next phase of HTAP.

401

402 **ACKNOWLEDGMENTS**

403 Dr Arlene Fiore (Columbia University) and the HTAP modeling community are
404 acknowledged for making the model and observational data available for the current analysis
405 (<http://www.htap.org/>) and for the openness to our investigation. Dr Frank Dentener (JRC) is
406 acknowledged for the valuable comments that greatly improved this manuscript. The authors
407 also thank Dr Brigitte Koffi (JRC) for having retrieved some of the data used in this paper.

408 REFERENCES

409

410 Anenberg S.C., J. J. West, A. M. Fiore, D. A. Jaffe, M.J. Prather, D. Bergmann, K. Cuvelier, F. J.
 411 Dentener, B. N. Duncan, Michael Gauss, Peter Hess, Jan Eiof Jonson, Alexandru Lupu, Ian
 412 A. MacKenzie, Elina Marmer, Rokjin J. Park, Michael G. Sanderson, Martin Schultz, Drew
 413 T. Shindell, Sophie Szopa, Marta Garcia Vivanco, Oliver Wild, Guang Zeng, (2009),
 414 Intercontinental Impacts of Ozone Pollution on Human Mortality, *Environ. Sc.Tech*, 43,
 415 6482–6487

416 Annan, J.D., Hargreaves, J.C., (2010), Reliability of the CMIP3 ensemble. *Geophys. Res. Lett.*, 37, p.
 417 L02703

418 Bishop, C.H., Abramowitz, G., (2013), Climate model dependence and the replicate Earth paradigm.
 419 *Clim. Dyn* 41, 885-900

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

422 Cover, T. T. (1974). The best two independent measures are not the two best, *IEEE Trans. System*
 423 *Man. and Cybernetics*, 4, 116–117.

424 Dentener F, T. Keating and H. Akimoto, (eds) , (2010), Hemispheric Transport of Airpollution, Part
 425 A, Ozone and Particulate Matter, Edited by F, Economic Commission for Europe, Air
 426 Pollution Studies, 17, ISBN, 978-92-1-117043-6, UNECE, Geneva.

427 Elashoff, J.D., Elashoff, R.M., Goldman, G.E., (1967), On the choice of variables in classification
 428 problems with dichotomous variables. *Biometrika* 54, pp. 668–670

429 Evans, J.P., Ji, F., Abramowitz, G., Ekstrom, M., (2013), Optimally choosing small ensemble
 430 members to produce robust climate simulations. *Environ. Res. Lett.* 8, 044050 (4pp).

431 Fiore, A. M., Dentener, F. J., wild, O., Cuvelier, C., Schultz, M. G., Hess, P., Textor, C., Schulz, M.,
 432 Doherty, R. M., Horowitz, L. W., MacKenzie, I. A., Sanderson, M. G., Shindell, D. T.,
 433 Stevenson, D. S., Szopa, S., Van Dingenen, R., Zeng, G., Atherton, C., Bergmann, D., Bey,
 434 I., Carmichael, G., Collins, W. J., Duncan, B. N., Faluvegi, G., Folberth, G., Gauss, M.,
 435 Gong, S., Hauglustaine, D., Holloway, T., Isaksen, I. S. A., Jacob, D. J., Jonson, J. E.,
 436 Kaminski, J. W., Keating, T. J., Lupu, A., Marmer, E., Montanaro, V., Park, R. J., Pitari, G.,
 437 Pringle, K. J., Pyle, J. A., Schroeder, S., Vivanco, M. G., Wind, P., Wojcik, G., Wu, S., and
 438 Zuber, A., (2009) Multimodel estimates of intercontinental source-receptor relationships for
 439 ozone pollution, *J. Geophys. Res.*, 114, D04301, doi:10.1029/2008JD010816.

440 Fortin, V., Abaza, M., Anctil, F., Turcotte, R., (2014). Why should ensemble spread match the RMSE
 441 of the ensemble mean?. *J.Hydrometeor* 15, 1708-1713.

