



## RESEARCH ARTICLE

10.1029/2021MS002974

# Neural Network Emulation of the Formation of Organic Aerosols Based on the Explicit GECKO-A Chemistry Model

## Key Points:

- Incorporation of explicit organic chemistry into 3D chemistry-simulations requires emulation
- We developed two types of neural network emulators for the Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere chemistry model
- The emulators produced accurate and stable simulations for three precursor species

## Supporting Information:

Supporting Information may be found in the online version of this article.

## Correspondence to:

A. Hodzic,  
alma@ucar.edu

## Citation:

Schreck, J. S., Becker, C., Gagne, D. J., Lawrence, K., Wang, S., Mouchel-Vallon, C., et al. (2022). Neural network emulation of the formation of organic aerosols based on the explicit GECKO-A chemistry model. *Journal of Advances in Modeling Earth Systems*, 14, e2021MS002974. <https://doi.org/10.1029/2021MS002974>

Received 26 DEC 2021  
Accepted 23 JUL 2022

John S. Schreck<sup>1</sup> , Charles Becker<sup>1</sup> , David John Gagne<sup>1</sup> , Keely Lawrence<sup>1</sup>,  
Siyuan Wang<sup>2,3</sup> , Camille Mouchel-Vallon<sup>4</sup>, Jinkyul Choi<sup>5</sup>, and Alma Hodzic<sup>1</sup> 

<sup>1</sup>National Center for Atmospheric Research (NCAR), Boulder, CO, USA, <sup>2</sup>Cooperative Institute for Research in Environmental Sciences (CIRES), University of Colorado, Boulder, CO, USA, <sup>3</sup>National Oceanic and Atmospheric Administration (NOAA), Chemical Sciences Laboratory (CSL), Boulder, CO, USA, <sup>4</sup>Laboratoire d'Aérodynamique, Université de Toulouse, CNRS, UPS, Toulouse, France, <sup>5</sup>Environmental Engineering Program, University of Colorado, Boulder, CO, USA

**Abstract** Secondary organic aerosols (SOA) are formed from oxidation of hundreds of volatile organic compounds (VOCs) emitted from anthropogenic and natural sources. Accurate predictions of this chemistry are key for air quality and climate studies due to the large contribution of organic aerosols to submicron aerosol mass. Currently, only explicit models, such as the Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A), can fully represent the chemical processing of thousands of organic species. However, their extreme computational cost prohibits their use in current chemistry-climate models, which rely on simplified empirical parameterizations to predict SOA concentrations. This study demonstrates that machine learning can accurately emulate SOA formation from an explicit chemistry model with an approximate error of 2%–8%, up to five days for several precursors and for potentially up to one month for recurrent neural network models, and with 100 to 100,000 times speedup over GECKO-A, making it computationally useable in a chemistry-climate model. We generated the training data using thousands of GECKO-A box simulations sampled from a broad range of initial environmental conditions, and focused on three representative SOA precursors: the oxidation by OH of two anthropogenic (toluene, dodecane), and the oxidation by O<sub>3</sub> of one biogenic VOC ( $\alpha$ -pinene). We compare several neural models and quantify their underlying uncertainty and robustness. These are promising results, suggesting that neural network models could be applied to predict SOA in chemistry-climate models, limited however to the range of environmental conditions that were considered in the training datasets.

**Plain Language Summary** Detailed and accurate representation of organic aerosol chemistry is needed to predict the effect of atmospheric aerosols formed from natural and anthropogenic sources on both human health and climate. Ideally, these complex representations of chemistry would be directly included within state-of-the-art weather and climate models to get a fully coupled system with meteorological and climatological feedback all over the globe. However, we are many years away from having the computational power needed to run such fully coupled large-scale simulations due to the complexity of organic chemistry, which involves hundreds of thousands of organic gaseous and particle species and chemical reactions. As a potential solution, we test an approach that uses a neural network to mimic the solution of an explicit representation of organic chemistry which would be computationally feasible to link with current air quality and climate models.

## 1. Introduction

Secondary organic aerosols (SOA) have been an active area of research over the past decades with the goal of improving their representation in air quality and climate models (Hodzic et al., 2016; Tsigaridis et al., 2014), which is essential for predicting their effect on human health (Mauderly & Chow, 2008) and their contribution to radiative forcing in the climate system (Boucher et al., 2013). The misrepresentation of SOA formation pathways in 3D models has led to a long-standing discrepancy between observed and modeled organic aerosol concentrations that has been reported from urban to remote regions (de Gouw, 2005; Hodzic et al., 2020). Unlike sulfate and other inorganic aerosols, which are made from a few dominant chemical pathways, SOAs result from the condensation of a very large number of partly oxidized gases. These gases are generated from the multi-generational oxidation of volatile organic compounds (VOCs) emitted from anthropogenic and natural sources. This complexity is not included in current 3D models that rely on simplified SOA parameterizations that have been developed and

© 2022 The Authors. Journal of Advances in Modeling Earth Systems published by Wiley Periodicals LLC on behalf of American Geophysical Union. This is an open access article under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

optimized based on laboratory measurements or ambient aircraft data (Hodzic & Jimenez, 2011; Ng et al., 2007). This empirical approach does not include the mechanistic understanding of processes leading to SOA formation, and the adequate sensitivity to environmental variables modulating SOA concentration.

Detailed chemistry models such as the widely used “near-explicit” Master Chemical Mechanism (MCM; Jenkin et al., 2003) or the “fully explicit” Generator of Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A; Aumont et al., 2005; Camredon et al., 2007) provide a mechanistic representation of the organic aerosol chemistry and relevant process, and lead to an improved agreement with ambient SOA measurements (Lee-Taylor et al., 2011; Mouchel-Vallon et al., 2020). Chemical mechanisms generated by GECKO-A typically include millions to tens of millions of reactions, and hundreds of thousands intermediate species (Aumont et al., 2005). MCM mechanisms are handwritten and much smaller than GECKO-A ones as they represent only 2–3 first generations of chemistry. Due to the remarkable computational cost, no air quality models or chemistry-climate models can afford to run with GECKO-A included in the foreseeable future. To our knowledge, only the study by (Li et al., 2015; Ying & Li, 2011) attempted to include the MCM organic chemistry mechanism into a regional 3D model and faced computational challenges.

Recent years have seen an increase in the number of applications in developing machine learning emulators using explicit/process-level models and implementing the trained emulators into large-scale models (Beucler et al., 2020; Brenowitz & Bretherton, 2018; Gettelman et al., 2021). Replacing complex processes with ML emulators have potential advantages by learning non-linear relationships that can represent underlying physical or chemical processes not captured in simple empirical characterizations, as well as multiple orders of magnitude speedups in computation when compared to fully coupled process-based models. However, maintaining both an acceptable level of accuracy and a system that remains numerically stable through an adequate amount of time with emulators remains challenging.

Current efforts in atmospheric chemistry emulation with machine learning (ML) have focused on inorganic gas-phase chemistry, such as ozone within GEOS-Chem (Keller & Evans, 2019; Kelp et al., 2018). Using random forest regression and neural network models they were able to reproduce the hourly concentration of 77 gaseous species predicted by the GEOS-Chem chemical mechanism, with a significantly reduced computational expense (250 times). However, the emulator for gas-phase chemistry was subject to runaway errors and numerical instability, as well as performance degradation on out-of-domain inputs. In a follow-up study (Kelp et al., 2020), used a neural network with a recurrent training approach, where a multi-time step loss function was used in conjunction with dimensionality reduction of the chemical system, that resulted in observed reduced error accumulation and provided greater stability. The use of recurrent neural networks (RNNs) was not reported in any of these studies.

Our primary aim in this work is to extend atmospheric chemistry emulation to organic aerosols for which current climate models do not currently account for. Additionally, this proof of concept study evaluates two different types of neural network architectures: (a) a feed-forward, fully connected network and (b) a RNN. Both models were designed to be feasibly integrated into current 3D transport models, that can provide fast and accurate predictions for organic aerosol concentrations. The RNN was chosen to determine if it could help to overcome numerical stability problems, as observed by others using fully connected model architectures (Brenowitz & Bretherton, 2018; Kelp et al., 2020), as they come equipped with feedback connections that can store information about previous events in the form of a latent vector, for example, RNNs possess memory (Hochreiter, 1991). As these architectures were developed to learn representations of sequential data to solve temporal problems, they offer the ability to use multi-length inputs. However, incorporating multi-length inputs into 3D models would require us to dissociate chemistry production from other processes (e.g., transport, removal). We address this with the development of a 1-step RNN that only requires a single time step of input, but relies on a separate simple neural network to initialize the hidden state vector for the very first time-step. A consequence of a 1-step approach is that the RNN will have a larger memory requirement relative to feed-forward, fully connected architectures, due to the presence of the hidden state.

The paper is organized as follows: Section 2 outlines the data generation, training, hyperparameter tuning, and evaluation procedures for the reference model and both neural network types. In Section 3, we characterize the two models' performance relative to the GECKO-A data sets for different precursor species, and compared to each other to assess the overall strengths and weaknesses of each model architecture. Both model types are also tested on data sets that help assess the ability of the models to generalize into new domains. Sections 4 and 5 provide

**Table 1**  
*Environmental Parameter Ranges Used in Generator For Explicit Chemistry and Kinetics of Organics in the Atmosphere Simulations*

Temperature	240–320 K	Uniform
Solar zenith angle (SZA)	0–90°	Uniform
Pre-existing aerosols	0.01–10 $\mu\text{g}/\text{m}^3$	Logarithmic
Ozone	1–150 ppb	Uniform
NO <sub>x</sub>	0.01–10 ppb	Logarithmic
OH	10 <sup>1</sup> –10 <sup>6</sup> molecules/cm <sup>3</sup>	Uniform

a brief discussion about the pros and cons of these different model architectures and the ongoing challenges regarding numerical stability, computational cost, and interpretability for emulating SOA production with ML.

## 2. Methods

### 2.1. Description of the Reference Model

To provide reference chemical mechanisms, we used the GECKO-A chemical generator (Aumont et al., 2005; Camredon et al., 2007), which describes in great details the chemical oxidation of organic compounds in the atmosphere. The resulting chemical mechanisms for each SOA precursor species are complete (down to the ultimate products CO<sub>2</sub> and H<sub>2</sub>O), and explicit by

