

Response to reviewers

We thank the two reviewers for their thorough and constructive reviews of this manuscript. We have incorporated most of the suggestions and believe the revised paper is substantially improved. Please find below the author's responses point-by-point. The original comments by the reviewer are in *blue italics*, our responses are in black. The revised manuscript with tracked changes is also appended.

Reviewer #1:

Major Comments

This study by Zhang et al. entitled “Unleashing the Potential of Geostationary Satellite Observations in Air Quality Forecasting Through Artificial Intelligence Techniques” presents a new machine-learning framework – GeoNet – that synthesizes geostationary observations of columnar NO₂ from the Geostationary Environment Monitoring Spectrometer (GEMS) with meteorological parameters to forecast surface-level NO₂ in East China. Overall, this study represents a significant advancement in surface-level pollution forecasting given its use of the unprecedented hourly data provided by GEMS. I believe that this manuscript is well-written and consistent; however, I have a few comments below.

We appreciate the reviewer for the positive comments. We have addressed the following concerns point-by-point below:

- 1. First, if possible, it would be useful to validate the GEMS observations using ground-based spectrometers (e.g., PGN) specifically for the study region and time period.*

We thank the reviewer for this suggestion! In our previous work (e.g., (Li et al., 2023)), we validated the GEMS NO₂ retrievals with the six ground-based MAX-DOAS distributed over the Jing-Jin-Ji region, Yangtze River Delta, Pearl River Delta, and Sichuan-Chongqing Basin in China. Generally, the correlation coefficient between GEMS and MAX-DOAS NO₂ retrieval ranges between 0.69-0.92 for different sites (see Fig. R1), indicating high consistency between both datasets. We have now discussed these results in L125-127, section 2.1:

Our previous validation results indicated that GEMS NO₂ retrievals generally agreed well with ground-based MAX-DOAS measurements from 6 sites in China, with correlation coefficients ranging between 0.69-0.92 (Li et al., 2023).

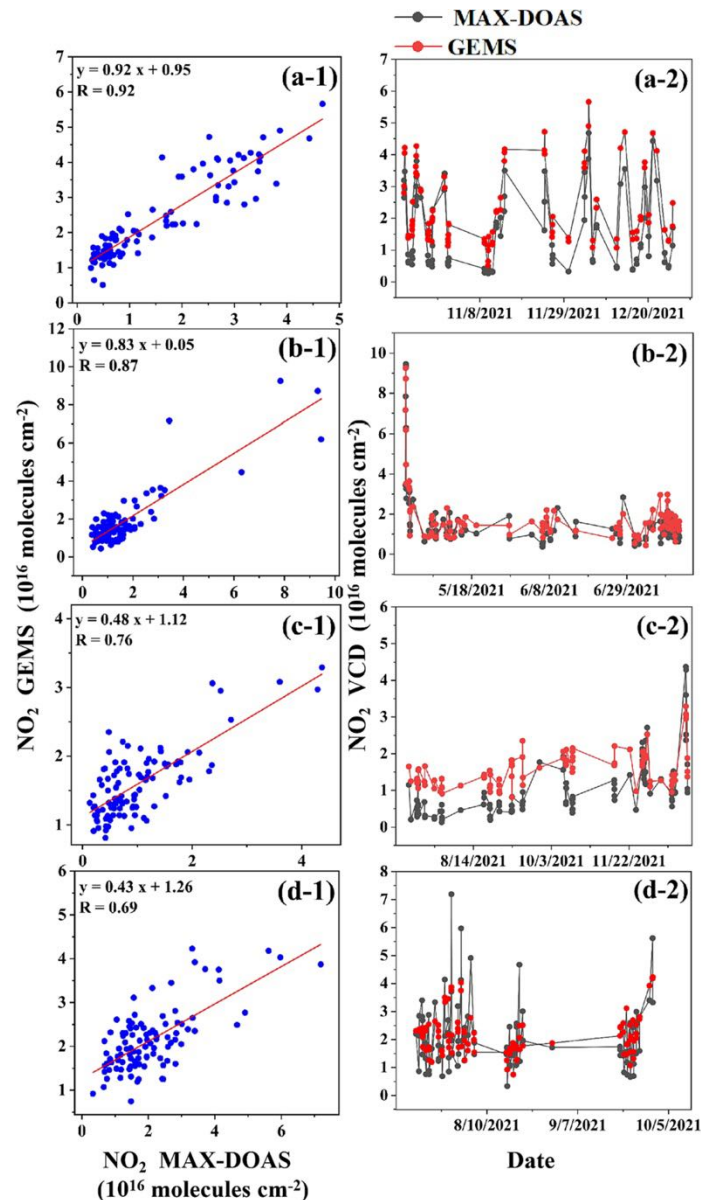


Fig. R1 (from Li et al., 2023). The comparison between GEMS and the MAX-DOAS NO_2 VCDs data from different stations. Panels (a-1), (b-1), (c-1), and (d-1) represent the comparisons with CAMS, GIG, HNU, and CQ stations, respectively. Panel (a-2), (b-2), (c-2), and (d-2) represent time series plots with CAMS, GIG, HNU, and CQ stations, respectively.

2. *Additionally, unless I missed it, I don't believe the time periods for model training and validation were ever stated; if this is the case they should be added to the main text.*

We thank the reviewer for this helpful reminder! We now add the following descriptions in L185-187, section 2.3:

We train the GeoNet model with input features during the whole year of 2021. The training datasets were randomly selected from 75% of the whole samples, while the remaining 25% were used as validation sets.

3. *Second, when investigating feature importance, it would be useful to also identify variability in the feature importance to uncover whether some components are more stable than others in GeoNet and to identify if the significance of geostationary observations is consistent across different days and seasons.*

We now investigated the temporal variations of these features' importance across different seasons. Fig. R2 is similar to Fig. 3a in the main text, but for the feature importance in Spring, Summer, Autumn, and Winter. We added the following discussions in L321-322:

The significance of the different input variables remained generally consistent across seasons, with minor variations (as shown in Fig. S12).

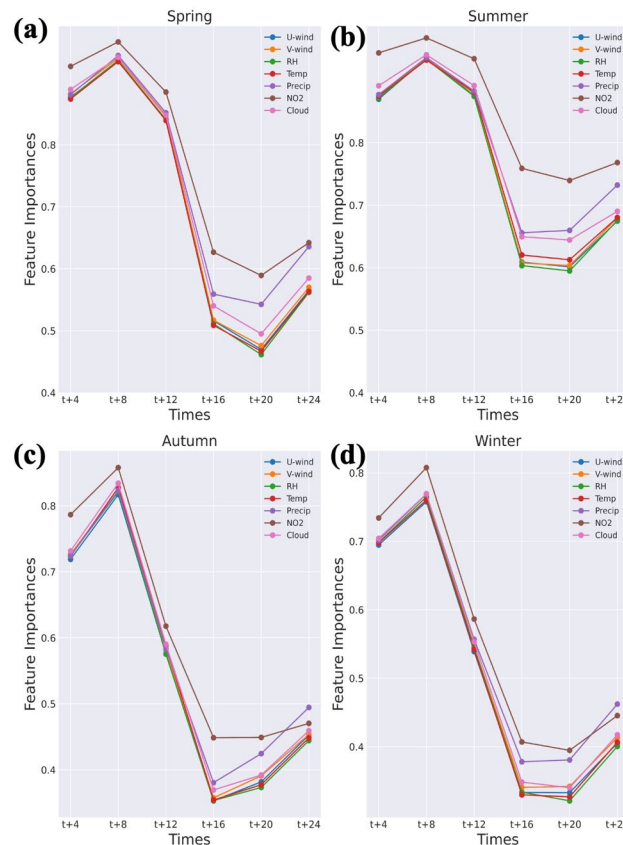


Fig. R2 (also moved into the Supplementary Information, Fig. S12). Similar to Fig 3a, but for different seasons, including Spring (a), Summer (b), Autumn (c), and Winter (d).

4. *Lastly, I suggest that the authors update their analysis in Figure 4 to include the GeoNet predictions regridded to the CAMS grid to identify how much of the improvement in predictions is attributable specifically to enhancements in spatial resolution.*

Thanks for this suggestion! We now re-grid the GeoNet predictions to the CAMS grid in Fig. R3 (and also replaced it with Fig. 5). It can be seen that GeoNet prediction results between the original (0.1°) and CAMS resolution (0.4°) are similar in spatial patterns overall.

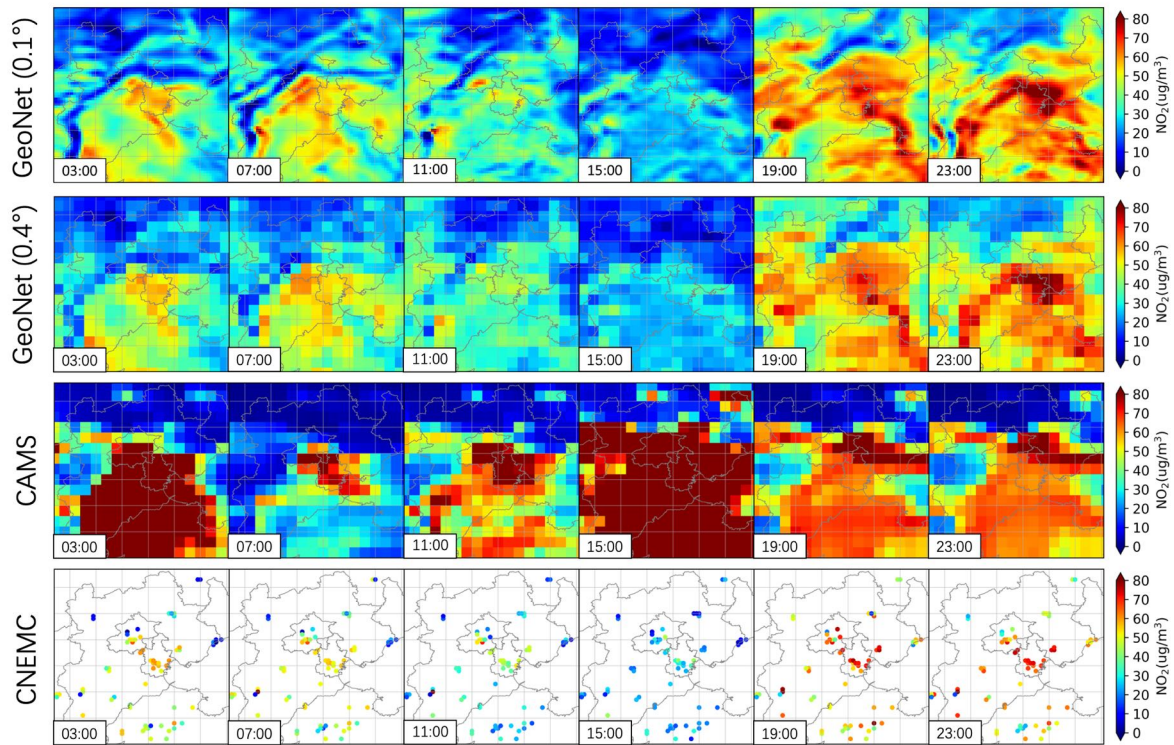


Fig. R3 (Fig. 5 in the revised manuscript). The spatial distribution comparisons of surface NO_2 concentration between (a) GeoNet prediction at the original resolution of 0.1° , (b) GeoNet prediction resampled to the CAMS resolution of 0.4° , (c) CAMS prediction, and (d) ground-based CNEMC site measurements. Note that the results are presented for different continuing local hours (labeled text in the subplot) on 23 November 2021.

