Implications of VOC Oxidation in Atmospheric Chemistry: Development of a Comprehensive AI Model for Predicting Reaction Rate Constants

Xin Zhang^{1,2#}, Jiaqi Luo^{1,2#}, Wenxiao Pan^{1,2}, Qiao Xue^{1,2}, Xian Liu^{1,2*}, Jianjie Fu^{1,2,3*}, Aiqian Zhang^{1,2,3}, Guibin Jiang^{1,2,3}

¹State Key Laboratory of Environmental Chemistry and Ecotoxicology, Research Center for Eco-Environmental Sciences, Chinese Academy of Sciences, Beijing 100085, P. R. China

² College of Resources and Environment, University of Chinese Academy of Sciences, Beijing 100190, P. R. China

³ School of Environment, Hangzhou Institute for Advanced Study, University of Chinese Academy of Sciences, Hangzhou 310012, P. R. China

#These authors contributed equally to this work and should be considered co-first authors.

Correspondence to: Xian Liu (xianliu@rcees.ac.cn); Jianjie Fu (jjfu@rcees.ac.cn)

Abstract. Volatile Organic Compounds (VOCs) significantly influence global atmospheric chemistry through oxidative reactions with oxidants. These reactions produce key precursors to the formation of atmospheric fine particulate matter (PM_{2.5}) and ozone (O₃), which in turn play a crucial role in regulating O₃ pollution and reducing PM_{2.5} concentrations. With the increasing diversity of VOCs, the need for advanced modeling techniques to accurately estimate the atmospheric oxidation reaction rate constants (k_i , where $i \in \{\text{*OHOH}, \text{*CICI}, \text{NO}_3, \text{ or O}_3\}$) has become more urgent. Here we introduce Vreact, a Siamese message passing neural networks (MPNN) architecture that jointly models VOC–oxidant reactivity. The model simultaneously predicts $\log_{10}k_i$ values and achieves a mean squared error (MSE) of 0.299 and a coefficient of determination (R²) of 0.941 on the internal test set. This framework overcomes the single-oxidant constraint of traditional models, enabling unified and scalable prediction of VOC oxidation kinetics across multiple oxidants. An interactive web tool (http://vreact.envwind.site:8001) is provided to facilitate non-expert access to reactivity screening. Vreact offers valuable insights into the formation and evolution of atmospheric pollutants, and serves as a critical resource for developing effective control and emission strategies, ultimately supporting global efforts to mitigate air pollution and improve public health.

25 1 Introduction

The rapid advancement in data-driven methodologies has revolutionized various fields, such as protein structure prediction (Abramson et al., 2024), molecular generation (Zhang et al., 2023), organic reaction prediction (Burés and Larrosa, 2023), and bioinformatics (Theodoris et al., 2023). Environmental challenges, particularly those associated with atmospheric chemistry and climate change (Chen et al., 2024; Kubečka et al., 2023; Qiu et al., 2023; Zhao et al., 2025), have also benefited from these innovations. As pollutants evolve under both anthropogenic and natural influences, the understanding of their chemical and physical properties has become increasingly vital for addressing global air quality and climate issues. Volatile Organic

Compounds (VOCs) are organic chemicals that readily vaporize at ambient temperature, contributing significantly to the complexity of atmospheric processes. Sources of VOCs are both natural and anthropogenic, with human activities such as industrial production, petrochemical processing, and vehicle exhaust contributing to the emission of a variety of VOCs. Additionally, biosphere sources, such as plants and forests, release compounds like isoprene and monoterpenes, which further complicate atmospheric VOC dynamics (Qin et al., 2021; Sindelarova et al., 2014). These highly reactive VOCs drives critical atmospheric reactions, such as the formation of ozone and secondary organic aerosols (SOA), and significantly contribute to environmental pollution. For instance, VOCs interact with nitrogen oxides (NO_x) and radicals to form tropospheric O₃ and SOA (Finlayson-Pitts and Pitts, 1997; Hallquist et al., 2009; Han et al., 2018; Zhang et al., 2020; Ziemann and Atkinson, 2012). The role of VOCs in the formation of secondary pollutants such as PM_{2.5} (Huang et al., 2014; Zhao et al., 2015) and O₃ is a growing concern due to the adverse impacts on human health (Kamarrudin et al., 2013), including respiratory diseases, cardiovascular conditions, and overall mortality. The dynamic interactions between VOCs and atmospheric oxidants determine the persistence and transformation of these pollutants, which in turn influence their contribution to global haze, photochemical smog, and acid deposition.

VOCs undergo degradation and removal from the troposphere through diverse mechanisms driven by atmospheric oxidants. During the daytime, *OHOH radicals serve as the primaryprimarily oxidants, facilitating rapid VOC oxidation. At night, however, the concentration of *OHOH decreases sharply due to the lack of photochemical reactions, shifting the dominant oxidation pathways to NO₃ radicals and O₃. The reaction rates of VOCs with *OHOH are approximately 30 times faster than those with NO₃ radicals, with NO₃ radicals, significantly influencing the spatial and temporal variation of the atmosphere's self-cleaning capacity and the formation of organic aerosols (Palmer et al., 2022; Zha et al., 2023). For example, regions with high isoprene concentrations often reflect differences in its reaction products and rates with *OHOH and NO_x rather than solely high emissions (Wells et al., 2020). Additionally, the structural diversity of VOCs determines their reaction mechanisms, influencing reaction rates. Highly ___reactive compounds such as alkenes, multi-substituted aromatics, and phenols exhibit higher reaction rates, whereas alkanes, alkyl nitrates, and ketones demonstrate relatively low reactivity (Ziemann and Atkinson, 2012). These variations underscore the significance of atmospheric oxidation reaction rates as key indicators of the persistence of organic pollutants in the atmosphere. Accurate assessment of these rates is essential for understanding the fate of VOCs, elucidating SOA formation processes, and addressing global challenges related to PM_{2.5} and ozone development.

Given their importance, accurately predicting the atmospheric oxidation rates of VOCs is critical for understanding their persistence, transformation, and contribution to secondary pollutant formation. Traditionally, such predictions have relied on experimental kinetic modeling methods and computational methods (e.g., quantum-chemistry (OC) and quantitative structure-

experimental kinetic modeling methods and computational methods (e.g., quantum-chemistry (QC) and quantitative structure-activity relationship (QSAR) approaches) (Basant and Gupta, 2018; Liu et al., 2021). Experimental methods involve tracking reactant and product concentrations using techniques like chemical ionization mass spectrometry (CIMS), followed by kinetic fitting to determine Arrhenius parameters (Logan, 1982; Wells et al., 1996). However, these methods are time-consuming and cover only a narrow subset of atmospheric VOCs. QC approaches combine ab initio or density-functional theory calculations with transition-state theory (TST), canonical or variational TST to obtain temperature-dependent rate constants (Canneaux et

al., 2014; Liu et al., 2021; Meana-Pañeda et al., 2024). While QC methods offer detailed mechanistic insight, their computational cost scales steeply with molecular size and conformational complexity, limiting routine application to large numbers of VOCs. However, traditional computational methods have shortcomings such as high computational complexity and low efficiency. As a more scalable alternative, QSAR model leverage molecular descriptors and statistical learning, and it has become one of the important methods for evaluating reaction rate constants. (Meana Pañeda et al., 2024)(Canneaux et al., 2014; Liu et al., 2021)QSAR models offer a scalable alternative by leveraging molecular descriptors and statistical learning. Early Previous Notable examples include AOPWINTM module integrated in US EPI SuiteTM software, which applies Partial Least Squares (PLS) regression to 109 gas-phase reaction with hydroxyl radicals (Atkinson, 1986, 1987; Kwok and Atkinson, 1995), and later expansions using a broader dataset (Öberg, 2005). Some models have also incorporated machine learning algorithms such as multiple linear regression (MLR) (Liu et al., 2020, 2022) for predicting reactions with NO₃ and •OHOH and artificial neural networks for predicting reactions with O₃ (Fatemi, 2006). Despite their utility, these models generally rely on predefined descriptors and are typically limited to reactions with a single type of oxidant, which constrains the scalability of the model. Recent advances in deep learning (DL), particularly graph neural networks (GNN), have improved molecular representation by learning features directly from molecular graphs. This enables more flexible and accurate prediction of chemical properties without requiring predefined descriptors. GNNs have been successfully applied in atmospheric chemistry and other fields tasks, such as in predicting vapor pressures with GC2NN (Krüger et al., 2025) and modeling reaction rate constants involving with *OHOH using GAT-GIN hybrid architectures (Huang et al., 2024). However, like traditional models. these GNN-based frameworks have been developed for single-molecule systems and thus fall short in capturing the complexity of multi-molecule reactions in real environments. In contrast, the atmosphere involves competing and sequential reactions between VOCs and multiple oxidants—OHOH, NO_X, OHOH, and O₃—depending on time of day, region, and chemical conditions. This multiplicity underscores the urgent need for models that can simultaneously learn and predict VOC reactivity across multiple oxidants. To meet this need, message passing neural networks (MPNN) offer a powerful framework (Gilmer et al., 2017). MPNNs propagate information across molecular graphs, capturing both atomic-level features and topological context. Extensions of MPNN, such as the communicative GraphRXN -(Li et al., 2023) and directed MPNN Chemprop (Heid et al., 2024), have shown promise in learning reactivity across multiple reactants. They extract the interaction features of chemical reactions in depth, rather than performing simple reactant concatenating. Yet, their application has largely focused on synthesis or materials chemistry, not atmospheric multiphase oxidation.

This study addresses this gap by proposing Vreact, a novel Siamese MPNN architecture capable of jointly modelling reactions between VOCs and four major atmospheric oxidants. Unlike previous models that treat each oxidant independently, Vreact processes VOC-oxidant pairs in a unified framework, it learns representations from the molecular graphs of VOCs and oxidants through the MPNN, and encodes their interactions via feature aggregation. This design enables the model to accept arbitrary VOC-oxidant combinations and simultaneously predict reaction rate constants k_i (where $i \in \{\text{*OHOH}, \text{*CICI}, \text{NO}_3, \text{or O}_3\}$). Compared to traditional and simple single-oxidant prediction models, Vreact shows significantly improved performance, achieving higher accuracy, stronger interpretability and wider scalabilityachieving higher accuracy and and

broader generalizability across multiple oxidants. Furthermore, based on the flexibility of the DL architecture, the designed the model's interaction module captures atomic-level interaction patterns, providing mechanistic insights into VOC oxidation process *via* interpretable interaction weight matrices. Applying Vreact to 447 atmospheric VOCs not included in the training data revealed a wide distribution of oxidation reactivities and confirmed that alkenes and aromatics exhibit higher reactivity, acting as key precursors for ozone and SOA formation.