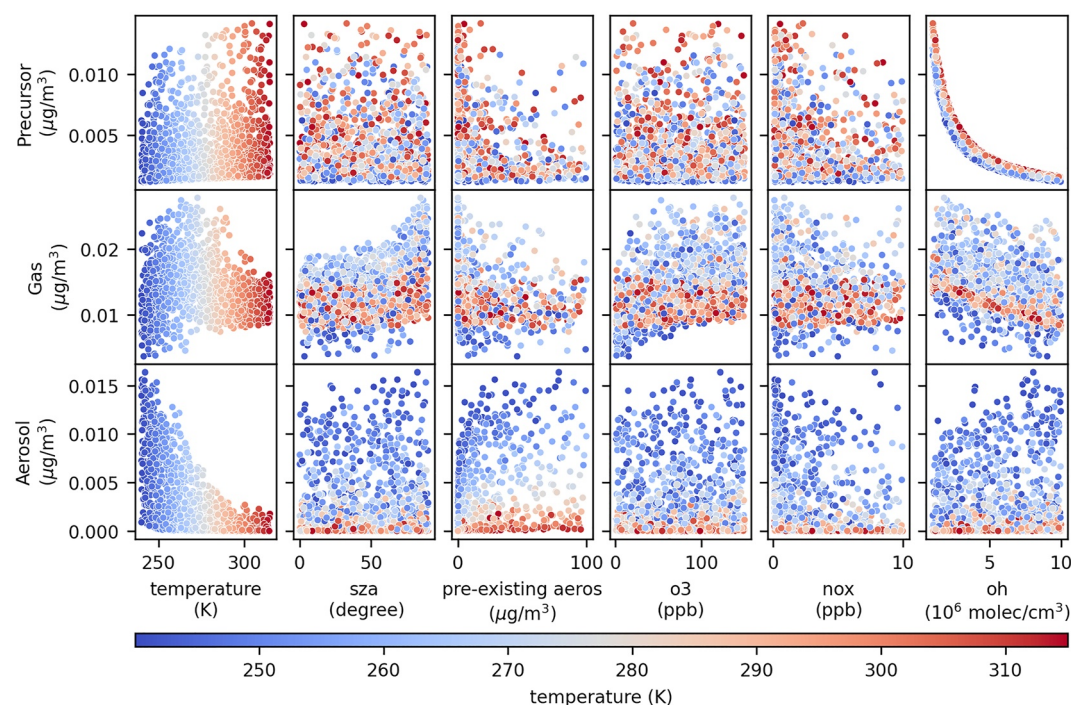
preserving knowledge of the molecular structures of all the intermediate compounds. Compared to other widely used semi-explicit chemical models such as MCM, GECKO-A can consider many generations of oxidative chemistry that is, 20 generations are considered here (Lee-Taylor et al., 2011). This has important implications for the formation of organic aerosols as SOA formation arises from a multitude of partly oxidized compounds, rather than from a few dominant molecules. In addition, the SOA formation timescale varies greatly for different precursors. For example, at ambient conditions, it takes only a few hours to form SOA from dodecane versus several days to form SOA from toluene (Hodzic et al., 2014). For the gas-particle partitioning of organic molecules, dynamic partitioning is used. GECKO-A assumes a bulk aerosol size distribution. It is reasonable to consider GECKO-A simulations as a benchmark for building an emulator for SOA chemistry given its reasonable agreement with observations shown for both comparisons with chamber measurements for example, for alkanes and alkenes compounds (La et al., 2016) and ambient measurements for example, during MIRAGE, BEACHON and Go-AMAZON (Lee-Taylor et al., 2011, 2015; Mouchel-Vallon et al., 2020).

### 2.2. GECKO-A Generated Training Datasets

We ran the GECKO-A model for a variety of initial conditions and physical parameters to generate the data set used to train the machine learning emulator. At this stage we focus on three representative SOA precursors: toluene, dodecane, and  $\alpha$ -pinene. Toluene and dodecane are emitted from a wide range of anthropogenic sources, while  $\alpha$ -pinene is one of the major SOA precursors emitted by vegetation. These compounds together contribute substantially to the global SOA burden. We generate chemical mechanisms and corresponding datasets for these precursors with their dominant oxidants, including the OH oxidation mechanisms of toluene, dodecane, and the ozonolysis mechanism of  $\alpha$ -pinene. Reactions of the precursor with nitrate radical (NO<sub>3</sub>) were not considered. Thus, the considered chemistry is mostly representative of daytime conditions. This is a reasonable assumption for the considered precursor species (Fry et al., 2014).

Based on our current understanding of atmospheric chemistry and the common chemistry-climate modeling frameworks, we identified the following six variables that are key to predicting SOA formation from VOC oxidation under tropospheric conditions: (a) temperature, (b) solar zenith angle, (c) pre-existing aerosol mass, (d) ozone concentrations, (e) nitrogen oxides (NO<sub>x</sub>) concentrations, and (f) OH radical concentrations. The range of variability considered for these parameters and the associated sampling scheme is summarized in Table 1 and illustrated in Figure 1. Additionally, a diurnal variation in temperature of an average 5° amplitude is used. In each training data set, we use the Latin Hypercube sampling approach to obtain two-thousand environmental input combinations. Temperature, solar zenith angle, ozone, and OH were all sampled uniformly, whereas pre-existing aerosol concentrations and NO<sub>x</sub> were sampled as a logarithmic distribution. The combination of these ranges is relevant for a wide range of tropospheric conditions, from remote to moderately polluted environments. These environmental variables, with the exception of temperature, are held constant in a given 5-day GECKO-A box model simulation in an attempt to coerce the ML models to generalize better and be more robust to over fitting.

In each data set, the initial concentrations of a given SOA precursor were set to an arbitrary low value of 10 ppt similar to previous studies (Hodzic et al., 2014, 2015; Lannuque et al., 2018). Although the tropospheric concentrations of the precursor can be higher in polluted regions, this low value is representative of the remote atmosphere, and was chosen so that the amount of aerosol produced from the given precursor is small compared to preexisting OA and will not impact the gas/particle partitioning, nor the overall photochemical reactivity.



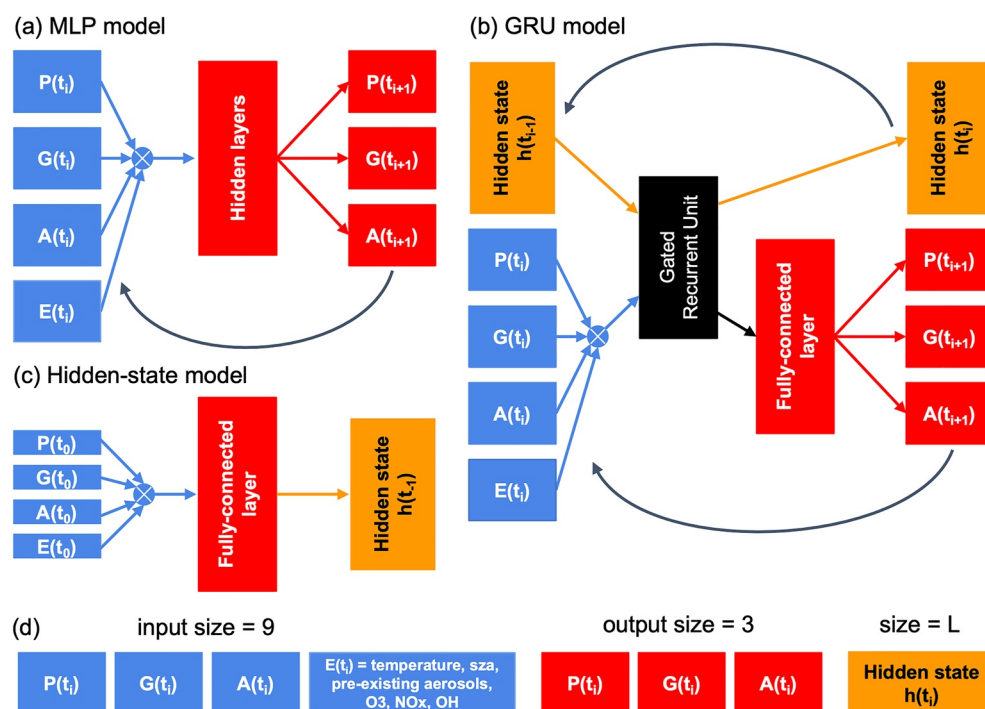
**Figure 1.** Training distributions (mean value through 5-day) of targets (rows) versus environmental variables (columns) for toluene.

Precursor spans several orders of magnitude in our 5-day GECKO-A simulations as it decays exponentially and is effectively consumed before the end of each simulation leading to the production of thousands of intermediate organic gases. During a given simulation, the levels of oxidants are kept at realistic concentrations representative of ambient atmospheric conditions as in Hodzic et al. (2015). As shown on Figure 1, complexities of the organic chemistry are illustrated by the wide variation of produced SOA mass with respect to each environmental variable. For example, significantly higher SOA mass concentrations are produced at colder temperatures.

As the precursor mass is exponentially distributed, before using the precursor data as input to the ML model, we transform the precursor values by taking the base-10 logarithm to avoid any stiffness in the system. Next, each input variable  $X_j$  in the training data, including chemical concentrations and environmental variables, were standardized independently into z-scores according to the formula  $X_j = (X_j - u)(X_j - s)^{-1}$ , where  $u$  and  $s$  are the mean and standard deviation of  $X_j$ . Standardization was chosen over other common transformations because many features are not normally distributed in the data sets for the three species. Hence the transformation recasts the values of the input and output channels into a format where the values of each variable are centered and have similar spread. This is especially important when computing the error on the model predictions against the training data values when the weights in the model are being updated, for example, to prevent the model from over-fitting on the quantities having the largest spread.

A total of 2,000 experiments were run in GECKO-A for each precursor species and output every 300 s (5 min) over the course of five days. This sets the model time-step prediction at 300 s. Thus, a total of 2,880,000 samples were generated per species. However, as the target variables are a subset of the input feature variables at the previous time step, we removed the first (final) time step of the output (input) variables from each experiment leaving a total of 2,878,000 samples.

Each GECKO-A experiment run should be independent from the others, thus we split the 2,000 total experiments into training (80%), validation (10%), and testing (10%) sets by run, so all time steps from each run are kept in one of the three sets. The training set is used to optimize the weights of an ML model, while the validation set is used to select the hyperparameters, or meta-settings describing the ML model architecture, such as the number of neurons in the hidden layer, that result in the best-performing model across the space of possible models (see



**Figure 2.** Model architectures showing the (a) multi-layer perceptron and (b) Gated Recurrent Unit models for predicting concentrations, and (c) the model for the initialization of the recurrent hidden state given an initial condition as input. In all three models, inputs at time  $t_i$  are colored blue (precursor, gas, and aerosol, and the environmental variables). The inputs are concatenated into a vector as illustrated by the cross-marked blue dot. The model outputs at time  $t_{i+1}$  are colored red in (a) and (b) for precursor, gas, and aerosol, and orange in (b) and (c) for the hidden state. In (a) and (b) when the models are used in box simulations, the black arrow represent using the output prediction from the model along with the environmental variables as input to the model for the next time step prediction. (d) The size of the input, output, and the hidden state (denoted  $L$ ). The input  $E(t_i)$  is of size 6.

below). The final testing set is not used to train or adjust the machine learning models and is only used in the final evaluations described herein.

### 2.3. Neural Network Models

Figure 2 shows several neural network models that we consider here as emulators for GECKO-A. In Figure 2a, a multi-layer perceptron (MLP) architecture is shown (referred to as the “MLP model”). The MLP model accepts as input the scaled values of the precursor, gas, aerosol, and environmental variables at time  $t_i$  (blue boxes in the figure). These quantities are concatenated into a vector, and denoted  $X(t_i) \equiv (P(t_i), G(t_i), A(t_i), E(t_i))$ . The model outputs the precursor, gas, and aerosol values at time  $t_{i+1} = t_i + \delta t = t_i + 300$  s. The MLP is an artificial neural network that contains, in addition to input and output layers, at least one hidden layer (horizontal red block in Figure 2a), and is the simplest neural network architecture available. Each hidden layer contains a set of perceptrons, which mathematically are linear regressions on the outputs of the previous layer followed by a non-linear transformation called the activation function. For our MLP, we use the Rectified Linear Unit (ReLU) activation function, which sets negative values to 0 and allows positive values to pass through unchanged. The final hidden layer is connected to the output layer, which is a linear regression on the hidden layer outputs. We tested using multiple hidden layers but found that a single hidden layer produced the lowest validation set error. The MLP model for each of the precursor species make future predictions for the values of the chemical quantities of interest, using only the current chemical state of the atmosphere. As such, any MLP model satisfies the Markovian condition, and possesses no memory about atmospheric chemical states visited in the past beyond the previous time step. GECKO-A uses concentrations from previous times steps to calculate the new concentration at  $t+1$ .

As GECKO-A generates sequential trajectories describing the evolution of species' concentrations over time, we have also applied a RNN model to investigate whether it may have an advantage over the MLP model due to its ability to utilize its memory about the past at  $\{t_{i-n}, \dots, t_{i-1}\}$  to make the next prediction at  $t_i$ . As such, a temporal model may be better equipped compared to the state-less MLP model to describe the changes occurring to the quantities over time, in addition to limiting and/or preventing runaway error propagation, but will also be slower due to the larger input and output requirements of RNNs. The input to an RNN is a sequence and a hidden state. The sequence elements could be scalars, vectors, or other higher-dimensional tensors. The first (or last) element in the sequence is ingested by the RNN along with a starting hidden state, producing an encoding of the element and a hidden state for the current encoding. The next element in the sequence is then used as input along with the current hidden state, which produces an encoding of the second element and another hidden state. This encoding represents not just the second element, but the elements that came before it, due to the fact that the model leverages its feedback connections to produce the encoding. This process continues until all of the elements in the sequence have been seen by the model, which produces a final encoding of the entire sequence, as well as a hidden state.