I have included line-specific comments below:

Minor Comments

- L53-54: While I agree with this statement, it should be mentioned that for air pollution forecasting to facilitate health benefits, infrastructure needs to be created that communicate risks and appropriate responses to risks to the public.*

We added the following lines to mention the necessity of implementing infrastructure to communicate air quality risks and public recommendations, in L50-52 of the revised manuscript:

Meanwhile, air quality health risk (AQHI) forecasts and corresponding public response recommendations need to be communicated to the public promptly through public facilities (Fino et al., 2021; Tang et al., 2024).

- L55: I think you can drop the second limited in this line.*

Removed.

- L75: Maybe it would be useful to give an example or two here (i.e., TROPOMI + OMI).*

Added “such as the Ozone Monitoring Instrument (OMI) and the TROPospheric Monitoring Instrument (TROPOMI)”.

4. *L78-81: Another limitation of the polar orbiting satellites that is worth mentioning is that typically (at least in the case of TROPOMI) the satellite observes at roughly the same time of day (early afternoon) which makes it difficult to predict concentrations at other times of the day with different meteorological (boundary layer height) and photochemical conditions.*

Thanks for your suggestion. We added the following sentences in L84-86:

However, these observations at a fixed daily overpass time could hardly support the prediction of atmospheric trace gas concentrations at other times of the day under different meteorological conditions.

5. *L92: It would be better to describe GEMS as having “unprecedented temporal and spatial resolution and coverage” as ground-level monitors can observe hourly NO₂ but are limited in time and aircraft remote-sensing can observe NO₂ at sub hourly resolution but over a limited temporal coverage (usually a few days or weeks). The resolution alone isn’t necessarily unique but rather than combined spatial + temporal resolution with extended spatial and temporal coverage.*

Corrected.

6. *L117-120: Were you able to validate these data for the study time period / domain? If possible, it may be useful to compare GEMS to ground-based spectrometers in the study domain to get an idea of performance.*

Please refer to the response to major comment #1.

7. *L207-208: I don’t think you need this sentence as it is already mentioned in the methods section.*

Removed.

8. *Figure 3: It would be interesting to present the variance of these different components as well in a). Are these importance values pretty consistent regardless of season and day, or do they vary substantially day to day?*

Please refer to the response to major comment #3.

9. *Figure 4: Have you assessed how much of the reductions in performance are attributable to resolution? If not, I suggest regridding the GeoNet prediction to the resolution of CAMS and comparing this “GeoNET_coarse” product to the observations to characterize how much of the improved performance is attributable to enhanced spatial resolution.*

Please refer to the response to major comment #4.

10. *Figure 5: The colorbar in a is not labeled, and throughout the font is small (especially in the yaxis of c and d), I suggest updating to improve readability.*

Thanks for your helpful suggestion! We added the figure label and increased the font as follows.

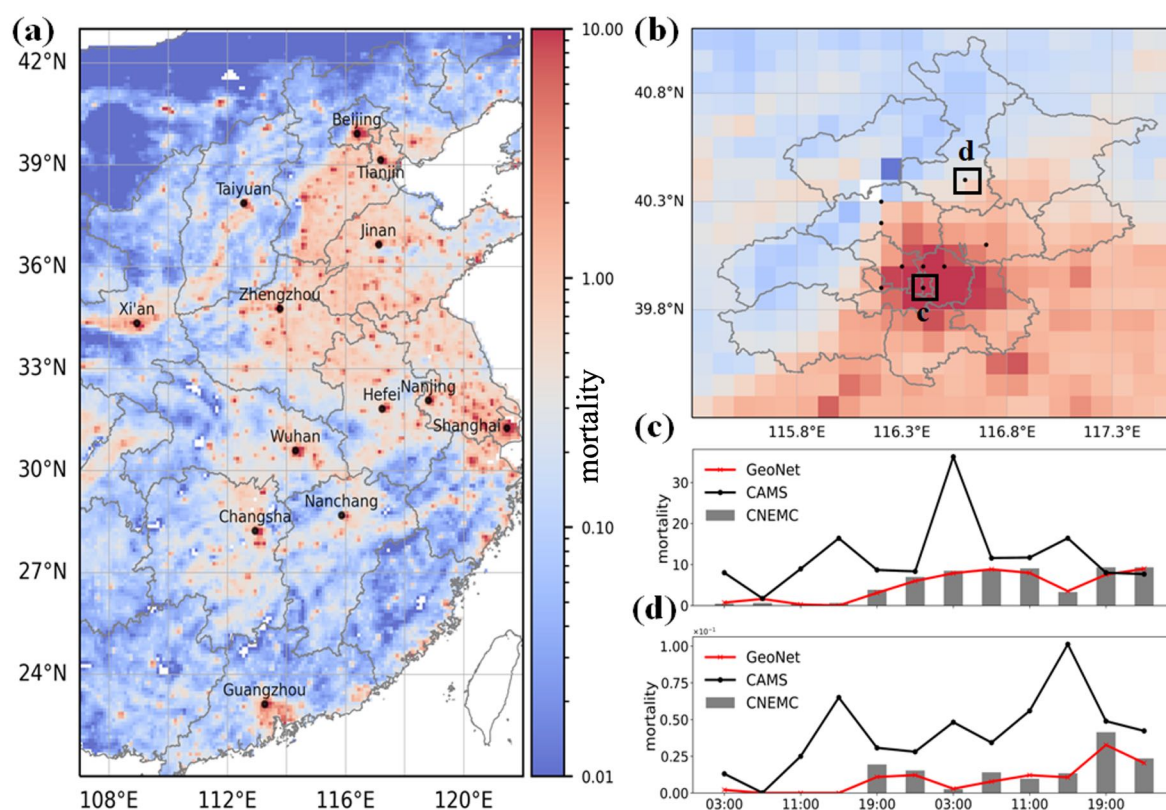


Fig R4 (Fig. 6 in the manuscript). Mortality risk of short-term NO₂ exposure based on the GeoNet prediction on November 23, 2021. **(a)** mean mortality due to the predicted NO₂ exposure in East China; **(b)** a zoom-in map over Beijing and its neighboring area; **(c)** and **(d)** are comparisons of mortality estimation over the Beijing urban and rural regions (the rectangle areas presented in **b**), respectively, based on different NO₂ exposure prediction among GeoNet, CAMS, and CNEMC.

11. L338-339: I don't believe the timeframe of this study was mentioned at all in the main text. What months / years was this prediction trained on and for what period was it validated?

Please refer to the response to major comment #2.

Reviewer #2:

The authors attempted to improve the short-term prediction of surface NO₂ at a high spatial and temporal resolution by taking advantage of the GEMS NO₂ products and a neural network model. They successfully forecasted full-coverage surface NO₂ for the next 24 hours and identified the critical role of GEMS NO₂. Their results demonstrate the potential application of the GEMS products in air quality prediction.

Overall, this is an important study and the results presented here will be useful for future applications of GEMS products as well as the geostationary satellite observations. I am happy to see its publication in due course. However, before that, I still have a few concerns or suggestions for the authors.

We appreciate the reviewer for the positive comments. We have addressed the following concerns point-by-point below:

- 1. I would suggest to move the model configuration and optimization into the main text. This will be very helpful for readers to understand the model.*

Thanks for your suggestion. We now moved the model configuration and optimization details to the L188-219 (Section 2.4) of the main text. Please refer to the details in the revised manuscript.

- 2. In the handling of missing data, the authors tried to set them to a fill value of zero. Is it reasonable? It looks reasonable to fill values of diurnal NO₂ climatology (e.g., the seasonal mean diurnal NO₂). In addition, as shown in Fig.2, it looks the three methods of handling missing data perform similarly in term of R² and RMSE. So I don't think it is necessary to highlight the "weakest" or "strongest" configuration.*

We greatly appreciate the reviewer's concerns regarding our method of handling missing data.

The reason we chose to use a direct fill value of zero is that during the experiments, we observed that this approach, combined with the model's feature extraction capability, could effectively mitigate the negative impact of missing data on prediction performance. Meanwhile, by incorporating cloud cover (mainly effectors of the missing satellite measurements), this method allows the model to fully leverage its potential, overcoming the issues caused by missing data. We acknowledge the reasonableness of using diurnal NO₂ climatology, as the reviewer suggested. However, the data discontinuities between satellite NO₂ measurements and gap-filled NO₂ climatology could lead to systematical biases in model predictions. Moreover, as indicated by the results, R² and RMSE metrics did not show significant differences across the various methods of missing data handling. Therefore, we believe our current approach is valid, with minimum effect on data discontinuities.

In general, we now provide a more detailed discussion of the different missing data handling techniques, clarifying the rationale for our choice and highlighting the positive role that cloud-cover data plays in enhancing the model's performance. Please refer to the detailed revisions in L256-263 of the manuscript:

The comparison results to the validation datasets indicate that the scenario using CAMS meteorology datasets and replacing missing satellite NO₂ data with fill-values (Fig. 2c), corresponds to a modest NO₂ prediction performance with R²=0.68 and RMSE=12.26 μg/m³. In contrast, the configuration scenario using ERA-5 reanalysis meteorology and imputing with WRF-Chem simulations (Fig. 2a), corresponds to the best prediction

performance of $R^2=0.69$ and $RMSE=11.88 \mu\text{g}/\text{m}^3$. This may indicate that the importance of satellite missing data imputation may be diminished by cloud mask inputs, especially since the model can extract informative features from spatial and temporal neighboring inputs.

3. I would also suggest to move Fig.S12 in the main text which shows the advantage of GEMS measurements.

Done. We moved Fig. S12 into the main text (Fig. 4).

4. The authors also show that the performance of GeoNet model degrades notably after t+16h. Is there any possible solution to overcome this short predicability?