442 Fry M.M., V. Naik, J. J. West, M. D. Schwarzkopf, A.M. Fiore, W.J. Collins, F.J. Dentener, D. T.
 443 Shindell, C. Atherton, D. Bergmann, B. N. Duncan, P. Hess, I. A. MacKenzie, E. Marmer,
 444 M. G. Schultz, S. Szopa, O.Wild, G Zeng, (2012), The influence of ozone precursor
 445 emissions from four world regions on tropospheric composition and radiative climate
 446 forcing, *J. Geophys. Res.*, 117, D7, doi:10.1029/2011JD017134.

447 Galmarini S., R. Bianconi, W. Klug, T. Mikkelsen, R. Addis, S. Andronopoulos, P. Astrup, A.
 448 Baklanov, J. Bartniki, J.C. Bartzis, R. Bellasio, F. Bompay, R. Buckley, M. Bouzom, H.

- 449 Champion, R. D'Amours, E. Davakis, H. Eleveld, G.T. Geertsema, H. Glaab, M. Kollax, M.
450 Ilvonen, A. Manning, U. Pechinger, C. Persson, E. Polreich, S. Potemski, M. Prodanova, J.
451 Saltbones, H. Slaper, M.A. Sofiev, D. Syrakov, J.H. Sørensen, L. Van der Auwera, I.
452 Valkama, R. Zelazny (2004). Ensemble dispersion forecasting—Part I: concept, approach
453 and indicators. *Atmos. Environ.*, 38 (28), pp. 4619–4632
- 454 Hamill, T.M., (2001), Interpretation of rank histograms for verifying ensemble forecasts. *Mon.*
455 *Weather Rev.*, 129 (3), 550–560
- 456 Jolliffe, I., (2002), *Principal component analysis*, Springer, 2nd edition.
- 457 Jonson J.E., A. Stohl, A.M. Fiore, P. Hess, S. Szopa, O. Wild, G. Zeng, F.J. Dentener, A. Lupu, M.G.
458 Schultz, B.N. Duncan, K. Sudo, P. Wind, M. Schulz, E. Marmer, C. Cuvelier, T.j. Keating,
459 A. Zuber, A. Valdebenito, V. Dorokhov, H. De Backer, J. Davies, G.H. Chen, B. Johnson,
460 and D.W. Tarasick, (2010), A multi-model analysis of vertical profiles, *Atmos. Chem. Phys.*
461 10, 5759-5783.
- 462 Knutti, R. , (2010), The end of model democracy?, *Climate Change*, 102, 395–404.
- 463 Knutti, R., Furrer, R., Tebaldi, C., Cermak, J., Meehl, G.A., (2010). Challenges in combining
464 projections from multiple climate models. *American Meteorological Society* 23, 2739-2758.
- 465 Kioutsioukis I. and S. Galmarini, (2014), De praeceptis ferendis: good practice in multi-model
466 ensembles, *Atmos. Chem. Phys. Discuss.*, 14, 15803-15865
- 467 Pennel, C., Reichler, T., (2011), On the effective numbers of climate models *J. Clim.*, 24, 2358–2367
- 468 Pirtle, Z., Meyer, R., Hamilton, A., (2010), What does it mean when climate models agree? A case for
469 assessing independence among general circulation models. *Environ. Sci. Policy*, 799, 351–
470 361
- 471 Potemski, S., Galmarini, S., (2009), Est modus in rebus: analytical properties of multi-model
472 ensembles. *Atmos. Chem. Phys.* (2009), pp. 9471–9489
- 473 Riccio, A., Ciaramella, A., Giunta, G., Galmarini, S., Solazzo, E., Potemski, S., 2012. On the
474 systematic reduction of data complexity in multi-model ensemble atmospheric dispersion
475 modelling. *Journal Geophysical Research* 117, D05314.
- 476 Solazzo, E., Bianconi, R., Vautard, R., Appel, K. W., Moran, M., D., Hogrefe, C., Bessagnet, B., 5
477 Brandt, J., Christensen, J. H., Chemel, C., Coll, I., van der Gon, H. D., Ferreira, J., Forkel,
478 R., Francis, X. V., Grell, G., Grossi, P., Hansen, A. B., Jericevic, A., Kraljevic, L., Miranda,
479 A. I., Nopmongcol, U., Pirovano, G., Prank, M., Riccio, A., Sartelet, K. N., Schaap, M.,
480 Silver, J. D., Sokhi, R. S., Vira, J., Werhahn, J., Wolke, R., Yarwood, G., Zhang, J., Rao, S.
481 T., and S. Galmarini, (2012): Ensemble modelling of surface level ozone in Europe and
482 North America in the context of AQMEI, *Atmos. Environ.*, 53, 60–74.
- 483 Solazzo, E., Riccio, A., Kioutsioukis, I., Galmarini, S., (2013), Pauci ex tanto numero: reducing
484 redundancy in multi-model ensembles. *Atmos. Chem. Phys.* 13, 8315–8333
- 485 Solazzo, E., Galmarini, S., (2014), The Fukushima-¹³⁷Cs deposition case study: properties of the
486 multi-model ensemble. *J. Environ. Radioact.* 139, 226-233
487 <http://dx.doi.org/10.1016/j.jenvrad.2014.02.017>.

