1 We appreciate the reviewer's valuable comments and constructive suggestions which

2 help us improve the quality of the manuscript. We have carefully revised the manuscript

3 according to these comments. Point-to-point responses are provided in the attachment.

4 The reviewers' comments are in black, our responses are in blue, and the quotes from

5 our manuscript are in italics.

6

7 **Reviewer #1**

8

9 [Comment]: How does the authors ensure the robustness of the model?

10 [Response]: We thank the reviewer for the valuable comment. We ensure the robustness

11 of the model from three aspects:

12 a) Model structure. Inspired by computer vision tasks, we adopt the batch-

13 normalization (Ioffe and Szegedy, 2015), dropout (Srivastava et al., 2014), L2

14 regularization (Zhang et al., 2016) to improve the generalization and robustness.

15 b) Early stop. When we train the NN-CTM, we split the data into train dataset and

16 validation dataset. As introduced in Sec. 3.1, we trained NN-CTM on the data of the

17 first 22 days in January, April, July, and October 2015 and tested it on the remaining

18 successive 8 days of each month. We stop the model training when the evaluation in

19 validation dataset does not improve within 1000 iterations.

20 c) Data augmentation. During training, we employ the noise injection, random

21 rescaling, random rotation method to avoid the overfitting in training dataset.

22 We have clarified the model robustness in the revised manuscript, as follows:

23 (Section 2.2, Paragraph 5) *"Model robustness. We ensure the robustness of the model*

24 *from three aspects: 1) Model structure. Inspired by computer vision tasks, we adopt the*

25 *batch-normalization (Ioffe and Szegedy, 2015), dropout (Srivastava et al., 2014), L2*

26 *regularization (Zhang et al., 2016) to improve the generalization and robustness. 2)*

27 *Early stop. When we train the NN-CTM, we split the data into train dataset and*

28 *validation dataset, and we stop the model training when the evaluation in validation*

29 *dataset does not improve within 1000 iterations. 3) Data augmentation. During*

30 *training, we employ the noise injection, random rescaling, random rotation method to*

31 *avoid the overfitting in training dataset."*

32

33 Reference:

34 Ioffe S, Szegedy C. Batch Normalization: Accelerating Deep Network Training by

35  Reducing Internal Covariate Shift. JMLR.org 2015.

36  Srivastava N, Hinton G, Krizhevsky A, Sutskever I, Salakhutdinov R. Dropout: A

37  Simple Way to Prevent Neural Networks from Overfitting. Journal of Machine

38  Learning Research 2014; 15: 1929-1958.

39  Zhang C, Be Ngio S, Hardt M, Recht B, Vinyals O. Understanding deep learning

40  requires rethinking generalization, 2016.

41

42  [Comment]: The authors use the observation data to update the emissions, however they

43  do not mention what happens in case more than one observation station is in the grid.

44  27 km × 27 km is a large grid size and hence would include many observation stations

45  in one grid. The averaged observed concentration of all stations if used won't serve the

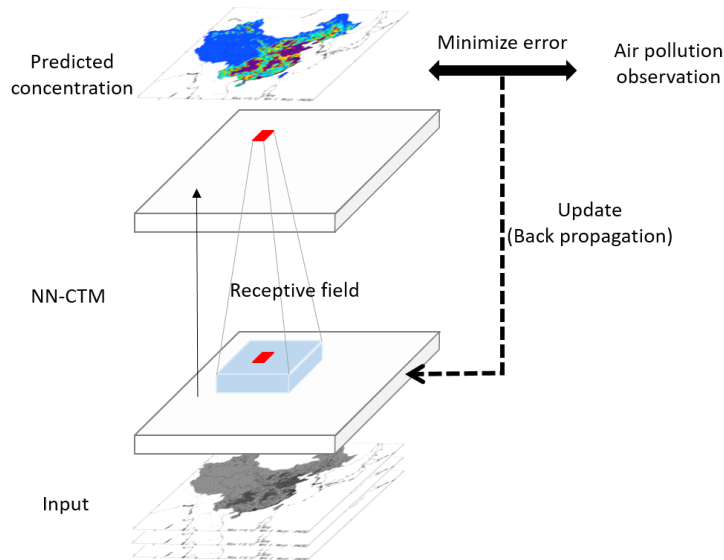46  purpose to accurately update the emissions at a station.

47  [Response]: We thank the reviewer for the valuable comment. As mentioned in Section

48  2.3, we use the average value in case of multiple observation stations in a 27 km $\times$ 27

49  km grid. We use the same processing method for observations when calculating MAE.

50  We focus on the emission estimation in one grid, which will be limited by the grid size.

51  If we want to get the higher resolution emission inventory estimation result (such as

52  focus on one typical region instead of whole China domain), we should use a finer-

53  grained emission inventory as the input. What's more, the lack of observation data in

54  some regions limits our updating, so we are more concerned about making good use of

55  existing observation data.