We appreciate the valuable feedback regarding the performance degradation of our short-term air quality forecasting model beyond t+16h. We acknowledge that this is a common phenomenon in many air quality prediction models, particularly due to the inherent uncertainties in meteorology, emissions, and chemical reactions that significantly affect longer-term forecasts. Although the accuracy of predictions after t+16h declines, the results still provide meaningful insights. In future work, we aim to improve long-term forecasting performance by adopting more sophisticated model architectures and incorporating enhanced observational data or physical constraints.

In response to this comment, we also revised the manuscript in the following aspects:

- (1) Explicitly discuss the reasons for the model's performance degradation beyond t+16h in the manuscript.
- (2) Outline future research directions, including the exploration of model frameworks with better temporal generalization and the integration of stronger observational and physical constraints to improve long-term predictions.

Please refer to L387-392 in the revised manuscript:

The variation of the model forecasting performance also shows that accurate prediction for longer time windows and heavy pollution events is still a major difficulty. This may be due to the high level of uncertainty in emissions and meteorology. In the future, a combination of higher resolution and more accurate multi-source data constraints, as well as machine learning models coupled with atmospheric physical mechanisms, may be needed to improve the existing forecasts.

5. As mentioned in my last comment, the authors also highlight the possible applications to other air pollutants. However, the chemistry and lifetime of other air pollutants might be very different from NO₂. For example, if the GEMS tropospheric ozone measurement is useful for the prediction of surface ozone? Some more detailed discussions on the possible application to other pollutants would be very insightful.

Thanks for your comment. We elaborated it in the L398-402 in the revised manuscript:

This work also has important implications for the prediction of near-surface O₃ and particulate matter. For example, the integration of using vertical O₃ profiles from the GEMS satellite, in particular near-surface layer concentrations, and their joint observations of important O₃ precursors including NO₂ and HCHO, is expected to significantly improve the uncertainty of existing estimates of near-surface air pollution.

6. *BTW, I am not sure whether it is necessary to include NO2 in the Title of this manuscript since the authors didn't talk too much about major air pollutants (i.e., PM2.5 and ozone)*

Thanks for your comment. Considering the broad implications of this research for forecasting other pollutants using Geostationary satellite measurement (also discussed in response to comment #5), we think it would be better to retain the current title.

References:

- (1) Fino, A., Vichi, F., Leonardi, C., & Mukhopadhyay, K. (2021). An overview of experiences made and tools used to inform the public on ambient air quality. *Atmosphere*, *12*(11), 1524.
- (2) Li, Y., Xing, C., Peng, H., Song, Y., Zhang, C., Xue, J., et al. (2023). Long-term observations of NO₂ using GEMS in China: Validations and regional transport. *Science of The Total Environment*, *904*, 166762.
- (3) Tang, K. T. J., Lin, C., Wang, Z., Pang, S. W., Wong, T.-W., Yu, I. T. S., et al. (2024). Update of Air Quality Health Index (AQHI) and harmonization of health protection and climate mitigation. *Atmospheric Environment*, *326*, 120473.

1 **Unleashing the Potential of Geostationary Satellite Observations in Air**
2 **Quality Forecasting Through Artificial Intelligence Techniques**

3 Chengxin Zhang¹, Xinhan Niu¹, Hongyu Wu², Zhipeng Ding², Ka Lok Chan³, Jhoon Kim⁴,
4 Thomas Wagner⁵, Cheng Liu^{1,6,7*}

5 ¹Department of Precision Machinery and Precision Instrumentation, University of Science and
6 Technology of China, Hefei, 230026, China

7 ²School of Environmental Science and Optoelectronic Technology, University of Science and
8 Technology of China, Hefei, 230026, China

9 ³Rutherford Appleton Laboratory Space, Harwell Oxford, United Kingdom

10 ⁴Department of Atmospheric Sciences, Yonsei University, Seoul, Republic of Korea

11 ⁵Satellite Remote Sensing Group, Max Planck Institute for Chemistry, Mainz, Germany

12 ⁶Key Laboratory of Environmental Optics and Technology, Anhui Institute of Optics and Fine
13 Mechanics, Chinese Academy of Sciences, Hefei, 230031, China

14 ⁷Key Laboratory of Precision Scientific Instrumentation of Anhui Higher Education Institutes,
15 University of Science and Technology of China, Hefei, 230026, China

16

17 *Correspondence: Cheng Liu (chliu81@ustc.edu.cn)

18

19

20 **Abstract.**

21 Air quality forecasting plays a critical role in mitigating air pollution. However, current
22 physics-based air pollution predictions encounter challenges in accuracy and spatiotemporal
23 resolution due to limitations in the understanding of atmospheric physical mechanisms,
24 observational constraints, and computational capacity. The world's first geostationary satellite
25 UV-Vis spectrometer, i.e., the Geostationary Environment Monitoring Spectrometer (GEMS),
26 offers hourly measurements of atmospheric trace gas pollutants at high spatial resolution over
27 East Asia. In this study, we successfully incorporate Geostationary satellite observations into
28 a neural network model (GeoNet) to forecast full-coverage surface nitrogen dioxide (NO₂)
29 concentrations over eastern China at 4-hour intervals for the next 24 hours. GeoNet leverages
30 spatiotemporal series of satellite NO₂ observations to capture the intricate relationships among
31 air quality, meteorology, and emissions in both temporal and spatial domains. Evaluation
32 against ground-based measurements demonstrates that GeoNet accurately predicts diurnal
33 variations and spatial distribution details of next-day NO₂ pollution, yielding the coefficient of
34 determination of 0.68 and root mean square of error of 12.31 μg/m³, significantly surpassing
35 traditional air quality model forecasts. The model's interpretability reveals that geostationary
36 satellite observations notably improve NO₂ forecast capability more than other input features,
37 especially over polluted regions. Our findings demonstrate the significant potential of
38 geostationary satellite observations in artificial intelligence-based air quality forecasting, with
39 implications for early warning of air pollution events and human health exposure.

40 **Keywords:** air quality forecast; deep learning; health impact; satellite remote sensing;
41 nitrogen dioxide;

42 1 Introduction

43 Since the industrial revolution, numerous countries worldwide have encountered severe
44 air pollution issues such as photochemical ozone smog and haze pollution (Hong et al., 2019),
45 which significantly affect human health, crop yields, and the global environment (Guarin et al.,
46 2024; Manisalidis et al., 2020; Sathe et al., 2021). Recent studies have shown that both long-
47 term and short-term exposure to air pollutants such as nitrogen dioxide (NO₂) can significantly
48 affect human health, especially the respiratory system (Meng et al., 2021). Accurate and high
49 spatial resolution predictions of air pollutant concentrations can provide critical information
50 for sensitive persons to mitigate health risks. [Meanwhile, air quality health risk \(AQHI\)](#)
51 [forecasts and corresponding public response recommendations need to be communicated to the](#)
52 [public promptly through public facilities \(Fino et al., 2021; Tang et al., 2024\)](#). In recent decades,
53 the advancement of atmospheric monitoring and modeling has enabled significant progress in
54 air quality forecasting based on our understanding of atmospheric physics and chemistry
55 (Peuch et al., 2022). Air pollution forecasting not only facilitates responses to environmental
56 health risks but also improves the accuracy of climate and weather simulations (Makar et al.,
57 2015). However, due to our still limited understanding of atmospheric mechanisms and
58 observational and emission constraints, existing air quality forecasts based on physical or
59 statistical models still face challenges in terms of temporal, spatial, and accuracy aspects
60 (Campbell et al., 2022; Zhong et al., 2021).

61 Artificial Intelligence (AI) technology has made breakthroughs in the field of Earth
62 science (Boukabara et al., 2020; Zhong et al., 2021), particularly excelling in addressing
63 complex problems that are challenging for traditional physical paradigms to simulate (Irrgang
64 et al., 2021), such as weather and climate forecasting (Andersson et al., 2021). Concerning
65 meteorological data, some large-scale deep learning models have surpassed the predictive
66 capabilities of existing numerical weather models to some extent, examples include Climax

Deleted: limited

68 (Nguyen et al., 2023), Pangu-Weather (Bi et al., 2023), and GraphCast (Lam et al., 2023).
69 Despite significant progress and impressive performance achieved in meteorological variables
70 forecasting with AI methods, there are still limitations in predicting atmospheric pollutant
71 compositions. Compared to meteorological parameters, the prediction of air pollutant
72 concentrations is affected by synoptic meteorology, chemistry, and anthropogenic emission
73 activities, usually with more complex driven mechanisms and associated uncertainties. Current
74 AI-based air quality forecasts often involve time series predictions at a limited number of
75 observation stations, rather than full-coverage predictions over the entire spatial domain (Du
76 et al., 2021). This is primarily due to the lack of effective air quality observations with high
77 temporal and spatial resolution simultaneously.

78 While past polar-orbiting satellite observations [such as the Ozone Monitoring Instrument](#)
79 [\(OMI\) and the TROPospheric Monitoring Instrument \(TROPOMI\)](#), have provided extensive
80 coverage of atmospheric pollutant distributions such as nitrogen dioxide (NO₂), sulfate dioxide
81 (SO₂), ozone (O₃), and aerosols, they are limited to once-daily overpasses and usually affected
82 by clouds (Chan et al., 2023; Van Geffen et al., 2022). This frequency usually exceeds the
83 chemical lifetimes of many reactive gas pollutants like NO₂, making it challenging to offer
84 effective observational constraints for machine learning short-term air quality forecasting
85 (Shah et al., 2020). [However, these observations at a fixed daily overpass time could hardly](#)
86 [support the prediction of atmospheric trace gas concentrations at other times of the day under](#)
87 [different meteorological conditions](#). In February 2020, the world's first geostationary satellite
88 payload for air pollution monitoring, the Geostationary Environment Monitoring Spectrometer
89 (GEMS), began to provide high-coverage and high-precision air quality observations at an
90 hourly rate for the East Asian region (J. Kim et al., 2020). The dynamic processes of air
91 pollutants including emission, transformation, and transport can be observed by the
92 geostationary satellite during the daytime. This monitoring capability may advance data-driven

93 air quality forecasting such as machine learning techniques by offering unprecedented
94 observational constraints with high spatial and temporal coverage. Recent observing system
95 simulation experiments (OSSE) indicate that assimilating trace gas observations by
96 geostationary satellites into chemical models can effectively improve surface ozone
97 simulations (Shu et al., 2023), nitrogen oxides (NO_x), and emission estimates (Hsu et al., 2024).