Figure 2b illustrates a model architecture that combines an RNN layer type called a Gated Recurrent Unit (GRU) (Chung et al., 2014b) with a fully connected layer to make chemical concentration predictions. We refer to this model as the GRU model. The GRU layer performs three key operations: filtering the contents of the hidden vector from the previous time, calculating a new hidden state from a combination of the filtered hidden state and new inputs, and finally calculating a new output. The GRU is similar to the well-known long-short term memory (LSTM) model (Hochreiter & Schmidhuber, 1997) but has fewer parameters to learn. Even so, the GRU often performs comparably to the LSTM in language modeling tasks (Chung et al., 2014a).

Similar to the MLP, the GRU model shown in Figure 2b illustrates the model at the  $t_i$  time step such that  $t_{i-1}$  came before it, accepting as input  $(P(t_i), G(t_i), A(t_i), E(t_i))$  at time  $t_i$  and hidden state  $h(t_{i-1})$ . The GRU layer produces an encoded representation of the input  $X(t_i)$  denoted  $Y^*(t_i)$ , and a hidden state  $h(t_i)$  for the current time, which has the same dimension as  $h(t_{i-1})$ . The GRU layer avoids the vanishing gradient problem by computing these quantities according to

$$\begin{aligned} z(t_i) &= \sigma(W_z h(t_{i-1}) + U_z X(t_i) + b_z) \\ r(t_i) &= \sigma(W_r h(t_{i-1}) + U_r X(t_i) + b_r) \\ c(t_i) &= \tanh(W_c (r_c \odot h(t_{i-1})) + U_c X(t_i) + b_c) \\ h(t_i) &= z(t_i) \odot h(t_{i-1}) + (1 - z(t_i)) \odot c(t_i) \\ Y^*(t_i) &= \text{softmax}(W_y h(t_i) + b_y) \end{aligned}$$

where the sigmoid function  $\sigma(x) = (1 + \exp^{-x})^{-1}$  projects input values to be within  $[0, 1]$ . The softmax function for a K-component  $z$  is  $\exp(z_k) / \sum_{j=1}^K \exp(z_j)$ . The tensors  $W_*$  and  $U_*$ , and bias terms  $b_*$ , contain the fit parameters that are updated through back-propagation during training. The  $\odot$  symbol refers to element-wise multiplication.

The quantity  $z(t_i)$  is the update gate, and  $r(t_i)$  is the reset gate, which determines how much information over past time steps to forget. The quantity  $c(t_i)$  represents the current memory content in the layer and utilizes the reset gate to store information from the past. The equation for the current hidden state  $h(t_i)$  uses the update gate to determine how much information from the current time step to collect and how much to collect from past time steps. The encoded information at the current time step,  $Y^*(t_i)$ , is computed using the current hidden state. Then, as Figure 2b illustrates, it is passed through a ReLU activation function (black arrow), and then through a fully connected layer (tall red box), which outputs the scalar precursor, gas, and aerosol values at the time  $t_{i+1}$ .

It is common that the input sequence to an RNN is longer than length one, in which case a hidden state can be immediately informed by the sequence (trajectory) to make the next prediction. Applied to GECKO-A trajectories, we are free to choose the length of the input sequence, which does not have to be fixed. Indeed, in our GECKO-A box model simulations concentrations of organic species are only undergoing chemistry processing, whereas in 3D models other processes (e.g., transport, dry and wet removal) are included. Thus, we must also consider the added complexity of incorporating RNNs into 3D transport models. For a RNN that uses a sequence

for input, each chemical variable needs to be transported and stored for the number of time steps needed as input, which could be memory intensive and programmatically challenging.

Here we describe a “1-step” approach, which means that when incorporated into a 3D climate simulation, the RNN is similar to the MLP in that only a single time step of input is required. The only difference being that a single hidden state is also input to the RNN, which can be understood mainly as a larger input compared to the MLP. Unlike the MLP, the RNN hidden state vector needs to be stored at every model grid cell to inform the calculation of the next time step. Depending on the size of the hidden state, this could create a large memory burden on the simulation but would be less disruptive to simulation codes than creating and managing multiple copies in time of every model field.

When there is no initial hidden state available, the initial condition  $X(t_0)$  is passed through a separate MLP, referred to as the hidden-state model, to obtain  $h(t_{-1})$ . Figure 2c illustrates that this model's architecture accepts as input values at some  $t_0$  for precursor, gas, aerosol, and the environmental variables ( $P(t_0)$ ,  $G(t_0)$ ,  $A(t_0)$ ,  $E(t_0)$ ), and outputs a hidden state  $h(t_{-1})$ . The hidden-state model contains one fully connected layer with a linear activation. The input size is equal to the length of  $X(t_0)$  while the output size is equal to the length of the hidden state used by the GRU.

#### 2.4. Training Procedure

Each MLP model is trained by first initializing an architecture and the trainable weights. The training data split is randomly shuffled removing the time sequence in the data. A fixed number of training data points is selected from the training data, called a batch, and then passed through the model to obtain a prediction for each point in the batch. The mean-absolute error (MAE), the training loss, is computed for the batch from the prediction. Using the loss, the weights of the model are updated accordingly using gradient descent with back-propagation (Rumelhart et al., 1986), and a pre-specified learning rate to reduce the error. This process is repeated until all of the training samples are passed through the model once, and is referred to as one epoch of model training. At the end of every epoch, the training data is randomly shuffled. This procedure is repeated for a prescribed number of epochs.

In order to train the GRU model and the hidden-state model, the input and output data for each experiment needs to be ordered by time, thus it is not shuffled along this coordinate, as is done with the MLP. The training procedure is then similar to how the model would be used in evaluation, and starts by setting the initial condition for an experiment along with the environmental variables as the initial input,  $X_0$  to the model. As there is no hidden state available at the beginning of a box simulation,  $X(t_0)$  is passed through the hidden state model to produce  $h(t_0)$ , then  $(X(t_0), h(t_{-1}))$  is passed through the GRU model to obtain the prediction  $Y(t_1)$  and  $h(t_0)$ . For the GRU model, we use the Huber formula as the loss function for the predicted chemical concentrations, which computes the mean-squared difference between the predicted output  $Y(t_1)$  and the known values produced by GECKO-A, when the difference is greater than a fixed cutoff value, and computes the MAE otherwise.

Next, the predicted output  $Y(t_1)$  is concatenated with the environmental variables one time step into the future to create the input to the model  $X(t_1)$ . This quantity is passed through the hidden state model to obtain  $h^*(t_0)$ , where the \* notation is used to distinguish this hidden state prediction from the one that the GRU model predicted. The mean absolute difference between  $h(t_0)$  and  $h^*(t_0)$  is computed. The total loss for the time step adds this quantity, multiplied by a loss weight, to the loss contributions computed for the chemical quantities. The loss weight is left as a parameter to be optimized (see below). The total loss for the time step is then used with a given learning rate value to update the adjustable parameters of both the GRU and the hidden state models in tandem. This procedure is repeated until the second-to-last time step in an experiment trajectory is used as input to the model.

During training of the GRU model, we select random experiments to create batches of data when computing the total loss at each time step. One epoch is defined as all training experiment trajectories passing through the model once, so the same data as with training the MLP model, except that the data is ordered by time (and randomized by experiment). At the end of every epoch, the model is put into evaluation mode and used to predict the trajectories for the validation experiments. The MAE is then computed between the model predictions and the validation experiments. After each epoch the validation MAE is used to measure improvement of the model predictions in two ways: (a) to anneal the learning rate if the models performance does not improve after some number of epochs, for example, it “over-fits” on the validation experiments, and (b) to stop the training entirely

once the model does not improve on the validation experiments after some number of epochs. We chose 3 and 7 epochs in (a) and (b) respectively.

## 2.5. Evaluation Procedure

The ability of MLP and GRU-based models to predict the time evolution of precursor, gas, and aerosol mass concentrations is evaluated by comparing the box model predictions against the benchmark values as produced by the GECKO-A model. Here each model is placed into evaluation mode, which disables any stochastic components such as the recurrent dropout used when training the GRU, and is used to make predictions on the hold-out validation set of data that was not used during the training to influence the weight updates in each model. For each validation experiment, a starting amount of precursor, gas, aerosol, and environmental variables at time  $t_0$  is passed through the MLP network to obtain the predicted quantities at the first time step  $t_1$ , where  $t_1 = t_0 + \delta t$ . These predictions are then used along with the environmental variables at the next time step as the next input to the model, and so forth for the length of the experiment (see Figure 2b).

### 2.5.1. Ensembles

Machine learning models must be stochastically initialized with random weights, as gradient descent requires variation amongst the weights to perform an initial adjustment that is not uniform across the weights. As a result, two models trained with the exact same basic architecture (without setting a random seed for initialization) will yield a different set of weights and biases, and thus can have a different result during evaluation. Generally, if a model has sufficient data and ample time for training, identical models will converge to similar values, and differences in performance may be small and/or negligible resulting in a robust model. However, for transport or propagation problems that require the input from a prior model prediction in order to make a future prediction, these small differences may accumulate through time and quickly become non-negligible. To further evaluate the robustness or sensitivity to the initialization process, we trained and evaluated 30 ensemble members for each precursor model.

## 2.6. Hyperparameter Optimization

At different stages in training an emulator model, from data post-processing to selecting an architecture, there are hyper parameters that need to be set that can affect the performance outcome of a trained model. They may include, for example, the learning rate used to update the model weights during training, or the size and number of the hidden layers in an MLP or GRU model. As the main objective is to minimize the difference between the model predictions and the test experiments in box simulations, we want to understand how the models performance depends on the hyper parameter choices. From such an understanding, an informed choice can be made in selecting potentially optimal parameter values.

To estimate such a dependency, we use the package Earth Computing Hyperparameter Optimization (ECHO) developed by the authors at NCAR (Schreck & Gagne, 2021), and perform hyper parameter optimization given an objective metric for the three species for both MLP and GRU models. The objective metric for both the MLP and GRU models is the box MAE on the validation holdout set. With the MLP model, a box simulation begins with the initial precursor concentration at  $t_0$ , while for the GRU model the MAE for box simulations is computed for a set of starting times  $\{t_0, t_p, \dots, t_j\}$  and added together, to also test the hidden-state model on different initial precursor amounts.

MLP and GRU models were optimized with ECHO for the three species, with the outcomes described in Section S1 in Supporting Information S1. Tables S1 and S2 in Supporting Information S1 list the best hyperparameters found in each optimization study for the two models. Using the best hyperparameter set, an ensemble of 30 models were trained, where each model had a different random weight initialization. See Section SII in Supporting Information S1 for more details.

Although it is rather trivial to train a model to output realistic predictions one time step ahead from the truth (primarily due to high auto-correlation), limiting cumbersome error accumulation when propagated through time is highly sensitive to model parameters in complex problems such as this. We found that efficient hyperparameter searches were crucial for finding models that could successfully stabilize and limit error accumulation through the length of the simulation. For further details, see Section SI in Supporting Information S1.



