

Regional prediction of ground-level ozone using a hybrid sequence-to-sequence deep learning approach



Hong-Wei Wang ^a, Xiao-Bing Li ^b, Dongsheng Wang ^a, Juanhao Zhao ^c, Hong-di He ^{a,*}, Zhong-Ren Peng ^{d, **}

^a Center for Intelligent Transportation Systems and Unmanned Aerial Systems Applications, State Key Laboratory of Ocean Engineering, School of Naval Architecture, Ocean and Civil Engineering, Shanghai Jiao Tong University, Shanghai, 200240, China

^b Institute for Environmental and Climate Research, Jinan University, Guangzhou, 510632, China

^c Department of Computer Science, Viterbi School of Engineering, University of Southern California, Los Angeles, CA, 90089, USA

^d International Center for Adaptation Planning and Design (iAdapt), School of Landscape Architecture and Planning, College of Design, Construction and Planning, University of Florida, P.O. Box 115706, Gainesville, FL, 32611-5706, USA

ARTICLE INFO

Article history:

Received 7 July 2019

Received in revised form

22 November 2019

Accepted 21 December 2019

Available online 27 December 2019

Handling editor: Kathleen Aviso

Keywords:

Ozone prediction

Spatiotemporal correlation

Air quality monitoring network

Sequence to sequence model

Deep learning

ABSTRACT

Ozone is one of the most important greenhouse gases and air pollutants in urban areas, and has significantly negative impacts both on the climate change and human health. In addition to alert the public for health concerns, accurate ozone prediction is also crucial to understand the formation mechanisms of surface ozone episodes, and has significant implications for making emission control strategies of ozone precursors such as methane, carbon monoxide, and volatile organic compounds. However, existing methods of ozone concentration prediction could not effectively capture temporal dependency, and most neglect spatial correlations. In this study, a hybrid sequence to sequence model embedded with the attention mechanism is proposed for predicting regional ground-level ozone concentration. In the proposed model, the inherent spatiotemporal correlations in air quality monitoring network are simultaneously extracted, learned and incorporated, and auxiliary air pollution and meteorological information are adaptively involved. The hourly data collected from 35 air quality monitoring stations and 651 meteorological girds in Beijing, China are used to validate the present model. The results demonstrate that the spatiotemporal correlations in the monitoring network present enormous advantages for the regional ozone prediction. Auxiliary data and time lags matching day-of-week or diurnal periods of ozone are confirmed to benefit the improvement of prediction accuracy. Monitoring stations in urban areas exhibit better prediction performances than stations in remote areas. The addressed model outperforms the baseline models, and is proven to have excellent performance in all monitoring station categories of Beijing and different months with significant disparity of ozone concentrations.

© 2019 Elsevier Ltd. All rights reserved.

1. Introduction

Tropospheric ozone is an important greenhouse gas and air pollutant, and has detrimental impacts both on climate change and human health. In recent years, great efforts have been made to reduce surface ozone levels by implementing stringent emission control measures of ozone precursors (NOx and VOCs)(Cheng et al., 2019), which is also closely associated with the reduction budget of carbon emissions around the world (Ye et al., 2019; Fasihi et al.,

2019; Yusuf et al., 2019). However, severe ozone pollution episodes still frequently occurred over urban areas, especially in developing countries like China (Su et al., 2017). In comparison to aerosol pollution, people usually do not give many concerns on ozone pollution issues due to insignificant sensory stimulations and take necessary measures in advance to reduce their exposure levels to ozone pollution (He et al., 2017). Thus, it is extremely important to provide accurate forecasts in urban areas to alert the public when ambient ozone concentration will exceed preset thresholds (Gradišar et al., 2016). Compared with the prediction of ozone concentration only at monitoring stations (Awang et al., 2015; Gao et al., 2018), predicting the regional ozone concentration suggests higher spatial resolution and directly reflects the

* Corresponding author.

** Corresponding author.

E-mail addresses: hongdihe@sjtu.edu.cn (H.-d. He), zpeng@ufl.edu (Z.-R. Peng).

spatial distribution change of ozone, which assists to the pollution control and the emission reduction of ozone. Therefore, conducting accurate regional prediction for ground-level ozone concentration is of greatly significance for both public health protection and emission reduction of greenhouse gases.