98 Here, based on the unprecedented temporal [and spatial](#) resolution [and coverage](#) of the
99 GEMS satellite (J. Kim et al., 2020), we incorporated Geostationary satellite remote sensing of
100 tropospheric NO₂ column densities (refer to section 4 for details) into a neural Network model
101 (GeoNet), to forecast full-coverage surface NO₂ concentration over the next day from the
102 current time t (i.e., $t+24h$). Compared with previous air quality forecasting based on the
103 simulation of atmospheric physics and chemistry possibly combined with data assimilation
104 approaches, GeoNet relies solely on geostationary satellite measurements and ancillary
105 meteorology data. GeoNet effectively addresses the complex nonlinear relationships between
106 future short-term air quality and current satellite observations, as well as temporally adjacent
107 meteorological variables (C. Zhang et al., 2022). The method employs satellite and
108 meteorological variables within the spatial vicinity of individual air quality monitoring sites as
109 input features, with site observations serving as labels for model training. The resulting model
110 achieves accurate and comprehensive air quality predictions across the entire domain over East
111 China, which is a significant achievement given that past machine learning technologies have
112 relied on only a few stations or polar-orbiting satellite observations.

113 **2 Materials and Methods**

114 **2.1 Geostationary satellite observations of atmospheric NO₂**

115 GEMS is the first UV-Vis spectrometer at a geostationary satellite orbit, measuring
116 atmospheric pollutants such as NO₂, SO₂, O₃, and HCHO over East Asia, at a spatial resolution
117 of $3.5 \text{ km} \times 7.5 \text{ km}$ at nadir and a temporal resolution of 1 hour during the daytime (J. Kim et

118 al., 2020). Based on the unique spectral absorption of trace gases, the atmospheric NO₂ column
119 can be retrieved in visible wavelengths from the spectra of back-scattered sunlight. The details
120 of the GEMS NO₂ retrieval can be found in the Algorithm Theoretical Basis Document
121 (available at <https://nesc.nier.go.kr/ko/html/satellite/doc/doc.do>, last access: June 1, 2023). In
122 this study, we used the tropospheric NO₂ column from the GEMS NO₂ version 2.0 product, as
123 well as the cloud fraction for each satellite ground pixel. Overall, GEMS NO₂ measurements
124 have a good correlation with ground-based remote sensing instruments, with correlation
125 coefficients (R) between 0.69-0.81, and root mean square of errors (RMSE) between 3.2-
126 4.9×10^{15} molecules/cm² (S. Kim et al., 2023). [Our previous validation results indicated that](#)
127 [GEMS NO₂ retrievals generally agreed well with ground-based MAX-DOAS measurements](#)
128 [from 6 sites in China, with correlation coefficients ranging between 0.69-0.92 \(Li et al., 2023\).](#)

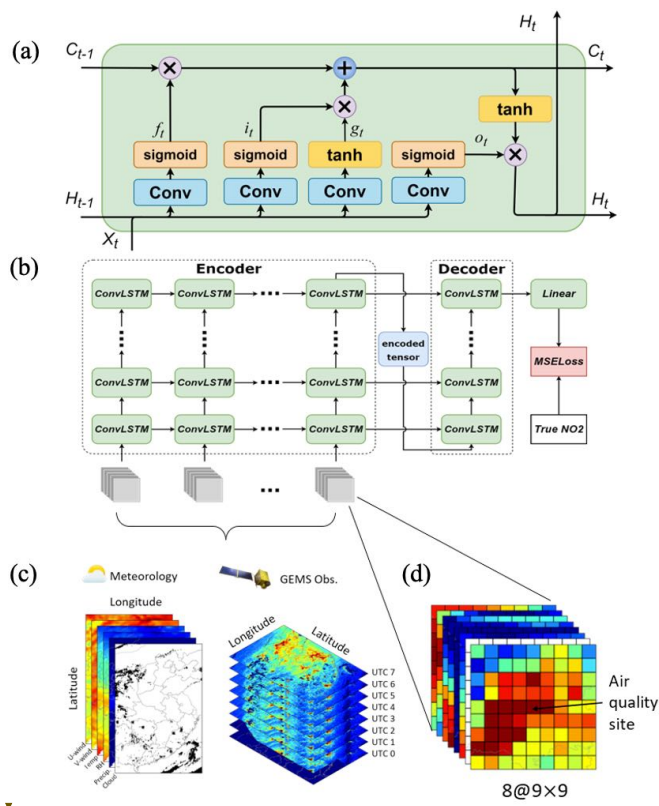
129 2.2 Ancillary datasets

130 Other input information including meteorological datasets is necessary to better constrain
131 the prediction of future NO₂ pollution. Here, both the ERA5 meteorology reanalysis (Hersbach
132 et al., 2020) and the CAMS forecast (Peuch et al., 2022) were used to provide meteorological
133 parameters such as zonal and meridional wind (U-wind and V-wind), temperature (Temp),
134 relative humidity (RH), and precipitation (Precip). In addition, the fraction of cloud cover
135 available from the satellite NO₂ datasets was also considered. To fill the missing gaps in the
136 satellite NO₂ measurements, we use both the NO₂ concentrations from the WRF-Chem model
137 (C. Zhang et al., 2022) and the CAMS forecast of atmospheric composition. Note that the
138 reanalysis datasets were typically updated with a week delay from real-time, while the forecast
139 datasets can provide future 7-day meteorology from the current time. Therefore, the latency of
140 input datasets would affect the operational prediction of the GeoNet model. Surface NO₂
141 measurements were used as the ground-truth label in the model training phase, available from

142 over 1000 national air quality sites by the China National Environmental Monitoring Centre
 143 (CNEMC) (Kong et al., 2021).

144 The preprocessing steps of model input datasets, including outlier detection, missing value
 145 handling, resampling, and normalization, are described in Supplementary Text S1.

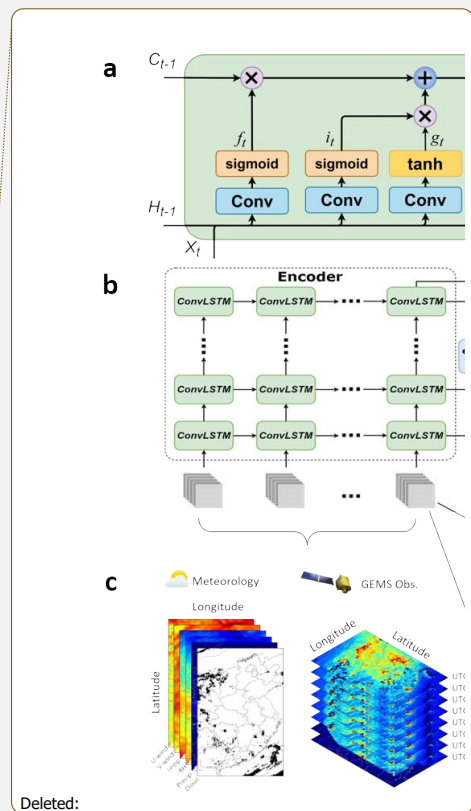
146 **2.3 The GeoNet model**



147

148 **Figure 1.** The framework of predicting surface NO₂ map based on Geostationary satellite measurements and
 149 a ConvLSTM neural network model (GeoNet). (a) the structure of the ConvLSTM block; (b) a diagram of
 150 GeoNet model structure with inputs and output; (c) an illustration of the model input parameters including
 151 meteorological variables and hourly NO₂ measurements by the Geostationary satellite; (d) the input data
 152 cube of different features for single training batch, which is centered at an air quality site.

153 Fig. 1 illustrates the structure and methodology of the artificial intelligence air quality
 154 forecasting model established in this study. Given the distinctive nature of spatiotemporal



Deleted:

156 sequence data for air quality, predictions must consider not only temporal relationships but also
 157 spatial correlations. The deep learning model employed in this research utilizes convolutional
 158 long short-term memory (ConvLSTM) as its kernel, a variant of the LSTM model designed for
 159 the time series forecasting (Lin et al., 2020). It incorporates a convolutional network structure
 160 to capture spatial features of three-dimensional inputs. Both input-to-state and state-to-state
 161 transitions involve convolutional structures. ConvLSTM determines the future state of a unit
 162 within a grid based on inputs from its local neighbors and past states, allowing it to effectively
 163 model the spatiotemporal dynamics of air quality. The ConvLSTM kernel structure employed
 164 in training is illustrated in Fig. 5a. Here, X_t represents the input at time t, H_t and H_{t-1} denote
 165 the outputs at times t and t-1, and C_t and C_{t-1} represent the states at times t and t-1. The
 166 computational process is as follows:

$$167 \quad i_t = \sigma(X_t * w_{xi} + H_{t-1} * w_{hi} + b_i) \quad (1)$$

$$168 \quad f_t = \sigma(X_t * w_{xf} + H_{t-1} * w_{hf} + b_f) \quad (2)$$

$$169 \quad o_t = \sigma(X_t * w_{xo} + H_{t-1} * w_{ho} + b_o) \quad (3)$$

$$170 \quad g_t = \tanh(X_t * w_{xg} + H_{t-1} * w_{hg} + b_g) \quad (4)$$

$$171 \quad C_t = f_t \times C_{t-1} + i_t \times g_t \quad (5)$$

$$172 \quad H_t = o_t \times \tanh(C_t) \quad (6)$$

173 Where the asterisk (*) represents the convolution operator, w is the convolution kernel, b is the
 174 offset, \tanh is the hyperbolic tangent function, and σ is the activation function of Sigmoid.