**Table 2**  
Table of Computed Metrics for Multilayer Perceptron (MLP) and Gated Recurrent Unit (GRU) Models, for Each of Toluene, Dodecane, and  $\alpha$ -Pinene

	Toluene		Dodecane		$\alpha$ -pinene	
	$R^2$ , $R_e^2$	RMSE, HD	$R^2$ , $R_e^2$	RMSE, HD	$R^2$ , $R_e^2$	RMSE, HD
MLP prec	0.972, 0.994	1.112, 0.763	0.778, 0.996	2.591, 0.240	-18.01, 0.765	20.56, 1.529
MLP gas	0.946, 0.978	1.389, 4.471	0.905, 0.930	7.740, 3.343	0.942, 0.966	3.973, 2.824
MLP aero	0.848, 0.870	1.264, 14.89	0.782, 0.785	10.20, 16.10	0.879, 0.888	3.813, 15.31
GRU prec	0.985, 0.997	0.890, 1.812	0.844, 0.967	3.496, 2.366	0.468, 0.942	4.367, 2.482
GRU gas	0.979, 0.987	0.871, 1.198	0.969, 0.980	4.479, 1.994	0.982, 0.988	2.223, 0.653
GRU aero	0.948, 0.974	0.731, 10.50	0.923, 0.952	6.076, 8.302	0.952, 0.968	2.393, 9.131

*Note.* The average  $R^2$  coefficient for ensemble members and that for the ensemble mean (denoted  $R_e^2$ ), root-mean squared errorD, and Hellinger distance are listed for the precursor, gas, and aerosol prediction tasks. All reported metrics for both models were computed using the testing hold-out set of experiments.

### 3. Results

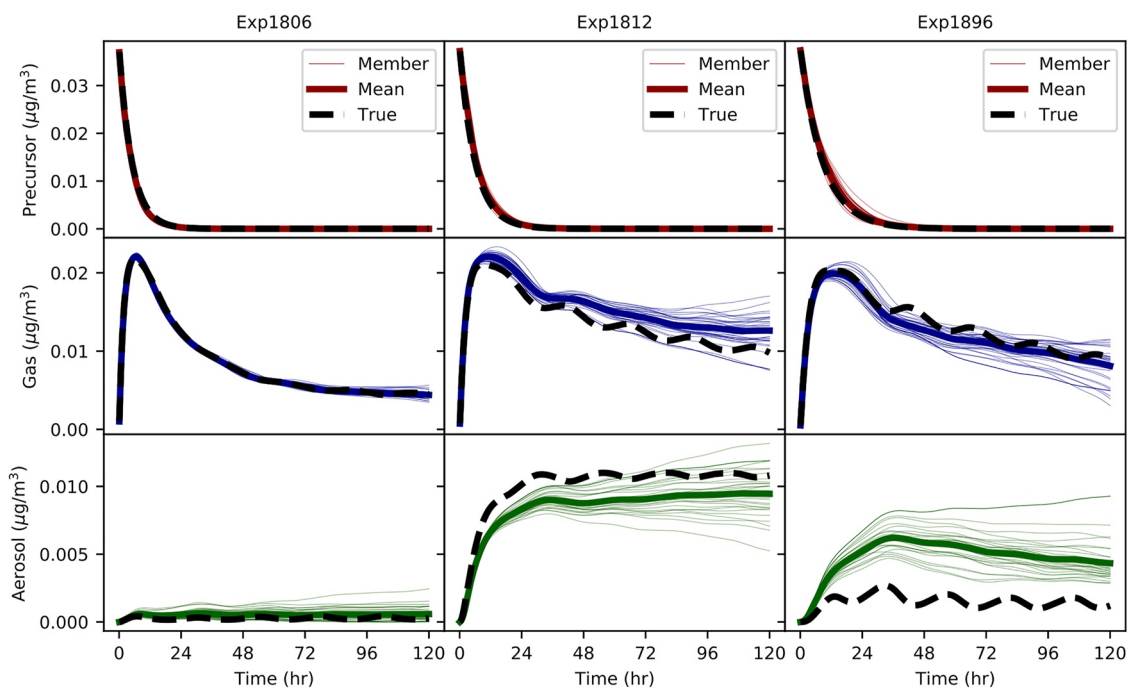
#### 3.1. Performance of Trained MLP and GRU Models

Table 2 lists the bulk validation performance metrics for the MLP and GRU models for toluene, dodecane, and  $\alpha$ -pinene. Each model was trained to predict precursor, gas, and aerosol quantities at 300 s (5 min) time steps. The metrics are the coefficient of determination ( $R^2$ ) defined by Equation B1, the Hellinger distance (HD) defined by Equation B2, and the root-mean squared error (RMSE), which are computed for each prediction task, and the number of unstable or runaway experiments observed. Here, an experiment is considered as unstable when predicted values exceed  $1 \mu\text{g}/\text{m}^3$  which corresponds to an unrealistic formation yield from 10 ppt of initial precursor. Unstable experiments were not used in the computed metrics reported below. It should be noted that both GECKO-A and the ML models assume a bulk size distribution for the representation of aerosols.

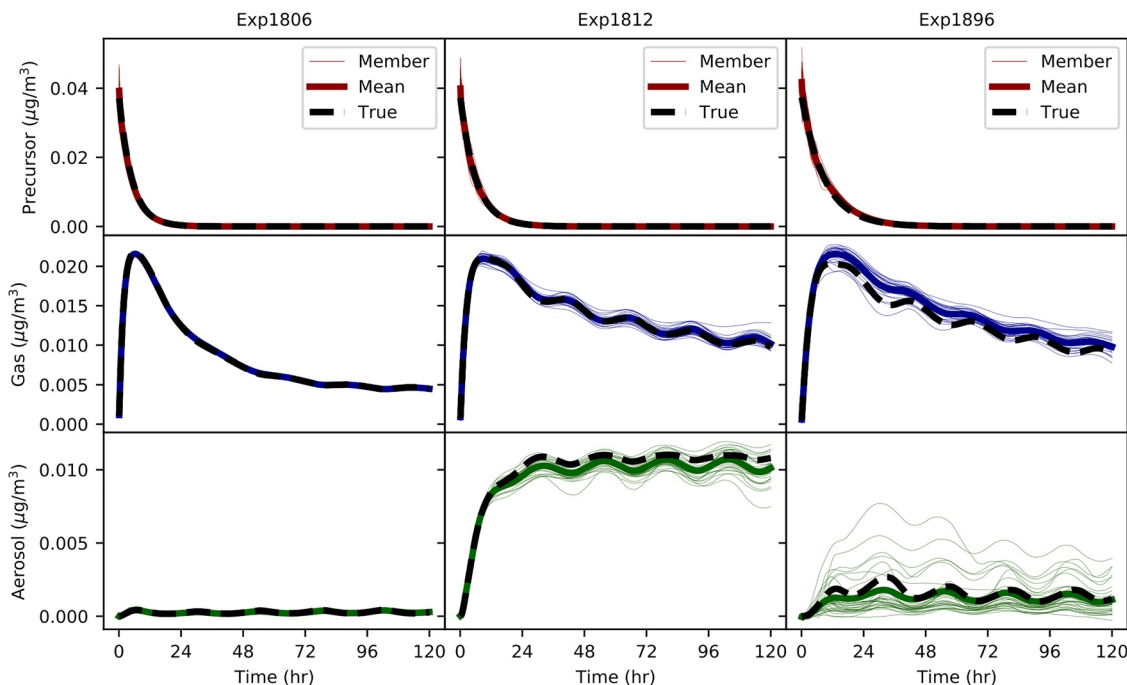
Figures 3 and 4 illustrate the MLP and GRU models' performance on reproducing the experimental data for three experiments selected from the test set of toluene experiments. Overall, they show that both the MLP and GRU models can predict experiment trajectories that resemble GECKO-A ones within a factor of two, which is within the uncertainty of the SOA predictions of GECKO-A (Valorso et al., 2011). In particular, the predicted concentrations for the ensemble mean matched closely with the GECKO-A trajectories for the three prediction tasks, for both models, having similar magnitudes and curvatures at different points in time. The predicted concentrations for precursor decayed exponentially in accordance with the training data, while the gas and aerosol concentrations generally showed qualitative agreement with the GECKO-A trajectories, with the GRU doing better at capturing the diurnal fluctuations. For toluene, Table 2 shows that all of the model prediction tasks led to  $R^2$  coefficients greater than about 0.95 for the individual ensemble members and greater than 0.97 for the ensemble mean. Additionally, the HDs for each task are all low for both MLP and GRU models, indicating that the temporal distributions predicted for different initial conditions matched closely with those generated by GECKO-A.

For the other precursor species, each model degraded in performance in different ways. The GRU model mainly did not perform as well at predicting precursor. For example, the  $R^2$  coefficient was 0.468 for  $\alpha$ -pinene compared with 0.985 for toluene. The HD for the other two species also modestly increased compared to that for toluene as did the RMSE. On gas and aerosol predictions, Table 2 indicates that the GRU performed about the same for all three species according to these three metrics, and that none of the predicted numerical values became unstable during the box simulations.

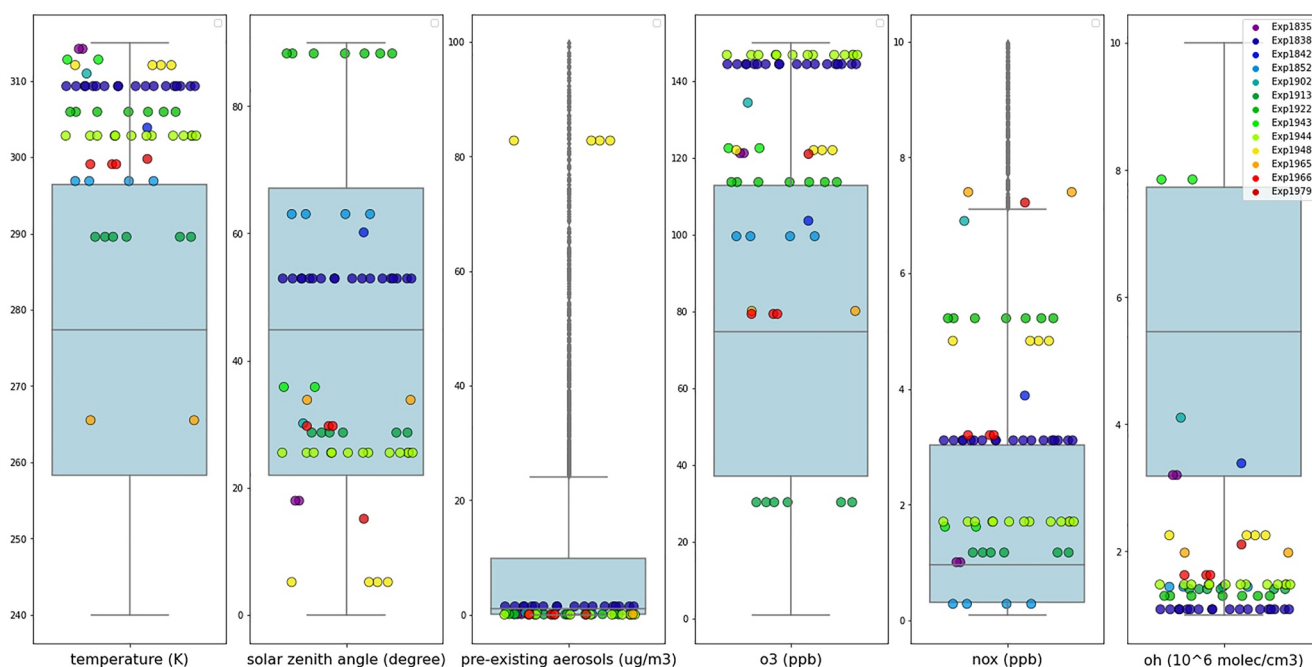
Table 2 shows the MLP model performance on precursor prediction was the best for toluene which had an  $R^2$  value of 0.972 and  $R_e^2$  of 0.994. However, the MLP ensemble members struggled by comparison for both dodecane and  $\alpha$ -pinene, where the  $R^2$  value was 0.778 and -18.055, respectively, with the later value being driven negative by a single ensemble member. By comparison,  $R_e^2$  was much higher, at 0.996 for dodecane and 0.765 for  $\alpha$ -pinene. On the gas and aerosol predictions, the MLP model was mostly consistent for the three species, with gas prediction performing better by comparison to aerosol prediction. Furthermore, out of 200 experiments that



**Figure 3.** Three toluene sample experiment trajectories for the multilayer perceptron model. Solid (thick) lines show the mean Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A) trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.