2 Methods and Data

105

110

115

120

125

130

2.1 Collection and Preprocessing of Reaction Rate Constant Dataset

The VOCs reaction rate constant dataset compiled by McGillen et al. is utilized in the study, which includes gas-phase reaction rate constants of natural atmospheric VOCs, halocarbons, and their degradation products with •OHOH, •ClCl, NO₃ radicals, and O₃, within a temperature range of 250-370K (McGillen et al., 2020). Under thermodynamic standard conditions at 298K, a total of 2802 gas-phase reaction rate constant data points were obtained, encompassing 1586 VOCs and 4 oxidants. This dataset includes k_i values for 1363 VOCs with •OHOH, 735 VOCs with •ClCl, 393 VOCs with NO₃ radicals, and 311 VOCs with O₃. Due to the wide range of reaction rate constants k_i in the dataset $(1.460 \times 10^{-21} \sim 7.550 \times 10^{-10} \text{cm}^3/(\text{molecule·s})$, S.D.=±1.040×10⁻¹⁰), the data were log-transformed to $\log_{10}k_i$ to reduce skewness and mitigate the influence of outliers on the model. To ensure a balanced distribution of each type of oxidant in the training, validation, and internal test sets, the dataset was divided using stratified random sampling into training, validation and internal test sets in an 8:1:1 ratio (Table S1)). Combinations of the same VOC with different oxidants may appear across the training, validation, and internal test sets.

2.2 Construction and Training of the Vreact Model

All VOCs and oxidant molecules were converted into graphs G(V, E) (Text S1). The generated molecular graph G includes ten types of atomic information for each non hydrogen atom, such as element type, chirality, and atomic hybridization type, as well as four types of bond information, including bond type and conjugation (Table S2). A Siamese MPNN architecture Vreact, was designed to simultaneously accept input features of VOCs and oxidant molecules (Fig. 1). The model takes the SMILES of VOCs and oxidants as input and primarily includes a VOC molecular graph representation layer and a MPNN layer, an oxidant molecular graph representation layer and MPNN layer, an interaction layer, and a prediction layer. The molecular graph G(V, E) encoding layers of VOCs and oxidants containing node feature matrix X and edge feature matrix X, which learn molecular properties through the MPNN layer (Gilmer et al., 2017). The MPNN forward propagation process consists of two phases: Message Passing Phase and Readout Phase and generates molecular feature tensors A for VOCs and B for oxidants. Subsequently, the interaction layer transforms the molecular features A of VOCs and B of oxidants into tensors A_I and B_I of the same shape and concatenates them into tensor B. Reaction rate constants are determined not only by the molecular structure of the reactants but also by the interactions between the reactants. The interaction feature tensor I is dot-

multiplied with A to obtain the oxidant-affected VOC feature tensor A'; similarly, it is dot-multiplied with B to obtain the VOC affected oxidant feature tensor B'. These operations embed the learned interaction features into the molecular structure features, providing a more comprehensive representation of the chemical reaction mechanisms between the two reactants. The prediction phase is composed of a pooling layer and three fully connected layers. The pooling layer uses the Set2 Set method to achieve global average pooling, and the fully connected layers map the input features to the final predicted values ($log_{10}k_i$). More details can be found in Text S2.

During model training, Adaptive Moment Estimation (Adam) (Kingma and Ba, 2017) was employed to address the fixed learning rates issue in traditional gradient descent methods. Adam adaptively adjusts the learning rate of each parameter using first order moment estimates (mean of the gradients) and second order moment estimates (exponentially moving average of the uncentered variance of the gradients), aiding in rapid model convergence. Bayesian optimization was utilized for hyperparameter tuning, which included the initial learning rate of the optimizer (*Ir*), batch size, L2 regularization parameter (*weight decay*), dropout rate (*p*), and MPNN time steps (*T*) (Text S3). After identifying the optimal hyperparameter combination (Table S3) on the validation set, and the best model was saved. The predictive performance of the model was assessed using Mean Squared Error (MSE), Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and coefficient of determination (R²) (Text S4). For more information on the model implementation, please refer to Text S5.

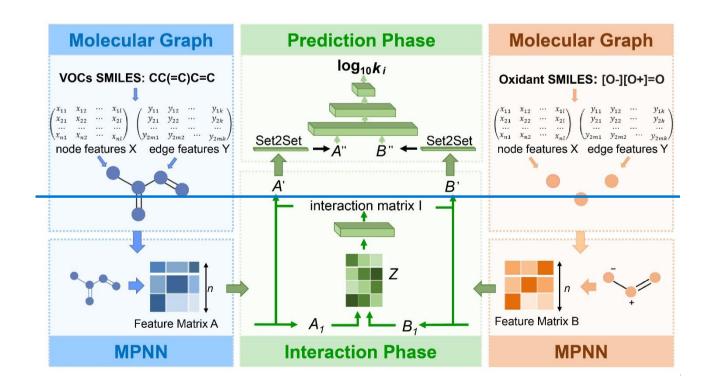


Figure 1. Schematic of the Vreact Architecture. SMILES of VOCs and oxidants are converted into molecular graphs, where nodes represent atoms and edges represent bonds. Atomic and bond features form matrices X and Y. Using a Siamese MPNN architecture, the Vreact model processes these features through separate MPNN layers for VOCs and oxidants. The final prediction layer outputs log₁₀k_i, incorporating both molecular and interaction features. 2.2 Construction and Training of the Vreact Model

All VOCs and oxidant molecules were converted into graphs G(V, E) (Text S1). The generated molecular graph G includes ten types of atomic information for each non-hydrogen atom, such as element type, chirality, and atomic hybridization type, as well as four types of bond information, including bond type and conjugation (Table S2). A Siamese MPNN architecture-Vreact, was designed to simultaneously accept input features of VOCs and oxidant molecules (Fig. 1). The model takes the SMILES of VOCs and oxidants as input and primarily includes a VOC molecular graph representation layer and a MPNN layer, an oxidant molecular graph representation layer and MPNN layer, an interaction layer, and a prediction layer. The molecular graph G(V, E) encoding layers of VOCs and oxidants containing node feature matrix X and edge feature matrix Y, which learn molecular properties through the MPNN layer (Gilmer et al., 2017). The MPNN forward propagation process consists of two phases: Message Passing Phase and Readout Phase and generates molecular feature tensors A for VOCs and B for oxidants. Subsequently, the interaction layer transforms the molecular features A of VOCs and B of oxidants into tensors A_I and B_I of the same shape and concatenates them into tensor Z. Reaction rate constants are determined not only by the molecular structure of the reactants but also by the interactions between the reactants. The interaction feature tensor I is dotmultiplied with A to obtain the oxidant-affected VOC feature tensor A'; similarly, it is dot-multiplied with B to obtain the VOC-affected oxidant feature tensor B'. These operations embed the learned interaction features into the molecular structure features, providing a more comprehensive representation of the chemical reaction mechanisms between the two reactants. The prediction phase is composed of a pooling layer and three fully connected layers. The pooling layer uses the Set2Set method to achieve global average pooling, and the fully connected layers map the input features to the final predicted values ($\log_{10}k_i$).

More details can be found in Text S2.

150

155

165

175

180

During model training, Adaptive Moment Estimation (Adam) (Kingma and Ba, 2017) was employed to address the fixed learning rates issue in traditional gradient descent methods. Adam adaptively adjusts the learning rate of each parameter using first-order moment estimates (mean of the gradients) and second-order moment estimates (exponentially moving average of the uncentered variance of the gradients), aiding in rapid model convergence. Bayesian optimization was utilized for hyperparameter tuning, which included the initial learning rate of the optimizer (*Ir*), batch size, L2 regularization parameter (*weight decay*), dropout rate (*p*), and MPNN time steps (*T*) (Text S3). During hyperparameter optimization, the hyperparameter combination that minimizes the Mean Squared Error (MSE) of the validation set was selected as the optimal hyperparameter combination, and the best model was saved (Table S3). The predictive performance of the model was assessed using MSE, Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and coefficient of determination (R²) (Text S4). For more information on the model implementation, please refer to Text S5.

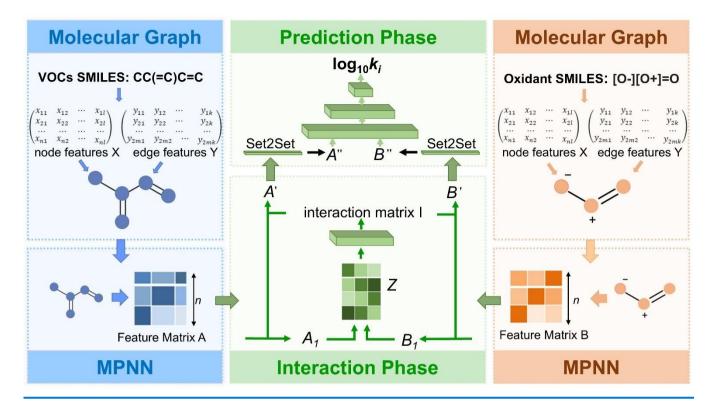


Figure 1. Schematic of the Vreact Architecture. SMILES of VOCs and oxidants are converted into molecular graphs, where nodes represent atoms and edges represent bonds. Atomic and bond features form matrices X and Y. Using a Siamese MPNN architecture, the Vreact model processes these features through separate MPNN layers for VOCs and oxidants. The final prediction layer outputs log₁₀k_i, incorporating both molecular and interaction features.

2.3 Clustering Analysis

185

190

Morgan fingerprints (radius 2, 1024 bits, generated using RDKit) was used as the molecular embeddings before clustering and visualization. To investigate VOC structural diversity and reactivity trends, two methods were applied: the Self-Organizing Map (SOM) (Kohonen, 2006) and the Uniform Manifold Approximation and Projection (UMAP). The SOM algorithm clustered VOCs into 100 structural groups (10×10 grid), using a sigma of 0.3 and learning rate of 0.5. The UMAP algorithm projected the high-dimensional fingerprint space into 2D for visualization, with the number of neighbors set to 50, minimum distance to 0.6, and metric as correlation.