Generally, ozone prediction models can be classified into two types, including deterministic models and empirical models (Li et al., 2017; Wen et al., 2019). Deterministic models refer to numerical simulation models that incorporate various physical and chemical mechanisms related to the emission, transformation, dispersion and transport of air pollutants (Sharma et al., 2016), such as the Weather Research Forecasting (WRF) model. In the WRF model, different combinations of macro-scale physical process parameterization schemes (e.g. planetary boundary layer schemes and land surface modelling schemes) have significant impacts on the accuracy of model results (Cai et al., 2016). Therefore, deterministic models are favorable for the analysis of air quality problems at large spatial scales, but have limited prediction accuracies at micro scales like urban areas (Gradišar et al., 2016). In empirical models, statistics and machine learning methods have been widely developed for general environmental predictions (e.g. ozone and aerosol pollution, renewable energy utilization and carbon emission reduction). Commonly used statistics models include Multiple Linear Regression (MLR)(Awang et al., 2015), Multivariate Adaptive Regression Splines (Srinivas et al., 2019) and Principal Component Analysis (PCA)(Huang et al., 2019). Recently, machine learning models are rapidly developed to improve the performances of general environmental predictions due to their nonlinear mapping ability, such as Random Forest (Feng et al., 2019; Srinivas et al., 2019), Bayesian Network (Mehdipour et al., 2018), Extreme Gradient Boosting (XGBoost)(Ma et al., 2019) and Support Vector Machine (SVM) (He et al., 2018; Mehdipour et al., 2018; Mehdipour and Memarianfar, 2019; Sumathi and Manivannan, 2019). These machine learning models take nonlinearity into account, and present strong abilities to capture temporal features of air pollutants and renewable resources.

As a form of machine learning, deep learning models are based on artificial neural networks with more complicated architectures, and have grown up to efficient approaches for handling multi-dimensional data due to flexible model structure, strong generalization, and powerful learning ability (Li et al., 2019; Hao et al., 2019). Recently, several deep learning methods have achieved outstanding performance in general environmental prediction issues (e.g., air pollution, carbon emission and renewable energy), such as Deep Neural Network (DNN)(Gao et al., 2018; Acheampong and Boateng, 2019), Stacked Autoencoder (SAE)(Bai et al., 2019), Recurrent Neural Network (RNN)(Feng et al., 2019), Long Short-Term Memory Recurrent Network (LSTM)(Li et al., 2017, Li et al., 2019; Zhou et al., 2019; Huang et al., 2019) and Gated Recurrent Network (GRU)(Athira et al., 2018). However, RNN-based models (i.e., LSTM and GRU) neglect neighboring factors in air quality monitoring network and are incapable of characterizing spatial correlations (Wen et al., 2019). Based on the idea of Convolutional Neural Network (CNN), CNN-LSTM (Wen et al., 2019) and Conv-LSTM (Shi et al., 2015) are employed to conduct temporal predictions considering spatial features in grid-structured data (e.g., images). These CNN-based models can effectively capture the spatiotemporal dependencies, but have limitations to handle long-term input data because of simply stacking multiple layers of LSTM units (Hao et al., 2019). In recent years, the emergence of Sequence to Sequence learning (Seq2Seq) (Cho et al., 2014) provides a more flexible and extendable framework on temporal modelling. However, the seq2seq model suffers from a potential issue that all historical data need to be compressed into an internal vector with fixed length (Bahdanau et al., 2015). This may limit the ability of the

seq2seq model coping with long sequences (Cho et al., 2014; Bahdanau et al., 2015). Thus, attention mechanism methods (Bahdanau et al., 2015) are proposed to resolve the drawbacks of simple Seq2Seq and adept at modelling long-range dependencies (Bahdanau et al., 2015; Hao et al., 2019), such as long-term processes of chemical and physical mechanisms in air pollution episodes (Liu et al., 2019).