**Figure 4.** Three toluene sample experiment trajectories for the gated recurrent unit model. Solid (thick) lines show the mean Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A) trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.



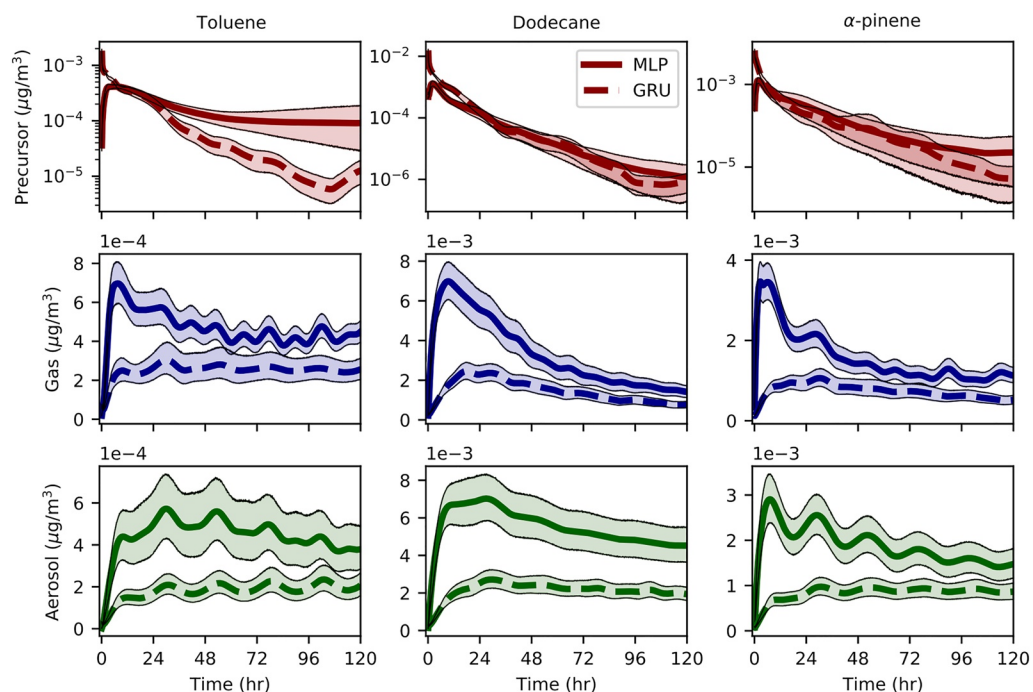
**Figure 5.** Unstable dodecane simulations. The blue box (IQR) and whisker plots ( $Q3/Q1 \pm IQR \times 1.5$ ) represent the training data initial conditions distribution for each input parameter. The colored dots represent where in the input space (vertical) the unstable experiment is. Different experiments are coded by color and different ensemble members for each experiment spread out horizontally.

went unstable during a box simulation there were 13 for dodecane and 2 for  $\alpha$ -pinene, despite the fact that the  $R^2$  score for dodecane remained high across the prediction tasks.

### 3.2. Model Stability

In the training process, we were not able to add additional constraints on the mass balance between the precursor, and its gaseous and aerosol products as the GECKO-A model does not keep track of the total mass. Ensuring mass and energy conservation can prevent some forms of model drift and instability (Yuval & O’Gorman, 2020), however the conservation does not guarantee numerical model stability (Gettelman et al., 2021) or a better performing model (Beucler et al., 2021). Other techniques, such as stochastic parameterization (Gagne et al., 2020), can prevent instability due to accumulating errors by sampling from the full distribution of possible tendencies. We also tried training the MLP models similarly to how the GRU was trained, as described in Section 2.4. However, we were unable to obtain trained MLP models with comparable performance as listed in Table 2.

Although only a very small percentage of total experiments were considered unstable (1.07% dodecane, 0.68%  $\alpha$ -pinene), it is useful to examine the environmental conditions for which they went unstable. Figure 5 shows all 64 of the 6,000 dodecane runs that went unstable in the environmental parameter space. A primary takeaway is that the vast majority of unstable runs happened in conditions where the pre-existing aerosols were extremely low, in conjunction with very low levels of the primary oxidation mechanism, OH. This intuitively makes sense, as the presence of existing aerosols is fundamental to allow the condensation of SOA, and it is difficult for the ML models to capture the very small production from GECKO-A as a result of tiny amounts of pre-existing aerosols. In a situation where this model were coupled to a 3D climate model, it could be circumvented by only running the neural network if a certain threshold of pre-existing aerosols were exceeded in a specific grid cell as tiny amounts would lead to a negligible amount of aerosol production. Experiment 1948 (in yellow) is the only experiment that had members become unstable when the pre-existing aerosol amount was relatively high ( $83 \mu\text{g}/\text{m}^3$ ), which could be due to three of the six inputs being near the tail of the distribution (Temperature, SZA, and pre-existing aerosols) which saw limited training samples from and likely made it difficult to learn. The unstable data for  $\alpha$ -pinene shows similar trends and can be seen in Figure C3. Toluene remained stable for all simulations.



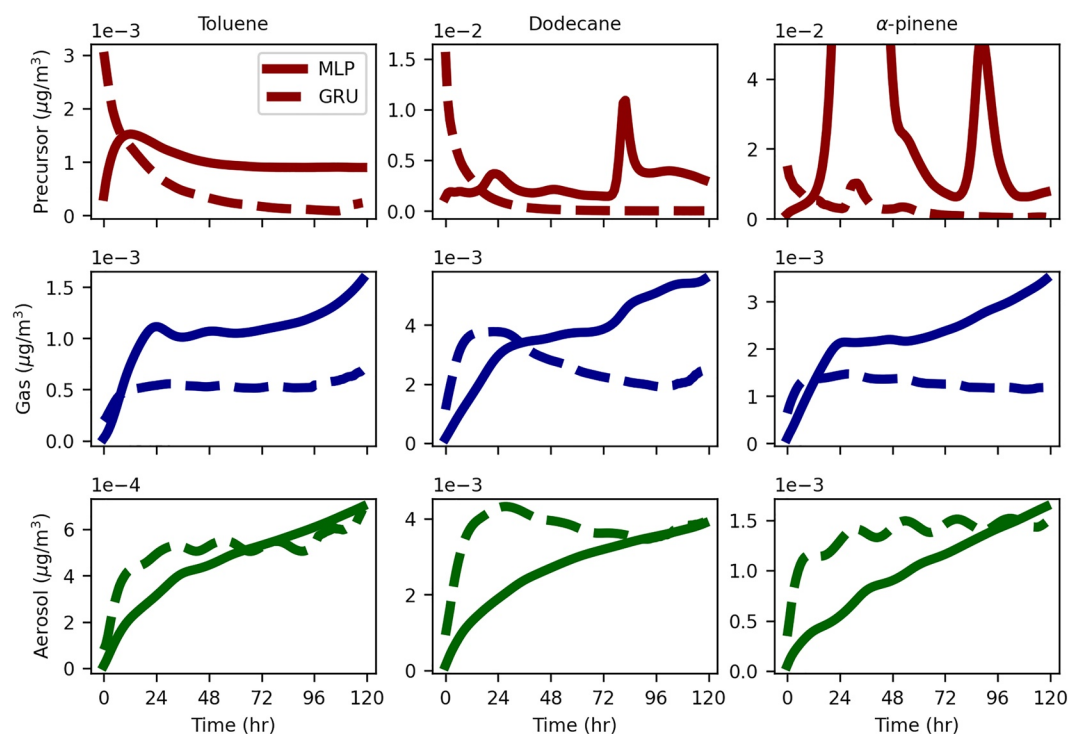
**Figure 6.** Ensemble bootstrapped continuously-ranked probability score through time. Multilayer perceptron (solid) and gated recurrent unit (dashed) with the shaded regions representing the 95% confidence interval.

Finally, we ran box simulations for both models for an extended period of time (up to 1 month) to test the stability for the toluene precursor. Figures C1 and C2 show the same experiments as in Figures 3 and 4 for MLP and GRU models, respectively, for the 30-day simulations. Figure C1 shows that the MLP model remains stable over the 5-day training period, but at later times the ensemble predictions for gas and aerosol quantities go negative and hence have become unphysical. The GRU ensemble prediction remains positive for the examples shown in Figure C2, however greater numbers of individual ensemble members are observed to make unphysical predictions at times beyond 5 days. This time range is compatible with the typical lifetime of organic aerosols in the atmosphere which is about a week (Hodzic et al., 2016).

### 3.3. Quantification of Model Variances

Figure 3 shows that the predicted aerosol quantities for the ensemble suite begins to diverge as the simulation time progresses, whereas the GRU variation (see Figure 4) appears to be roughly constant or improving as time progresses. In order to capture how well the models are predicting time-dependent quantities, we also computed the bootstrapped continuously ranked probability score (CRPS) in Figure 6 and the mean standard deviation (Figure 7) between the different species across all 30 ensemble members. Figure 6 shows the CRPS (lines) and the 95% bootstrap confidence interval (shaded areas) changing in time for the two model types. For the three species, the MLP model has a lower CRPS on all predictions at early times, then the GRU at later times. Figure 6 shows the two models' CRPS values for precursor crossing typically within one simulation day, with the CRPS for the GRU starting out relatively high compared to the MLP, but then quickly declining. Except at the earliest simulation times, the GRU had a lower CRPS as well as a smaller confidence interval on the gas and aerosol prediction tasks and stayed flat or declined.

Figure 7 shows a notable difference in the variation among ensemble members between the MLP and GRU, specifically the slope of the gas and aerosol trajectories. The steep positive slope of the MLP demonstrates the inherent growing uncertainty in the model itself as it progresses further from the starting condition, which is also seen in the ensemble spread on Figure 3. The GRU has a much flatter trajectory, especially in the later time steps due to the short-term memory of the trajectory being encoded into the hidden state, which could potentially make it much more suitable for maintaining stability in much longer running simulations. Additionally, if it is



**Figure 7.** Mean ensemble standard deviation across all validation experiments as a function of simulation hour.

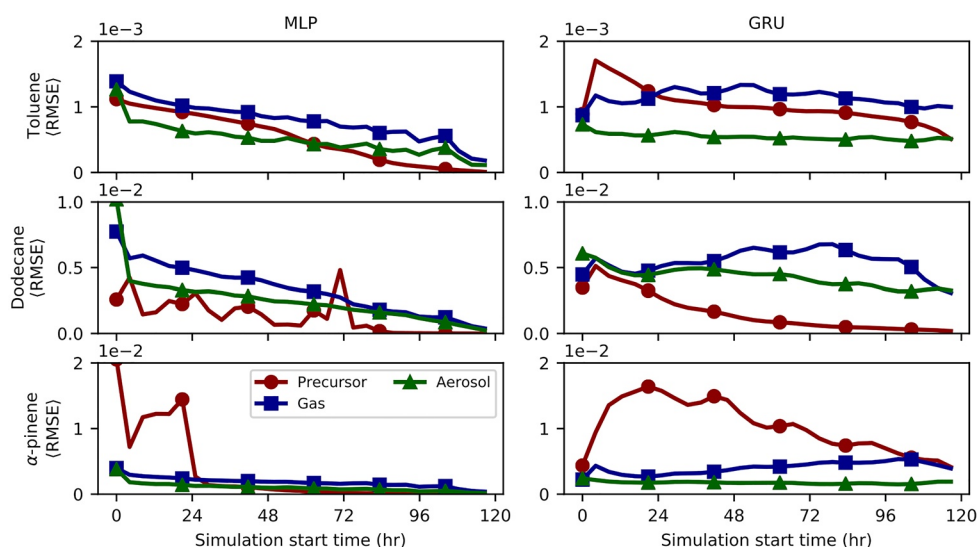
computationally feasible, one could choose to run the ensemble suite and take the mean to increase accuracy. With this approach, the mean absolute percentage errors for the GRU are less than 2% for the gas and aerosol partitions, and approximately 3%–8% for the MLP.