175 The model primarily consists of three components: an encoder, a decoder, and fully
 176 connected layers. Tropospheric NO₂ observations from the GEMS satellite for different
 177 consecutive hours within a day, along with corresponding meteorological forecast field data,
 178 serve as input features for model training. The encoder processes the spatiotemporal sequences
 179 of input features for the preceding 8 hours (t-7h, t-6h, ..., t), which are then decoded by the
 180 decoder. The final output, representing NO₂ concentrations at 4-hour intervals for the next 24

181 hours (t+4h, t+8h, t+12h,..., t+24h), is produced through fully connected layers. The loss
182 function of mean squared error (MSE) is calculated by comparing the model output with the
183 actual values from station observations, and the model undergoes iterative training. In the
184 training task for a single station sample, the model utilizes continuous and distinct hourly
185 dynamic images of all variables within the spatiotemporal vicinity of the station as input (see
186 Fig. 1c-d). This effectively considers the intricate correlations in time and space between air
187 quality, satellite observations, and meteorological input features. [We train the GeoNet model](#)
188 [with input features during the whole year of 2021. The training datasets were randomly selected](#)
189 [from 75% of the whole samples, while the remaining 25% were used as validation sets.](#)

190 **2.4 The model configuration and optimization**

191 [The model configurations and hyperparameters such as the optimizer, loss function, L1 or](#)
192 [L2 regularization, dropout, training steps, and epochs can make a difference to the model](#)
193 [performance including the prediction accuracy and generalizability. ~~The performance metrics~~](#)
194 [such as the coefficient of determination \(\$R^2\$ \), root mean square of error \(RMSE\), mean absolute](#)
195 [error \(MAE\), and mean absolute percentage error \(MAPE\), were used to diagnose the model](#)
196 [\(see definition in Supplementary Text S2\). ~~Thus, several scenarios of model hyperparameters~~](#)
197 [have been tested during the model training phase. The model accuracy on validation datasets](#)
198 [and the learning rate curve were used to diagnose the model hyperparameters. The model](#)
199 [parameters mainly include the number of layers and the dimensions of the hidden layers, both](#)
200 [control the model's capacity. If the model capacity is relatively small, underfitting may occur;](#)
201 [overfitting may exist if it is too large. Therefore, selecting an appropriate model capacity is](#)
202 [crucial for improving model performance. During the pre-training process, the model is trained](#)
203 [by combining different numbers of layers and dimensions of the hidden layers. The Mean](#)
204 [Squared Error \(MSE\) Loss is recorded for each training iteration, and a heatmap is generated](#)
205 [as shown in Fig. S2. From the heatmap, it can be observed that when the number of layers is 2](#)

Moved down [1]: The performance metrics such as the coefficient of determination (R^2), root mean square of error (RMSE), mean absolute error (MAE), and mean absolute percentage error (MAPE), were used to diagnose the model (see definition in Supplementary Text S2). T

Deleted: The model configuration and optimization are also described in detail in Supplementary Text S2

Moved (insertion) [1]

Deleted: 2.4

214 [and the dimension of the hidden layer is 256, the model achieves the minimum MSE Loss. Fig.](#)
215 [S3 shows the sensitivity test results of model loss varying with different batch size settings,](#)
216 [indicating that a batch size of 64 is optimal. Based on the model's MSE loss under different](#)
217 [hyperparameter configurations, the best-fitting model can be selected.](#)

218 [The Adam optimization algorithm controls the learning rate, which can design](#)
219 [independent adaptive learning rates for different parameters. The three initialization parameters](#)
220 [\$\epsilon\$, \$\rho_1\$, and \$\rho_2\$ of the Adam algorithm are set to be 0.0001, 0.9, and 0.99, respectively. For the](#)
221 [epoch, its size is controlled by the early stop method. The early stop method monitors the](#)
222 [change of the model's loss function on the validation set during the training process and stops](#)
223 [the model training immediately when the validation loss of the model starts to become larger.](#)
224 [Due to the fluctuation of the loss function, a threshold \$p\$ is set for the early stopping method in](#)
225 [practice, and when the validation loss of the model becomes large for \$p\$ consecutive epochs,](#)
226 [the model is rolled back to the lowest validation loss and the training is stopped, and the](#)
227 [threshold \$p\$ is set to 10 in this paper. Fig. S4 shows a typical learning curve of the MSE loss in](#)
228 [training and validation data sets for different learning steps in training an optimal model. Such](#)
229 [diagnostics can be used to avoid the model overfitting.](#)

230 **[2.5 The importance of the model input feature](#)**

231 [Permutation feature importance is a technique used to assess the significance of each input](#)
232 [feature in a machine-learning model \(Altmann et al., 2010\). The core idea is to evaluate the](#)
233 [impact of each feature on model performance by randomly shuffling its values and observing](#)
234 [the resulting change in the model's accuracy. In this study, for each input feature of the GeoNet,](#)
235 [we iteratively shuffle its value independently while keeping other features unchanged, and then](#)
236 [observe the model prediction on the modified input. The difference in the model prediction](#)
237 [performance between using the original and shuffling input quantifies the feature's importance.](#)
238 [Here, we measure the relative importance of each input feature using the metric of \$1-R^2\$, due](#)

239 to its good standardized and indicative ability (C. Zhang et al., 2022). Generally, a larger
240 performance drop indicates greater importance, as the model heavily relies on that feature for
241 predictions. Conversely, smaller drops or increases suggest the feature may be less crucial or
242 redundant. By permuting the input feature array based on the different spatial and temporal
243 domains, we can gain a deeper understanding of how feature importance varies spatially and
244 temporally. For example, the relative importance of one meteorology variable may vary with
245 different diurnal, weekly, and monthly cycles, revealing the variability of its impact on the
246 predicted NO₂ levels.

247 **3 Results and Discussion**

248 **3.1 Model performance**

249 Based on the GeoNet model and necessary input data (refer to section 2), we have
250 achieved preliminary predictions of near-surface NO₂ concentration with full spatial coverage
251 and a spatial resolution of 0.1 degrees over eastern China, at four-hour intervals over the next
252 24 hours. In this study, we first tested the impact of using reanalysis and forecast meteorology
253 datasets and filling in missing values in satellite observation data on the model predictions. The
254 reanalysis datasets usually have higher precision than the forecast. Previous studies revealed
255 that the accuracy of the information on meteorology and chemical composition significantly
256 affects the performance of machine learning models in estimating air pollutant concentrations
257 (Wang et al., 2024; Zuo et al., 2023). Due to the shielding effect of clouds, a considerable
258 proportion of missing values may even exist in satellite NO₂ observations. Recent air quality
259 big-data research usually requires the gap-filling of missing satellite data before inputting it
260 into the machine learning model, either by spatial interpolation or regression techniques (M.
261 Kim et al., 2021). We tested three methods for handling missing data, such as setting them to
262 a fill value of zero, or replacing them by real-time CAMS simulated NO₂, or WRF-Chem
263 simulated NO₂ results (not real-time, but with higher precision).

Deleted: 01

Deleted: Meteorological data sources included ERA5 reanalysis meteorology datasets with a latency of one week, and CAMS forecast meteorology data for the upcoming 7 days.

268 The comparison results to the validation datasets indicate that the scenario using CAMS
269 meteorology datasets and replacing missing satellite NO₂ data with fill-values (Fig. 2c),
270 corresponds to a modest NO₂ prediction performance with R²=0.68 and RMSE=12.26 μg/m³.
271 In contrast, the configuration scenario using ERA-5 reanalysis meteorology and imputing with
272 WRF-Chem simulations (Fig. 2a), corresponds to the best prediction performance of R²=0.69
273 and RMSE=11.88 μg/m³. This may indicate that the importance of satellite missing data
274 imputation may be diminished by cloud mask inputs, especially since the model can extract
275 informative features from spatial and temporal neighboring inputs. To compromise between
276 the performance of real-time and accuracy, we selected the configuration scenario of using
277 CAMS meteorology and imputing with CAMS NO₂ (Fig. 2d) for subsequent discussion and
278 operational forecasting, with an R²=0.68 and RMSE=12.31 μg/m³. In summary, the use of
279 higher-precision meteorology and filling missing NO₂ data enhances the model's prediction
280 accuracy on the validation dataset, but to a rather limited extent. This suggests that, unlike
281 previous machine learning techniques, GeoNet can effectively adapt to three-dimensional
282 inputs of varying accuracy and type, fully explore the spatiotemporal correlation of data
283 features, and demonstrate strong model generalization capabilities.

Deleted: (10% of whole datasets, randomly sampled)

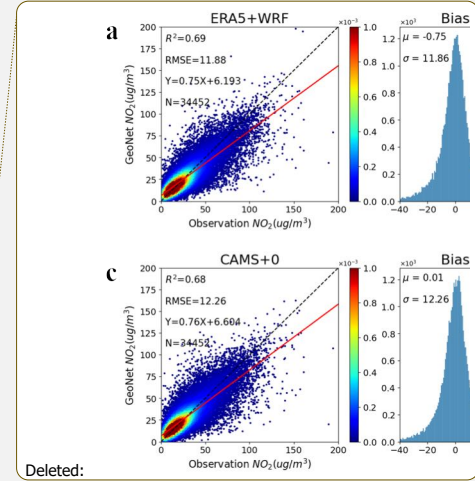
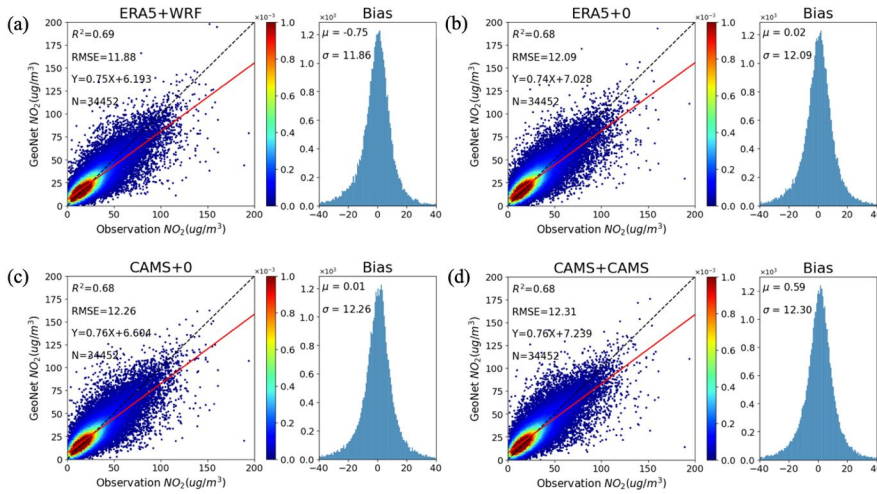
Deleted: with the "weakest" input, i.e.,

Deleted: the

Deleted: metrics of

Deleted: strongest

Deleted: ,



290

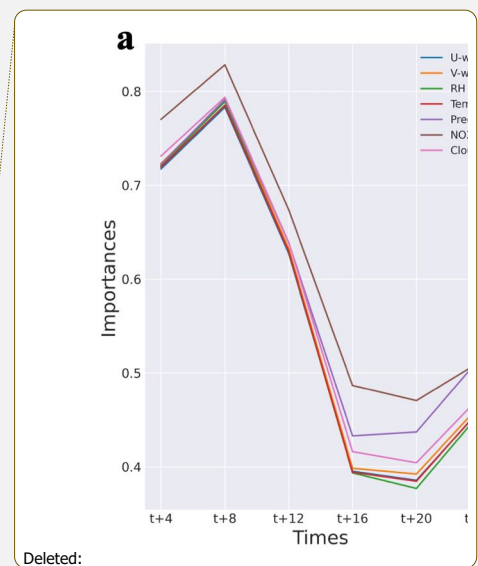
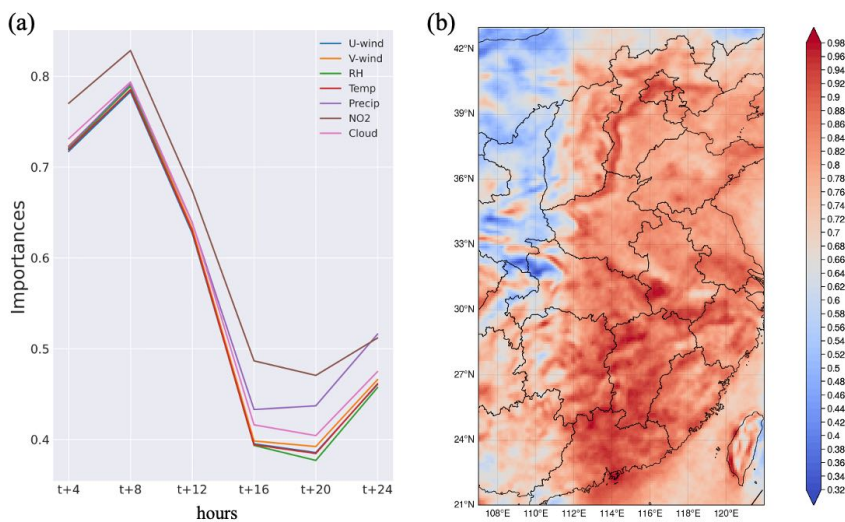
291 **Figure 2.** The GeoNet prediction performance of the surface NO₂ concentration compared to the validation
 292 samples, based on different input datasets of meteorology and atmospheric composition: (a) use ERA5
 293 meteorology and fill satellite measurement gaps with WRF-Chem simulated NO₂; (b) use ERA5
 294 meteorology and NO₂ fill-value of zero for over gaps; (c) use CAMS meteorology and NO₂ fill-value of zero
 295 for gaps; (d) use CAMS meteorology and CAMS NO₂. The left plot shows the scatter comparisons between
 296 GeoNet predictions and site observations, while the right plot shows the bias distribution between the two.

