

Review result of “Estimation of PM_{2.5} Concentration in China Using Linear Hybrid Machine Learning Model.” (AMT-2021-64) by Song et al.

Response to RC1:

referee’s comments are given in blue,

our responses are given in red.

RC1: The submitted article develops a method to estimate PM_{2.5} values over China using a linear combination of three machine learning model. The innovative of this approach is the method to have an ensemble PM_{2.5} data from multiple machine learning model outputs. The research method is solid, and the results are convincing.

Response: We would like to thank the editor and referee for carefully reading the manuscript and providing detailed and constructive comments, which have helped a lot in improving the manuscript. We quote each comment below, followed by our response.

RC1: The background of the research does not cover all of the most recent machine learning produced PM_{2.5} products over China and provide convincing reason of why this approach is superior to the rest products. The big advantage of using AHI is the high temporal data (sub-hourly), however, the results section does not reflect this advantage.

Response: Due to the early start of this study, the latest research progress

was not quoted when writing the research background. To make up for these deficiencies, we will add 18 references to the manuscript. These references are listed at the page 8-10 of this document.

The advantage that AHI can provide high temporal resolution data is also discussed, but for some reasons it was not included in the previous version of the manuscript. In the revised manuscript we have added this content.

The results are shown in the figure below.

Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 Local Time. The model R^2 ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R^2 exceeded 0.7 at other times, indicating that the model had a good performance. The optimal performance time is 13:00, R^2 is 0.88. According to the results, the hourly differences in model performance were significant.

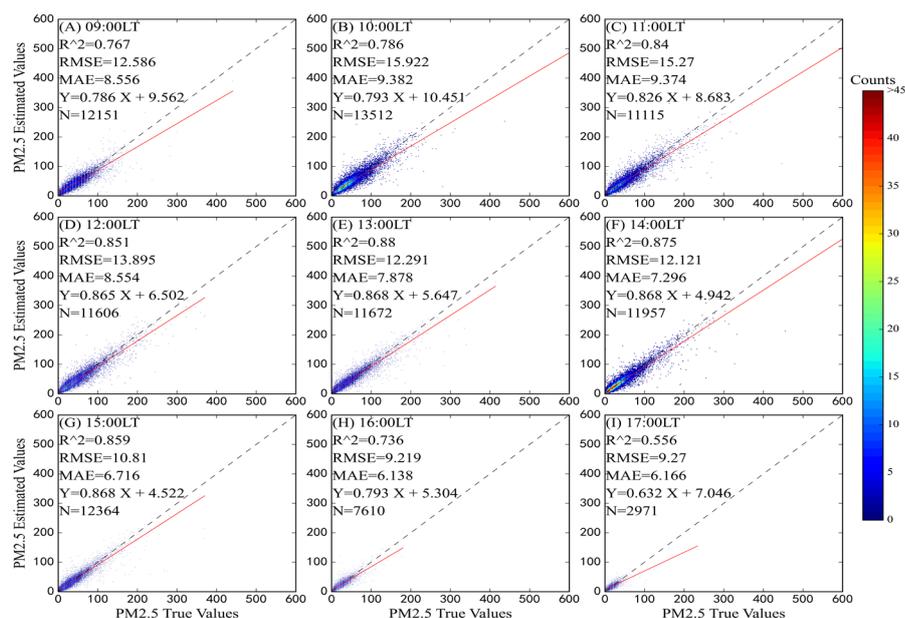


Figure 6 Hourly validation of model performance

The temporal distribution of PM_{2.5} is shown in Figure 10, The PM_{2.5} concentration began to rise from 9:00, and peaked at 55.65 $\mu\text{g}/\text{m}^3$ between 10:00 and 11:00 every day. After that, it maintained a high concentration until 15:00, and began to decrease. In the most polluted areas of China, the peak concentration of PM_{2.5} can reach 85.05 $\mu\text{g}/\text{m}^3$, while the peak in the less polluted areas is only about 40 $\mu\text{g}/\text{m}^3$. On a national scale, daily PM_{2.5} concentrations fluctuate little.

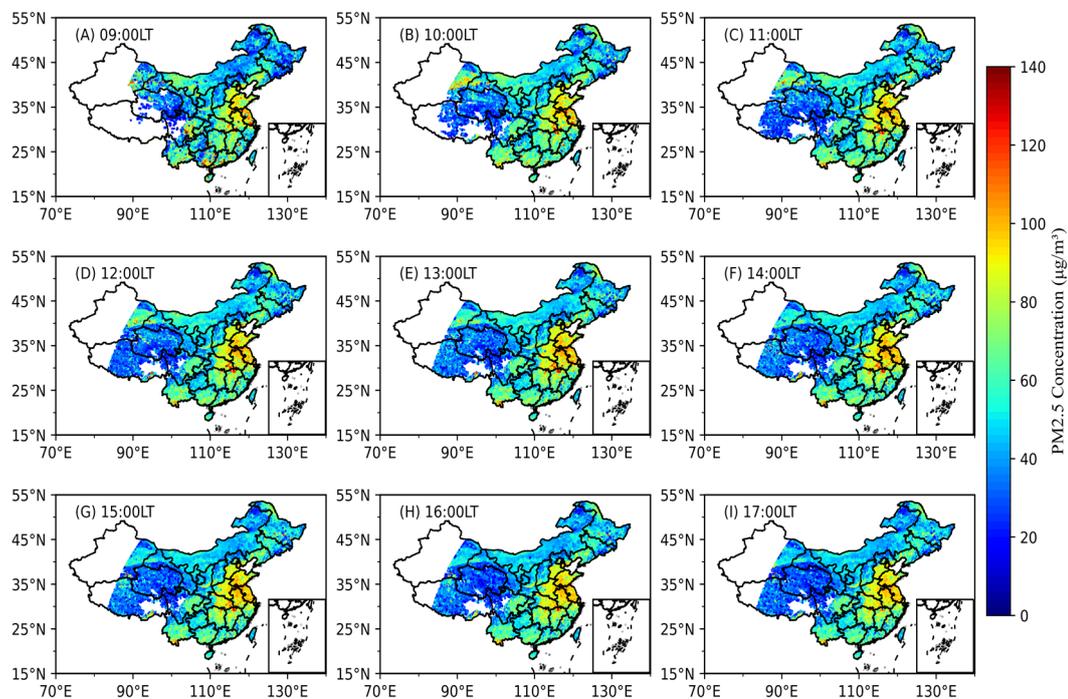


Figure 10 Hourly distribution of PM_{2.5} in China in 2019

RC1: The most contribution of this study is the linear hybrid ML model. However, the paper does not explain details of this procedure. For example, why using linear combination, and how are the coefficients are determined? Instead of a simple regression, complexed error evaluations of individual

ML PM2.5 data may provide insights on a better way of combining these model outputs.

Response: Wolpert et al. (1992) pointed out that the combination of multiple models can improve the robustness and generalization ability of the model. In other words, machine learning models can be integrated in the same way as multi-mode ensemble forecasting. Thus, we could further improve the accuracy of the fitting by hybrid model.

The coefficient is determined by multiple linear regression model. Firstly, we use three sub-models to calculate the predicted value under the corresponding model. Then, multiple linear regressions are performed between the calculated predicted values and the label values in the original data. Finally, the output coefficients and intercepts of the multiple linear regression model are taken as the parameters of the **RGD-LHMLM**.

RC1: The parameter impotency is listed but no further explanation of parameter selection is mentioned.

Response: We mainly used feature importance to analyze the contribution of different parameters to the model. This can provide an explanation of the interpretability of the model. The selection of parameters is mainly based on the variable information provided in some references. Finally, these characteristics we screened are all physical quantities that have a certain influence on PM_{2.5}, such as AOD, boundary layer height, relative humidity, population density.

RC1: Bias analysis as functions of other influence factors is needed to better understand the uncertainties in PM_{2.5} product.

Response: We use formula (5) and formula (6) to calculate the value of the Bias and the generalization error of the Bias (GEB). It is generally believed that when we take the generalization error, the Bias must be expressed in the form of a square. The average GEB between estimated PM_{2.5} based on the RGD-LHMLM and measured PM_{2.5} are shown in Table 1.

The results show that the average GEB of the mixed model is smaller, and the deviation between the predicted data and the label data is lower.

$$\text{Bias} = y_{\text{label},i} - y_{\text{predict},i} \quad (5)$$

$$\text{GEB} = \frac{\sum_{i=1}^N (y_{\text{label},i} - y_{\text{predict},i})^2}{N} \quad (6)$$

Table 1 Comparison of model accuracy

Model	Fitting				Validation			
	R ²	RMSE	MAE	GEB	R ²	RMSE	MAE	GEB
RF	0.95	6.99	4.05	114.19	0.79	14.89	9.33	208.97
GBRT	0.96	6.87	4.52	110.00	0.81	14.09	9.18	198.65
DNN	0.97	5.03	3.49	59.16	0.80	14.45	9.06	221.86
RGD-LHMLM	0.98	4.39	3.00	44.97	0.84	12.92	8.01	166.95

Then the bias of the mixed model in different PM_{2.5} concentration ranges was analyzed. As shown in the figure below: The average bias of the mixed model in different PM_{2.5} concentration ranges was analyzed, and the result is shown in the figure 4. when the PM_{2.5} concentration is less than 60 µg/m³,

the average bias of the model is less than 0. As the $PM_{2.5}$ concentration increases, the model deviation gradually increases. In other words, when the $PM_{2.5}$ concentration is small, the predicted value of the model will generally overestimate $PM_{2.5}$, and when the $PM_{2.5}$ further increases, it will underestimate the $PM_{2.5}$ concentration.