In this study, we aim to propose a hybrid sequence to sequence model embedded with the attention mechanism (HSA-Net) to predict the ozone concentration at a city scale. In order to fully consider the spatiotemporal correlations, the neighboring factors in the monitoring network are adaptively learned by the proposed model, and the spatial information of the entire monitoring network is well extracted and integrated. The dataset collected in Beijing, China is used for training, validating and testing the proposed model, and the auxiliary air pollution data (SO_2 , NO_2 , $\text{PM}_{2.5}$, PM_{10} , and CO), meteorological data and weather forecast data are also employed. The novelties of the presented study are summarized as follows: (1) The advantage of incorporating spatiotemporal correlations in the model is quantitatively explored and evaluated. (2) The spatiotemporal dependencies of air quality monitoring network are simultaneously captured. (3) The auxiliary air pollution and meteorological data are integrated into the model. (4) The ozone prediction performances of the proposed model are estimated monthly and geographically by different categories of monitoring stations.

2. Data and methods

The details of data are described in this section, then we introduce the methodology of RNN-series methods, Seq2Seq model and attention mechanism, and finally the experiment settings and network architecture are presented.

2.1. Data description

In this study, Beijing, the capital of China, is selected as the study area. The research area of interest approximately ranges from E 115° to E 118° and N 39° to N 41° . Hourly mean concentrations of six regulatory air pollutants (in China) including O_3 ($\mu\text{g}/\text{m}^3$), SO_2 ($\mu\text{g}/\text{m}^3$), NO_2 ($\mu\text{g}/\text{m}^3$), $\text{PM}_{2.5}$ ($\mu\text{g}/\text{m}^3$), PM_{10} ($\mu\text{g}/\text{m}^3$), and CO (mg/m^3) are obtained from 35 air quality monitoring stations from January 1, 2017 to May 31, 2018. Hourly gridding meteorological data (21×31 points) in the same period are processed by the Weather Research and Forecasting (WRF) model, and have a grid spacing of 5 km. Meteorological parameters including temperature, air pressure, relative humidity, wind speed, and wind direction are selected as the principal meteorological features due to their close relationships with the change of ground-level ozone concentration (Li et al., 2017; Wen et al., 2019; Bai et al., 2019). Fig. 1 shows the geographical locations of the air quality monitoring stations. The monitoring stations are officially divided into four categories according to their locations and surrounding environments. The stations labelled by 1–12, 13 to 17, 18 to 28 and 29 to 35 refer to urban, traffic, rural and cross reference monitoring stations, respectively. Urban and rural monitoring stations (Cheng et al., 2019) are used to assess the average condition and variation of air quality in urban and suburban environments, respectively; Traffic monitoring stations are used to monitor the impact of road traffic sources on ambient air quality; Cross reference monitoring stations are used to characterize the regional environmental background level and reflect the transmission of pollution within the area.

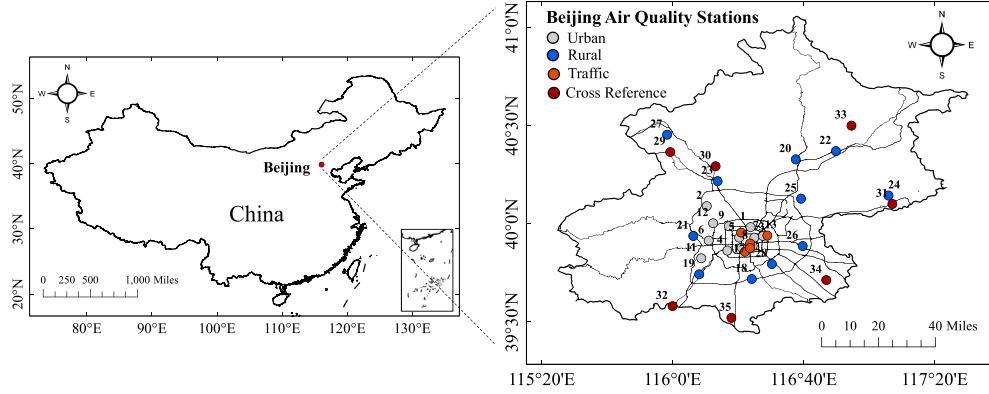


Fig. 1. Geographical locations of the air quality monitoring stations in Beijing.