297 Figs. S5-S8 provide an overview of the major metrics (e.g., R², RMSE, MAE, and MPE)
 298 of GeoNet prediction performance varying with prediction hours from t+4h to t+24h in
 299 different months. The results indicate that the model exhibits a higher correlation in NO₂
 300 forecast during the spring and winter seasons compared to the summer, while the RMSE errors
 301 show the opposite trend. This could be attributed to much higher NO₂ pollution levels in winter
 302 months. Additionally, GeoNet's NO₂ prediction errors gradually increase during the next 24
 303 hours, particularly after t+20h. This is primarily due to the short lifetime of atmospheric NO₂,
 304 leading to a diminishing constraint from historical observational data on future NO₂ predictions.
 305 Similar phenomena are also observed in machine learning or model-assisted weather forecasts
 306 (Andersson et al., 2021).

308 To assess the GeoNet model's performance for short-term pollution events, we compared
309 it with near-surface NO₂ from CAMS forecasts, and in situ observations from CNEMC ground
310 stations. Fig. S9 illustrates the daily time series of t+4h NO₂ from GeoNet, CAMS, and
311 CNEMC for three typical sites in Beijing, Shanghai, and Guangzhou in 2021. As shown from
312 the plot, NO₂ predictions by both GeoNet and CAMS generally agreed with the variation trends
313 of CNEMC measurement. However, CAMS forecasts systematically overestimate the surface
314 NO₂ concentration by 100%, possibly resulting from the biases in the NO_x emission inventory
315 (Douros et al., 2023). Compared to CAMS, the GeoNet prediction closely aligns with the
316 ground-truth observations at CNEMC sites over eastern China, with an overall R² > 0.5 and
317 mean bias < 5 μg/m³ for polluted regions (see Fig. S10 and S11, respectively).

318 **3.2 Main factors in NO₂ forecast and their implications**

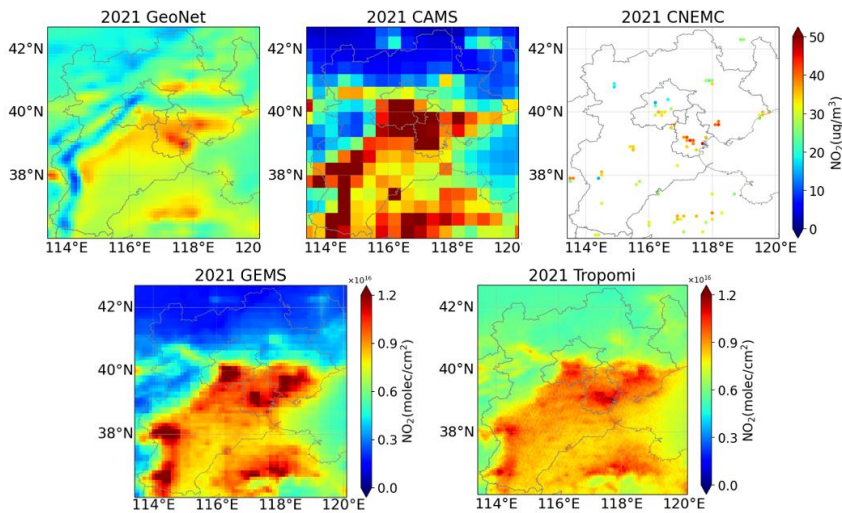
319 Previous physics-based numeric models of air quality prediction, e.g., the CAMS global
320 forecast model and the regional WRF-CMAQ model (Kuhn et al., 2024; Kumar et al., 2021;
321 Liu et al., 2023), can simulate the atmospheric physical and chemical processes (such as
322 advection, diffusion, deposition, and chemical reactions) by solving the atmospheric equations.
323 Recent data assimilation techniques further take real-time monitoring data from satellite and
324 ground-based platforms as model constraints to better predict air quality variables (Antje Inness
325 et al., 2022). Compared with physics-based models, "black-box" models such as the deep
326 learning technique usually lack interpretability and explainability (Q.-s. Zhang & Zhu, 2018).
327 This hinders the understanding and implications for predicting air quality variables such as
328 NO₂. Here, we measure the relative importance of each input feature on the NO₂ forecast
329 accuracy, by iteratively permuting the input array and observing its influences on the model
330 prediction.



331
 332 **Figure 3.** (a) The overall relative importance of different input features such as wind, surface pressure,
 333 satellite NO₂, and cloud mask, in GeoNet NO₂ forecast, varying with different hour steps from t+4h to t+24h.
 334 (b) The spatial distribution of the relative importance of satellite NO₂ measurements in the GeoNet NO₂
 335 forecast in 2021.

336 Fig. 3a presents the relative importance (1-R²) of different input features varying with
 337 prediction hour steps from t+4h to t+24h. The geostationary satellite NO₂ measurements play
 338 the highest role in predicting surface NO₂ levels of the next day, although it degrades after t+8h.
 339 Other meteorological input features also show a major impact on NO₂ prediction performance.
 340 [The significance of the different input variables remained generally consistent across seasons,](#)
 341 [with minor variations \(as shown in Fig. S12\).](#) By permutating the input array for each ground
 342 pixel, Fig. 3b derived the spatial distribution of the relative importance of geostationary satellite
 343 NO₂ in the predicting performance. Overall, satellite NO₂ has a higher impact in densely
 344 populated areas experiencing severe air pollution, such as the Pearl River Delta, Yangtze River
 345 Delta, and Jianghuai Plain, than in western China. Such results highlight the underappreciated
 346 role of satellite NO₂ measurements with high spatial and temporal coverage in air pollution
 347 forecasts.

348 **3.3 NO₂ pollution episodes and health exposure forecast**



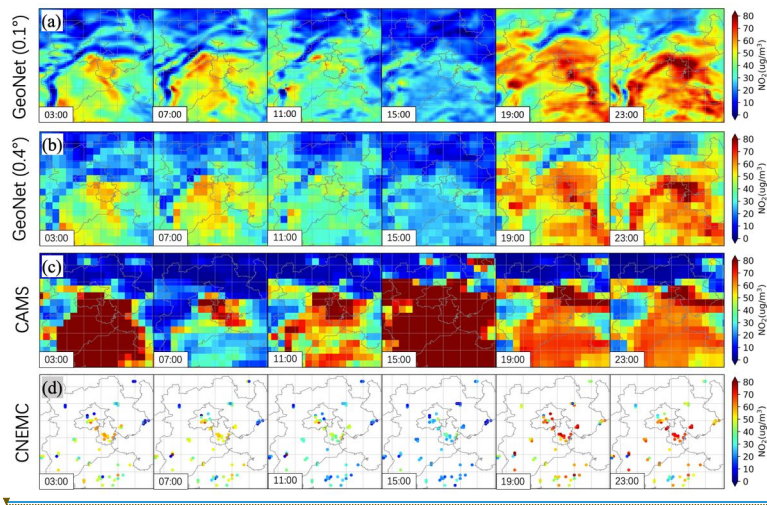
350 [Figure 4.](#) The comparisons of annual surface NO₂ concentrations from GeoNet, CAMS, and CNEMC,
 351 respectively, (in the top panel), as well as the tropospheric NO₂ column observations from GEMS and
 352 TROPOMI over East China in 2021 (in the bottom panel).
 353

354 Beyond its prediction accuracy, GeoNet exhibits a pronounced advantage in spatial
 355 coverage and resolution, allowing for capturing finer-scale details in the pollutant distribution.
 356 Illustrated in Fig. 4, GeoNet demonstrates remarkable performance in predicting spatial
 357 nuances of NO₂ pollution, particularly when contrasted with ground-based and satellite
 358 observations. During a typical winter NO₂ pollution event (as shown in Fig. 5), GeoNet
 359 accurately simulates a significant decrease in concentrations at 11:00 and 15:00, probably led
 360 by intense photochemical activity in the daytime, coincident with ground-based observations.
 361 It also outperforms CAMS in predicting NO₂ variations throughout the day. The GeoNet model
 362 also retains the distributional differences in NO₂ concentrations between urban and rural areas,
 363 consistent with emission source characteristics and satellite observations. The suboptimal
 364 performance of CAMS predictions can be attributed to insufficient observational constraints
 365 and the use of outdated emission inventories (Douros et al., 2023). In the European region, the
 366 assimilation of TROPOMI observations into CAMS forecasts significantly improves the

Deleted: S12

Deleted: 4

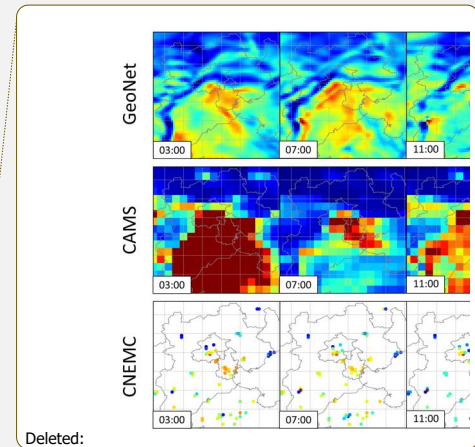
369 simulation accuracy of near-surface NO₂ concentrations and tropospheric column densities (A.
 370 Inness et al., 2019). Neural network methods, similar to GeoNet, could be used to correct and
 371 downscale forecast results by existing models (Baghanam et al., 2024). This approach holds
 372 promise for achieving operational air quality forecasts that balance efficiency and accuracy.