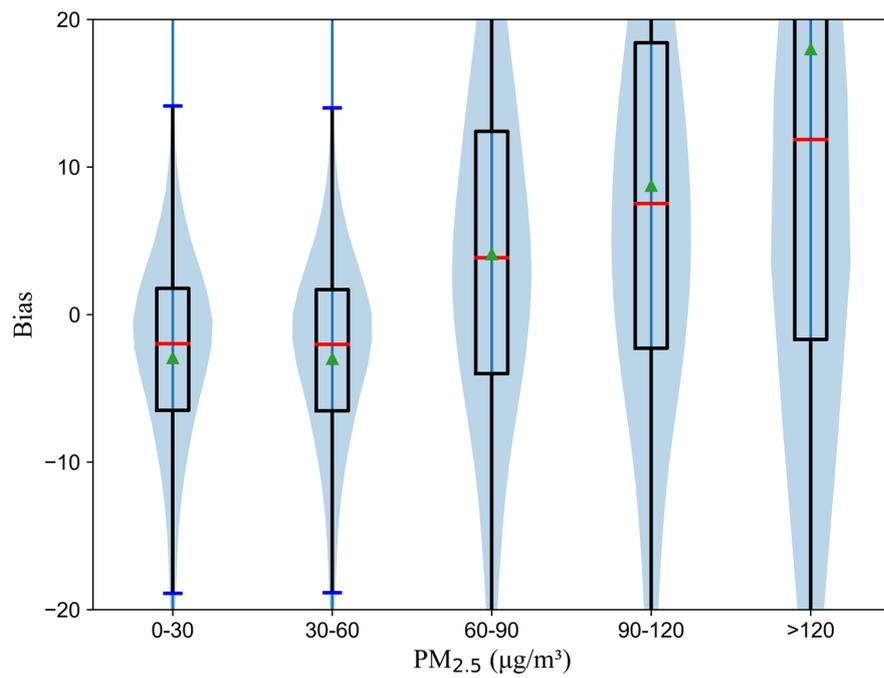


Figure 4 Bias between model predicted values and label values

We have compared other studies with our own and listed the results in Table 1:

Table 1

Model	R ²	RMSE	MAE	Reference
Stacking model	0.85	17.3	10.5	(Chen et al., 2019)
Two-stage random forests (YRD)	0.86	12.4	/	(Tang et al., 2019)
LME (BTH)	0.86	24.5	14.2	(Wang et al., 2017)
GTWR	0.78	20.10	/	(Xue et al., 2020b)
STLG	0.85	13.62	8.49	(Wei et al., 2021a)
RGD-LHMLM	0.84	12.92	8.01	This paper

References

(Yin et al., 2021;Xue et al., 2021;Wei et al., 2021b;Mao et al., 2021;Chen et al., 2021;Xue et al., 2020a;Wei et al., 2020;Wei et al., 2019a;Zeng et al., 2021;Guo et al., 2021;Qin et al., 2017;Li et al., 2016;Zheng et al., 2017;Gui et al., 2019;Shen et al., 2016;Li et al., 2021;Yang et al., 2016;Gui et al., 2020) (China, 2012;Yoshida et al., 2018)(Zhang et al., 2019)(Wei et al., 2019b)

Chen, B. J., You, S. X., Ye, Y., Fu, Y. Y., Ye, Z. R., Deng, J. S., Wang, K., and Hong, Y.: An interpretable

self-adaptive deep neural network for estimating daily spatially-continuous PM_{2.5} concentrations across China, *Sci Total Environ*, 768, <https://doi.org/10.1016/j.scitotenv.2020.144724>, 2021.

Chen, J. P., Yin, J. H., Zang, L., Zhang, T. X., and Zhao, M. D.: Stacking machine learning model for estimating hourly PM_{2.5} in China based on Himawari 8 aerosol optical depth data, *Sci Total Environ*, 697, <https://doi.org/10.1016/j.scitotenv.2019.134021>, 2019.

China: Ambient air quality standards. GB 3095-2012., China Environmental Science Press, Beijing, 2012.

Gui, K., Che, H. Z., Wang, Y. Q., Wang, H., Zhang, L., Zhao, H. J., Zheng, Y., Sun, T. Z., and Zhang, X. Y.: Satellite-derived PM_{2.5} concentration trends over Eastern China from 1998 to 2016: Relationships to emissions and meteorological parameters, *Environ Pollut*, 247, 1125-1133, <https://doi.org/10.1016/j.envpol.2019.01.056>, 2019.

Gui, K., Che, H. Z., Zeng, Z. L., Wang, Y. Q., Zhai, S. X., Wang, Z. M., Luo, M., Zhang, L., Liao, T. T., Zhao, H. J., Li, L., Zheng, Y., and Zhang, X. Y.: Construction of a virtual PM_{2.5} observation network in China based on high-density surface meteorological observations using the Extreme Gradient Boosting model, *Environ Int*, 141, <https://doi.org/10.1016/j.envint.2020.105801>, 2020.

Guo, B., Zhang, D. M., Pei, L., Su, Y., Wang, X. X., Bian, Y., Zhang, D. H., Yao, W. Q., Zhou, Z. X., and Guo, L. Y.: Estimating PM_{2.5} concentrations via random forest method using satellite, auxiliary, and ground-level station dataset at multiple temporal scales across China in 2017, *Sci Total Environ*, 778, <https://doi.org/https://doi.org/10.1016/j.scitotenv.2021.146288>, 2021.

Li, H. M., Yang, Y., Wang, H. L., Li, B. J., Wang, P. Y., Li, J. D., and Liao, H.: Constructing a spatiotemporally coherent long-term PM_{2.5} concentration dataset over China during 1980-2019 using a machine learning approach, *Sci Total Environ*, 765, <https://doi.org/10.1016/j.scitotenv.2020.144263>, 2021.

Li, Z. Q., Zhang, Y., Shao, J., Li, B. S., Hong, J., Liu, D., Li, D. H., Wei, P., Li, W., Li, L., Zhang, F. X., Guo, J., Deng, Q., Wang, B. X., Cui, C. L., Zhang, W. C., Wang, Z. Z., Lv, Y., Xu, H., Chen, X. F., Li, L., and Qie, L. L.: Remote sensing of atmospheric particulate mass of dry PM_{2.5} near the ground: Method validation using ground-based measurements, *Remote Sens Environ*, 173, 59-68, <https://doi.org/10.1016/j.rse.2015.11.019>, 2016.

Mao, F. Y., Hong, J., Min, Q. L., Gong, W., Zang, L., and Yin, J. H.: Estimating hourly full-coverage PM_{2.5} over China based on TOA reflectance data from the Fengyun-4A satellite, *Environ Pollut*, 270, <https://doi.org/10.1016/j.envpol.2020.116119>, 2021.

Qin, K., Wang, L. Y., Wu, L. X., Xu, J., Rao, L. L., Letu, H., Shi, T. W., and Wang, R. F.: A campaign for investigating aerosol optical properties during winter hazes over Shijiazhuang, China, *Atmos Res*, 198, 113-122, <https://doi.org/10.1016/j.atmosres.2017.08.018>, 2017.

Shen, Z. X., Cao, J. J., Zhang, L. M., Zhang, Q., Huang, R. J., Liu, S. X., Zhao, Z. Z., Zhu, C. S., Lei, Y. L., Xu, H. M., and Zheng, C. L.: Retrieving historical ambient PM_{2.5} concentrations using existing visibility measurements in Xi'an, Northwest China, *Atmos Environ*, 126, 15-20, <https://doi.org/10.1016/j.atmosenv.2015.11.040>, 2016.

Tang, D., Liu, D. R., Tang, Y. L., Seyler, B. C., Deng, X. F., and Zhan, Y.: Comparison of GOCI and Himawari-8 aerosol optical depth for deriving full-coverage hourly PM_{2.5} across the Yangtze River Delta, *Atmos Environ*, 217, <https://doi.org/10.1016/j.atmosenv.2019.116973>, 2019.

Wang, W., Mao, F. Y., Du, L., Pan, Z. X., Gong, W., and Fang, S. H.: Deriving Hourly PM_{2.5} Concentrations from Himawari-8 AODs over Beijing-Tianjin-Hebei in China, *Remote Sens-Basel*, 9, <https://doi.org/10.3390/rs9080858>, 2017.

Wei, J., Huang, W., Li, Z. Q., Xue, W. H., Peng, Y. R., Sun, L., and Cribb, M.: Estimating 1-km-resolution PM_{2.5} concentrations across China using the space-time random forest approach, *Remote Sens Environ*,

231,<https://doi.org/10.1016/j.rse.2019.111221>, 2019a.

Wei, J., Li, Z., Sun, L., Peng, Y., Zhang, Z., Li, Z., Su, T., Feng, L., Cai, Z., and Wu, H.: Evaluation and uncertainty estimate of next-generation geostationary meteorological Himawari-8/AHI aerosol products, *Sci Total Environ*, 692, 879-891,<https://doi.org/10.1016/j.scitotenv.2019.07.326>, 2019b.

Wei, J., Li, Z. Q., Cribb, M., Huang, W., Xue, W. H., Sun, L., Guo, J. P., Peng, Y. R., Li, J., Lyapustin, A., Liu, L., Wu, H., and Song, Y. M.: Improved 1 km resolution PM_{2.5} estimates across China using enhanced space-time extremely randomized trees, *Atmos Chem Phys*, 20, 3273-3289,<https://doi.org/10.5194/acp-20-3273-2020>, 2020.