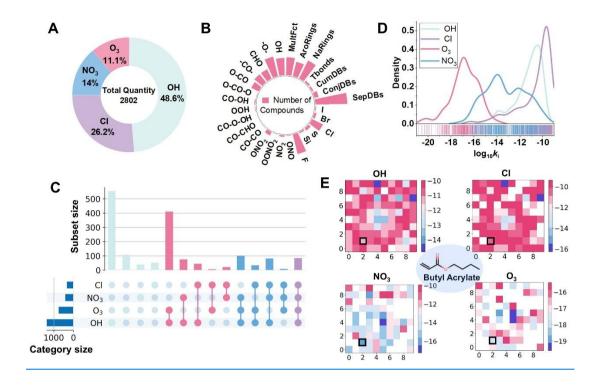
3 Results and Discussion

195 3.1 Analysis of VOC and Oxidant Reaction Data Distribution and Characteristics

The categories and distribution characteristics of VOC and oxidant reaction data are first explored in the study, which includes $\log_{10}k_i$ data for 1586 VOCs with \bullet OHOH, \bullet ClCl, NO₃, and O₃ (Fig. 2A). The dataset contains the most data for \bullet OHOH,

accounting for 48.64% of the total, as *OHOH plays a crucial role in the atmosphere, rapidly reacting with organic pollutants and dominating their removal process. The remaining data points are for *CICI (26.23%), NO₃ (14.03%), and O₃ (11.1%) in descending order of data quantity. O₃ is primarily produced through photochemical reactions involving NO_x and VOCs, while NO₃, as the principal nighttime atmospheric oxidant, significantly contributes to the oxidation and removal of trace gases. The dataset encompasses VOCs with diverse chemical structures, including 22 molecular motifsfunctional groups such as double bonds, esters, benzene rings, and halogen atoms (F, S, Cl, Br, and I) (Fig. 2B). This extensive chemical structure space facilitates the model's ability to learn more structural features and enhances its generalization capability.

Moreover, although there is some overlap in the reactions of the four oxidants with VOCs, each oxidant also has specific VOC reactions (Fig. 2C). There are 747 VOCs with k_i data for only one oxidant and 839 VOCs with k_i data for multiple oxidants, of which 81 VOCs have data for all four oxidants. For example, isoprene can react with *OHOH, NO3, and *CICI through hydrogen abstraction reactions, and undergo addition reactions with O3 via its unsaturated double bonds. Furthermore, the four oxidants exhibit different $\log_{10}k_i$ value distribution with VOCs due to differences in chemical structures and reactivity (Fig. 2D). *OHOH, due to its high oxidation potential, usually reacts quickly with VOCs via hydrogen abstraction, with $\log_{10}k_i$ concentrated in the range of -14.000 to -10.000. In contrast, O3 typically undergoes slower addition reactions with unsaturated bonds in reactants (Ziemann and Atkinson, 2012), with $\log_{10}k_i$ ranging from -20.836 to -13.721. NO3 can participate in both hydrogen abstraction and addition reactions, resulting in a wider range of $\log_{10}k_i$ values. The diverse reaction rates of these oxidants maintain the composition and oxidative state of aerosols in the atmosphere, but the uneven distribution of their values makes predicting k_i more challenging. Even for the same oxidant, VOCs with different structures exhibit varied reaction rates in gas-phase oxidation reactions. For example, NO3 reacts very slowly with aromatic rings, with a k_i value of 3.900×10⁻¹⁶ cm³/(molecule·s) for xylene. In contrast, NO3 can rapidly abstract hydrogen from hydroxyl groups, with a k_i value of up to 1.72×10^{-10} cm³/(molecule·s) for 3-methylcatechol.



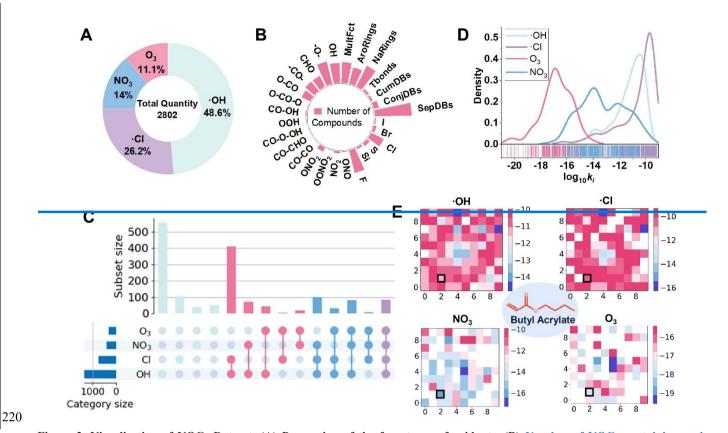


Figure 2. Visualization of VOCs Dataset. (A) Proportion of the four types of oxidants. (B) Number of VOCs containing each molecular motiffunctional group. MultFct: multifunctional; AroRings: aromatic rings; NaRings: non-aromatic rings; Tbonds: triple bonds; CumDBs: cumulated double bonds; ConjDBs: conjugated double bonds; SepDBs: separated double bonds. (C) Number of VOCs that can undergo oxidation reactions with the four oxidants. (D) Distribution of $log_{10}k_i$ values for the four oxidants. (E) Heatmap of reaction rate constants based on VOCs clustering, where each grid represents a cluster of structurally similar VOCs. The color gradient indicates the $log_{10}k_i$ values, with red indicating higher $log_{10}k_i$ values (faster reaction rates), blue indicating lower $log_{10}k_i$ values (slower reaction rates), and white indicating the absence of $log_{10}k_i$ data for that cluster. The cluster containing butyl acrylate are enclosed within the black box.

Furthermore, the same VOCs show different reaction rates with different oxidants. The SOM algorithm was used to explore the relationship between VOC structural variation and $\log_{10}k_i$. Each grid in Fig. 2E represents a VOC cluster, and the color gradient indicates reactivity (average $\log_{10}k_i$ values) for the corresponding oxidants. By comparing $\log_{10}k_i$ values across clusters, oxidant-specific reactivity patterns can be assessed. To explore the relationship between structural differences of VOCs and reaction rates, the study employed the Self Organizing Map (SOM) algorithm (Kohonen, 2006) to visualize $\log_{10}k_i$ values. Based on the Morgan fingerprint similarity of VOCs, the VOCs were clustered into 100 groups, each containing VOCs with similar molecular structures. Each grid in Fig. 2E represents a cluster of VOCs, and the color gradient indicates the $\log_{10}k_i$ values of their reactions with the corresponding oxidants. By comparing the $\log_{10}k_i$ values of the same VOCs with four oxidants, the relationship between structural features and reaction rates for each oxidant can be evaluated. For example, butyl acrylate (CAS RN.141-32-2) reacts slowly with NO₃ radicals and O₃, mainly due to the unsaturated addition reactions through the

carbon-carbon double bond, where the ester group in the molecular structure produces an electron-withdrawing effect, reducing the electron density in the π bond and thus lowering the reaction rate (Gai et al., 2009; Wang et al., 2010). In contrast, it reacts faster with •OHOH and •CICI through hydrogen abstraction rather than addition (Le Calvé et al., 1997; Ohta, 1984; Wang et al., 2018). This demonstrates that the dataset, which includes various oxidants and VOCs, exhibits diverse log₁₀k_i values. The overall log₁₀k_i values differ significantly between different oxidants. This diverse dataset enables the model to learn the reaction information between VOCs and different oxidants, thereby improving model performance and prediction accuracy.

3.2 Performance Evaluation of Vreact Model

245

250

255

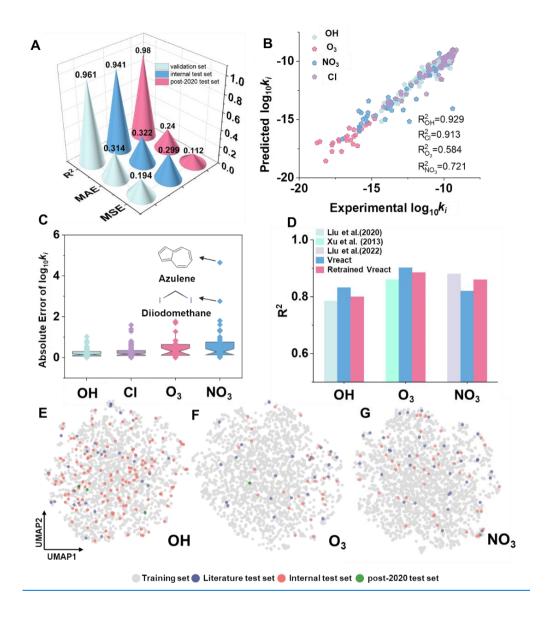
260

265

The Siamese MPNN architecture of the_-Vreact captures both molecular features of VOCs and oxidants as well as their interaction dynamics simultaneously. During hyperparameter optimization, the set of hyperparameters that minimized MSE on the validation set was selected. After training for 46 epochs (Fig. S1),_-Vreact achieved robust predictive performance on the validation set, with R^2 of 0.961, MSE of 0.194 and MAE of 0.314 for $\log_{10}k_i$ (Fig. 3A). On the internal test set, the model achieved R^2 of 0.941, MSE of 0.299 and MAE of 0.322 for $\log_{10}k_i$ (Fig. 3A), indicating robust predictive capability and excellent generalization ability for unseen VOC-oxidant combinations. The small MAE difference between the validation set and internal test sets, despite a larger difference in MSE, indicates that MSE is more sensitive to outliers or large errors, while MAE directly reflects the average absolute prediction error. Although the R^2 on the internal test set is slightly lower than on the validation set, this minor discrepancy does not affect the model's robust predictive ability. The result on the internal test set is available in Table S4.