- 488 Talagrand, O., Vautard, R., and Strauss B., (1998), Evaluation of probabilistic prediction systems,
489 paper presented at aa seminar on predictability, Eur. cent. For Medium Weather Forecasting,
490 Reading (UK).
- 491 Vautard, R., Moran, M. D., Solazzo, E., Gilliam, R. C., Matthias, V., Bianconi, R., Chemel, C.,
492 Ferreira, J., Geyer, B., Hansen, A. B., Jericevic, A., Prank, M., Segers, A., Silver, J. D.,
493 Werhahn, J., Wolke, R., Rao, S. T., and Galmarini, S., (2012), Evaluation of the
494 meteorological forcing used for AQMEII air quality simulations, *Atmos. Environ.*, 53, 15–37
- 495 Weigel, A.P., Liniger, M.A., Appenzeller, C., (2008). Can multi-model combination really enhance
496 skill of probabilistic ensemble forecast? *Q.J.R. Meteorolo. Soc.* 134, 241-260.
- 497

498 **Table 1.** Number of effective models N_{eff} for the sub-regions object of the analysis (with reference to Figure 2
 499 of Fiore et al (2009) top panel, based on $corr(di,dj)$). $nrec$ is the number of surface receptors used for evaluation

500

Sub-region	N_{eff}
EU Mediterranean region (nrec=6)	4.0
EU central region 0-1 km (nrec=24)	3.1
EU central region 1-2 km (nrec=11)	3.5
NE-USA (nrec=13)	1.9
SW USA (nrec=5)	1.8
SE USA (nrec=6)	1.9
Great Lakes USA (nrec=8)	2.0
Mountainous USA (nrec=10)	1.8
Japan EANET (nrec=10)	2.6