2.2. Hybrid sequence to sequence model with attention (HSA-Net)

In this section, we first elaborate the methodologies involved in the HSA-Net model, from LSTM, GRU, Bidirectional LSTM (Bi-LSTM) to Seq2Seq architecture with attention mechanism, then introduce the structure details of the HSA-Net model, as shown in Fig. 2. In

temporal modelling, the encoder is built by Bi-LSTM network that can read input data in both forward and backward directions, and the decoder is developed by GRU. The combination of the Bi-LSTM encoder and the GRU decoder is widely-used in Seq2Seq learning (Bahdanau et al., 2015; Liu et al., 2019). In spatial modelling, the Deep Neural Network (DNN) is implemented for extracting and

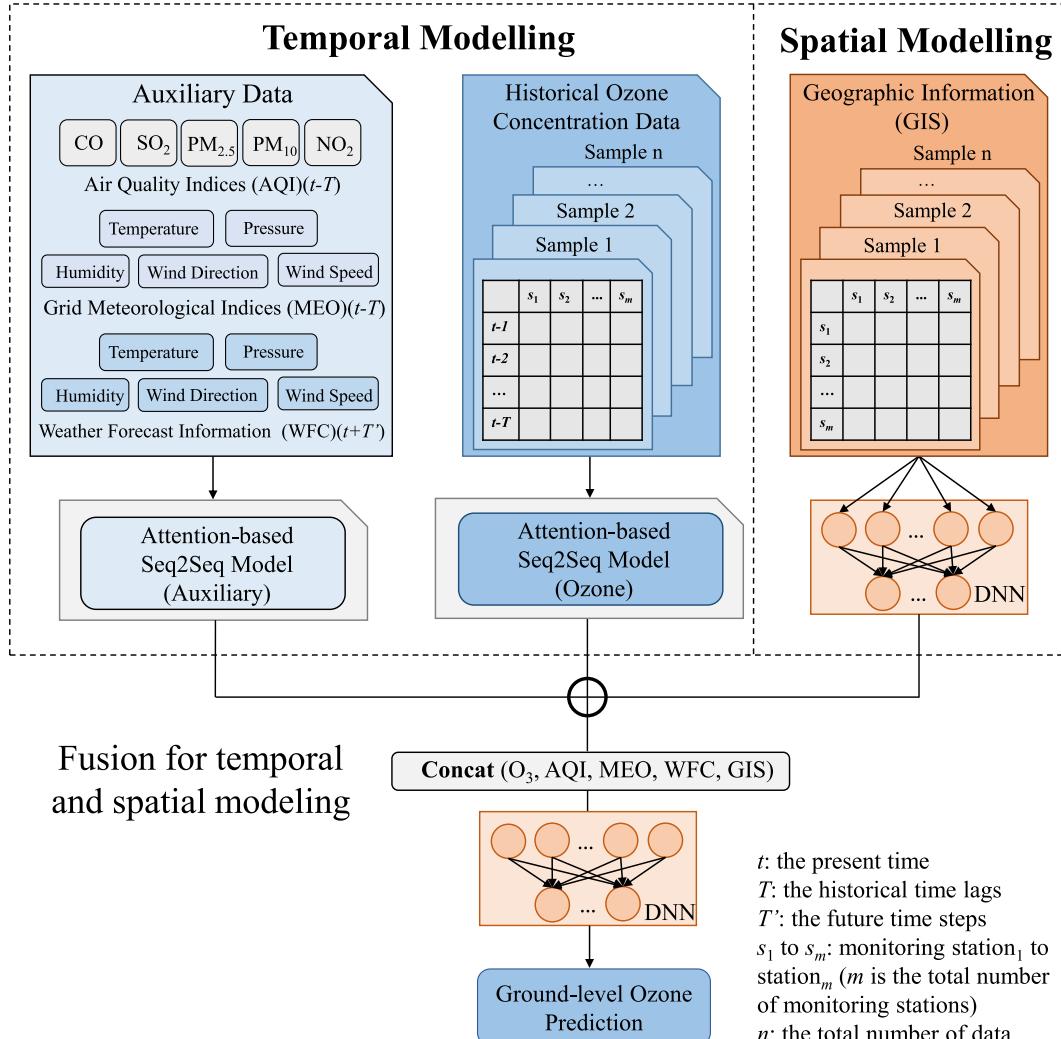


Fig. 2. The framework of the hybrid seq2seq model with attention (HSA-Net).