To explore the predictive performance of the Vreact model for different types of oxidants, we evaluated the prediction performance for $\bullet OHOH$, $\bullet ClCl$, O₃, and NO₃ separately. The regression fit of predicted $\log_{10}k_i$ values versus experimental values for the four oxidants (Fig. 3B) shows that O₃ and NO₃ have higher dispersion compared to $\bullet OHOH$ and $\bullet ClCl$. The R² values for the reactions of the four oxidants, in descending order, are $\bullet OHOH > \bullet ClCl > NO_3 > O_3$, with $\bullet OHOH$ and $\bullet ClCl$ having R² values of 0.92942 and 0.9136, respectively. The prediction performance for NO₃ radicals and O₃ is comparatively lower, with R² values below 0.800.

The •OHOH dataset is the most abundant and concentratedbalanced, while data amount of O₃ and NO₃ was relatively small, and the model can't fully capture the reaction features, leading to prediction bias. whereas the log₁₀k_i values for In addition, the log₁₀k_i values for NO₃ are highly dispersed, which may cause the model to have difficulty capturing all patterns and relationships for NO₃, also reducing the prediction performance. Additionally, the order of the size of R² is consistent with the order of the data volume of the four oxidant datasets. This indicates that the amount of data is also an important factor affecting the prediction performance of reaction rate constants, and that more available data help the model to fully capture reaction features, leading to prediction biases.



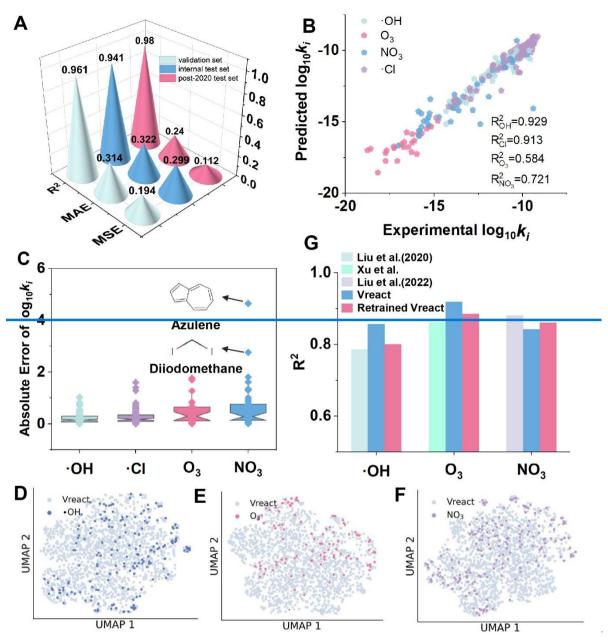


Figure 3. Evaluation and comparison of the predictive performance of the Vreact model. (A) MSE, MAE, R² of Vreact (trained on the McGillen et al. dataset) on the validation set, internal test set, and external post-2020 test set. (B) R² values for log₁₀k_i predictions of four oxidants' reactions in the internal test set. (C) Distribution of AE between predicted and experimental log₁₀k_i values for the four oxidants in the internal test set. (D) R² comparison among previously published single-oxidant models, the original Vreact (evaluated on cleaned literature test sets), and Retrained Vreact (trained and tested using the same original splits as the literature) highlighting adaptability. -(DE-FG) The chemical spatial distribution of VOCs in the *OHOH, O₃, and NO₃ datasets used in this study and prior literatures. (G) R² comparison among previously published single-oxidant models, the original Vreact (evaluated on literature test set), and Retrained Vreact (trained and tested using the same splits as the literature models) highlighting adaptability.

The Absolute Error (AE) between the predicted and experimental $\log_{10}k_i$ values for the four types of oxidants are presented in Fig. 3C. The median AE for *OHOH is 0.149, while O₃ and NO₃ exhibit median AEs of 0.301 and 0.287, respectively, which are slightly higher than that of *OHOH. Overall, 84% of the AE values for O₃ and NO₃ are within 1. As depicted in the Fig. 3C, individual outliers in AE contribute to the increased RMSE and MAE for O₃ and NO₃, and the consequent decrease in R². For example, the AE for the reaction of NO₃ with azulene ($C_{10}H_8$) is 4.653. Azulene, an aromatic hydrocarbon composed of a seven-membered ring fused to a five-membered ring, is an isomer of naphthalene ($C_{10}H_8$). NO₃, as electrophilic reagents, tend to attack regions with higher electron density. Compared to naphthalene, the electron density distribution of azulene is uneven, with certain regions having high electron density that may facilitate effective interactions with NO₃. Additionally, the structure of azulene may reduce steric hindrance, allowing NO₃ radicals easier access to reaction sites (Atkinson et al., 1992), resulting in a higher reaction rate constant and increasing the model's prediction difficulty. Similarly, the predicted $\log_{10}k_i$ value for the reaction of NO₃ with diiodomethane (CH₂I₂) is significantly lower than the true value (AE=2.763). This discrepancy may be attributed to the limited representation of iodine-containing VOCs in the dataset, with only iodomethane (CH₃I) and iodoethane (C₂H₃I) having k_i values in the training and validation sets. This limited data prevents the model from fully learning the reaction characteristics of iodine-containing compounds, resulting in a larger prediction error for diiodomethane with NO₃ radicals.

3.3 Comparation with Single-Oxidant Prediction Models

280

285

290

Most existing machine learning models for predicting VOC reaction rates constants are tailored for individual oxidants, 295 limiting their applicability to complex atmospheric systems involving multiple oxidants. In contrast, the Siamese MPNN architecture of the Vreact enables simultaneous learning of molecular features and interaction patterns across different VOCoxidant pairs within a unified framework. To benchmark Vreact against previously published single-oxidant QSAR/ML models, we selected three top-performing models developed under 298K conditions: Liu et al. (2020) for *OHOH (training/test = 144/36training set:144/ test set:36180 data points), Xu et al. (2013) for O₃ (training set:60/test set:3595 data points), and 300 Liu et al. (2022) for NO₃ radicals (1training set:151/test set:38189 data points). Prior to evaluation, we UMAP was applied Uniform Manifold Approximation and Projection (UMAP) to reduce the dimensionality of the Morgan molecular fingerprints to visualize the chemical space of both the comparison literature datasets and the Vreact training set (Fig.s. 3D F S2). The observed structural overlap confirms that Vreact's dataset spans a broad and diverse chemical space. Given that our study used different data than those reported in the literature, we employed two strategies for comparison. First, the pre-trained 305 Vreact model (trained on the McGillen dataset) was directly applied to the literature test sets from the literature to evaluate extrapolation performance. To ensure a fair comparison, overlapping data points between the literature test sets and the McGillen training set were removed However, there are some duplicates in the data points of the literature test set and the McGillen training set, and an unfair comparison can be made if duplicate data exist. So, we removed the duplicate data to construct cleaned literature test sets for evaluation (2 of 38 for NO₃, 13 of 35 for O₃, and 6 of 36 for OHAmong 38 NO₃ data points, 2 are in the training set; among 35 O₃ data points, 13 are in the training set; and among 36. OH data points, 6 are in the

<u>training set</u>). Second, Vreact was retrained on each literature dataset using their original train/test splits (Retrained Vreact), allowing a direct comparison with published models on original literature test sets (Retrained Vreact).

315

320

325

330

335

340

As shown in Fig. 363D₇, both the original Vreact model and its retrained version consistently outperformed the single-oxidant models from Liu et al. (2022) and Xu et al. (2013) on the *OHOH and O₃ literature test sets, achieving higher R² values and demonstrating superior regression fits between predicted and experimental values. These results highlight the capability of the Vreact architecture—whether trained on a broad multi-oxidant dataset or fine-tuned on smaller single-oxidant datasets—to effectively learn structural features of VOCs and oxidants and capture complex molecular interactions through its Siamese MPNN framework. Notably, Vreact shows opposite performance trends for OH and O₃ between the internal and literature test set. To understand this, UMAP was applied to project compounds from the training, internal, and literature test sets into a shared chemical space. As shown in Fig. 3E, the internal OH test set overlaps well with the training data, leading to consistently strong performance. In contrast, the literature OH set is sparse and scattered near the dataset boundaries. Despite this, Vreact still achieves a high R², demonstrating good generalization. For O₃ (Fig. 3F), the internal test set lies farther from the dense training distribution, contributing to lower R². Meanwhile, the literature O₃ set is better aligned with the training data, resulting in higher prediction accuracy. For NO₃ (Fig. 3G), both internal and literature sets show similar distributions, and the model achieves comparable R² values (~0.815). AAlthough Vreact underperforms slightly compared to the original single-oxidant model, retraining on the literature data improves performance. This suggests that multi-oxidant training may introduce some noise but does not significantly compromise prediction accuracy. Notably, the outcomes for OH and O3 exhibit contrasting results between the literature test set and the internal test set. Use UMAP to map the compounds in the training set, literature test set and the internal test set into the same space simultaneously, which was used to analyse this difference. As shown in Fig 3E, the OH internal test set has the largest amount of data and shares similar distribution features with the training set. Furthermore, most of data points don't exceed the coverage of the training set. In contrast, the literature test set for OH is characterized by sparse data points with a scattered distribution, predominantly located at the dataset boundaries. However, since the model already performs excellently on OH, it can also achieve a high R² on the literature test set, only losing some precision. For O₃ (Fig. 3F), a distinct distributional boundary exists between the internal and literature test sets. The lower region (internal test set) is marked by sparse features and low overlap with the training set, where individual samples deviate from the overall distribution, leading to reduced R². Conversely, the upper region (literature test set) aligns with the training set's distribution, vielding the model can predict better. As a result, the differences between OH and O₂ on the literature dataset and the internal dataset emerge. For the NO₂ prediction task, the R² of the model on the literature test set was 0.815, which was similar to the result on the internal test set. The similar spatial distribution patterns of the two test sets likely contributed to the similarity of their prediction results (Fig. 3G). However, Vreact performed poorer than the original literature results on the literature test set, and the retrained Vreact model showed an improvement in R2. This suggests that a unified dataset containing multiple exidents may introduce additional noise during training, affecting the model's ability to learn key interaction features between NO2 and VOCs. Nonetheless, the noise only causes some loss of prediction accuracy, which is acceptable.

For the NO₃ prediction task, the original Vreact model performed poorly, but the retrained Vreact model showed an improvement in R². This suggests that a unified dataset containing multiple oxidants may introduce additional noise during training, affecting the model's ability to learn key interaction features between NO₃ and VOCs. Nonetheless, the model's R² on the literature test set still reached 0.842, indicating only a slight loss in predictive accuracy, which is acceptable.