### 3.4. Performance Dependency on Initial Conditions

Next the models' performance dependency on different initial conditions was probed by selecting initial values  $X(t_N)$  for some  $t_N$  in a GECKO-A experiment trajectory as the initial starting point in a box simulation. For the GRU model,  $X(t_N)$  is initially passed through the hidden-state model to obtain a starting hidden state. As the experiments contain 1,440 total time steps, box simulations were left to run for 1,440 -  $N$  time steps once the initial time was selected. As the observed instabilities discussed above occurred at a range of different time steps after the box simulation was first started, we wanted to check each models' stability over as many possible time steps as there was data available. This means that box simulations started at earlier times in experiments will run for more time steps compared to those which started at later times in the experiments. Figure 8 shows the average RMSE skill score for the ensemble members for MLP and GRU models versus the initial box simulation start time. Similar plots for the  $R^2$  coefficient and HD are shown in Figures S4 and S5 in Supporting Information S1, respectively.

Figure 8 shows that the RMSE for the precursor task is usually higher at the earlier start times, then declines at later times when the precursor is chemically negligible, for all three species. The GRU precursor RMSE also increases early on before gradually declining. For  $\alpha$ -pinene in particular, the GRU was also observed to go unstable at earlier start times. However, by comparison to the MLP model for all initial starting times, the total number of experiments having gone unstable was significantly less. Overall, the RMSE is lower for toluene compared to both dodecane and  $\alpha$ -pinene by nearly a factor of ten.

On the gas prediction task, the GRU model typically performed comparable or better than the MLP at the earliest start times (longer box simulations), while at the later start times (shorter box simulations) the MLP had lower RMSE scores for precursor and gas prediction. We observed in some experiments the GRU struggling at later start times to reproduce the GECKO-A predicted gas values as accurately as the MLP, in particular, when those



**Figure 8.** The average root-mean squared error versus the initial box simulation start time, computed from the 200 test experiments for the three species considered. The multilayer perceptron and gated recurrent unit results are shown in panels on the left and right, respectively.

concentration values were very small compared to the initial experiment precursor amount, but the model predictions remained stable at these times.

### 3.5. Stabilization Through Fewer Prediction Targets

The results above indicate that predicting the evolution of the precursor's concentrations is the most difficult task of the three that we considered, especially for the MLP models. As both MLP and GRU models always have to predict finite quantities of precursor at early times before mass moves into the other two phases at later times, small precursor prediction inaccuracies can lead to numerically inaccurate predictions for gas and aerosol quantities, as well as lead to the observed experiments having gone unstable, as reported above. Furthermore, when the precursor value is small it would be set to zero in a 3D model, whereas both GRU and MLP models were trained to predict all values in the training data, even extremely small values observed in many experiments.

In the reference model, the precursor decays exponentially from its initial concentrations, at different rates depending on the environmental conditions and the species, and could be estimated using other heuristic models (such as a linear regression model), or directly calculated within the chemical model. Additionally, when the precursor is exhausted there is effectively no need to predict that quantity in the 3D simulation (as it would be set to zero). Dimensionality reduction is a common technique that can lead to performance improvements (Keller & Evans, 2019; Whitehouse et al., 2004). Thus, we consider MLP and GRU models that only perform gas and aerosol prediction, and not precursor, to probe whether not predicting precursor will improve model performance and stability on these prediction tasks. The inputs to the model and all other architecture choices remains the same, just that the output layer size is size 2 rather than 3. Both model types were optimized using ECHO, and 30 ensemble members were trained using the parameters from the best study.

Table 3 shows the same metrics as in Table 2 for the MLP and GRU models tasked with gas and aerosol predictions only. An overall improvement in scores is seen when predicting the precursor concentrations was not a model task. The  $R^2$  coefficients were comparable for toluene, but higher for the GRU on the other two species. Comparing Figure 8 and Figure S8 in Supporting Information S1 shows that the MLP and GRU both improved at gas and aerosol prediction with overall lower RMSE scores, especially at later times when the precursor had been exhausted in the experiments. Furthermore, Figure S3 in Supporting Information S1 shows that the  $R^2$  was above 0.95 for gas and aerosol tasks at all start times, while it was modestly lower for the MLP. The other notable difference for the MLP is that all model runs remained stable.

**Table 3**  
Table of Computed Metrics for Multilayer Perceptron (MLP) and Gated Recurrent Unit (GRU) Models Which Are Tasked With Prediction of Gas and Aerosol for Each of Toluene, Dodecane, and  $\alpha$ -Pinene

	Toluene		Dodecane		$\alpha$ -pinene	
	$R^2$ , $R_e^2$	RMSE, HD	$R^2$ , $R_e^2$	RMSE, HD	$R^2$ , $R_e^2$	RMSE, HD
MLP gas	0.958, 0.975	1.223, 2.227	0.803, 0.875	10.38, 9.326	0.904, 0.955	5.040, 5.029
MLP aero	0.904, 0.926	1.002, 4.025	0.764, 0.811	10.38, 24.79	0.794, 0.878	4.918, 33.21
GRU gas	0.978, 0.985	0.881, 14.794	0.981, 0.988	3.472, 1.414	0.980, 0.986	2.330, 0.981
GRU aero	0.968, 0.984	0.566, 22.168	0.980, 0.990	3.075, 9.698	0.972, 0.983	1.804, 10.48

Note. The average  $R^2$  coefficient and average Hellinger distance are listed for the two prediction tasks. All reported metrics for both models were computed using the testing set of experiments.

### 3.6. GECKO-A Emulator Evaluation With External Datasets

The performance abilities of both MLP and GRU models were tested by expanding the data sets to include additional simulations, (a) for 10 times (X10, = 100 ppt) and 100 times (X100, = 1 ppb) higher initial concentrations of the precursor, which are more representative of somewhat polluted atmospheric conditions, and (b) for simulating the diurnal variation in the precursor levels, that was not present in the original data sets which did not include the daily variability on the emissions.

#### 3.6.1. Model Performance on Increased Precursor Concentrations

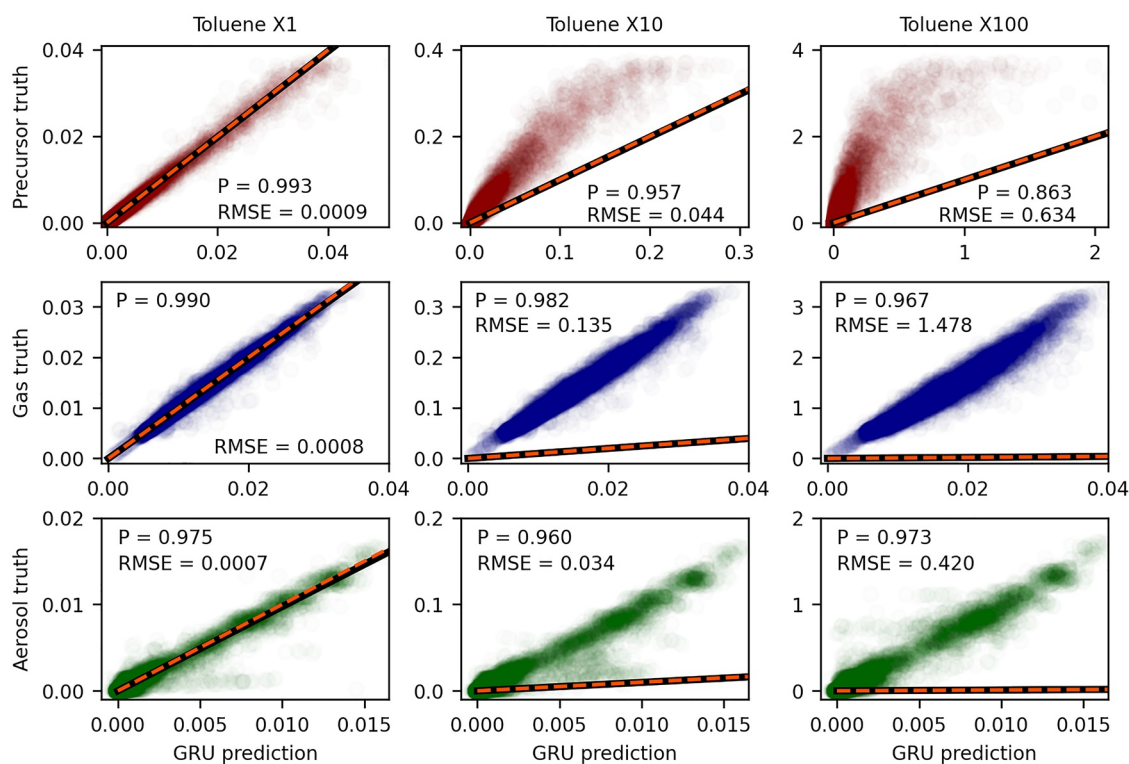
The simulations performed to create X10 and X100 data sets were carried out identically compared with the reference simulations starting at 10 ppt precursor concentrations, except that the initial precursor concentrations were increased by a factor of 10 (X10) and 100 (X100). The new data sets were then split into train, validation, and test data sets just as before, then transformed using the fitted scaling transformations on the original data sets. Then, the X10 and X100 test data sets were passed through MLP and GRU models that were trained on the original data set containing the smaller initial value of the precursor.

Figure 9 shows the predictions of the GRU model on the reference test set of experiments (left column), and the expanded X10 and X100 test data sets (middle and right columns, respectively). The Pearson coefficient, defined by Equation B3, which measures the correlation between two data sets, and MAE for each prediction task are listed in the sub-panels. The figure shows that the GRU trained on the smaller initial precursor concentrations made predictions on the X10 and X100 data sets that correlated strongly with the true values for precursor, gas, and aerosol, as is seen by high values of the Pearson coefficients for the different prediction tasks, but the MAE for each task increased by orders of magnitude with larger starting precursor concentrations.

The figure also clearly indicates that the GRU model under-predicted the true values for gas and aerosol by approximately 1 and 2 orders of magnitude for the X10 and X100 data sets, respectively. For the precursor prediction task, the predicted decay times were significantly shorter compared to that observed in the GECKO-A experiments. Overall, similar performance declines were observed for the MLP model (results not shown). Models which did not have the precursor prediction task did better by comparison but overall performance still declined. These results indicate that the neural models cannot be extrapolated outside of the training data sets. This poses a real challenge for 3D model applications given the wide range of precursor's concentrations in the atmosphere going from very clean conditions in the remote regions, and upper troposphere to polluted conditions found that is, in urban or fire plumes.

#### 3.6.2. Evaluation With Varying Environmental Conditions

Lastly, the models' performance was tested on 36 experiments run for toluene that simulated daily varying conditions for five days. Like for the training data set, the precursor's initial concentration was set to 10 ppt. Initial temperature, pre-existing aerosol seed, ozone and NO<sub>x</sub> were randomly selected in the same ranges as the training data set (Table 1). CO mixing ratio was initialized to 100 ppb. Relative humidity was held constant to a random value picked in the 50%–80% range. The latitude was also randomly selected in the 80°S–80°N range. Contrary to the training data set, after initialization, all chemical concentrations were free to evolve with the diurnal cycle to simulate a realistic atmospheric degradation of toluene and the subsequent organic aerosol formation.