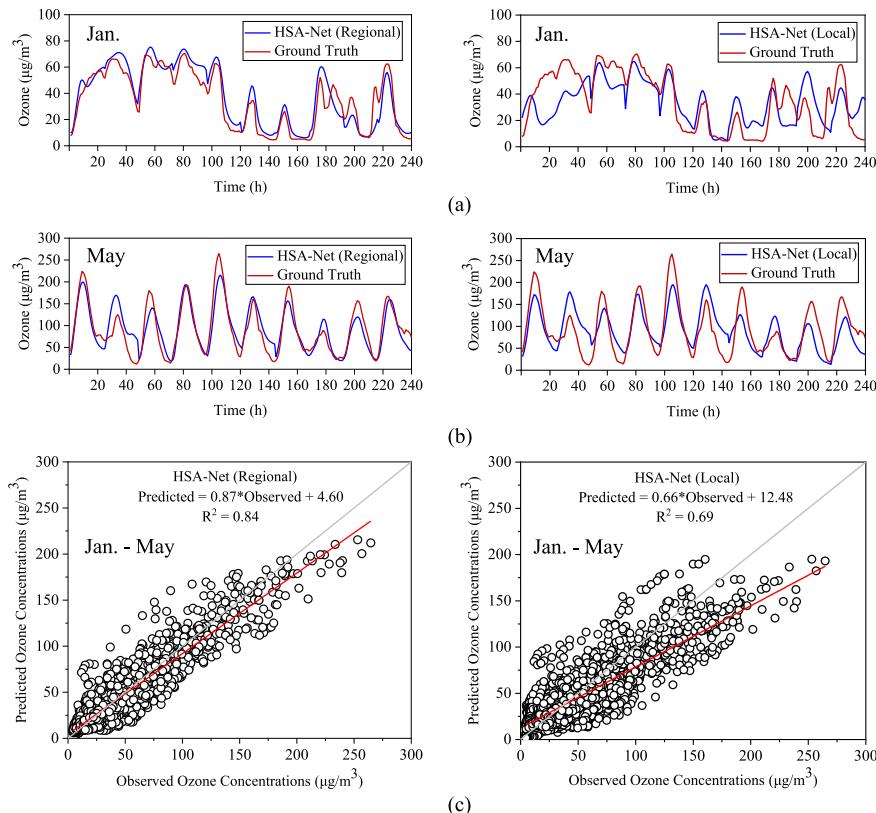


Fig. 5. Observed and predicted ozone concentrations of the regional pattern and the local pattern in the test dataset: (a) time series in January, (b) time series in May, (c) scatter plots from January to May.

increase from January to May in the test dataset. The forecasting results indicate that the regional model outperforms the local model both in the months with lower ozone concentration (January) and higher ozone concentration (May). In Fig. 5 (c), the R^2 values between predictions and observations demonstrate that 84% and 69% of the explained variance are captured by the regional pattern and the local pattern, respectively. The scatter plots also demonstrate that predictions are more consistent with the observations in the regional pattern. In Fig. 6 (a), the values of RMSE and r indicate that the regional model has dominated superior in prediction accuracy for most of the monitoring stations, except the station 12. In Fig. 6 (b), the regional model wins the comparisons in terms of the indices of RMSE and r for the four categories of monitoring stations. The results demonstrate that the adoption of spatiotemporal correlation can notably improve the prediction performance.

3.2. Performance evaluation of ground-level ozone predictions

In this section, the prediction performances of the HSA-Net model are carefully estimated. we first compare the proposed model with baseline methods, then evaluate the ability of multi-scale predictions, and finally examine the effects of auxiliary input variables.

3.2.1. Multi-scale predictions

This experiment evaluates the predictive ability of the HSA-Net model when the predictions of different lengths are required. In order to confirm the key historical period that significantly affects the prediction of ozone concentration, the influence of different input time lags is examined. We estimate the accuracy of

predictions when the historical time lags are set as 24, 48, 120 and 168 h, respectively. In Table 2, when the input time lags increase from 24 h (1 day) to 168 h (7 days), prediction accuracy first drops and then grows up. However, the running time notably increases from 2026.57 s to 8518.81 s. This result indicates that the time lags corresponding to the diurnal variation period (24 h) of ozone concentrations (Su et al., 2017) and the weekend effects (168 h) of ozone concentrations caused by traffic activities (Gao, 2007) can achieve a better prediction performance, but longer input data result in longer computing time.