3.4 Mechanism Insights Through Interaction Analysis

355

360

365

The interaction layer of the Vreact model can elucidate the atomic interaction mechanisms between VOCs and oxidants. The interaction matrix, sized $n_1 \times n_2$, where n_1 represents the number of non-hydrogen atoms in the VOC molecule and n_2 represents the number of non-hydrogen atoms in the oxidant molecule. Mapping these interaction coefficients onto the molecular structure highlights key atoms that determine the reaction rate.

To exemplify this mechanism, we analysed specific cases. 2-methyl-4-penten-2-ol is an unsaturated oxygenated volatile organic compound (OVOC) that constitutes a significant proportion of the atmospheric VOCs, primarily sourced from industrial solvents used in ink and jet ink manufacturing (Li et al., 2021). As shown in Fig. 4A, the interaction coefficient for the distal unsaturated carbon atoms is the highest during the reaction with O₃, indicating these are likely the reaction sites for O₃ attack. It is inferred that O₃ adds to the unsaturated carbon-carbon double bond through an addition reaction, forming primary ozonides (POZs). These POZs are unstable intermediates that rapidly cleave to produce carbonyl compounds and carbon-based radicals, which further rearrange to form secondary ozonides (SOZs). The SOZs and their reaction products are precursors of SOA. Another example is γ-caprolactone (GCL), a five-membered ring ester used in perfumes, which rapidly reacts and degrades with •OHOH upon entering the atmosphere. Interaction weight analysis shows that the carbon atom linked to the ethyl group contributes most to GCL's oxidative degradation by •OHOH (Fig. 4B), suggesting that •OHOH initially attacks this carbon atom, abstracting a H atom to form a carbon radical. Previous studies indicate that the reactivity of carbons adjacent to the oxygen atom in lactones is particularly significant in reactions with •OHOH, especially when alkyl substituents are attached to this carbon, which enhances its reactivity (Barnes et al., 2014).

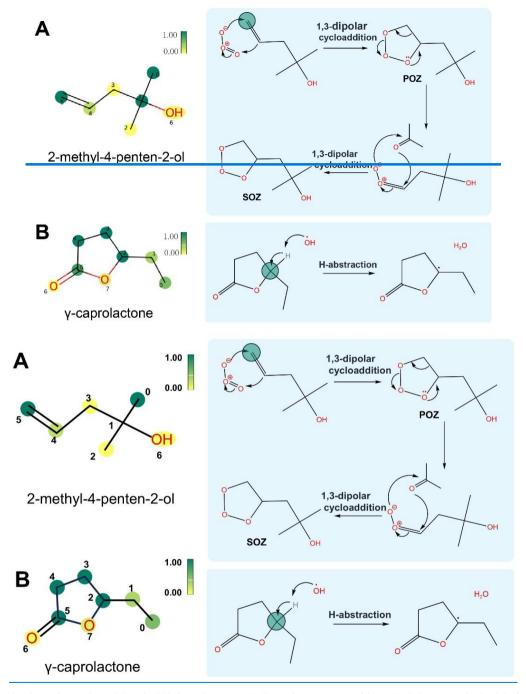


Figure 4. Visualization of atomic weights in VOC molecules. (A) Reaction process of 2-methyl-4-penten-2-ol with O₃. (B) Reaction process of γ-caprolactone with •ΟΗΟΗ. The darker the highlighted color of the atom, the stronger its interaction in the gas-phase oxidation reaction.

3.5 Evaluating Extrapolation Ability and Prioritizing VOCs for Environmental Impact

375

380

385

390

To further validate the extrapolation capability and generalization performance of the Vreact model, developed using a dataset compiled up to the year 2020 (Baptista et al., 2021; Joudan et al., 2022; Li et al., 2021), additional k_i data from experimentally measured VOCs and oxidants published after 2020 were collected as an external test set (post-2020 test set) (Table 1). The prediction results showed that the AE between the experimental $\log_{10}k_i$ and the predicted values was within 1, with the reaction rate constant prediction for γ -heptalactone and \bullet OHOH exhibiting the smallest prediction error. The AE for γ -heptalactone with \bullet OHOH was only 0.005, and the overall MAE was 0.240, with an MSE of 0.112 and an R² of 0.98 (Fig. 3A shown in red). The results indicate that the Vreact can accurately predict the atmospheric oxidation reaction rate constants of unknown VOCs, demonstrating its potential application in addressing complex atmospheric chemistry issues involving the interactions between VOCs and oxidants.

Table 1. The prediction results on the post-2020 test set.

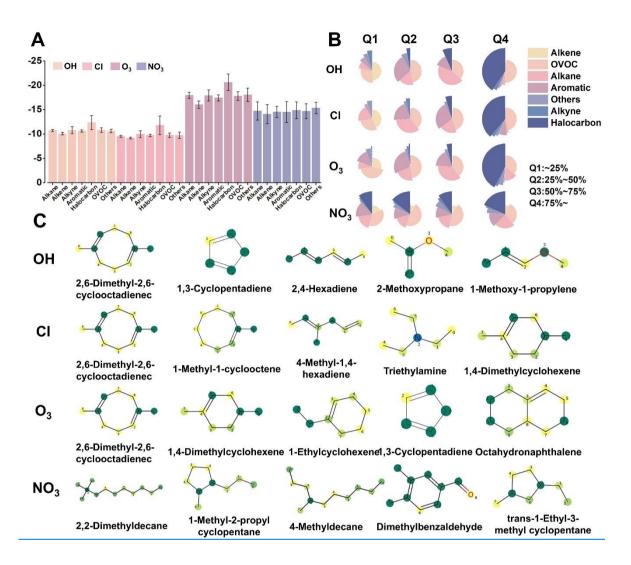
| VOC name | Chemical structure | Oxidant | Experimental log ₁₀ k _i | Predicted log ₁₀ k _i | AE | Ref. |
|----------------------------|--|----------------|--|--|-------|-------------------------|
| 2-methyl-4- penten-2-ol | HO | O ₃ | -17.370 | -16.712 | 0.658 | (Li et al., 2021) |
| γ-caprolactone | | <u>+OHOH</u> | -11.194 | -11.209 | 0.015 | (Baptista et al., 2021) |
| | | <u> •ClCl</u> | -9.886 | -10.149 | 0.263 | (Baptista et al., 2021) |
| γ-heptalactone | | <u>•OHOH</u> | -11.056 | -11.051 | 0.005 | (Baptista et al., 2021) |
| | | <u> +ClCl</u> | -9.770 | -9.943 | 0.173 | (Baptista et al., 2021) |
| FESOH | F ₃ C F F F F F F F F F F F F F F F F F F F | <u>•OHOH</u> | -11.377 | -11.876 | 0.499 | (Joudan et al., 2022) |
| | | <u>+ClCl</u> | -10.824 | -10.759 | 0.065 | (Joudan et al., 2022) |

FESOH: 2- (1,1,2-trifluoro-2-heptafluoropropyloxy-ethylsulfanyl)-ethanol; AE: absolute error

Despite the identification of hundreds of VOC species, the environmental behavior of most VOCs in the atmosphere and their potential contributions to particulate matter formation and ozone increase remain largely unclear. To address this gap, we employed the Vreact model to evaluate the atmospheric oxidation reaction rate constants of a broad spectrum of VOCs. Molecular structures for 447 VOCs with unknown atmospheric oxidation k_i values were collected from previous research, which evaluated more than 500 Chinese domestic source profiles, including literature and field measurements (Sha et al., 2021) (Table S5). After excluding VOCs already included in the Vreact dataset, 296, 339, 416, and 369 data points for $\frac{1}{2}$

*CICI, O₃, and NO₃ were retained, respectively. The prediction results indicated that, although the oxidation reaction rates of VOCs in the atmosphere vary (Fig. 5A), the differences in $\log_{10}k_i$ values are primarily influenced by the type of oxidant, with smaller variations in $\log_{10}k_i$ values observed for different VOCs reacting with the same oxidant. Among these, reactions with *OHOH and *CICI were the fastest, consistent with the results from the McGillen dataset analysis used in the modeling (Fig. 2D). Additionally, the changes in the proportion of VOC types within different reaction rate intervals (Fig. 5B) demonstrated that the composition of VOC types varied with reaction rates. Halocarbons exhibited relatively slower reaction rates, while alkenes and aromatics reacted relatively quickly, and oxygenated compounds showed a more uniform rate distribution. Consequently, areas with high emissions of alkenes and aromatics will produce more reaction products per unit time, providing precursors for O₃ and SOA formation (Gao et al., 2021).

The top five VOCs with the fastest reaction rates with •OHOH, •ClCl, O3, and NO3 were further examined in the study (Fig. 5C). Among these, 2,6-Dimethyl-2,6-cyclooctadiene (CAS RN: 3760-14-3) is a volatile compound with an irritating odor, exhibiting the fastest reaction rates with •OHOH, •ClCl, and O3. Additionally, 1,3-cyclopentadiene (CAS RN: 542-92-7) and 1,4-Dimethylcyclohexene (CAS RN: 70688-47-0) also showed high reaction rates with O3, •ClCl, and •OHOH, likely due to the presence of double bonds and cyclic structures in these molecules. The carbon atoms in the double bonds and those connected to methyl groups generally have high reactivity. Therefore, it could be inferred that these VOCs, or VOCs with similar structures, may significantly contribute to the formation of fine particulate matter and the increase in ozone in the atmosphere.



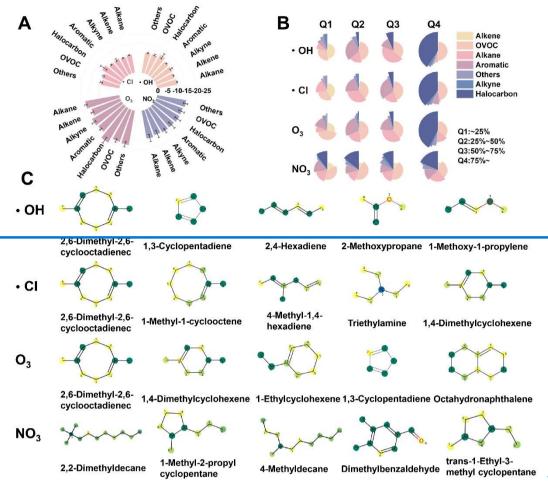


Figure 5. Predicted reaction rate constants for VOCs atmospheric oxidation reactions. (A) Predicted mean $\log_{10}k_i$ values for different types of VOCs. (B) Distribution of VOC types ranked by predicted reaction rates, divided into quartiles: the fastest 25% (Q1), 25%-50% (Q2), 50%-75% (Q3), and the slowest 25% (Q4). (C) Molecular structures of VOCs with the fastest reaction rates with the four oxidants.