501

502

503

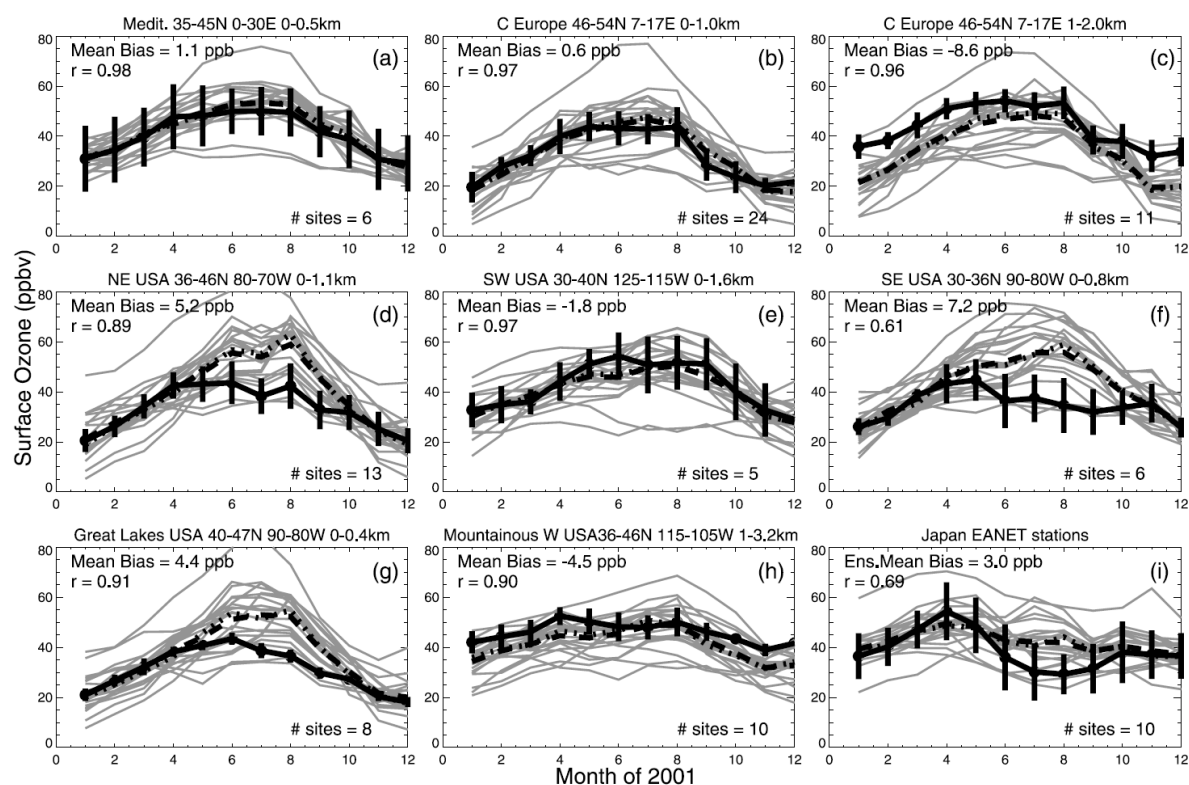
504 **Table 2.** RMSE-ranking and scores of the reduced MM ensemble mean for the sub-regions object of the
 505 analysis (RMSE: Roor-Mean-Square-Error; PCC: Pearson Correlation Coefficient; σ : ratio of the modelled to
 506 the observed standard deviation)