The multi-scale prediction performance is assessed under the scenarios that the lengths of prediction are set as 6, 12, 24, and 48 h, respectively. In Table 3, the result indicates that the prediction ability drops continuously from 17.75 to 23.23 in RMSE when the lengths of prediction increase from 6 h to 48 h. Meanwhile, the time cost rises notably from 692.45 s to 3498.88 s. The result demonstrates that longer-term predictions are more difficult and consume much more computing time.

3.2.2. Comparison of models

Two comparisons under the 24-h predictions are separately conducted for further performance estimation of the HSA-Net model in ozone predictions. The first comparison is implemented between the HSA-Net model and conventional models within the hybrid architecture. The second one is conducted between the HSA-Net model and the models without the hybrid architecture (e.g. Seq2Seq, Liu et al., 2019; LSTM, Li et al., 2017; GRU, Athira et al., 2018). All the results are presented in Table 4.

In the first comparison, the HSA-Net model has the optimal performance (22.08 in RMSE), and the hybrid Seq2Seq model and the hybrid Bi-LSTM model follow behind closely. There is an

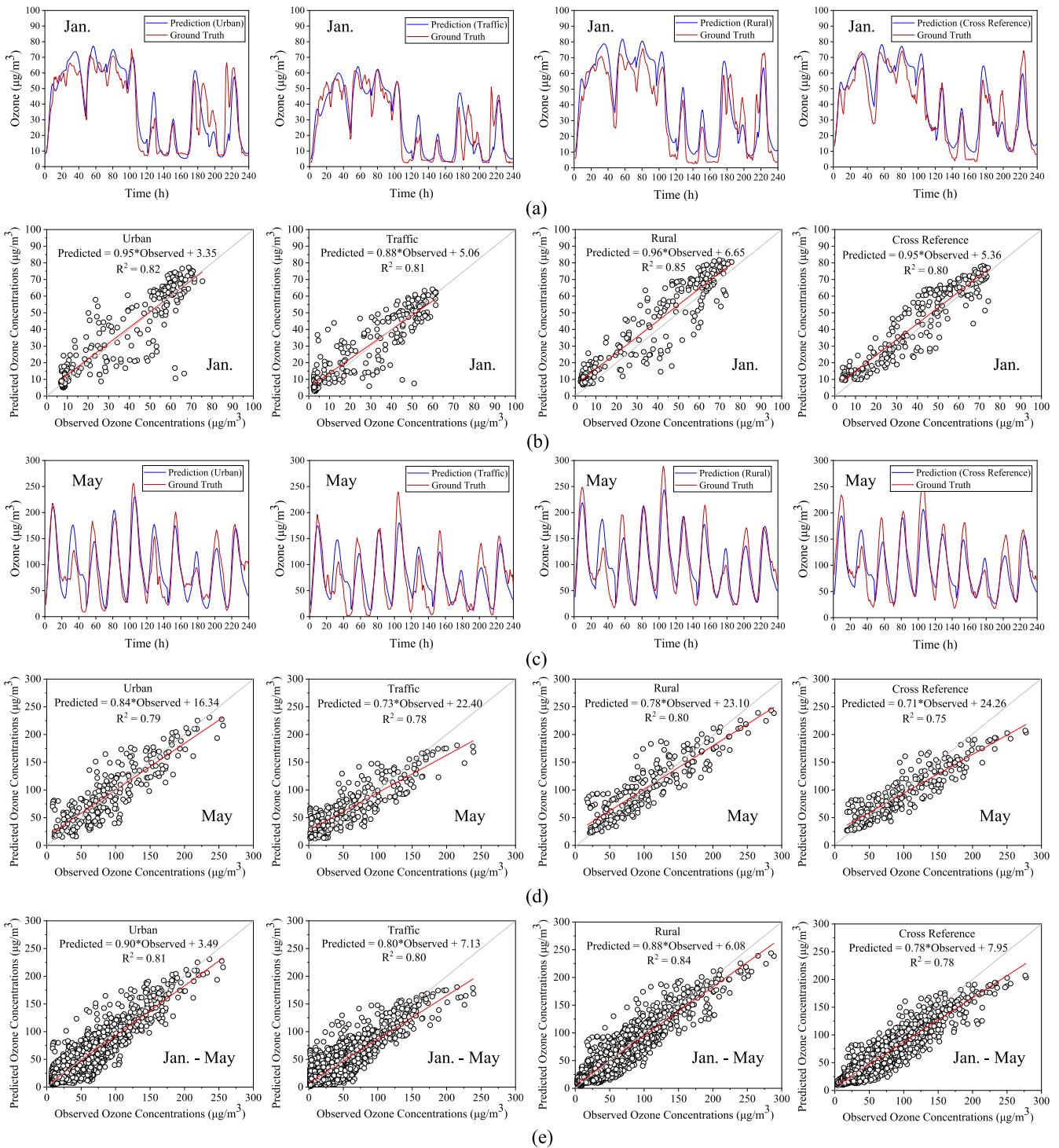