Wei, J., Li, Z., Pinker, R. T., Wang, J., Sun, L., Xue, W., Li, R., and Cribb, M.: Himawari-8-derived diurnal variations in ground-level PM_{2.5} pollution across China using the fast space-time Light Gradient Boosting Machine (LightGBM), *Atmos. Chem. Phys.*, 21, 7863-7880,<https://doi.org/10.5194/acp-21-7863-2021>, 2021a.

Wei, J., Li, Z. Q., Lyapustin, A., Sun, L., Peng, Y. R., Xue, W. H., Su, T. N., and Cribb, M.: Reconstructing 1-km-resolution high-quality PM_{2.5} data records from 2000 to 2018 in China: spatiotemporal variations and policy implications, *Remote Sens Environ*, 252,<https://doi.org/10.1016/j.rse.2020.112136>, 2021b.

Xue, W. H., Zhang, J., Zhong, C., Ji, D. Y., and Huang, W.: Satellite-derived spatiotemporal PM_{2.5} concentrations and variations from 2006 to 2017 in China, *Sci Total Environ*, 712,<https://doi.org/10.1016/j.scitotenv.2019.134577>, 2020a.

Xue, W. H., Zhang, J., Zhong, C., Li, X. Y., and Wei, J.: Spatiotemporal PM_{2.5} variations and its response to the industrial structure from 2000 to 2018 in the Beijing-Tianjin-Hebei region, *J Clean Prod*, 279,<https://doi.org/10.1016/j.jclepro.2020.123742>, 2021.

Xue, Y., Li, Y., Guang, J., Tugui, A., She, L., Qin, K., Fan, C., Che, Y. H., Xie, Y. Q., Wen, Y. N., and Wang, Z. X.: Hourly PM_{2.5} Estimation over Central and Eastern China Based on Himawari-8 Data, *Remote Sens-Basel*, 12,<https://doi.org/10.3390/rs12050855>, 2020b.

Yang, Y., Liao, H., and Lou, S.: Increase in winter haze over eastern China in recent decades: Roles of variations in meteorological parameters and anthropogenic emissions, *J Geophys Res-Atmos*, 121, 13050-13065,<https://doi.org/10.1002/2016jd025136>, 2016.

Yin, J. H., Mao, F. Y., Zang, L., Chen, J. P., Lu, X., and Hong, J.: Retrieving PM_{2.5} with high spatio-temporal coverage by TOA reflectance of Himawari-8, *Atmospheric Pollution Research*, 12, 14-20,<https://doi.org/10.1016/j.apr.2021.02.007>, 2021.

Yoshida, M., Kikuchi, M., Nagao, T. M., Murakami, H., Nomaki, T., and Higurashi, A.: Common Retrieval of Aerosol Properties for Imaging Satellite Sensors, *Journal of the Meteorological Society of Japan. Ser. II*, 96B, 193-209,<https://doi.org/10.2151/jmsj.2018-039>, 2018.

Zeng, Z. L., Gui, K., Wang, Z. M., Luo, M., Geng, H., Ge, E. J., An, J. C., Song, X. Y., Ning, G. C., Zhai, S. X., and Liu, H. Z.: Estimating hourly surface PM_{2.5} concentrations across China from high-density meteorological observations by machine learning, *Atmos Res*, 254,<https://doi.org/10.1016/j.atmosres.2021.105516>, 2021.

Zheng, C. W., Zhao, C. F., Zhu, Y. N., Wang, Y., Shi, X. Q., Wu, X. L., Chen, T. M., Wu, F., and Qiu, Y. M.: Analysis of influential factors for the relationship between PM_{2.5} and AOD in Beijing, *Atmos Chem Phys*, 17, 13473-13489,<https://doi.org/10.5194/acp-17-13473-2017>, 2017.

Response to RC2:

referee's comments are given in blue,

our responses are given in red.

RC2: The study by Song et al. presents a linear hybrid machine learning model to estimate regional $PM_{2.5}$ distributions from Himawari-8 AOD observations. In the manuscript, the authors stated that the proposed RGD-LHMLM method outperforms than three conventional machine learning methods and can perform accurate estimations.

Response: We would like to thank the editor and referee for carefully reading the manuscript and providing detailed and constructive comments, which have helped a lot in improving the manuscript. We quote each comment below, followed by our response.

RC2: The paper does not provide enough evidence to support the major conclusions. The proposed method does not have generality in terms of target period as the training relies fully on the Himawari-8 AOD data over 2019. What about for the $PM_{2.5}$ estimation in some other years? To have a completely new training? Since the authors did not perform any $PM_{2.5}$ estimation for other years, I'd like to ask whether the training data already includes all possible cases between satellite AOD and ground $PM_{2.5}$. Even if by including more satellite AOD datasets over a longer period, it can still be questionable whether the selected training data are considered to

be representative.

Response: Our research is mainly based on two decision tree models and a neural network model to build a semi-explanatory estimation model. This semi-explanatory nature is mainly reflected in the analysis of the feature importance. In other words, deep learning models are often seen as black boxes with low interpretability. Therefore, we want to use the feature importance obtained by the decision tree model and the computational power of deep learning to build the semi-explanatory estimation model. Since DNN has the highest weight coefficient in the final hybrid model, we believe that this assumption has been realized to a certain extent.

Given factors such as climate change and human controls, the data from just one year cannot represent all possible scenarios between AOD and PM_{2.5}. However, the monthly and hourly variations contained in the data are very significant, and the number of samples retrieved from this data also meets the requirements of machine learning. So, we believe that one year's datasets can provide better training for the model; on the other hand, the Himawari-8 data was updated when we started this study. Based on the core thesis of this research and the above two reasons, we have selected the Himawari-8 AOD of 2019 for training.

In future research, we will extend the time period to study the change trend of PM_{2.5} on a long time scale.

RC2: Section 2: Please include information about data quality of all

datasets used for training (e.g., satellite AOD, ground-based data, meteorological data). The current training assumes that Himawari-8 AOD and ground PM_{2.5} data are true values, which in reality, is not true. Thus, please discuss how much impact of their data quality on the model performance in a quantitative way, i.e., what is the error propagation of these training data?

Response: Ground PM_{2.5} can be observed by two methods. The first is an automatic analysis method including trace element oscillation balance method or β -ray attenuation method. The other is manual gravimetric method (HJ618). The observed data are calibrated and quality-controlled according to national standards GB 3095-2012 (China's National Ambient air quality standards)(China, 2012).

Himawari-8 AOD is obtained by an aerosol retrieval algorithm based on Lambertian-surface-assumed developed by Yoshida et al. (2018).

Himawari-8 AOD was compared with the AOD data of AERONET (Aerosol Robotic Network)(Zhang et al., 2019), the results show that they are consistent ($R^2=0.75$), RMSE and MAE were 0.39 and 0.21, respectively(Wei et al., 2019). In the study, we selected AOD with strict cloud screening, that is, AOD data with low uncertainty.

Uncertainty estimation of ERA5 data has described in detail in the following website: <https://confluence.ecmwf.int/display/CKB/ERA5%3A+uncertainty+estimation>.

To sum up, the data we used have been quality-controlled and can

represent the real situation to some extent. As commented by Referee #2, we have added bias analysis.

There is an irretrievable error between the AOD or $PM_{2.5}$ and its true value. As shown in figure 4, the average bias of the mixed model in different $PM_{2.5}$ concentration ranges was analyzed, and the result is shown in the figure 4. when the $PM_{2.5}$ concentration is less than $60 \mu g/m^3$, the average bias of the model is less than 0. As the $PM_{2.5}$ concentration increases, the model deviation gradually increases. In other words, when the $PM_{2.5}$ concentration is small, the predicted value of the model will generally overestimate $PM_{2.5}$, and when the $PM_{2.5}$ further increases, it will underestimate the $PM_{2.5}$ concentration.

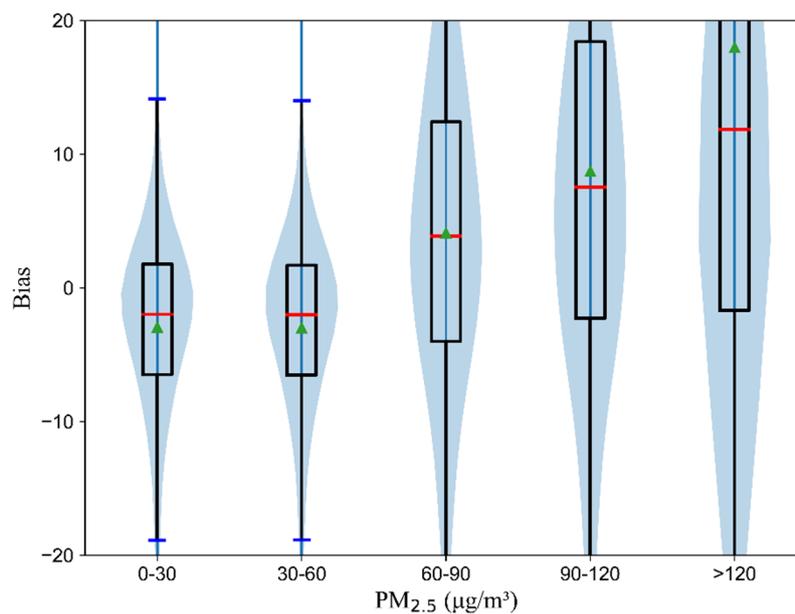


Figure 4 Bias between model predicted values and label values

In the machine learning algorithm, the error of the model will be corrected continuously according to the label value during the training.