415 4 Concluding Concluding

410

420

In response to growing concerns about atmospheric pollution and its impact on human health and climate, this study introduces Vreact, a deep learning model designed to predict oxidation rate constants for VOCs with multiple oxidants (OH, Cl, NO₃, O₃). Vreact demonstrates strong overall performance (MSE=0.299, R²=0.941 on internal test data) and provides mechanistic insights by capturing atomic-level interaction patterns through a Siamese MPNN framework. Its predictive accuracy varies by oxidant, reflecting the availability and diversity of training data. The model achieves high accuracy for OH (R²=0.929, n=1363) and Cl (R²=0.913, n=735), supporting robust application in daytime oxidation modeling. In contrast, lower performance is observed for NO₃ (R²=0.721, n=393) and O₃ (R²=0.584, n=311), pointing to challenges in modeling oxidants with fewer data

and more complex mechanisms. This underscores the importance of expanding high-quality experimental datasets to improve generalization, particularly for underrepresented oxidants and VOC classes.

Vreact supports high-throughput screening for emission inventories and atmospheric reactivity assessments. Its applications span VOC prioritization, emission control planning, and kinetic mechanism development, offering actionable insights for environmental policy and modeling. An interactive web interface (http://vreact.envwind.site:8001) (Fig. S3) enhances accessibility for researchers and policymakers. Further improvements in NO₃ and O₃ predictions will expand its utility in nighttime chemistry and secondary aerosol formation scenarios. In response to growing concerns about atmospheric pollution and its impact on human health and climate, this study introduces Vreact, a deep learning model designed to predict oxidation rate constants for VOCs with multiple oxidants (OH, Cl, NO₃, O₃). Vreact demonstrates strong overall performance (MSE = 0.299, R² = 0.941 on internal test data) and provides mechanistic insights by capturing atomic level interaction patterns through a Siamese MPNN framework. Its predictive accuracy varies by oxidant, reflecting the availability and diversity of training data. The model achieves high accuracy for OH (R² = 0.929, n = 1363) and Cl (R² = 0.913, n = 735), supporting robust application in daytime oxidation modeling. In contrast, lower performance is observed for NO₃ (R² = 0.721, n = 393) and O₃ (R² = 0.584, n = 311), pointing to challenges in modeling oxidants with fewer data and more complex mechanisms. This underscores the importance of expanding high quality experimental datasets to improve generalization, particularly for underrepresented oxidants and VOC classes.

Vreact supports high throughput screening for emission inventories and atmospheric reactivity assessments. Its applications span VOC prioritization, emission control planning, and kinetic mechanism development, offering actionable insights for environmental policy and modeling. An interactive web interface (http://vreact.envwind.site:8001) (Fig. S3) enhances accessibility for researchers and policymakers. Further improvements in NO₃ and O₃ predictions will expand its utility in nighttime chemistry and secondary acrosol formation scenarios.

Given the increasing complexity of atmospheric pollution and its global impacts on human health and climate, these advancements in predictive modeling offer valuable resources for addressing worldwide air quality challenges. Understanding the oxidation rates of VOCs is crucial for evaluating their impact on atmospheric chemistry and air pollution. In this study, Vreact, a deep learning based model, is introduced to predict VOC oxidation reaction rate constants with multiple oxidants (OH, Cl, NOs and Os) simultaneously while offering mechanistic insights into VOC oxidation by analyzing atomic level interaction patterns. Vreact achieves robust predictive performance, with MSE of 0.299 and R² of 0.941 on internal test data. By incorporating a broad range of VOC structures and oxidant interactions, Vreact enhances its generalizability, allowing for large scale screening of previously uncharacterized VOC oxidation rates. Additionally, an interactive web based tool (http://vreact.envwind.site:8001) is provided (Fig. S3), which facilitates VOC oxidation rate predictions for non-experts. This tool significantly improves accessibility, enabling researchers, policymakers, and environmental agencies to assess VOC reactivity and prioritize mitigation efforts effectively. Furthermore, the study found that the model's reliable model accuracy

depends on the quantity and quality of available experimental kinetic data, which vary significantly among oxidants. Vreact demonstrates superior accuracy for OH (R²=0.929) and Cl (R²=0.913), where extensive datasets (1363 and 735 VOCs) facilitate robust learning of structure reactivity relationships. In contrast, predictions for NO₂ (R²=0.721) and O₂ (R²=0.584) are less precise, reflecting smaller data volumes (393 and 311 VOCs) and higher complexity in reaction mechanisms. The difference underscores the critical role of data availability in achieving reliable model performance, particularly for oxidants with sparse experimental coverage or complex reaction mechanisms, which requires further and adequate data collection. Considering the current accuracy of Vreact and the discussion and analysis of this study, some potential applications and improvements can be amenable to an outlook. Vreact has excellent overall performance and supports high throughout reactivity assessment, thus enabling it to identify highly reactive VOCs for prioritizing emission control. In addition, the high accuracy of Vreact in OH and Cl data demonstrates its application advantages in daytime atmospheric simulation, allowing model to integrate into regional air quality models to simulate VOCs degradation nathways under daytime conditions. However, the model needs further improvement for simulating nighttime reaction processes (NO₃ and O₃). Improved data coverage for these exidents is essential for refining forecasts of necturnal aerosol formation. Due to limited training data, which may not cover all VOCs, especially those with complex structures or containing halogen and sulfur groups. As a result, VOC classes that are underrepresented in the dataset, such as jodine containing compounds, may exhibit prediction errors, which also highlights the need for further data collection to make the model more broadly applicable. Additionally, while Vreact captures essential molecular interactions, biases may arise from the existing reaction rate constants datasets, especially when reaction conditions or mechanisms differ from those used in training. Extrapolating reaction rate constants for uncharacterized VOCs presents another challenge. While the model shows strong generalization capabilities, its accuracy may decrease for highly reactive or structurally unique compounds. To address these challenges, future work should focus on integrating highthroughput quantum chemical calculations and automated experimental validation to augment existing datasets, especially for data poor exidents and functional groups. Optimizing Vreact's inference speed and coupling it with atmospheric chemistry models could further enable real time simulations of air quality scenarios, enhancing its applicability for regulatory decision-

making.

460

465

470

475

480

485

490

Understanding the oxidation rates of VOCs is crucial for evaluating their impact on atmospheric chemistry and air pollution. In this study, Vreact, a deep learning based model, is introduced to predict VOC oxidation reaction rate constants with multiple oxidants simultaneously. The model demonstrates high predictive accuracy while offering mechanistic insights into VOC oxidation by analyzing atomic level interaction patterns. By incorporating a broad range of VOC structures and oxidant interactions, Vreact enhances its generalizability, allowing for large scale screening of previously uncharacterized VOC oxidation rates. Additionally, an interactive web based tool (http://vreact.envwind.site:8001) is provided (Fig. S2), which facilitates VOC oxidation rate predictions for non-experts. This tool significantly improves accessibility, enabling researchers, policymakers, and environmental agencies to assess VOC reactivity and prioritize mitigation efforts effectively. Given the increasing complexity of atmospheric pollution and its global impacts on human health and climate, these advancements in predictive modeling offer valuable resources for addressing worldwide air quality challenges.

However, Vreact has several limitations. The model's performance depends on the availability and quality of experimental kinetic data. The training dataset primarily relies on measured reaction rate constants, which may not cover all VOCs, especially those with complex structures or containing halogen and sulfur groups. As a result, VOC classes that are underrepresented in the dataset, such as iodine-containing compounds, may exhibit prediction errors, highlighting the need for further data collection. Additionally, while Vreact captures essential molecular interactions, biases may arise from the existing reaction rate constants datasets, especially when reaction conditions or mechanisms differ from those used in training. Extrapolating reaction rate constants for uncharacterized VOCs presents another challenge. While the model shows strong generalization capabilities, its accuracy may decrease for highly reactive or structurally unique compounds. To improve predictions, future work could integrate high throughput quantum chemical calculations and automated experimental validation. Optimizing inference speed and integrating Vreact with atmospheric chemistry models could enhance its applicability in real time air quality simulations.

Data and Code Availability

495

500

510

The code and datasets used and/or analyzedanalysed during the current study are available at https://github.com/Luo-Jiaqi/Vreact and supplemental information.

505 Supplementary Material

Detailed information about the learning curve of the Vreact training process (Figure S1); The chemical spatial distribution of VOCs in the OH, O3, and NO3 datasets used in this study and prior literatures (Figure S2); User interface of the web platform for predicting VOC reaction rate constants using the Vreact model (Figure \$2\$3); Graph representation of molecular structures (Text S1); MPNN message passing and readout phases for molecular graphs (Text S2); Regularization and early stopping techniques in the Vreact model training (Text S3); Model performance evaluation metrics (Text S4); Implementation of the Vreact model (Text S5); Distribution of VOCs reactions with atmospheric oxidants across datasets (Table S1); Atomic features and bond features used in molecular graph representation (Table S2); Hyperparameter search space and optimal settings for the Vreact model (Table S3);Experimental and predicted $\log_{10}k_i$ values for VOCs on the internal test dataset (Table S4); 447 real-world atmospheric VOCs (Table S5).

515 Author Contributions

Methodology, Investigation, Formal analysis, Data curation, Visualization, Writing-original draft, X.Z. and J.Q.L; Resources, Conceptualization, Software, Writing-review & editing, Supervision, Funding acquisition, J.J.F and X.L.; Software,

Validation, Writing-review & editing, W.X.P and Q.X.; Software, Funding acquisition, Writing-review & editing, A.Q.Z; Resources, Supervision, G.B.J.

520 Competing Interests

The authors declare no competing interests.