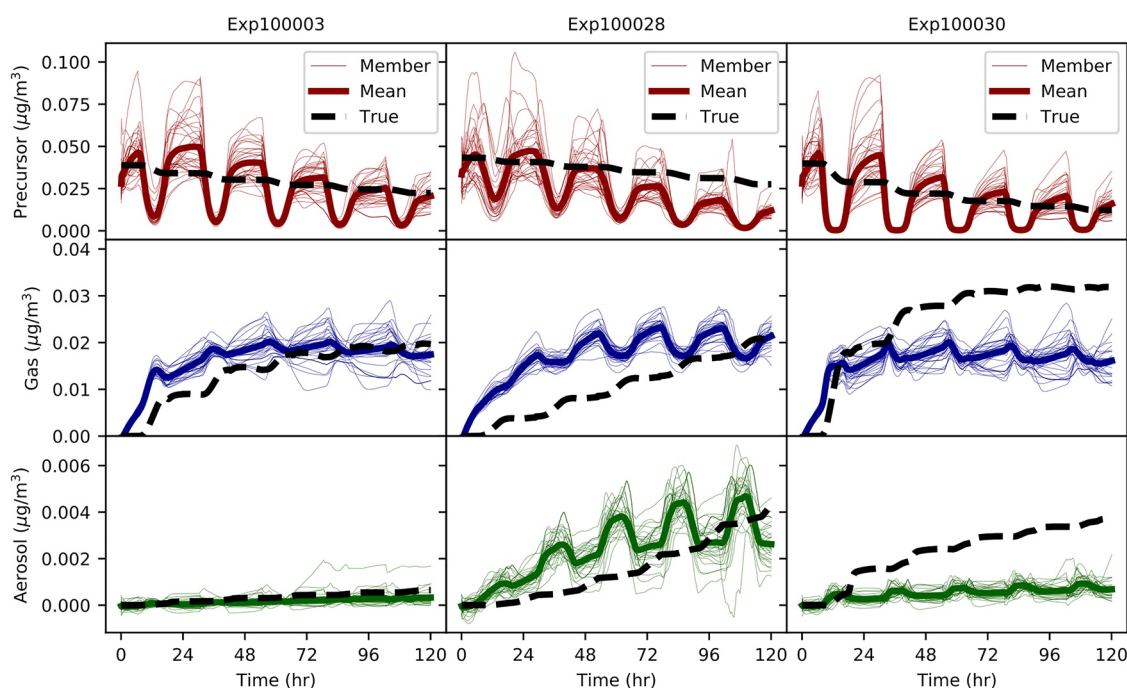
**Figure 9.** Scatter plots comparing the GECKO-A value (y-axis) against the gated recurrent unit predicted value (x-axis) for models trained on X1 data and applied to test sets with 10X and 100X higher initial concentrations. The dashed orange line shows  $y = x$ , while the solid black line shows the linear relationship for models trained on the respective  $X^*$  data set.

Because the experiments simulated a diurnal cycle and started at midnight, box simulations were performed with MLP and GRU models at different starting times in the experiments to assess the impact of training the models on daytime oxidation only. Three example experiment trajectories are shown in Figure 10 for the 3-task GRU model, for the full 5-day box simulations which all began at midnight (the same examples for the MLP model are shown in Figure C4). The simulations performed at other starting times covered a shorter 1-day window. Figure 11 shows the average  $R^2$  coefficient for these shorter simulations for both the 3- and 2-task MLP and GRU models.

Figure 10 shows the predicted concentrations for simulations starting at midnight are notably different compared with those in Figure 4. In particular, the GRU model seems to have captured some of the diurnal changes, where oxidation appears to proceed during the day, but not at night. Drastic changes in the predicted precursor amounts are observed during day-night transition periods. However, the predicted concentration values are clearly not in agreement with the true values. Although it is less obvious, close inspection of the MLP predicted precursor values in Figure C4 shows that it too responded to the diurnal variation in the extended experiments, but with poor numerical accuracy. Both MLP and GRU models predicted that all experiments remained stable at all simulation times.

Figure 11a shows that over a 1-day simulation window, the performance still dropped relative to the experiments where the environmental variables were held constant for toluene. The MLP model had comparably high  $R^2$  score for precursor prediction across the start times, but gas and aerosol performance was lower by comparison. The periodic response of the GRU to the diurnal signal in Figure 11a (b) is indicated by the sign-change in the average  $R^2$  value for predicted precursor, which goes negative when box simulations were started during day-time hours, while those started overnight stayed positive. The GRUs performance on gas and aerosol prediction also peaked for simulations that started during the middle of the day-time, and was poorest by comparison for those started late at night. Figure 11b shows that MLP and GRU models, which were only tasked with predicting gas and aerosol, performed mostly similar to the 3-task models, with notable gas performance improvement for the 2-task GRU.

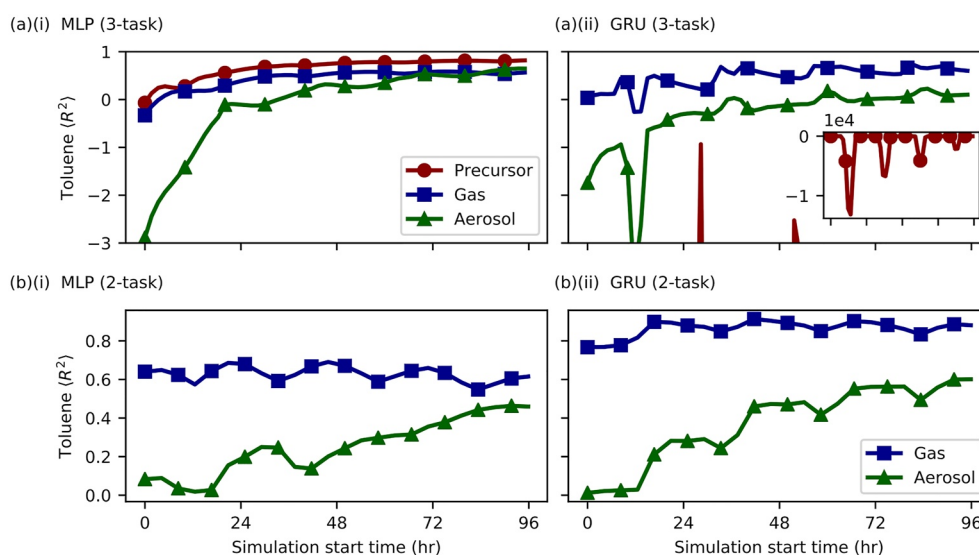




**Figure 10.** Examples of gated recurrent unit box simulations where select environmental variables were allowed to vary with time.

### 3.7. Computational Performance of Emulator Models and GECKO-A

In addition to being able to reproduce reasonably well the evolution of concentrations of organic compounds on the test data sets for the three species, the MLP and GRU emulators also led to significant computational gains. Table 4 lists estimates for the time required by GECKO-A, MLP, and GRU models to advance one time step, for example, 300 seconds (5 minutes) of simulation time, for the three precursor species (see Appendix A for more details). For toluene, GECKO-A requires 0.9 s and is about 78 and 244 times faster than dodecane and  $\alpha$ -pinene, respectively. By comparison, both MLP and GRU models require about the same time for the three precursor



**Figure 11.** The average  $R^2$  coefficient versus the initial box simulation start time computed from 36 experiments for toluene. (a) The results for the three-task multilayer perceptron (MLP) and gated recurrent unit (GRU) models are shown in (a) and (b), respectively. (b) The same quantities as in (a) except for the two-task MLP and GRU models. All box simulations ran for 24 hr.

**Table 4**

*The Average Time Each Model Required to Advance 300 s of Simulation Time*

Model	Toluene		Dodecane		$\alpha$ -pinene	
GECKO-A	0.9 s	1	71 s	1	220 s	1
MLP CPU	2.1 ms	430	0.8 ms	$8.88 \times 10^4$	1.6 ms	$1.38 \times 10^5$
MLP GPU	0.08 ms	11,250	0.07 ms	$1.01 \times 10^6$	0.08 ms	$2.75 \times 10^6$
GRU CPU	3.1 ms	290	3.2 ms	$2.22 \times 10^4$	3.3 ms	$6.67 \times 10^4$
GRU GPU	0.38 ms	2,368	0.38 ms	$1.87 \times 10^5$	0.38 ms	$5.79 \times 10^5$

*Note.* For each species, the first column shows the time step in seconds while the second column shows the ratio of the Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere time step to each model time step.

species on the CPU, typically a few microseconds, with the MLP faster than the GRU by up to a factor of five. Thus, for toluene both neural network models could be expected to perform hundreds of times faster, while for  $\alpha$ -pinene the expected speed-up could be up to 4–6 orders of magnitude faster than the explicit model.

### 3.8. Incorporation of ML Models Into 3D Models

A model trained on data with 300 s time step temporal resolution can be incorporated into a 3D model, for instance with a 30 min time step, by running the emulator simulation for 30 min. If a smaller time step is needed, the 5-min data set could be extrapolated to 2.5 or even 1 min, but for shorter times a new model may need to be trained on a data set with the desired temporal resolution. A GRU emulator will also require additional memory due to the presence of the hidden state for all grid cells. We estimate the memory requirement of the GRU in a 3D model to be around 10 gigabytes with Equation A1, which does not take into account, for example, any compression algorithms or parallelization with MPIs. The requirement could also be reduced, for instance, by pursuing smaller hidden layers during hyperparameter optimization with high overall performance (note that we limited the GRU to a single hidden layer during the optimization, which may have encouraged a larger hidden state size over a model with more, but smaller hidden layers). Note also that here the GRU is trained only to emulate a complex chemical solver, independent from other processes modulating aerosol concentrations such as transport, deposition, etc. It could be implemented as is into a 3D model with an operator splitting method: transport, advection, etc, that are solved with traditional ODE solvers and the chemistry is solved with the GRU. This is, for instance, the approach adopted by Keller and Evans (2019) for solving ozone chemistry with random forests in GEOS-Chem and by Kelp et al. (2022) with MLPs.

## 4. Discussion

In general, the comparative differences seen between the MLP and 1-step GRU applications show the advantages of a RNN emulator in a few key areas. First, its overall accuracy, especially at integrated time steps further away from its starting conditions, is notably higher than the MLP model. Furthermore, the encoding of a hidden state which can represent the trajectory of each input, the key feature of a recurrent network architecture, appears to help constrain model uncertainty and ultimately, numerical stability. Most neural network applications for atmospheric chemistry have not yet begun to examine such model uncertainties. Our use of training a suite of ensemble members using the exact same architecture for each model, and only initializing the weights differently prior to training, provides some evidence that model uncertainty can be sensitive to, and better constrained by, certain model types. We also note that our recurrent model remains numerically stable for all species and for most starting initial conditions, which is not true for a small percentage of MLP member/experiment combinations. This insight may not have been detected without the inspection of an ensemble suite, as most of the MLP models remained stable. Additionally, examination of the conditions that produced unstable runs were highly correlated with near zero values of pre-existing aerosols and the relative oxidation mechanism for each precursor.

Maintaining numerical stability with the use of emulators for atmospheric chemistry and other atmospheric parameterizations is a known issue and initial steps have been taken to address it (Brenowitz & Bretherton, 2018; Kelp et al., 2020, 2022). These recent studies found some performance improvements by using a “recurrent training” scheme, where a model was rolled out in time for  $n$  time steps during training, and a loss was calculated on

the sum of  $n$  time steps, instead of a single time step. However, the models used in these studies were not RNNs, as the network architectures were that of an MLP (Brenowitz & Bretherton, 2018) and an encoder/decoder framework (Kelp et al., 2020, 2022), which only utilized feed-forward connections. Rather, training these models relied on the multi-time step loss function as a means to update the model weights using a sequence instead of a single length input. Our GRU model provides an alternative approach, by rolling out the model to the end of the training experiment and calculating the loss at successive single time steps, the feedback connections' memory of the trajectory through  $t - 1$  is simply used as input at  $t$  along side the current values of the precursor, gas and aerosol.

RNNs have not yet been thoroughly explored in 3D atmospheric modeling, although there have been applications in other earth systems areas including hydrology (Ardabili et al., 2019; Kratzert et al., 2018), earthquake magnitude prediction (Mousavi & Beroza, 2020), rain-runoff (Boulmaiz et al., 2020) and wind velocity forecasting (Irrgang et al., 2020), as well as vegetation growth estimation (Reddy & Prasad, 2018). One reason for this might be the lower dimensional nature of many of these models, which would be computationally less burdensome to put into production as opposed to integrating a model that requires multiple time steps of input into a full 3D climate or weather model. For this reason, we have developed a method that still only requires one time step of input but maintains the advantage of having an encoded memory of past time steps. A small disadvantage of our framework is that it does require an additional model to predict the initial hidden state prior to running the GRU. However, if the community ultimately finds that it is computationally and programmatically feasible to couple large recurrent networks into full 3D transport models, investigation of training recurrent models with multiple time steps of input would be a recommended pathway.