As is known to all, the data calculated by the model are mainly related to the factors with high feature importance. In this model, the factor with the highest importance of feature is AOD. That is to say, when there is data error in AOD, it will be transmitted to the forecast result, and when there is data error in PM_{2.5}, it will interfere with the error correction of the model. Based on the above discussion, we believe that the errors in the model are mainly caused by the errors of AOD and PM_{2.5} when the pollution is relatively serious. In the case of low PM_{2.5} concentration, this error transfer phenomenon is relatively less.

RC2: These machine learning based models are sort of “black boxes”, which means that it would seem unclear what a physical relationship between input and output are learned, particularly to readers who are not familiar with PM_{2.5} estimation. I would suggest to reformulate the beginning of Section 3 by adding mathematical explanation for such context.

Response: It is a good suggestion. We will add the mathematical expression of the sub-model in the revised manuscript.

$$\begin{aligned}
 PM_{2.5i,j} = & AOD_{i,j} + BLH_{i,j} + RH_{i,j} + TM_{i,j} + LL_{i,j} + LH_{i,j} + SP_{i,j} \quad (1) \\
 & + RAIN_{i,j} + U_{10i,j} + V_{10i,j} + PD_{i,j} + HEIGHT_{i,j} + LON_{i,j} \\
 & + LAT_{i,j} + MONTH_{i,j} + HOUR_{i,j}
 \end{aligned}$$

Formula (1) is applicable to RF, GBRT and DNN. Where PM_{2.5i,j} is the PM_{2.5} at time i on station j.

RC2: Section 3: Please specify explicitly the input/output of the training(s).

Response: The input is 16 features including AOD (aerosol optical depth), surface relative humidity (RH, expressed as a percentage), air temperature at a height of 2 m (TM, expressed as K), Wind speed (U10, V10, in m/s), surface pressure (SP, in Pa), boundary layer height (BLH, in m) and cumulative precipitation (RAIN, in m) at 10 m above the ground, high and low vegetation index (LH, LL), ground elevation data (DEM), population density data (PD), longitude, latitude, month and hour.

The output is PM2.5 concentrations.

RC2: Section 3: Please describe in detail the linear combination of the three optimal sub-models.

Response: The coefficient is determined by multiple linear regression model. Firstly, we use three sub-models to calculate the predicted value under the corresponding model. Then, multiple linear regressions are performed between the calculated predicted values and the label values in the original data. Finally, the output coefficients and intercepts of the multiple linear regression model are taken as the parameters of the **RGD-LHMLM**.

RC2: Page 8, Line 13: According to Table 1, I do not notice any “significant” improvement from an individual sub-model to a linear-mixed model. I would prefer to say slightly improved, as can be seen also from

Figure 3.

Response: We have revised the description in the revised manuscript.

RC2: Section 4: The current manuscript only discusses the monthly performance of the linear-mixed model. But as far as I know, the usage of geostationary data such as Himawari-8, is especially beneficial to improving the understanding of daily variation of PM_{2.5}. If this study focuses solely on the monthly/seasonal variation, why not use MODIS AOD data over a longer period?

The advantage that AHI can provide high temporal resolution data is also discussed, but for some reasons it was not included in the previous version of the manuscript. In the revised manuscript we have added this content. The results are shown in the figure below.

Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 Local Time. The model R² ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R² exceeded 0.7 at other times, indicating that the model had a good performance. The optimal performance time is 13:00, R² is 0.88. According to the results, the hourly differences in model performance were significant.

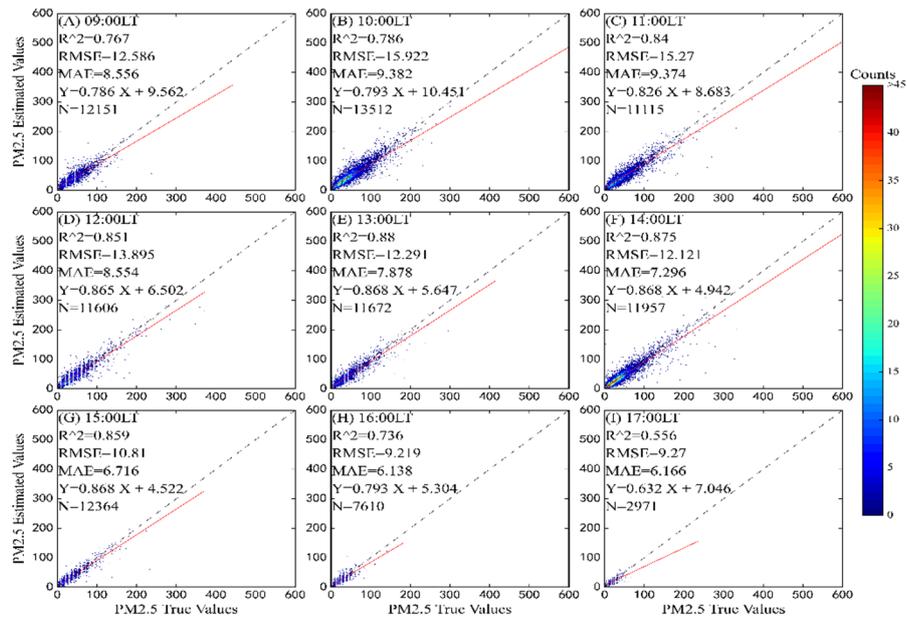


Figure 6 Hourly validation of model performance

The temporal distribution of PM_{2.5} is shown in Figure 10, The PM_{2.5} concentration began to rise from 9:00, and peaked at 55.65 $\mu\text{g}/\text{m}^3$ between 10:00 and 11:00 every day. After that, it maintained a high concentration until 15:00, and began to decrease. In the most polluted areas of China, the peak concentration of PM_{2.5} can reach 85.05 $\mu\text{g}/\text{m}^3$, while the peak in the less polluted areas is only about 40 $\mu\text{g}/\text{m}^3$. On a national scale, daily PM_{2.5} concentrations fluctuate little.

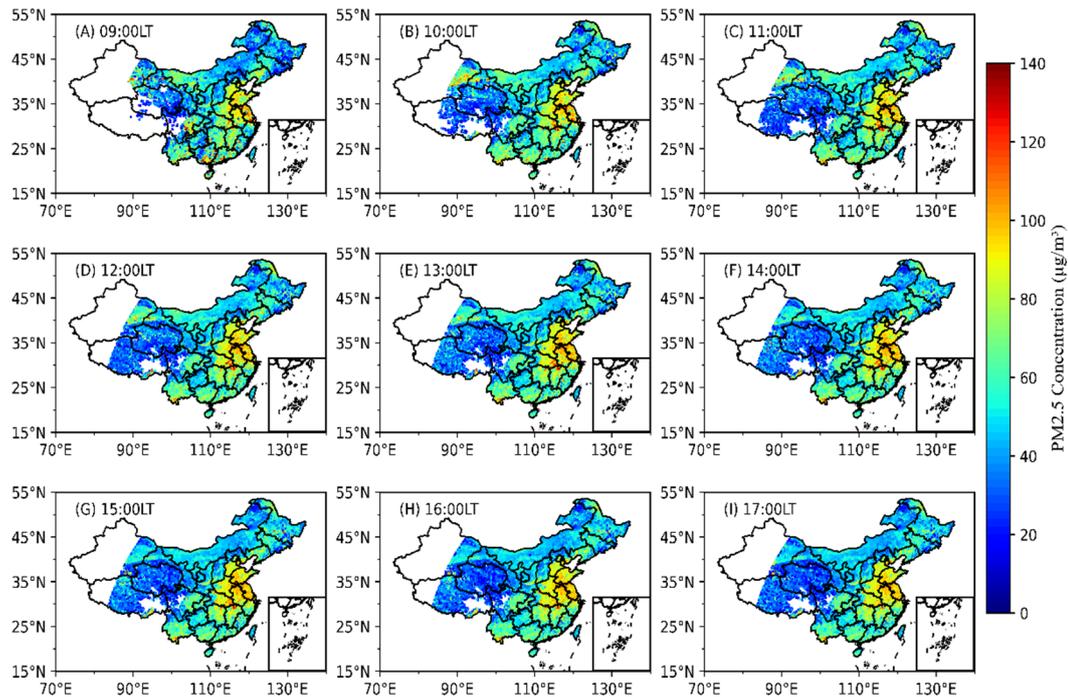


Figure 10 Hourly distribution of PM_{2.5} in China in 2019

RC3: Figure 5: It seems that the estimated PM_{2.5} are in general lower than the “true” values. Is this underestimation pattern related to Himawari-8 data? Please expand the relevant discussion.

Response: That's a very good question. As we all know, AOD is the integral of the aerosol extinction coefficient from the surface to the top of the atmosphere, and PM_{2.5} is small aerosol particles close to the surface which could float in the atmosphere for long period. Thus, PM_{2.5} contributes a significant portion of AOD, and the correlation between AOD and PM_{2.5} has a strong spatial and temporal variation (Ma et al., 2016; Xu et al., 2021). Combined with the feature importance of AOD and the above content, We believe that AOD has a very important influence on the model prediction values. In some studies,

however, Himawari-8 AOD has been found to be underestimated (Zang et al., 2018). Therefore, we believe that the underestimation of PM2.5 is closely related to the value of AOD. But, we need to note that the impact of meteorological parameters on the relationship between PM2.5 and AOD cannot be ignored (Gupta et al., 2006). So, the underestimated PM2.5 predicted value is greatly related to the influence of AOD, but the influence of meteorological factors should also be considered.