Fig. 7. Observed and predicted ozone concentrations of the HSA-Net model in the urban, traffic, rural and cross reference monitoring stations: (a) time series in January, (b) scatter plots in January, (c) time series in May, (d) scatter plots in May, (e) scatter plots from January to May.

monitoring stations are located in the downtown area of Beijing, where the ground-level ozone concentrations present much stronger diurnal variations due to the impacts of traffic emissions (Gao, 2007). This phenomenon benefits the deep learning models to learn the pattern of ozone generating and dispersing processes, and to achieve better performance. Cross reference monitoring stations are located in remote areas, and the diurnal variation

patterns of ozone concentrations at these monitoring stations are affected by local photochemical production and transport of ozone from the downtown Beijing (Cheng et al., 2019). Thus, the special temporal patterns in such stations are quite challenging to predict.

Fig. 9 presents the spatial distribution of mean daily ozone concentrations and mean 24-h prediction performance in January, 2018, May, 2018 and the whole test dataset. The observed ozone

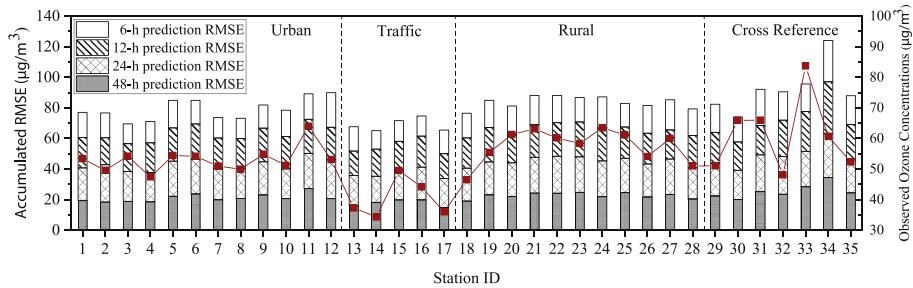


Fig. 8. Accumulative RMSE ($\mu\text{g}/\text{m}^3$) of the HSA-Net prediction (bars) and observed mean ozone concentrations (red line) for the 35 monitoring stations in the test dataset.