Financial Support

This research was supported by the Strategic Priority Research Program of the Chinese Academy of Sciences XDB0750100; the project of National Natural Science Foundation of China grant numbers 22193053, 22276197, 22022611 and 92143301; and the Youth Innovation Promotion Association of CAS grant number Y2022020.

References

- Abramson, J., Adler, J., Dunger, J., Evans, R., Green, T., Pritzel, A., Ronneberger, O., Willmore, L., Ballard, A. J., Bambrick, J., Bodenstein, S. W., Evans, D. A., Hung, C.-C., O'Neill, M., Reiman, D., Tunyasuvunakool, K., Wu, Z., Žemgulytė, A., Arvaniti, E., Beattie, C., Bertolli, O., Bridgland, A., Cherepanov, A., Congreve, M., Cowen-Rivers, A. I., Cowie, A., Figurnov, M., Fuchs, F. B., Gladman, H., Jain, R., Khan, Y. A., Low, C. M. R., Perlin, K., Potapenko, A., Savy, P., Singh, S., Stecula, A., Thillaisundaram, A., Tong, C., Yakneen, S., Zhong, E. D., Zielinski, M., Žídek, A., Bapst, V., Kohli, P., Jaderberg, M., Hassabis, D., and Jumper, J. M.: Accurate structure prediction of biomolecular interactions with AlphaFold 3, Nature, 630, 493–500, https://doi.org/10.1038/s41586-024-07487-w, 2024.
- Atkinson, R.: Kinetics and mechanisms of the gas-phase reactions of the hydroxyl radical with organic compounds under atmospheric conditions, Chem. Rev., 86, 69–201, https://doi.org/10.1021/cr00071a004, 1986.
 - Atkinson, R.: A structure-activity relationship for the estimation of rate constants for the gas-phase reactions of OH radicals with organic compounds, Inter. J. Chem. Kinet., 19, 799–828, https://doi.org/10.1002/kin.550190903, 1987.
 - Atkinson, R., Arey, J., and Aschmann, S. M.: Gas-phase reactions of azulene with OH and NO3 radicals and O3 at 298 ± 2 K, Inter. J. Chem. Kinet., 24, 467–480, https://doi.org/10.1002/kin.550240507, 1992.
- 540 Baptista, A., Gibilisco, R. G., Wiesen, P., and Teruel, M. A.: FTIR kinetic study of the reactions of γ-caprolactone and γ-heptalactone initiated by Cl and OH radicals at 298 K and atmospheric pressure, Chem. Phys. Lett., 765, 138313, https://doi.org/10.1016/j.cplett.2020.138313, 2021.
- Barnes, I., Kirschbaum, S., and Simmie, J. M.: Combined Experimental and Theoretical Study of the Reactivity of γ-Butyroand Related Lactones, with the OH Radical at Room Temperature, J. Phys. Chem. A, 118, 5013–5019, https://doi.org/10.1021/jp502489k, 2014.
 - Basant, N. and Gupta, S.: Multi-target QSPR modeling for simultaneous prediction of multiple gas-phase kinetic rate constants of diverse chemicals, Atmos. Environ., 177, 166–174, https://doi.org/10.1016/j.atmosenv.2017.11.028, 2018.

- Burés, J. and Larrosa, I.: Organic reaction mechanism classification using machine learning, Nature, 613, 689-695, https://doi.org/10.1038/s41586-022-05639-4, 2023.
- Canneaux, S., Bohr, F., and Henon, E.: KiSThelP: A program to predict thermodynamic properties and rate constants from quantum chemistry results†, J. Comput. Chem., 35, 82–93, https://doi.org/10.1002/icc.23470, 2014.
 - Chen, X., Ma, W., Zheng, F., Wang, Z., Hua, C., Li, Y., Wu, J., Li, B., Jiang, J., Yan, C., Petäjä, T., Bianchi, F., Kerminen, V.-M., Worsnop, D. R., Liu, Y., Xia, M., and Kulmala, M.: Identifying Driving Factors of Atmospheric N2O5 with Machine Learning, Environ. Sci. Technol., 58, 11568–11577, https://doi.org/10.1021/acs.est.4c00651, 2024.
- Fatemi, M. H.: Prediction of ozone tropospheric degradation rate constant of organic compounds by using artificial neural networks, Anal. Chim. Acta, 556, 355–363, https://doi.org/10.1016/j.aca.2005.09.033, 2006.
 - Finlayson-Pitts, B. J. and Pitts, J. N.: Tropospheric Air Pollution: Ozone, Airborne Toxics, Polycyclic Aromatic Hydrocarbons, and Particles, Science, 276, 1045–1051, https://doi.org/10.1126/science.276.5315.1045, 1997.
- Gai, Y., Ge, M., and Wang, W.: Rate constants for the gas phase reaction of ozone with *n*-butyl acrylate and ethyl methacrylate, 560 Chem. Phys. Lett., 473, 57–60, https://doi.org/10.1016/j.cplett.2009.03.070, 2009.
 - Gao, Y., Li, M., Wan, X., Zhao, X., Wu, Y., Liu, X., and Li, X.: Important contributions of alkenes and aromatics to VOCs emissions, chemistry and secondary pollutants formation at an industrial site of central eastern China, Atmos. Environ., 244, 117927, https://doi.org/10.1016/j.atmosenv.2020.117927, 2021.
- Gilmer, J., Schoenholz, S. S., Riley, P. F., Vinyals, O., and Dahl, G. E.: Neural message passing for Quantum chemistry, in: Proceedings of the 34th International Conference on Machine Learning Volume 70, Sydney, NSW, Australia, 1263–1272, 2017.
 - Hallquist, M., Wenger, J. C., Baltensperger, U., Rudich, Y., Simpson, D., Claeys, M., Dommen, J., Donahue, N. M., George, C., Goldstein, A. H., Hamilton, J. F., Herrmann, H., Hoffmann, T., Iinuma, Y., Jang, M., Jenkin, M. E., Jimenez, J. L., Kiendler-Scharr, A., Maenhaut, W., McFiggans, G., Mentel, T. F., Monod, A., Prévôt, A. S. H., Seinfeld, J. H., Surratt, J. D., Szmigielski,
- R., and Wildt, J.: The formation, properties and impact of secondary organic aerosol: Current and emerging issues, Atmos. Chem. Phys., 9, 5155–5236, https://doi.org/10.5194/acp-9-5155-2009, 2009.
 - Han, D., Gao, S., Fu, Q., Cheng, J., Chen, X., Xu, H., Liang, S., Zhou, Y., and Ma, Y.: Do volatile organic compounds (VOCs) emitted from petrochemical industries affect regional PM2.5?, Atmos. Res., 209, 123–130, https://doi.org/10.1016/j.atmosres.2018.04.002, 2018.
- Heid, E., Greenman, K. P., Chung, Y., Li, S.-C., Graff, D. E., Vermeire, F. H., Wu, H., Green, W. H., and McGill, C. J.: Chemprop: A Machine Learning Package for Chemical Property Prediction, J. Chem. Inf. Model., 64, 9–17, https://doi.org/10.1021/acs.jcim.3c01250, 2024.
- Huang, R.-J., Zhang, Y., Bozzetti, C., Ho, K.-F., Cao, J.-J., Han, Y., Daellenbach, K. R., Slowik, J. G., Platt, S. M., Canonaco, F., Zotter, P., Wolf, R., Pieber, S. M., Bruns, E. A., Crippa, M., Ciarelli, G., Piazzalunga, A., Schwikowski, M., Abbaszade,
 G., Schnelle-Kreis, J., Zimmermann, R., An, Z., Szidat, S., Baltensperger, U., Haddad, I. E., and Prévôt, A. S. H.: High secondary aerosol contribution to particulate pollution during haze events in China, Nature, 514, 218–222, https://doi.org/10.1038/nature13774, 2014.
- Huang, Z., Yu, J., He, W., Yu, J., Deng, S., Yang, C., Zhu, W., and Shao, X.: AI-enhanced chemical paradigm: From molecular graphs to accurate prediction and mechanism, Journal of Hazardous Materials, 465, 133355, https://doi.org/10.1016/j.jhazmat.2023.133355, 2024.

- Joudan, S., Orlando, J. J., Tyndall, G. S., Furlani, T. C., Young, C. J., and Mabury, S. A.: Atmospheric Fate of a New Polyfluoroalkyl Building Block, C3F7OCHFCF2SCH2CH2OH, Environ. Sci. Technol., 56, 6027–6035, https://doi.org/10.1021/acs.est.0c07584, 2022.
- Kamarrudin, N., Zulkafli, N. H., Sikirman, A., Mahayuddin, N. M., Sigau, B. A., Ku Hamid, K. H., and Akhbar, S.: Concentration and toxicological study on sanitary landfill gases at drilling point closed cell, in: 2013 IEEE Business Engineering and Industrial Applications Colloquium (BEIAC), 2013 IEEE Business Engineering and Industrial Applications Colloquium (BEIAC), 333–338, https://doi.org/10.1109/BEIAC.2013.6560142, 2013.
 - Kingma, D. P. and Ba, J.: Adam: A Method for Stochastic Optimization, http://arxiv.org/abs/1412.6980, 29 January 2017.
- Kohonen, T.: Self-organizing neural projections, Neural Networks, 19, 723–733, https://doi.org/10.1016/j.neunet.2006.05.001, 2006.
 - Krüger, M., Galeazzo, T., Eremets, I., Schmidt, B., Pöschl, U., Shiraiwa, M., and Berkemeier, T.: Improved vapor pressure predictions using group contribution-assisted graph convolutional neural networks (GC²NN), EGUsphere, 1–22, https://doi.org/10.5194/egusphere-2025-1191, 2025.
- Kubečka, J., Knattrup, Y., Engsvang, M., Jensen, A. B., Ayoubi, D., Wu, H., Christiansen, O., and Elm, J.: Current and future machine learning approaches for modeling atmospheric cluster formation, Nat. Comput. Sci., 3, 495–503, https://doi.org/10.1038/s43588-023-00435-0, 2023.
 - Kwok, E. S. C. and Atkinson, R.: Estimation of hydroxyl radical reaction rate constants for gas-phase organic compounds using a structure-reactivity relationship: An update, Atmos. Environ., 29, 1685–1695, https://doi.org/10.1016/1352-2310(95)00069-B, 1995.
- Le Calvé, Stéphane, Le Bras, G., and Mellouki, A.: Temperature Dependence for the Rate Coefficients of the Reactions of the OH Radical with a Series of Formates, J. Phys. Chem. A, 101, 5489–5493, https://doi.org/10.1021/jp970554x, 1997.
 - Li, B., Su, S., Zhu, C., Lin, J., Hu, X., Su, L., Yu, Z., Liao, K., and Chen, H.: A deep learning framework for accurate reaction prediction and its application on high-throughput experimentation data, Journal of Cheminformatics, 15, 72, https://doi.org/10.1186/s13321-023-00732-w, 2023.
- 610 Li, W., Dan, G., Chen, M., Chen, Y., Wang, Z., Zhao, Y., Wang, F., Li, F., Tong, S., and Ge, M.: The gas-phase reaction kinetics of different structure of unsaturated alcohols and ketones with O3, Atmos. Environ., 254, 118394, https://doi.org/10.1016/j.atmosenv.2021.118394, 2021.
- Liu, Y., Cheng, Z., Liu, S., Tan, Y., Yuan, T., Yu, X., and Shen, Z.: Quantitative structure activity relationship (QSAR) modelling of the degradability rate constant of volatile organic compounds (VOCs) by OH radicals in atmosphere, Sci. Total. Environ., 729, 138871, https://doi.org/10.1016/j.scitotenv.2020.138871, 2020.
 - Liu, Y., Liu, S., Cheng, Z., Tan, Y., Gao, X., Shen, Z., and Yuan, T.: Predicting the rate constants of volatile organic compounds (VOCs) with ozone reaction at different temperatures, Environ. Pollut., 273, 116502, https://doi.org/10.1016/j.envpol.2021.116502, 2021.
- Liu, Y., Cheng, Z., Liu, S., Ren, Y., Yuan, T., Zhang, X., Fan, M., and Shen, Z.: A quantitative structure activity relationship (QSAR) model for predicting the rate constant of the reaction between VOCs and NO3 radicals, Chem. Eng. J., 448, 136413, https://doi.org/10.1016/j.cej.2022.136413, 2022.