RC3: Figure 6: Please include importance of input parameters to DNN as well.

Response: As is answered in the first question, the feature importance of deep learning is difficult to obtain, and we only use the strong computational power of DNN to build the model. The DNN input is the same as the tree model, and the importance of the features in the tree model can explain which features are more important. In future research, we will study how to obtain the feature importance of DNN, and isolate them for analysis.

RC3: Section 4: An error characterization of model estimation is missing. Please discuss (quantitatively if possible) error contributions of the input parameters (at least including dominant error sources) to the final output.

Response: That's a tremendously good suggestion. We believe that the greater the importance of a feature in a model, the greater its contribution to the error of the model when there is an error. Perhaps this is not a

sufficient explanation. In future studies, we will try to discuss the error contribution of input parameters to the model.

RC2: Page 15, Line 19: Any examples of “other satellite data”? If other satellite observations are considered, how do you optimize the model training, as the current training is only based on Himawari-8 data.

Response: Some studies used “other satellite data”, such as FY-4A(Mao et al., 2021), MODIS(Wei et al., 2021b), GOIC(Tang et al., 2019) and VIIRS(Yao et al., 2019).

“If other satellite are considered”, I have two understandings. If it means not using Himawari-8 AOD data but using other satellite data for training, then the optimization process of the model is no different with Himawari-8. If this means using both Himawari-8 AOD data and other satellite data for training, then I think it's best to merge the two AOD datasets. In other words, the two kinds of AOD data are unified into one kind of integrated AOD data through linear regression or other algorithms. There are two benefits to doing this: firstly, The integrated AOD data can improve the data coverage to the surface; secondly, reducing the number of features can reduce the training time of the model and improve the efficiency.

We fully agree with the Referee #2’s opinion, and our follow-up work will be done through multi-satellite data fusion.

We have compared other studies with our own and listed the results in Table 1:

Table 1

Model	R ²	RMSE	MAE	Reference
Stacking model	0.85	17.3	10.5	(Chen et al., 2019)
Two-stage random forests (YRD)	0.86	12.4	/	(Tang et al., 2019)
LME (BTH)	0.86	24.5	14.2	(Wang et al., 2017)
GTWR	0.78	20.10	/	(Xue et al., 2020)
STLG	0.85	13.62	8.49	(Wei et al., 2021a)
RGD-LHMLM	0.84	12.92	8.01	This paper

Reference

Chen, J. P., Yin, J. H., Zang, L., Zhang, T. X., and Zhao, M. D.: Stacking machine learning model for estimating hourly PM_{2.5} in China based on Himawari 8 aerosol optical depth data, *Sci Total Environ*, 697, <https://doi.org/10.1016/j.scitotenv.2019.134021>, 2019.

China: Ambient air quality standards. GB 3095-2012., China Environmental Science Press, Beijing, 2012.

Gupta, P., Christopher, S. A., Wang, J., Gehrig, R., Lee, Y., and Kumar, N.: Satellite remote sensing of

particulate matter and air quality assessment over global cities, *Atmos Environ*, 40, 5880-5892, <https://doi.org/10.1016/j.atmosenv.2006.03.016>, 2006.

Ma, X. Y., Wang, J. Y., Yu, F. Q., Jia, H. L., and Hu, Y. N.: Can MODIS AOD be employed to derive PM_{2.5} in Beijing-Tianjin-Hebei over China?, *Atmos Res*, 181, 250-256, <https://doi.org/10.1016/j.atmosres.2016.06.018>, 2016.

Mao, F., Hong, J., Min, Q., Gong, W., Zang, L., and Yin, J.: Estimating hourly full-coverage PM_{2.5} over China based on TOA reflectance data from the Fengyun-4A satellite, *Environ Pollut*, 270, 116119, <https://doi.org/10.1016/j.envpol.2020.116119>, 2021.

Tang, D., Liu, D. R., Tang, Y. L., Seyler, B. C., Deng, X. F., and Zhan, Y.: Comparison of GOCI and Himawari-8 aerosol optical depth for deriving full-coverage hourly PM_{2.5} across the Yangtze River Delta, *Atmos Environ*, 217, <https://doi.org/10.1016/j.atmosenv.2019.116973>, 2019.

Wang, W., Mao, F. Y., Du, L., Pan, Z. X., Gong, W., and Fang, S. H.: Deriving Hourly PM_{2.5} Concentrations from Himawari-8 AODs over Beijing-Tianjin-Hebei in China, *Remote Sens-Basel*, 9, <https://doi.org/10.3390/rs9080858>, 2017.

Wei, J., Li, Z., Sun, L., Peng, Y., Zhang, Z., Li, Z., Su, T., Feng, L., Cai, Z., and Wu, H.: Evaluation and uncertainty estimate of next-generation geostationary meteorological Himawari-8/AHI aerosol products, *Sci Total Environ*, 692, 879-891, <https://doi.org/10.1016/j.scitotenv.2019.07.326>, 2019.

Wei, J., Li, Z., Pinker, R. T., Wang, J., Sun, L., Xue, W., Li, R., and Cribb, M.: Himawari-8-derived diurnal variations in ground-level PM_{2.5} pollution across China using the fast space-time Light Gradient Boosting Machine (LightGBM), *Atmos. Chem. Phys.*, 21, 7863-7880, <https://doi.org/10.5194/acp-21-7863-2021>, 2021a.

Wei, J., Li, Z. Q., Lyapustin, A., Sun, L., Peng, Y. R., Xue, W. H., Su, T. N., and Cribb, M.: Reconstructing 1-km-resolution high-quality PM_{2.5} data records from 2000 to 2018 in China: spatiotemporal variations and policy implications, *Remote Sens Environ*, 252, <https://doi.org/10.1016/j.rse.2020.112136>, 2021b.

Xu, Q. Q., Chen, X. L., Yang, S. B., Tang, L. L., and Dong, J. D.: Spatiotemporal relationship between Himawari-8 hourly columnar aerosol optical depth (AOD) and ground-level PM_{2.5} mass concentration in mainland China, *Sci Total Environ*, 765, <https://doi.org/10.1016/j.scitotenv.2020.144241>, 2021.

Xue, Y., Li, Y., Guang, J., Tugui, A., She, L., Qin, K., Fan, C., Che, Y. H., Xie, Y. Q., Wen, Y. N., and Wang, Z. X.: Hourly PM_{2.5} Estimation over Central and Eastern China Based on Himawari-8 Data, *Remote Sens-Basel*, 12, <https://doi.org/10.3390/rs12050855>, 2020.

Yao, F., Wu, J., Li, W., and Peng, J.: A spatially structured adaptive two-stage model for retrieving ground-level PM_{2.5} concentrations from VIIRS AOD in China, *ISPRS Journal of Photogrammetry and Remote Sensing*, 151, 263-276, <https://doi.org/10.1016/j.isprsjprs.2019.03.011>, 2019.

Yoshida, M., Kikuchi, M., Nagao, T. M., Murakami, H., Nomaki, T., and Higurashi, A.: Common Retrieval of Aerosol Properties for Imaging Satellite Sensors, *Journal of the Meteorological Society of Japan. Ser. II*, 96B, 193-209, <https://doi.org/10.2151/jmsj.2018-039>, 2018.

Zang, L., Mao, F., Guo, J., Gong, W., Wang, W., and Pan, Z.: Estimating hourly PM₁ concentrations from Himawari-8 aerosol optical depth in China, *Environ Pollut*, 241, 654-663, <https://doi.org/10.1016/j.envpol.2018.05.100>, 2018.

Zhang, Z., Wu, W., Fan, M., Tao, M., Wei, J., Jin, J., Tan, Y., and Wang, Q.: Validation of Himawari-8 aerosol optical depth retrievals over China, *Atmos Environ*, 199, 32-44, <https://doi.org/10.1016/j.atmosenv.2018.11.024>, 2019.

Response to RC3:

referee's comments are given in blue,

our responses are given in red.

RC3: The authors presented a new perspective to derive hourly PM_{2.5} concentrations from Himawari-8 satellite in China by combining different AI methods. This study is overall good, and the results are generally well presented.

Response: We would like to thank the editor and referee for carefully reading the manuscript and providing detailed and constructive comments, which have helped a lot in improving the manuscript. We quote each comment below, followed by our response.

RC3: My first concern is that the authors used all the data samples collected at the same locations having ground-based measurements using the cross-validation method, but the PM_{2.5} predictions are not evaluated at locations where ground-based measurements are unavailable. Thus, I suggest adding an additional validation to test the spatial prediction ability of your model based on the monitoring stations using the cross-validation method.

Response: We strongly agree with the comment. We have added the additional validation based on the monitoring stations. The results are shown in Fig. 3 (E), with a decrease in accuracy. In future studies, therefore,

we should add better spatial predictor features.