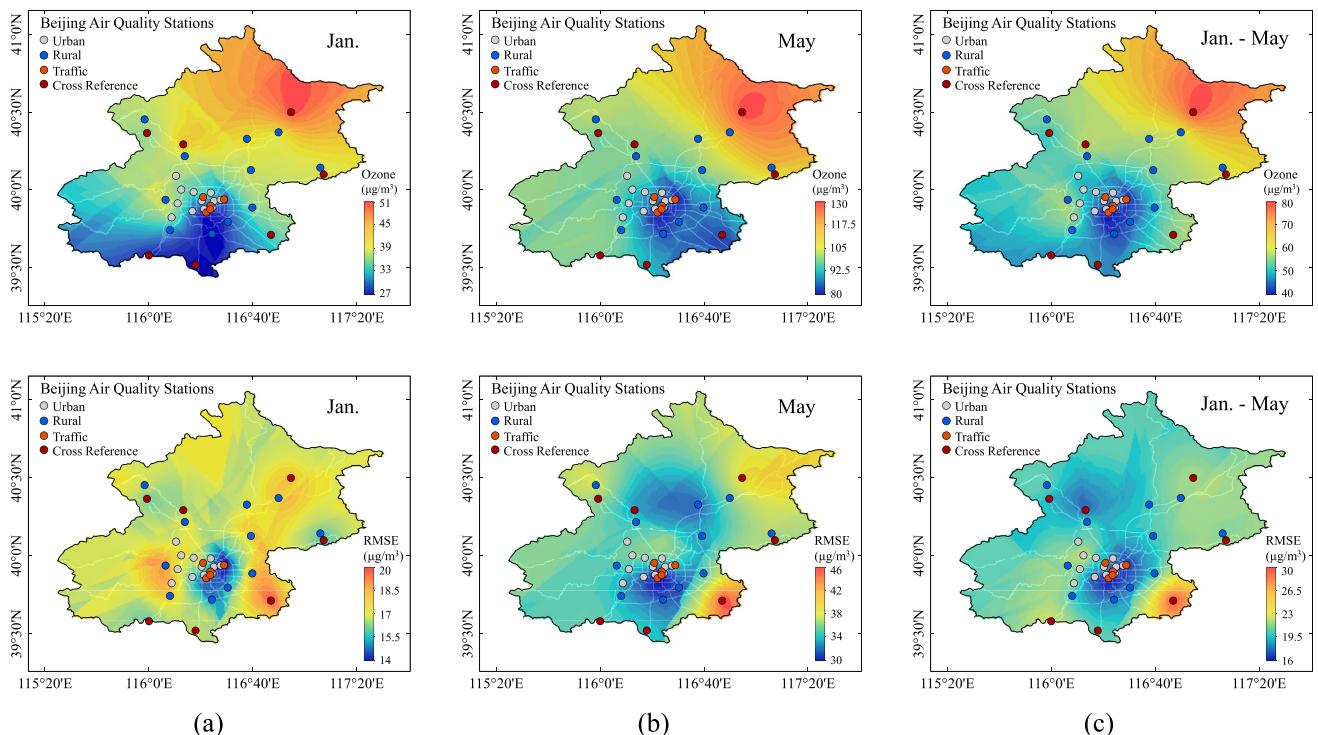


Fig. 9. Spatial distributions of observed mean daily ozone concentrations and mean RMSE ($\mu\text{g}/\text{m}^3$) of 24-h predicted ozone concentrations in the HSA-Net model for the 35 monitoring stations: (a) in January, (b) in May, (c) from January to May.

concentrations grow up notably from January to May, and are relatively lower over the urban areas. The spatial distribution patterns of ozone over Beijing are predominantly determined by meteorological and orographic factors (Cheng et al., 2019). Southeast, southwest and south winds prevailed over the period of January through May, 2018. The prevailing winds could blow ozone produced over urban areas to downwind areas in the north of Beijing. In addition, mountains are located to the north and northwest of Beijing, leading to evident accumulation of ozone to the north of Beijing. The RMSE of the predicted ozone concentrations also presents a significant increase from January to May. Additionally, the values of RMSE over the central urban area keep lower than the remote areas, and even show a more stable distribution in May. This result demonstrates that the prediction performance in strong spatial-correlated stations (such as traffic stations) is better than that in scattered stations (such as cross reference stations).

4. Conclusions

In this study, we propose a novel hybrid sequence to sequence

model embedded with attention mechanism (HSA-Net) to predict the ground-level ozone concentration, fully considering the spatiotemporal correlations in the air quality monitoring network. The Attention-based Seq2Seq model and the deep neural network are both employed to simultaneously capture the spatiotemporal dependencies. In order to improve prediction performance, the neighboring factors and the spatial information in the air quality monitoring network are adaptively incorporated and learned in the model. The auxiliary air pollution and meteorological data are also integrated, which are related to the change of ground-level ozone concentration but were neglected in previous studies. The experiment results demonstrate that the proposed model can effectively predict the regional ground-level ozone concentration and outperform the competitors in terms of accuracy. Several meaningful findings are summarized as follows:

- (1) Considering the spatiotemporal correlations in the air quality monitoring network presents enormous advantages for regional ozone prediction.
- (2) The supplement of auxiliary data can notably improve the prediction ability.