- Logan, S. R.: The origin and status of the Arrhenius equation, J. Chem. Educ., 59, 279, https://doi.org/10.1021/ed059p279, 1982.
- McGillen, M. R., Carter, W. P. L., Mellouki, A., Orlando, J. J., Picquet-Varrault, B., and Wallington, T. J.: Database for the kinetics of the gas-phase atmospheric reactions of organic compounds, Earth Sys. Sci. Data, 12, 1203–1216, https://doi.org/10.5194/essd-12-1203-2020, 2020.
- Meana-Pañeda, R., Zheng, J., Bao, J. L., Zhang, S., Lynch, B. J., Corchado, J. C., Chuang, Y.-Y., Fast, P. L., Hu, W.-P., Liu, Y.-P., Lynch, G. C., Nguyen, K. A., Jackels, C. F., Fernández-Ramos, A., Ellingson, B. A., Melissas, V. S., Villà, J., Rossi, I., Coitiño, E. L., Pu, J., Albu, T. V., Zhang, R. M., Xu, X., Ratkiewicz, A., Steckler, R., Garrett, B. C., Isaacson, A. D., and Truhlar, D. G.: Polyrate 2023: A computer program for the calculation of chemical reaction rates for polyatomics. New version announcement, Computer Physics Communications, 294, 108933, https://doi.org/10.1016/j.cpc.2023.108933, 2024.
 - Öberg, T.: A QSAR for the hydroxyl radical reaction rate constant: validation, domain of application, and prediction, Atmos. Environ., 39, 2189–2200, https://doi.org/10.1016/j.atmosenv.2005.01.007, 2005.
- Ohta, T.: Rate constants for the reactions of OH radicals with alkyl substituted olefins, Inter. J. Chem. Kinet., 16, 879–886, https://doi.org/10.1002/kin.550160708, 1984.
 - Palmer, P. I., Marvin, M. R., Siddans, R., Kerridge, B. J., and Moore, D. P.: Nocturnal survival of isoprene linked to formation of upper tropospheric organic aerosol, Science, 375, 562–566, https://doi.org/10.1126/science.abg4506, 2022.
- Qin, J., Wang, X., Yang, Y., Qin, Y., Shi, S., Xu, P., Chen, R., Zhou, X., Tan, J., and Wang, X.: Source apportionment of VOCs in a typical medium-sized city in North China Plain and implications on control policy, J. Environ. Sci., 107, 26–37, https://doi.org/10.1016/j.jes.2020.10.005, 2021.
 - Qiu, Y., Feng, J., Zhang, Z., Zhao, X., Li, Z., Ma, Z., Liu, R., and Zhu, J.: Regional aerosol forecasts based on deep learning and numerical weather prediction, npj Clim. Atmos. Sci., 6, 71, https://doi.org/10.1038/s41612-023-00397-0, 2023.
- Sha, Q., Zhu, M., Huang, H., Wang, Y., Huang, Z., Zhang, X., Tang, M., Lu, M., Chen, C., Shi, B., Chen, Z., Wu, L., Zhong, Z., Li, C., Xu, Y., Yu, F., Jia, G., Liao, S., Cui, X., Liu, J., and Zheng, J.: A newly integrated dataset of volatile organic compounds (VOCs) source profiles and implications for the future development of VOCs profiles in China, Sci. Total. Environ., 793, 148348, https://doi.org/10.1016/j.scitotenv.2021.148348, 2021.
 - Sindelarova, K., Granier, C., Bouarar, I., Guenther, A., Tilmes, S., Stavrakou, T., Müller, J.-F., Kuhn, U., Stefani, P., and Knorr, W.: Global data set of biogenic VOC emissions calculated by the MEGAN model over the last 30 years, Atmos. Chem. Phys., 14, 9317–9341, https://doi.org/10.5194/acp-14-9317-2014, 2014.
- Theodoris, C. V., Xiao, L., Chopra, A., Chaffin, M. D., Al Sayed, Z. R., Hill, M. C., Mantineo, H., Brydon, E. M., Zeng, Z., Liu, X. S., and Ellinor, P. T.: Transfer learning enables predictions in network biology, Nature, 618, 616–624, https://doi.org/10.1038/s41586-023-06139-9, 2023.
- Wang, K., Ge, M., and Wang, W.: Kinetics of the gas-phase reactions of NO3 radicals with ethyl acrylate, n-butyl acrylate, methyl methacrylate and ethyl methacrylate, Atmos. Environ., 44, 1847–1850, https://doi.org/10.1016/j.atmosenv.2010.02.039, 2010.
 - Wang, S., Du, L., Zhu, J., Tsona, N. T., Liu, S., Wang, Y., Ge, M., and Wang, W.: Gas-Phase Oxidation of Allyl Acetate by O3, OH, Cl, and NO3: Reaction Kinetics and Mechanism, J. Phys. Chem. A, 122, 1600–1611, https://doi.org/10.1021/acs.jpca.7b10599, 2018.

- Wells, K. C., Millet, D. B., Payne, V. H., Deventer, M. J., Bates, K. H., de Gouw, J. A., Graus, M., Warneke, C., Wisthaler, A., and Fuentes, J. D.: Satellite isoprene retrievals constrain emissions and atmospheric oxidation, Nature, 585, 225–233, https://doi.org/10.1038/s41586-020-2664-3, 2020.
 - Wells, R., Baxley, S., and Williams, D.: Rate constants and atmospheric transformations of Air Force VOCs, in: Advanced Technologies for Environmental Monitoring and Remediation, Advanced Technologies for Environmental Monitoring and Remediation, 153–160, https://doi.org/10.1117/12.259768, 1996.
- Ku, Y., Yu, X., and Zhang, S.: QSAR models of reaction rate constants of alkenes with ozone and hydroxyl radical, J. Braz. Chem. Soc., 24, 1781–1788, https://doi.org/10.5935/0103-5053.20130223, 2013.
- Zha, Q., Aliaga, D., Krejci, R., Sinclair, V. A., Wu, C., Ciarelli, G., Scholz, W., Heikkinen, L., Partoll, E., Gramlich, Y., Huang, W., Leiminger, M., Enroth, J., Peräkylä, O., Cai, R., Chen, X., Koenig, A. M., Velarde, F., Moreno, I., Petäjä, T., Artaxo, P., Laj, P., Hansel, A., Carbone, S., Kulmala, M., Andrade, M., Worsnop, D., Mohr, C., and Bianchi, F.: Oxidized organic molecules in the tropical free troposphere over Amazonia, Natl. Sci. Rev., 11, nwad138, https://doi.org/10.1093/nsr/nwad138,

2023.

- Zhang, O., Zhang, J., Jin, J., Zhang, X., Hu, R., Shen, C., Cao, H., Du, H., Kang, Y., Deng, Y., Liu, F., Chen, G., Hsieh, C.-Y., and Hou, T.: ResGen is a pocket-aware 3D molecular generation model based on parallel multiscale modelling, Nat. Mach. Intell., 5, 1020–1030, https://doi.org/10.1038/s42256-023-00712-7, 2023.
- Zhang, X., Gao, S., Fu, Q., Han, D., Chen, X., Fu, S., Huang, X., and Cheng, J.: Impact of VOCs emission from iron and steel industry on regional O3 and PM2.5 pollutions, Environ. Sci. Pollut. Res., 27, 28853–28866, https://doi.org/10.1007/s11356-020-09218-w, 2020.
- Zhao, M., Qiao, T., Huang, Z., Zhu, M., Xu, W., Xiu, G., Tao, J., and Lee, S.: Comparison of ionic and carbonaceous compositions of PM2.5 in 2009 and 2012 in Shanghai, China, Sci. Total. Environ., 536, 695–703, https://doi.org/10.1016/j.scitotenv.2015.07.100, 2015.
 - Zhao, Y., Zheng, B., Saunois, M., Ciais, P., Hegglin, M. I., Lu, S., Li, Y., and Bousquet, P.: Air pollution modulates trends and variability of the global methane budget, Nature, 642, 369–375, https://doi.org/10.1038/s41586-025-09004-z, 2025.
 - Ziemann, P. J. and Atkinson, R.: Kinetics, products, and mechanisms of secondary organic aerosol formation, Chem. Soc. Rev., 41, 6582–6605, https://doi.org/10.1039/C2CS35122F, 2012.