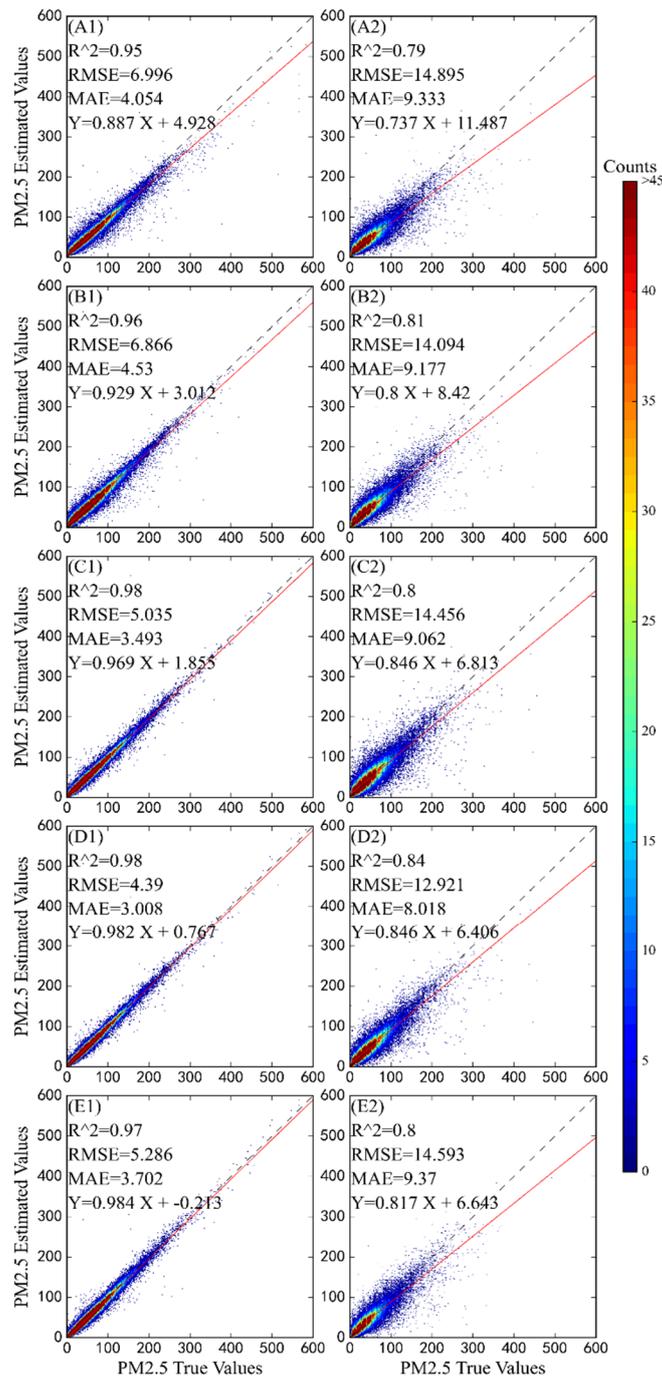


Figure 3 Accuracy of model Fitting and Validation (A: RF, B: GBRT, C: DNN, D: RGD-LHMLM (Based on sample), E: RGD-LHMLM (Based on site))

RC3: My other concern is that the purpose of this study is to derive hourly PM_{2.5} concentrations from geostationary satellites. However, the spatial analysis is performed on a monthly scale (Section 4.3), which will largely

reduce the sense of the current study. Thus, it is suggested to add more analysis on PM diurnal variations across China.

The advantage that AHI can provide high temporal resolution data is also discussed, but for some reasons it was not included in the previous version of the manuscript. In the revised manuscript we have added this content.

The results are shown in the figure below.

Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 Local Time. The model R^2 ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R^2 exceeded 0.7 at other times, indicating that the model had a good performance. The optimal performance time is 13:00, R^2 is 0.88. According to the results, the hourly differences in model performance were significant.

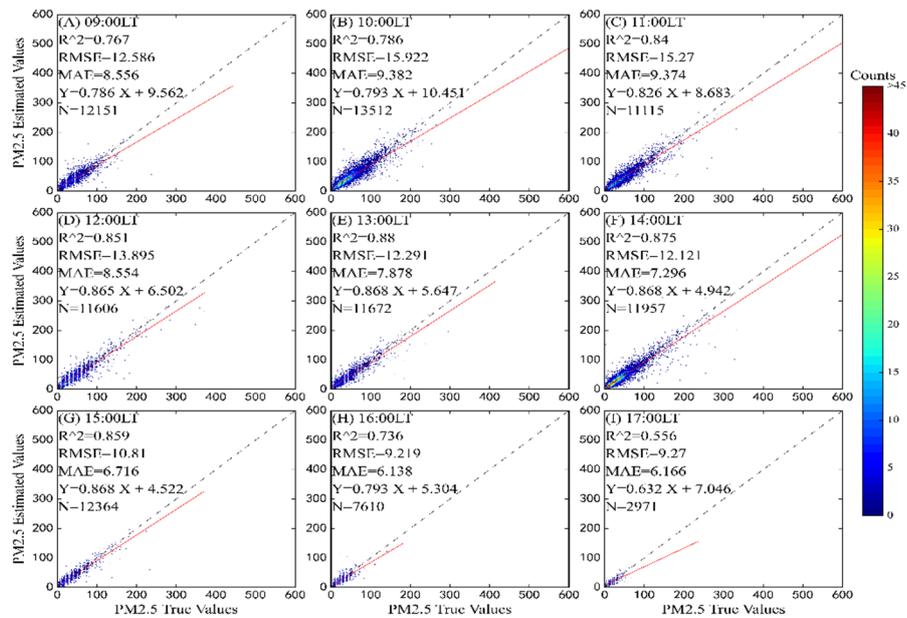


Figure 6 Hourly validation of model performance

The temporal distribution of $PM_{2.5}$ is shown in Figure 10, The $PM_{2.5}$

concentration began to rise from 9:00, and peaked at $55.65\mu\text{g}/\text{m}^3$ between 10:00 and 11:00 every day. After that, it maintained a high concentration until 15:00, and began to decrease. In the most polluted areas of China, the peak concentration of $\text{PM}_{2.5}$ can reach $85.05\mu\text{g}/\text{m}^3$, while the peak in the less polluted areas is only about $40\mu\text{g}/\text{m}^3$. On a national scale, daily $\text{PM}_{2.5}$ concentrations fluctuate little.

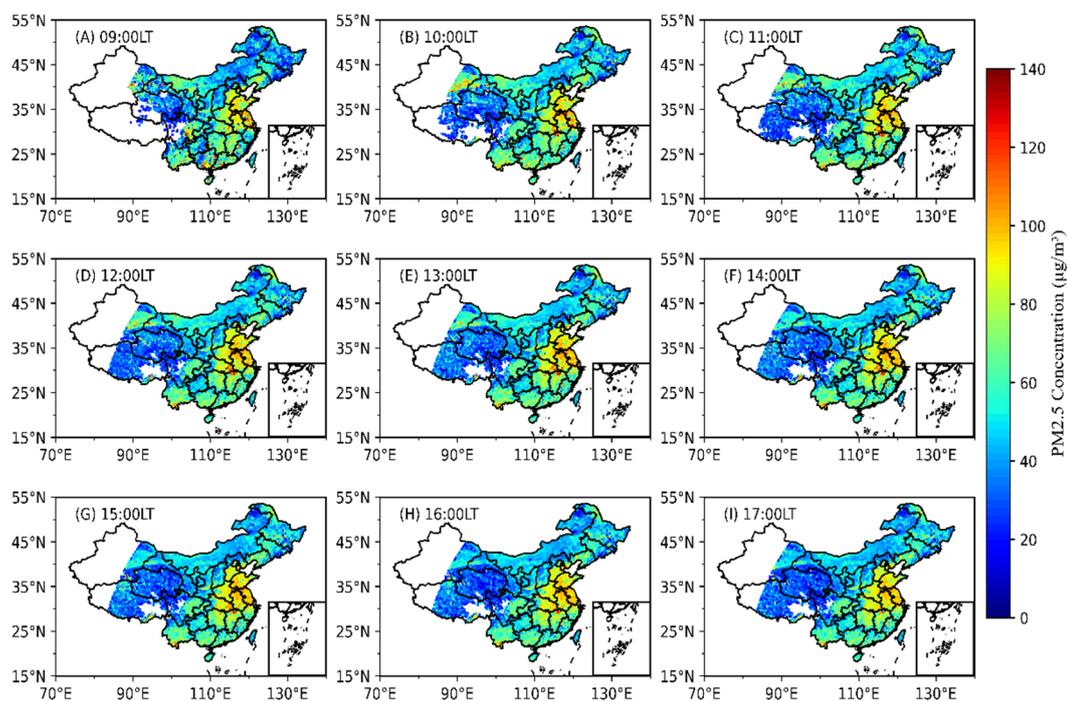


Figure 10 Hourly distribution of $\text{PM}_{2.5}$ in China in 2019

RC3: The authors are suggested to update the literature by summarizing more recent studies on $\text{PM}_{2.5}$ estimations using sun-synchronous and geostationary orbit satellites, especially those focusing on the whole of China. Below references may help you found more information on various recent studies to help enrich your study.

Section 2.2: Line 15, Reference for Himawari-8 aerosol algorithm is

needed.

Line 17: Below references provide a more comprehensive evaluation of Himawari-8 aerosol products in China.

Section 2.3: Reference for ERA5 reanalysis is needed.

References for these traditional ML or DL methods are needed.

Response: Many thanks for the references that were provided to our paper.

We have included it in the revised manuscript.

RC3: Lines 5-9: It is not clear to me how to determine the weight coefficients, and could you add more descriptions?

Response: The coefficient is determined by multiple linear regression model. Firstly, we use three sub-models to calculate the predicted value under the corresponding model. Then, multiple linear regressions are performed between the calculated predicted values and the label values in the original data. Finally, the output coefficients and intercepts of the multiple linear regression model are taken as the parameters of the weight coefficients.

RC3: Section 4.2.2: How about the accuracy of PM2.5 estimations for different hours?