507

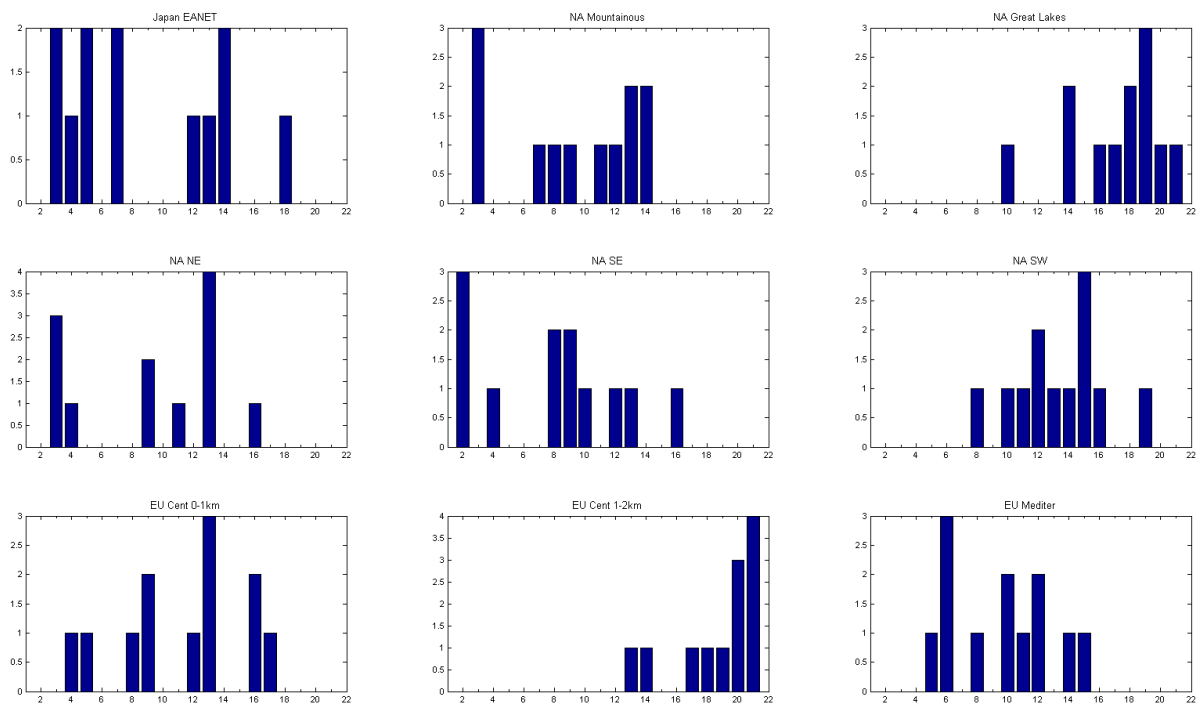
Domain	Ranking of the MinRMSE combination	score
EU central 0-1 km	1,15,19	RMSE=1.69 (2.65) PCC=0.98 (0.96) σ =0.99 (1.10)
EU central 1-2 km	7,17,18	RMSE=3.35 (9.2) PCC=0.98 (0.95) σ =1.03 (1.25)
EU mediterr	4,6,13,15,19	RMSE=0.76 (1.44) PCC=0.99 (0.98) σ =1.0 (1.13)
NA SW	8,10,11,15	RMSE=2.0 (2.9) PCC =0.95 (0.96) σ =0.87 (0.86)
NA SE	1,2,4,8	RMSE=3.61 (10.27) PCC=0.77 (0.62) σ =0.83 (1.81)
NA NE	3,5,6,7	RMSE=3.01 (7.8) PCC=0.93 (0.90) σ =0.90 (1.56)
NA Mountain	1,5,12	RMSE=1.53 (5.33) PCC=0.93 (0.90) σ =1.04 (1.44)
NA Lakes	1,5,6	RMSE=1.89 (6.58) PCC=0.97 (0.91) σ =1.03 (1.45)
Japan EANET	12,15	RMSE=3.11 (5.70) PCC=0.96 (0.79) σ =0.66 (0.51)