373

374 **Figure 5.** The spatial distribution comparisons of surface NO₂ concentration between (a) GeoNet prediction
 375 at the original resolution of 0.1°, (b) GeoNet prediction resampled to the CAMS resolution of 0.4°, (c) CAMS
 376 prediction, and (d) ground-based CNEMC site measurements. Note that the results are presented for different
 377 continuing local hours (labeled text in the subplot) on 23 November 2021.

378 In this study, we used a simplified linearized risk model for the short-term NO₂ exposure
 379 (Meng et al., 2021; C. Zhang et al., 2022) to calculate the distribution of all-cause mortality
 380 risks based on GeoNet NO₂ predictions (see Fig. 6). Short-term NO₂ exposure leads to
 381 remarkable regional differences in all-cause mortality, which are mainly concentrated in highly
 382 polluted and densely populated urban areas. For both urban and suburban locations in Beijing
 383 (see Fig. 6c-d), GeoNet-based NO₂ pollution exposure predictions are more consistent with
 384 actual in situ observations than the CAMS forecasts. Current air quality health indices
 385 forecasting based on limited station data has significant gaps, making it difficult to meet the
 386 refined needs for different populations in urban, suburban, and rural areas. Integrating GeoNet



Deleted:

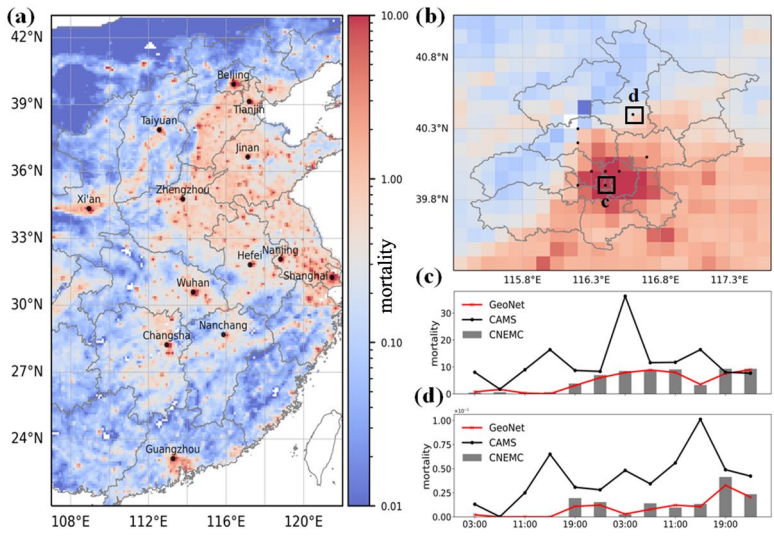
Deleted: 4

Deleted: c

Deleted: 5

Deleted: 5c

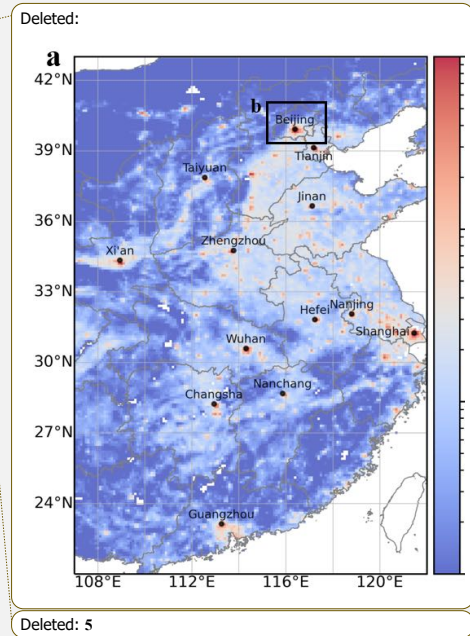
392 forecasts based on hourly geostationary satellite observations can support spatially
 393 comprehensive and fine-scale air quality health risk prediction. This, in turn, guides managing
 394 the risks of air pollution exposure-related diseases in sensitive populations and communities.



395
 396 **Figure 6.** Mortality risk of short-term NO₂ exposure based on the GeoNet prediction on November 23, 2021.
 397 (a) mean mortality due to the predicted NO₂ exposure in East China; (b) a zoom-in map over Beijing and its
 398 neighboring area; (c) and (d) are comparisons of mortality estimation over the Beijing urban and rural
 399 regions (the rectangle areas presented in b), respectively, based on different NO₂ exposure prediction among
 400 GeoNet, CAMS, and CNEMC.

401 **4 Conclusion**

402 The GeoNet model utilizes the unprecedented hourly air quality observations from
 403 geostationary satellites and resolves nonlinear associations in spatiotemporal proximity across
 404 multiple data sources. It achieves seamless short-term regional air quality predictions,
 405 exhibiting significant performance advantages over existing machine-learning air quality
 406 prediction models. To strike a balance between real-time and accuracy requirements, we
 407 evaluated the impact of using reanalysis- and forecast-based meteorology datasets, as well as
 408 imputing the missing values of satellite NO₂. The findings reveal that the GeoNet model
 409 demonstrates robust generalization across diverse datasets, with minimal fluctuations in



413 prediction performance. Overall, the model achieves an RMSE of 12.31 $\mu\text{g}/\text{m}^3$ and an R^2 of
414 0.68 in predicting NO_2 concentrations every 4 hours for the next 24 hours. However, validation
415 accuracy notably diminishes after t+16h within the next 24 hours, with stronger predictive
416 correlations observed in seasons characterized by severe pollution, such as spring and winter,
417 compared to summer. The variation of the model forecasting performance also shows that
418 accurate prediction for longer time windows and heavy pollution events is still a major
419 difficulty. This may be due to the high level of uncertainty in emissions and meteorology. In
420 the future, a combination of higher resolution and more accurate multi-source data constraints,
421 as well as machine learning models coupled with atmospheric physical mechanisms, may be
422 needed to improve the existing forecasts.

423 Compared to traditional chemical model forecasts and data assimilation predictions, the
424 GeoNet model handles various data sources, including meteorological simulations and air
425 quality observations, and more accurately captures spatial intricacies of air pollution evolution.
426 The GeoNet framework elucidated in this study forecasts short-term near-surface NO_2
427 concentrations and demonstrates transferable learning potentials for predicting other pollutants.

428 This work also has important implications for the prediction of near-surface O_3 and particulate
429 matter. For example, the integration of using vertical O_3 profiles from the GEMS satellite, in
430 particular near-surface layer concentrations, and their joint observations of important O_3
431 precursors including NO_2 and HCHO , is expected to significantly improve the uncertainty of
432 existing estimates of near-surface air pollution. This study underscores the pivotal role of next-
433 generation stationary satellite observations of air pollution constituents in air quality
434 forecasting, with the potential to advance operational air quality forecasting and mitigate
435 associated health risks by integrating machine learning technologies.

436

Deleted: In our forthcoming endeavors, we aim to enhance

Deleted: predictability

Deleted: ozone

Deleted: pollution by integrating

Deleted: tropospheric ozone and its

Deleted: like

443 **Data and code availability.** The GEMS NO₂ v2.0 data is available from the National Institute
444 of Environmental Research (NIER) of South Korea (<https://nesc.nier.go.kr/en/html/index.do>,
445 last access: December 10, 2023). We downloaded the NO₂ measurements from the CNEMC
446 real-time air quality platform (<https://air.cnemc.cn:18007/>, last access: Jun 8, 2023). ERA-5
447 reanalysis meteorological data is obtained from the European Center for Medium-Range
448 Weather Forecasts (<https://climate.copernicus.eu/climate-reanalysis>, last access: December 8,
449 2023). CAMS forecast of meteorological and atmospheric NO₂ datasets are retrieved from the
450 CAMS Atmosphere Data Store (<https://ads.atmosphere.copernicus.eu/>, last access: December
451 8, 2023). The source codes of the GeoNet model, surface NO₂ prediction, and necessary input
452 data can be obtained from Chengxin Zhang (zcx2011@ustc.edu.cn) upon reasonable request.
453

454 **Contributions:** C.Z. implemented the GeoNet model and analyzed the data. C.L. supervised
455 the study. C.Z. wrote the manuscript with input from all co-authors.
456

457 **Competing interests:** The contact author has declared that none of the authors has any
458 competing interests.
459

460 **Acknowledgments.** This study was supported by the National Natural Science Foundation of
461 China (Nos. 42225504, 62305322, and 42375120), the National Key Research and
462 Development Program of China (Nos. 2022YFC3700100 and 2023YFC3706104), [the](#)
463 [Fundamental Research Funds for the Central Universities \(Nos. YD2090002021 and](#)
464 [WK2090000038\)](#) and the New Cornerstone Science Foundation through the XPLOER
465 PRIZE (2023-1033).
466

467 **References**

- 468 Altmann, A., Tolosi, L., Sander, O., & Lengauer, T. (2010). Permutation importance: a
469 corrected feature importance measure. *Bioinformatics*, 26(10), 1340-1347.
470 <https://www.ncbi.nlm.nih.gov/pubmed/20385727>
- 471 Andersson, T. R., Hosking, J. S., Pérez-Ortiz, M., Paige, B., Elliott, A., Russell, C., et al.
472 (2021). Seasonal Arctic sea ice forecasting with probabilistic deep learning. *Nature*
473 *Communications*, 12(1), 5124.
- 474 Baghanam, A. H., Nourani, V., Bejani, M., Pourali, H., Kantoush, S. A., & Zhang, Y. (2024).
475 A systematic review of predictor screening methods for downscaling of numerical
476 climate models. *Earth-Science Reviews*, 104773.
- 477 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X., & Tian, Q. (2023). Accurate medium-range
478 global weather forecasting with 3D neural networks. *Nature*, 1-6.
- 479 Boukabara, S.-A., Krasnopolsky, V., Penny, S. G., Stewart, J. Q., McGovern, A., Hall, D., et
480 al. (2020). Outlook for exploiting artificial intelligence in the earth and environmental
481 sciences. *Bulletin of the American Meteorological Society*, 1-53.
- 482 Campbell, P. C., Tang, Y., Lee, P., Baker, B., Tong, D., Saylor, R., et al. (2022).
483 Development and evaluation of an advanced National Air Quality Forecasting
484 Capability using the NOAA Global Forecast System version 16. *Geoscientific Model*
485 *Development*, 15(8), 3281-3313.
- 486 Chan, K. L., Valks, P., Heue, K.-P., Lutz, R., Hedelt, P., Loyola, D., et al. (2023). Global
487 Ozone Monitoring Experiment-2 (GOME-2) daily and monthly level-3 products of
488 atmospheric trace gas columns. *Earth System Science Data*, 15(4), 1831-1870.
- 489 Douros, J., Eskes, H., van Geffen, J., Boersma, K. F., Compennolle, S., Pinaridi, G., et al.
490 (2023). Comparing Sentinel-5P TROPOMI NO₂ column observations with the

491 CAMS regional air quality ensemble. *Geoscientific Model Development*, 16(2), 509-
492 534.

493 Du, S., Li, T., Yang, Y., & Horng, S. J. (2021). Deep Air Quality Forecasting Using Hybrid
494 Deep Learning Framework. *IEEE Transactions on Knowledge and Data Engineering*,
495 33(6), 2412-2424.