56

57  [Comment]: The entire premise of the model depends on availability of observation

58  data, what happens if data is very sparsely available e.g. say out of 4 neighboring grids

59  only one has observation data how are the emissions in other 3 grids updated?

60  [Response]: We thank the reviewer for the valuable comment. When we train the NN-

61  CTM, the long short term memory (LSTM) block is employed to capture the temporal

62  information, and the convolution (U-Net) is employed to capture the spatial information

63  (e.g. the emission inventory, meteorological information, and geographic information

64  of its neighbor grid). That is to say, in NN-CTM, the convolution neural network will

65  capture the surrounding grids' information within the receptive field, and we make a

66  detailed introduction about the receptive field in the answer of next comment which

67  represents the transmission between different grids. Therefore, as shown in Fig. 1, if

68  only the red gird has observation data, the surrounding blue grids' emission inventory

69  within the receptive field will also be updated. At the same time, the grids with a longer
70  distance will have a lower update weight. In extreme circumstances, if we have no
71  observation data, our method will not work as we have no more information to adjust
72  the emission inventory. If the observation data is denser, the emission inventory
73  estimation is more accurate as it can consider more observation data.
74



76  Figure 1: The visualization of neighbor emission update.
77

78  We have clarified the relation between observation data and emission inventory in the
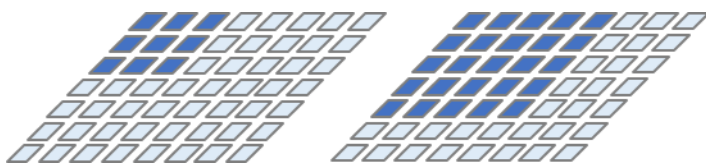79  revised manuscript, as follows:
80  (Section 2.3, Paragraph 1) *"The observation data will help update the surrounding*
81  *grids' emission inventory within the receptive field. However, in extreme circumstances,*
82  *if we have no observation data, our method will not work as we have no more*
83  *information to adjust the emission inventory. If the observation data is denser, the*
84  *emission inventory estimation is more accurate as it can consider more observation*
85  *data."*
86

87  [Comment]: Does the deep learning process consider the impact of transmission
88  between different grids? The authors are suggested to explain this point in detail.
89  [Response]: We thank the reviewer for the valuable suggestion to improve the quality
90  of the paper. The deep learning process has considered the impact of transmission
91  between different grids. The NN-CTM, which refers to U-Net branch in particular,
92  employs the convolution neural network to utilize neighbor information effectively. We
93  visualize a demo case of $3 \times 3$ convolution and $5 \times 5$ convolution in Fig. 2. In U-Net, the

stacked of convolution can get the neighbor information with a bigger receptive field (e.g. stacking 5×5 convolution and 5×5 convolution can get a 9×9 convolution), the non-linear function (P-RELU) is employed to improve model fitting with nearly zero extra computational cost and little overfitting risk, and the batch normalization and dropout are employed to enhance the robustness of the model. We calculate that the receptive field of our model is 38×38 grid. In other words, the predicted pollutant concentration is related to its surrounding 38×38 grid's information, which represents the transmission between different grids. Meanwhile, the closer the distance, the greater the contribution.



Figure 2: The visualization of convolution neural network (left: 3*3 kernel size, right: 5*5 kernel size).

We have clarified the impact of transmission between different grids in the revised manuscript, as follows:

(Section 2.2, Paragraph 3) *"In U-Net, the stacked of convolution can get the neighbor information with a bigger receptive field (e.g. stacking 5×5 convolution and 5×5 convolution can get a 9×9 convolution), the non-linear function (P-RELU) is employed to improve model fitting with nearly zero extra computational cost and little overfitting risk, and the batch normalization and dropout are employed to enhance the robustness of the model. We calculate that the receptive field of our model is 38×38 grid. In other words, the predicted pollutant concentration is related to its surrounding 38×38 grid's information, which represents the transmission between different grids. Meanwhile, the closer the distance, the greater the contribution."*

[Comment]: Lines 24-26, Abstract. Please be specific on the simulation year and the emission inventory you applied.

[Response]: We apologize for missing this information, and we have added year 2015 in Abstract.

[Comment]: Line 310, Page 15. I suggest the authors add more description for Figure 8, such as explaining why the performance of using the new emission inventory

127 worsened at some sites.

128 [Response]: We appreciate the reviewer for the valuable suggestion. We have added

129 more explanations accordingly as follows:

130 (Section 3.4, Paragraph 2) *"The model performance of most stations has been improved,*

131 *and a small number of stations with worsen performance show the link between*

132 *compound pollutants. For example, stations with larger deviations between $PM_{2.5}$*

133 *simulation results and observations tend to have greatly improved $O_3$ performance, and*

134 *vice versa."*

135

136 [Comment]: The language of the manuscript needs to be further polished.

137 [Response]: We thank the reviewer for the comment, and we have further polished the

138 manuscript and checked grammar carefully. All modifications will be marked in the

139 revised manuscript.