## _Reply to the 1<sup>st</sup> reviewer_

_Reply to the 1<sup>st</sup> reviewer_

We thank the referee for the positive and helpful comments that have improved the manuscript. They have all been taken on board and addressed in the revised version of our manuscript.

### _General Comments:_

_The paper provides an exhaustive analysis of the multi-model ensembling in air quality problems, using the data from AQMEII exercises. The authors give clear theoretical introduction based on various error decomposition, and in the following sections compare the predictive skill of three ensemble products with well-defined mathematical properties, namely: - the arithmetic mean of the entire ensemble, - the arithmetic mean of an ensemble subset, linked to the error decompositions, - the weighted mean of the entire ensemble, linked to the analytical optimization. For the selection of ensemble subsets the authors consider several clustering methods. In the analysis different indices and skills are used – the choice seems to be sufficient for the purpose, however – to some extent – this is a question of taste (but "de gustibus non est disputandum"). The analysis is based on AQMEII dataset, which is an appropriate and representative for this purpose. The authors raised the important issues of ensemble training and predictability – this part seems to be particularly valuable. Of course drawing any final conclusions from the analysis relying on large but one dataset is uncertain, nevertheless some reasonable hints have been formulated. The paper is a step forward towards better understanding of how to build good ensembles. Specific and technical comments are included in supplement file._

We do appreciate the positive comments.

### _Specific Comments:_

_• Page 8 lines 17-18 (remark in brackets): in principle the models can have different distributions. No such an assumption is needed to obtain formulas in Table 1. Actually the only assumption made is that the models are treated as random variables with known variances (distributions doesn't have to be known and can vary from model to model)._

We have removed the remark as it generated confusion. It was mistakenly pointing to the statistical distribution while the intention was to emphasize the randomness of the distribution in the i.i.d. sense (independent identically distributed).

_• Page 11 line 2: while selecting the subset theoretically optimal sequence is obtained if the models are ordered starting from the one with the smallest variance and the ensemble is built by adding step by step the next with smallest variance. This is however theoretical as it works for independent models – nevertheless one can consider this also as a possible approach. This procedure could be extended to the case of correlated models by making use of eigenvalue analysis._

Indeed, this can be seen in the Example. We have included another paragraph to summarize the effects of the various perturbations.

_"mme: its RMSE is reduced, compared to the i.i.d. case, if within the sample exist few members with lower variance or negative correlation."_

Indeed, the stabilization of the statistics required a 60-day hourly time-series. The results demonstrated the superiority of the analytically optimized full ensemble at all available monitoring stations, in predictive mode. Using shorter periods, the statistics behind the weights were not robust. On the other hand, the robustness of the ensemble subset in dynamic mode is due to (a) the persistence of ozone levels and (b) the successful modelling of its extremes by only few members.

Definitely, we were comparing a $1^{st}$ order bias correction that only removes the systematic errors with a $2^{nd}$ order correction that also adjusts the spread. We rephrased the terms in the manuscript to the general term 'bias correction'.

 **Technical Remarks:**

Done as suggested.

Done as suggested.

Done as suggested.