Response: Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 Local Time. The model R^2 ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R^2 exceeded 0.7 at other times,

indicating that the model had a good performance. The optimal performance time is 13:00, R^2 is 0.88. According to the results, the hourly differences in model performance were significant.

RC3: Page 11, Lines 12-15, Page 12, and Page 13, Lines 1-4: May move to a new separate Discussion section.

Response: This is a very good comment, and we have adjusted it in the revised manuscript. The part pointed out by the Referee #3 has been taken as a separate subsection.

RC3: How about your model compared with those developed in previous studies using the Himawari-8 AOD products in China?

Response: We have compared other studies with our own and listed the results in Table 1:

Table 1

Model	R^2	RMSE	MAE	Reference
Stacking model	0.85	17.3	10.5	(Chen et al., 2019)
Two-stage random forests (YRD)	0.86	12.4	/	(Tang et al., 2019)
LME (BTH)	0.86	24.5	14.2	(Wang et al., 2017)
GTWR	0.78	20.10	/	(Xue et al., 2020)
STLG	0.85	13.62	8.49	(Wei et al., 2021)
RGD-LHMLM	0.84	12.92	8.01	This paper

According to the result of the table 1, the accuracy of our model is similar

to other models, both of which can better complete the estimation of PM_{2.5}.

References

- Chen, J. P., Yin, J. H., Zang, L., Zhang, T. X., and Zhao, M. D.: Stacking machine learning model for estimating hourly PM_{2.5} in China based on Himawari 8 aerosol optical depth data, *Sci Total Environ*, 697, <https://doi.org/10.1016/j.scitotenv.2019.134021>, 2019.
- Tang, D., Liu, D. R., Tang, Y. L., Seyler, B. C., Deng, X. F., and Zhan, Y.: Comparison of GOCI and Himawari-8 aerosol optical depth for deriving full-coverage hourly PM_{2.5} across the Yangtze River Delta, *Atmos Environ*, 217, <https://doi.org/10.1016/j.atmosenv.2019.116973>, 2019.
- Wang, W., Mao, F. Y., Du, L., Pan, Z. X., Gong, W., and Fang, S. H.: Deriving Hourly PM_{2.5} Concentrations from Himawari-8 AODs over Beijing-Tianjin-Hebei in China, *Remote Sens-Basel*, 9, <https://doi.org/10.3390/rs9080858>, 2017.
- Wei, J., Li, Z., Pinker, R. T., Wang, J., Sun, L., Xue, W., Li, R., and Cribb, M.: Himawari-8-derived diurnal variations in ground-level PM_{2.5} pollution across China using the fast space-time Light Gradient Boosting Machine (LightGBM), *Atmos. Chem. Phys.*, 21, 7863-7880, <https://doi.org/10.5194/acp-21-7863-2021>, 2021.
- Xue, Y., Li, Y., Guang, J., Tugui, A., She, L., Qin, K., Fan, C., Che, Y. H., Xie, Y. Q., Wen, Y. N., and Wang, Z. X.: Hourly PM_{2.5} Estimation over Central and Eastern China Based on Himawari-8 Data, *Remote Sens-Basel*, 12, <https://doi.org/10.3390/rs12050855>, 2020.

List of all relevant changes made in the manuscript

Page1, line25-26: *The period from 10:00 to 15:00 every day is the best time for model inversion, also at this time the pollution is high.*

Page3, line2-4: *Using interpretable self-adaptive deep neural network, Chen et al. (2021) estimated daily spatially-continuous PM_{2.5} concentrations across China, and analyzed the contribution of various characteristics to the PM_{2.5} model.*

Page3, line12-13: *Due to its excellent performance, some scholars use Himawari-8 data to estimate ground PM_{2.5}(Wei et al., 2021a).*

Page3, line18-19: *Yin et al. (2021) used Himawari-8 hourly TOA data to estimate ground PM_{2.5} in China, improved data coverage area.*

Page4, line1: *In addition, there are some novel algorithms such as STET (Wei et al., 2021b) and STRF (Wei et al., 2019a) that are also used for PM_{2.5} inversion research.*

Page4, line13-15: *The PM_{2.5} datasets are calibrated and quality-controlled according to national standards GB 3095-2012 (China's National Ambient air quality standards) (China, 2012).*

Page5, line5-6: *Himawari-8 AOD is obtained by an aerosol retrieval algorithm based on Lambertian surface-assumed developed by Yoshida et al. (2018).*

Page5, line8-11: *In previous studies (Zang et al., 2018), Himawari-8 AOD was compared with the AOD data of AERONET (Aerosol Robotic Network) in China and achieved good performance (Zhang et al., 2019c), so that the results show that they are consistent ($R^2=0.75$), RMSE and MAE were achieved 0.39 and 0.21, respectively (Wei et al., 2019b).*

Page5, line13-14: *In the study, we selected AOD with strict cloud screening, that is, AOD data with low uncertainty.*

Page5, line23-25: *Uncertainty estimation of ERA5 data has described in detail in the following website: <https://confluence.ecmwf.int/display/CKB/ERA5%3A+uncertainty+estimation>.*

Page7, line15-19: *After data processing, RF, GBRT, and DNN are used for modeling.*

$$\begin{aligned} PM_{2.5i,j} = & AOD_{i,j} + BLH_{i,j} + RH_{i,j} + TM_{i,j} + LL_{i,j} + LH_{i,j} + SP_{i,j} \\ & + RAIN_{i,j} + U_{10i,j} + V_{10i,j} + PD_{i,j} + HEIGHT_{i,j} + LON_{i,j} + LAT_{i,j} \\ & + MONTH_{i,j} + HOUR_{i,j} \end{aligned} \quad (1)$$

Formula (1) is applicable to RF, GBRT and DNN, where $PM_{2.5i,j}$ is the PM_{2.5} at time i on station j .

Page8, line2-14: *the bias (Bias, is the difference between the predicted values and the true values, so that models with larger bias performed worse), and the GME (generalization error of the bias, It is generally*

believed that bias should be expressed as a square when using generalization error). The calculation formula of each indicator is shown as follows:

$$R^2 = 1 - \frac{SS_{res}}{SS_{tot}} \quad (2)$$

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2} \quad (4)$$

$$Bias = \frac{\sum_{i=1}^N \hat{y}_i - y_i}{N} \quad (5)$$

$$GEB = \frac{\sum_{i=1}^N (\hat{y}_i - y_i)^2}{N} \quad (6)$$

Where \hat{y}_i represents the predicted value, y_i shows the true value, SS_{res} denotes the error between the regression data and the mean value, SS_{tot} represents the error between the real data and the mean value, and the mean value is the mean value of the true value.

Page9, line14-16: Among all the models, the deviation generalization error of the linear mixed model is also the lowest, indicating that the difference between the results obtained by this model and the real value is the least.

Page10, line1-3: The accuracy of the model decreased in the site-based validation, in which the R_2 and RMSE values are 0.8 and 14.59 $\mu\text{g}/\text{m}^3$, respectively.

Page11, line2-7:

4.2.1 Bias analysis of Model

The average bias of the mixed model in different $\text{PM}_{2.5}$ concentration ranges was analyzed, and the result is shown in figure 4. When the $\text{PM}_{2.5}$ concentration is less than 60 $\mu\text{g}/\text{m}^3$, the average bias of the model is less than 0. As the $\text{PM}_{2.5}$ concentration increases, the model deviation gradually increases. In other words, when the $\text{PM}_{2.5}$ concentration is small, the predicted value of the model will generally overestimate $\text{PM}_{2.5}$, and when the $\text{PM}_{2.5}$ further increases, it will underestimate the $\text{PM}_{2.5}$ concentration.

Page12, line16- Page13, line2:

4.2.3 Time-Scale Model Performance Analysis

Figure 6 shows the scatterplot fitted with the inversion results of the mixed model from 9:00-17:00 local Time. The model R_2 ranged from 0.556 to 0.88 at different times. Except for 17:00 when the model had the worst performance, the model R_2 exceeded 0.7 at other times, indicating that the model had a good performance. The optimal performance time is 13:00, 1 and R_2 is 0.88. According to the results, the hourly differences in model performance were significant.

Page 16, line13-17: The temporal distribution of $\text{PM}_{2.5}$ is shown in Figure 10, The $\text{PM}_{2.5}$ concentration began to rise from 9:00, and peaked at 55.65 $\mu\text{g}/\text{m}^3$ between 10:00 and 11:00 every day. After that, it maintained a high concentration until 15:00; and began to decrease. In the most polluted areas of

China, the peak concentration of PM_{2.5} can reach 85.05µg/m³, while the peak in the less polluted areas is only about 40µg/m³. On a national scale, daily PM_{2.5} concentrations fluctuates slightly.