496 [Fino, A., Vichi, F., Leonardi, C., & Mukhopadhyay, K. \(2021\). An overview of experiences
497 made and tools used to inform the public on ambient air quality. *Atmosphere*, 12\(11\),
498 1524.](#)

499 Guarin, J. R., Jägermeyr, J., Ainsworth, E. A., Oliveira, F. A., Asseng, S., Boote, K., et al.
500 (2024). Modeling the effects of tropospheric ozone on the growth and yield of global
501 staple crops with DSSAT v4. 8.0. *Geoscientific Model Development*, 17(7), 2547-
502 2567.

503 Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., et al.
504 (2020). The ERA5 global reanalysis. *Quarterly Journal of the Royal Meteorological
505 Society*, 146(730), 1999-2049.
506 <https://rmets.onlinelibrary.wiley.com/doi/abs/10.1002/qj.3803>

507 Hong, C., Zhang, Q., Zhang, Y., Davis, S. J., Tong, D., Zheng, Y., et al. (2019). Impacts of
508 climate change on future air quality and human health in China. *Proceedings of the
509 National Academy of Sciences*, 116(35), 17193-17200.

510 Hsu, C. H., Henze, D. K., Mizzi, A. P., González Abad, G., He, J., Harkins, C., et al. (2024).
511 An Observing System Simulation Experiment Analysis of How Well Geostationary
512 Satellite Trace-Gas Observations Constrain NOx Emissions in the US. *Journal of
513 Geophysical Research: Atmospheres*, 129(2), e2023JD039323.

514 Inness, A., Aben, I., Ades, M., Borsdorff, T., Flemming, J., Jones, L., et al. (2022).
515 Assimilation of S5P/TROPOMI carbon monoxide data with the global CAMS near-
516 real-time system. *Atmospheric Chemistry and Physics*, 22(21), 14355-14376.

517 Inness, A., Flemming, J., Heue, K. P., Lerot, C., Loyola, D., Ribas, R., et al. (2019).
518 Monitoring and assimilation tests with TROPOMI data in the CAMS system: near-
519 real-time total column ozone. *Atmospheric Chemistry and Physics*, 19(6), 3939-3962.
520 <Go to ISI>://WOS:000462793200001

521 Irrgang, C., Boers, N., Sonnewald, M., Barnes, E. A., Kadow, C., Staneva, J., & Saynisch-
522 Wagner, J. (2021). Towards neural Earth system modelling by integrating artificial
523 intelligence in Earth system science. *Nature Machine Intelligence*, 3(8), 667-674.

524 Kim, J., Jeong, U., Ahn, M.-H., Kim, J. H., Park, R. J., Lee, H., et al. (2020). New era of air
525 quality monitoring from space: Geostationary Environment Monitoring Spectrometer
526 (GEMS). *Bulletin of the American Meteorological Society*, 101(1), E1-E22.

527 Kim, M., Brunner, D., & Kuhlmann, G. (2021). Importance of satellite observations for high-
528 resolution mapping of near-surface NO2 by machine learning. *Remote Sensing of
529 Environment*, 264, 112573. <Go to ISI>://WOS:000688451300002

530 Kim, S., Kim, D., Hong, H., Chang, L.-S., Lee, H., Kim, D.-R., et al. (2023). First-time
531 comparison between NO₂ vertical columns from Geostationary Environmental
532 Monitoring Spectrometer (GEMS) and Pandora measurements. *Atmospheric
533 Measurement Techniques*, 16(16), 3959-3972.

534 Kong, L., Tang, X., Zhu, J., Wang, Z. F., Li, J. J., Wu, H. J., et al. (2021). A 6-year-long
535 (2013-2018) high-resolution air quality reanalysis dataset in China based on the
536 assimilation of surface observations from CNEMC. *Earth System Science Data*,
537 13(2), 529-570. <Go to ISI>://WOS:000622997600001

538 Kuhn, L., Beirle, S., Kumar, V., Osipov, S., Pozzer, A., Bösch, T., et al. (2024). On the
539 influence of vertical mixing, boundary layer schemes, and temporal emission profiles

540 on tropospheric NO₂ in WRF-Chem—comparisons to in situ, satellite, and MAX-
541 DOAS observations. *Atmospheric Chemistry and Physics*, 24(1), 185-217.

542 Kumar, V., Remmers, J., Beirle, S., Fallmann, J., Kerkweg, A., Lelieveld, J., et al. (2021).
543 Evaluation of the coupled high-resolution atmospheric chemistry model system
544 MECO (n) using in situ and MAX-DOAS NO₂ measurements. *Atmospheric
545 Measurement Techniques*, 14(7), 5241-5269.

546 Lam, R., Sanchez-Gonzalez, A., Willson, M., Wirnsberger, P., Fortunato, M., Alet, F., et al.
547 (2023). Learning skillful medium-range global weather forecasting. *Science*,
548 382(6677), 1416-1421.

549 [Li, Y., Xing, C., Peng, H., Song, Y., Zhang, C., Xue, J., et al. \(2023\). Long-term observations
550 of NO₂ using GEMS in China: Validations and regional transport. *Science of The
551 Total Environment*, 904, 166762.](#)

552 Lin, Z., Li, M., Zheng, Z., Cheng, Y., & Yuan, C. (2020). *Self-attention convlstm for
553 spatiotemporal prediction*. Paper presented at the Proceedings of the AAAI
554 conference on artificial intelligence.

555 Liu, C., Wu, C., Kang, X., Zhang, H., Fang, Q., Su, Y., et al. (2023). Evaluation of the
556 prediction performance of air quality numerical forecast models in Shenzhen.
557 *Atmospheric Environment*, 314, 120058.
558 <https://www.sciencedirect.com/science/article/pii/S1352231023004843>

559 Makar, P., Gong, W., Milbrandt, J., Hogrefe, C., Zhang, Y., Curci, G., et al. (2015).
560 Feedbacks between air pollution and weather, Part 1: Effects on weather. *Atmospheric
561 Environment*, 115, 442-469.

562 Manisalidis, I., Stavropoulou, E., Stavropoulos, A., & Bezirtzoglou, E. (2020).
563 Environmental and health impacts of air pollution: a review. *Frontiers in public
564 health*, 14.

565 Meng, X., Liu, C., Chen, R., Sera, F., Vicedo-Cabrera, A. M., Milojevic, A., et al. (2021).
566 Short term associations of ambient nitrogen dioxide with daily total, cardiovascular,
567 and respiratory mortality: multilocation analysis in 398 cities. *bmj*, 372.

568 Nguyen, T., Brandstetter, J., Kapoor, A., Gupta, J. K., & Grover, A. (2023). ClimaX: A
569 foundation model for weather and climate. *arXiv preprint arXiv:2301.10343*.

570 Peuch, V.-H., Engelen, R., Rixen, M., Dee, D., Flemming, J., Suttie, M., et al. (2022). The
571 Copernicus Atmosphere Monitoring Service: From Research to Operations. *Bulletin
572 of the American Meteorological Society*, 103(12), E2650-E2668.

573 Sathe, Y., Gupta, P., Bawase, M., Lamsal, L., Patadia, F., & Thipse, S. (2021). Surface and
574 satellite observations of air pollution in India during COVID-19 lockdown:
575 Implication to air quality. *Sustainable cities and society*, 66, 102688.

576 Shah, V., Jacob, D. J., Li, K., Silvern, R. F., Zhai, S., Liu, M., et al. (2020). Effect of
577 changing NO_x lifetime on the seasonality and long-term trends of satellite-observed
578 tropospheric NO₂ columns over China. *Atmospheric Chemistry and Physics*, 20(3),
579 1483-1495.

580 Shu, L., Zhu, L., Bak, J., Zoogman, P., Han, H., Liu, S., et al. (2023). Improving ozone
581 simulations in Asia via multisource data assimilation: results from an observing
582 system simulation experiment with GEMS geostationary satellite observations.
583 *Atmospheric Chemistry and Physics*, 23(6), 3731-3748.

584 [Tang, K. T. J., Lin, C., Wang, Z., Pang, S. W., Wong, T.-W., Yu, I. T. S., et al. \(2024\).
585 Update of Air Quality Health Index \(AQHI\) and harmonization of health protection
586 and climate mitigation. *Atmospheric Environment*, 326, 120473.](#)

587 Van Geffen, J., Eskes, H., Compernelle, S., Pinardi, G., Verhoelst, T., Lambert, J.-C., et al.
588 (2022). Sentinel-5P TROPOMI NO₂ retrieval: impact of version v2. 2 improvements

589 and comparisons with OMI and ground-based data. *Atmospheric Measurement*
590 *Techniques*, 15(7), 2037-2060.

591 Wang, S., Zhang, M., Gao, Y., Wang, P., Fu, Q., & Zhang, H. (2024). Diagnosing drivers of
592 PM 2.5 simulation biases in China from meteorology, chemical composition, and
593 emission sources using an efficient machine learning method. *Geoscientific Model*
594 *Development*, 17(9), 3617-3629.

595 Zhang, C., Liu, C., Li, B., Zhao, F., & Zhao, C. (2022). Spatiotemporal neural network for
596 estimating surface NO2 concentrations over north China and their human health
597 impact. *Environmental Pollution*, 119510.

598 Zhang, Q.-s., & Zhu, S.-C. (2018). Visual interpretability for deep learning: a survey.
599 *Frontiers of Information Technology & Electronic Engineering*, 19(1), 27-39.

600 Zhong, S., Zhang, K., Bagheri, M., Burken, J. G., Gu, A., Li, B., et al. (2021). Machine
601 learning: new ideas and tools in environmental science and engineering.
602 *Environmental Science & Technology*, 55(19), 12741-12754.

603 Zuo, C., Chen, J., Zhang, Y., Jiang, Y., Liu, M., Liu, H., et al. (2023). Evaluation of four
604 meteorological reanalysis datasets for satellite-based PM2. 5 retrieval over China.
605 *Atmospheric Environment*, 305, 119795.

606

607