508

509

510
511

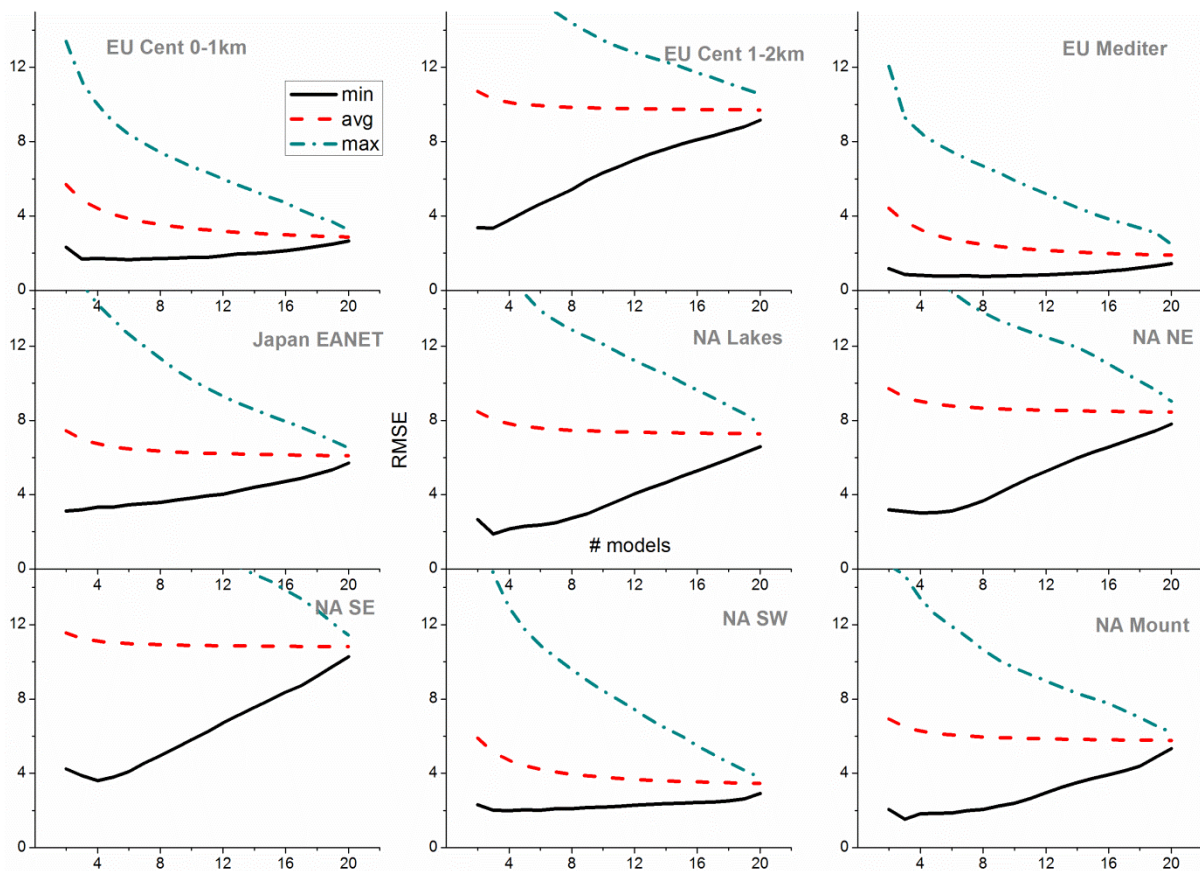
512 **Fig 1:** From Fiore et al. (2009): Monthly mean surface O₃ concentrations (ppb).
 513 Observed values (black circles) represent the average of all sites falling within the given latitude, longitude, and
 514 altitude boundaries and denoted by the symbols in Figure 1; vertical black lines depict the standard deviation
 515 across the sites. Monthly mean O₃ in the surface layer of the SR1 simulations from the 21 models are first
 516 sampled at the model grid cells containing the observational sites and then averaged within subregions (gray
 517 lines); these spatial averages from each model are used to determine the multimodel ensemble median (black
 518 dotted line) and mean (black dashed line). Observations are from CASTNET (<http://www.epa.gov/castnet/>) in
 519 the United States, from EMEP (<http://www.nilu.no/projects/ccc/emepdata.html>) in Europe, and from EANET
 520 (<http://www.eanet.cc/eanet.html>) in Japan.
 521



522

523 **Fig.2** Ranked histogram for the nine sub-regions subject to MM ensemble evaluation

524



525

526

527 **Fig 3** Maximum (dash-dot), average (dashed), and minimum (continuous line) RMSE for all subsets of MM
528 combinations and for the nine sub-regions subject to MM ensemble evaluation.

529

530