Page18, line16-17: *The diurnal variation of the model inversion effect is also obvious, and the 11:00-14:00 model usually has better performance.*

Table 1 Comparison of model accuracy

Model	Fitting				Validation			
	R ²	RMSE	MAE	GEB	R ²	RMSE	MAE	GEB
RF	0.95	6.99	4.05	114.19	0.79	14.89	9.33	208.97
GBRT	0.96	6.87	4.52	110.00	0.81	14.09	9.18	198.65
DNN	0.97	5.03	3.49	59.16	0.80	14.45	9.06	221.86
RGD-LHMLM	0.98	4.39	3.00	44.97	0.84	12.92	8.01	166.95

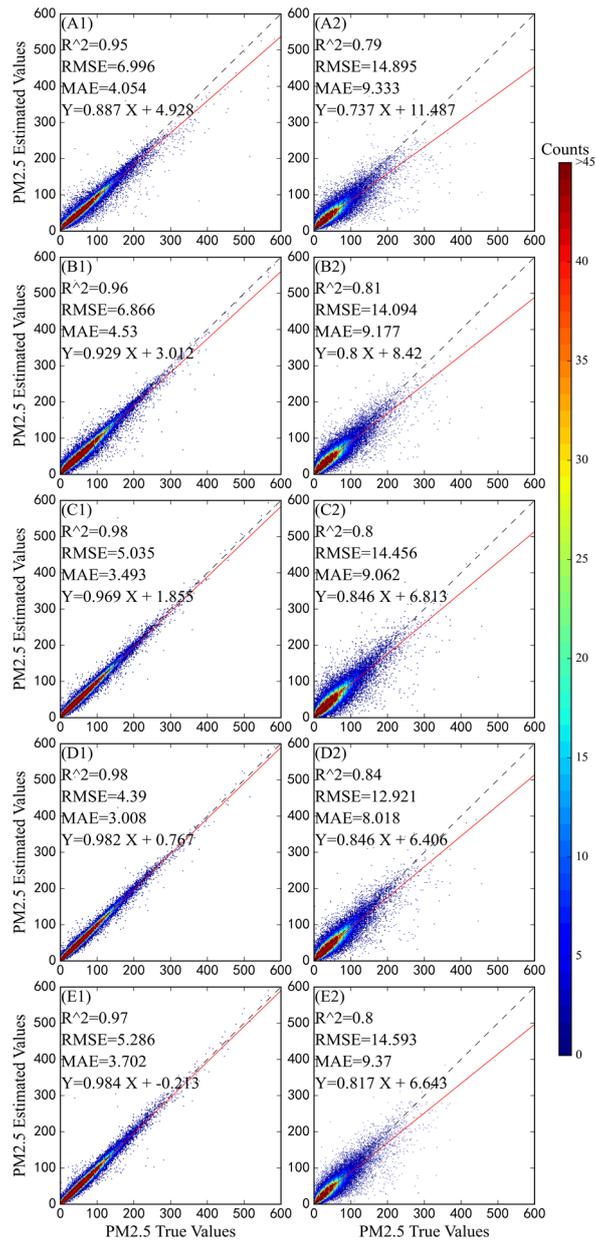


Figure 3 Accuracy of model Fitting and Validation (A: RF, B: GBRT, C: DNN, D: RGD-LHMLM (Based on sample), E: RGD-LHMLM (Based on site))

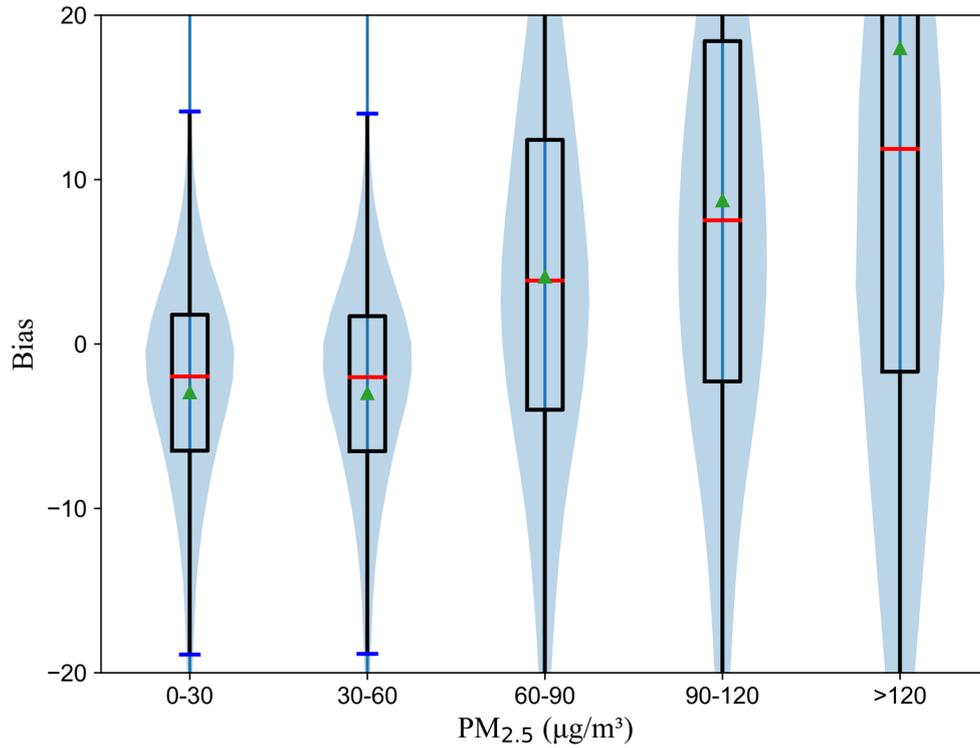


Figure 4 Model bias at different PM_{2.5} concentrations

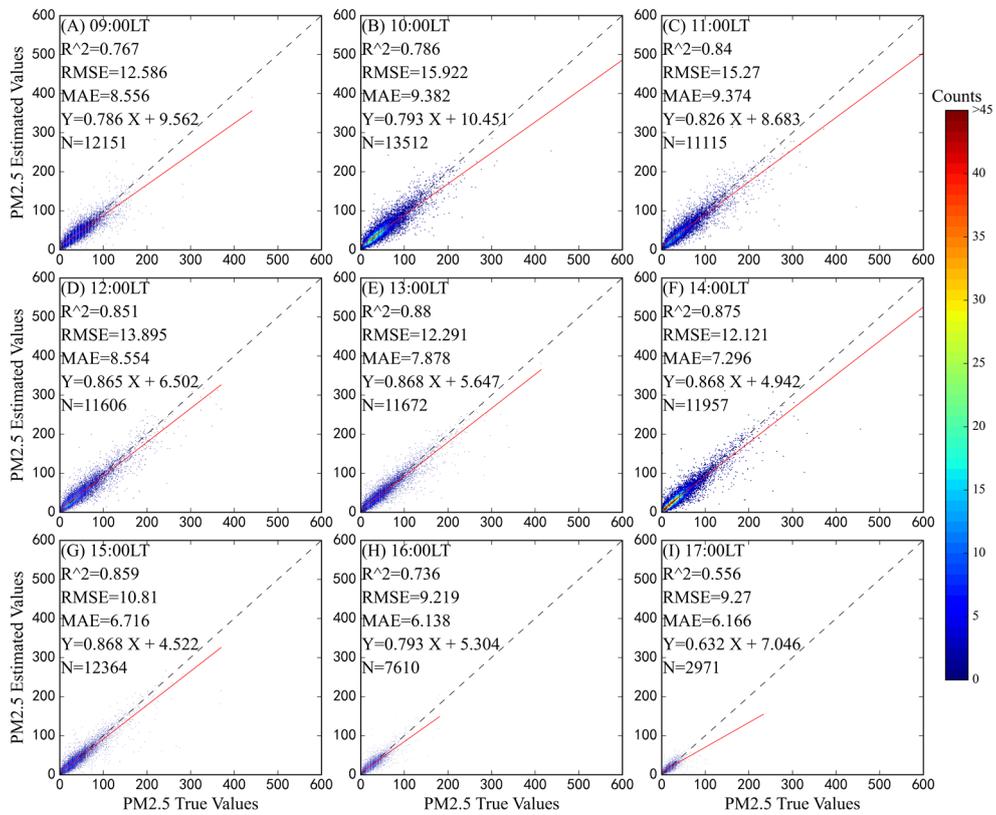


Figure 6 Hourly model performance fitting scatter diagram in 2019

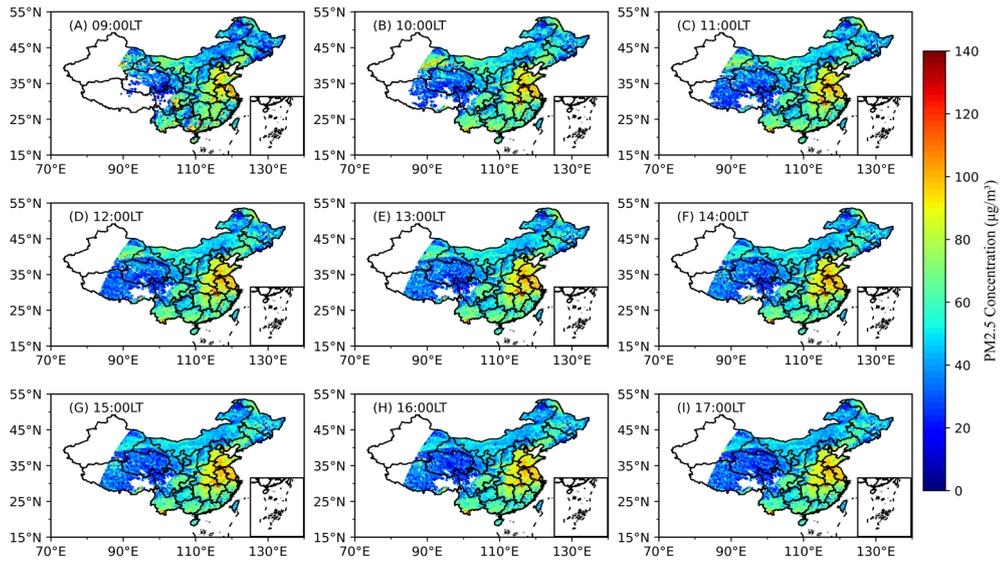


Figure 10 Monthly distribution of PM_{2.5} concentration in China in 2019