Although there are clear benefits demonstrated from use of a recurrent network, there are computational limitations. The hidden state increases the input needed for each prediction from 9 for the MLP model to 1,000 for the GRU. This is not problematic for 1D validation efforts, but would become too memory intensive if this model were integrated into 3D simulations. A smaller GRU hidden state is possible but may result in drops in performance. If directly shrinking the vector is not feasible, lossy compression of the vector with principal component analysis or an autoencoder may balance a smaller performance loss with slightly more computation. For this reason, despite the lower performance metrics, we still find value in simplified neural networks such as the MLP if they can still approximate a solution within the given tolerance. Additionally, some performance could be sacrificed for a smaller GRU model (see Figure S3 in Supporting Information S1).

Our results demonstrate some ML generalization challenges involving the selection and training of neural networks on experiments with both small initial precursor concentrations and select static environmental variables. For example, low precursor concentrations of 10 ppt were chosen primarily to limit the influence of a single precursor on the photochemical reactivity, and gas/particle partitioning in GECKO-A. However, this had a large impact on the generalizability outside the training range (Figure 9). If we were to implement our models into a 3D climate model, they would need to be trained on a larger range of precursor values that are also representative of more polluted atmospheric conditions. Additionally, environmental variables outside of temperature were held constant in an effort to help the models generalize better, but we observed the model predictions deviating from the external data sets they were tested on. While there did not appear to be any direct evidence of over-fitting to the training data, an open question remains of how to properly configure the reference box models to provide data to best capture the physical relationships in a complex chemical system. One speculation is that the GRU could generalize more effectively than observed here by having varying environmental fields within an experiment, to better parameterize the models feedback connections. Many other generalization questions also remain, such as the inclusion of night chemistry, as well as reactions between species originated from various precursors. Our NN models were built for use in bulk aerosol models (such as e.g., GEOS-Chem) based on the GECKO-A bulk size distribution, and will need to be expanded to account for other aerosol size representations.

To our knowledge, this is the first neural network emulation of organic atmospheric chemistry. As a result, there are many areas that warrant further exploration: (a) coupling both the MLP and GRU models to a 3D chemistry-climate model, such as WRF- or GEOS-Chem, to better understand their successes and shortcomings, (b) further quantification of the underlying uncertainties in model predictions to determine whether the error sources originate from the data or the model architecture choices, or both, (c) testing of different data sets, training regimes, and model architectures to better generalize across different chemical regimes (such as daytime vs. nighttime chemistry), (d) incorporating physical constraints into the model architecture or training procedure as a means for constraining model outputs, for example, the total mass or the number of C atoms needs to be

conserved, and (e) utilizing explainable and interpretable methodologies to better understand what the model has learned, and what it is using to drive its predictions.

## 5. Conclusions

In summary, we have developed two types of neural network emulators for the GECKO-A organic aerosol chemistry model and evaluated the emulators on a broad set of chemical precursors, atmospheric conditions, and edge cases for machine learning emulator robustness. Both the MLP and GRU neural networks can accurately replicate the general time evolution of complex chemical processes. The recurrent connections and persistent hidden state of the GRU enabled it to produce more stable and longer-running box simulations with better fidelity than the MLP, which only incorporated information from the previous time step. An extensive hyperparameter optimization search enabled the strong performance of both ML models. Performance and stability issues were further improved with both types of models by only predicting gas and aerosol and assuming precursor changes are prescribed.

The machine learning emulators in this paper were designed as a proof-of-concept for an emulator that can operate within a full 3D numerical simulation. The results show significant progress toward 3D simulation implementation but also highlight key limitations that need to be addressed further. The current training and emulation modeling pipeline should mechanically be feasible to run within a 3D simulation. The MLP and GRU both produce significant computational speedups compared with GECKO-A and should not require significant processing overhead to run within the model. Our GRU implementation will require a large but still manageable amount of memory overhead to store the hidden state between integration timesteps but does not require the storing of multiple timesteps like other implementations of RNNs. On the other hand, the robustness of the emulator is very sensitive to the composition of the training data set and does not generalize to precursor concentrations and temporal variations in precursor outside the bounds of the training data. Future emulation datasets should ensure that the initial, boundary, and forcing conditions for their runs are representative of the variations expected within a full 3D simulation.

## Appendix A: Time Step Comparison and GRU Memory Requirements

Table 4 compares the average speed in which GECKO-A, MLP, and GRU models take to advance 300 s of simulation time. GECKO-A is written in Fortran, the MLP is written in python using Tensorflow, and the GRU is written in python using PyTorch. For the GRU and MLP models, 200 box simulations were run out to 5 days, and the average over all the individual estimates was taken. The standard deviation for both MLP and GRU models was about 5% the average value reported in Table S1 of Supporting Information S1. Of the GPU estimates, approximately 70% of the reported value was time spent moving input and output data to and from the GPU. The neural network models were evaluated on NCAR's casper supercomputer, on a node that contained an 18-core 2.3-GHz Intel Xeon Gold 6,140 processors (CPU) and a NVIDIA Tesla V100 32 GB graphics cards (GPU). Each model utilized a single core and all of the GPU memory. GECKO-A is written in Fortran, and the box model simulations were performed on the NCAR's supercomputer, running on single cores using 2.3-GHz Intel Xeon E5-2697V4 (Broadwell) processors (CPU).

The memory requirement of the GRU due to the hidden state is estimated as

$$\begin{aligned} \text{Memory requirement} &= HSS * HGC * VLT * BN \\ &= 10^3 * 64,800 * 20 * 64 \text{ bits} \\ &\approx 10 \text{ Gb of memory} \end{aligned} \tag{A1}$$

where HSS is the GRU hidden state size, HGC is the number of horizontal grid cells, VLT is the number of vertical levels in the troposphere, and BN is the number of bits per neuron.

## Appendix B: Metrics Definitions

The coefficient of determination,  $R^2$ , is defined as

$$R^2 \equiv 1 - R_{ss}/T_{ss}$$

$$R_{ss} = \sum_i (y_i - f_i)^2$$

$$T_{ss} = \sum_i (y_i - \bar{y})^2$$
(B1)

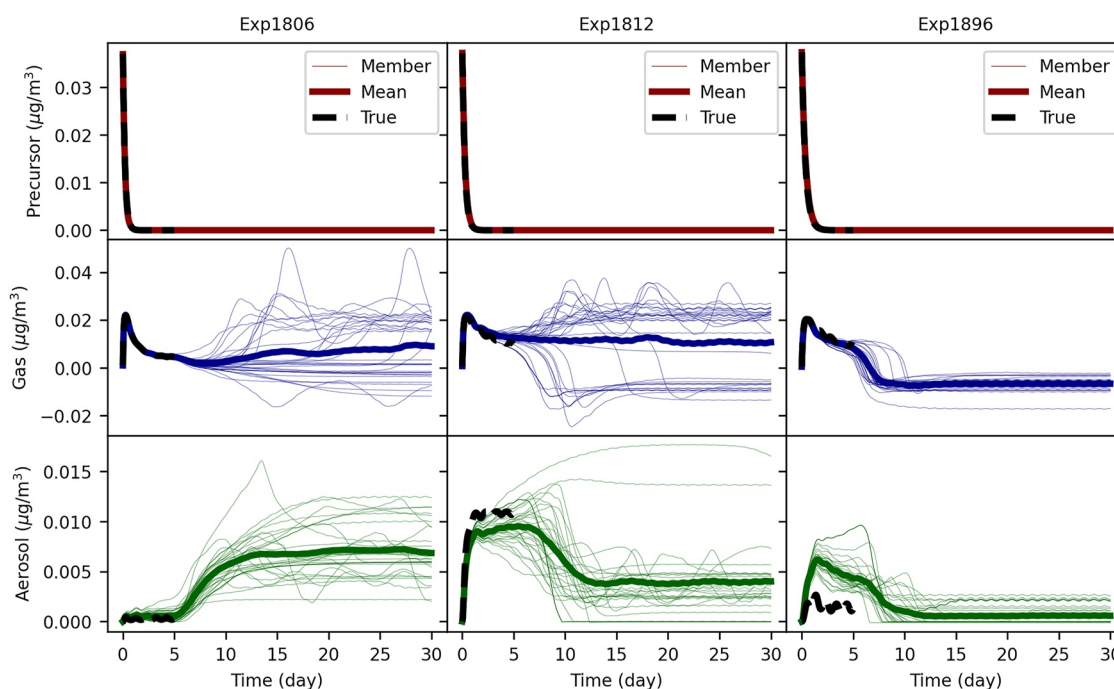
where  $y$  and  $f$  are the temporal true and the predicted values for trajectories, respectively, and  $\bar{y}$  is the mean value. The HD is defined as

$$HD \equiv \sqrt{\sum_i (\sqrt{y_i} - \sqrt{f_i})^2}$$
(B2)

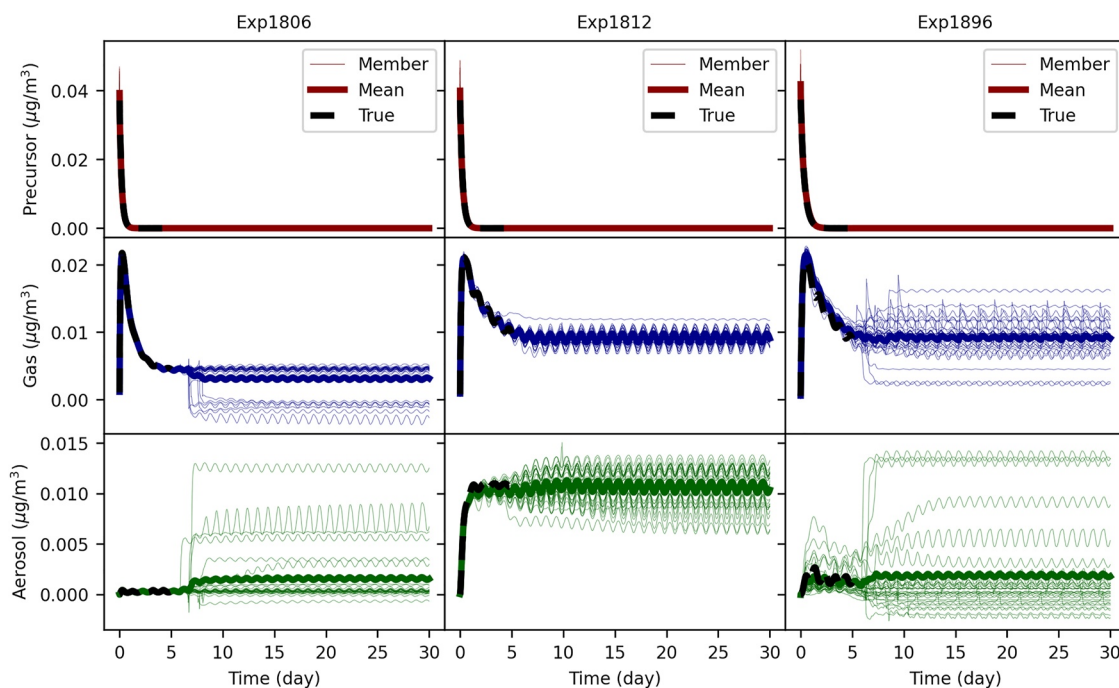
and the Pearson coefficient ( $P$ ) as the covariance between  $y$  and  $f$  as given by

$$P \equiv \frac{\sum_i (y_i - \bar{y})(f_i - \bar{f})}{\sqrt{\sum_i (y_i - \bar{y})^2} \sqrt{\sum_i (f_i - \bar{f})^2}}$$
(B3)

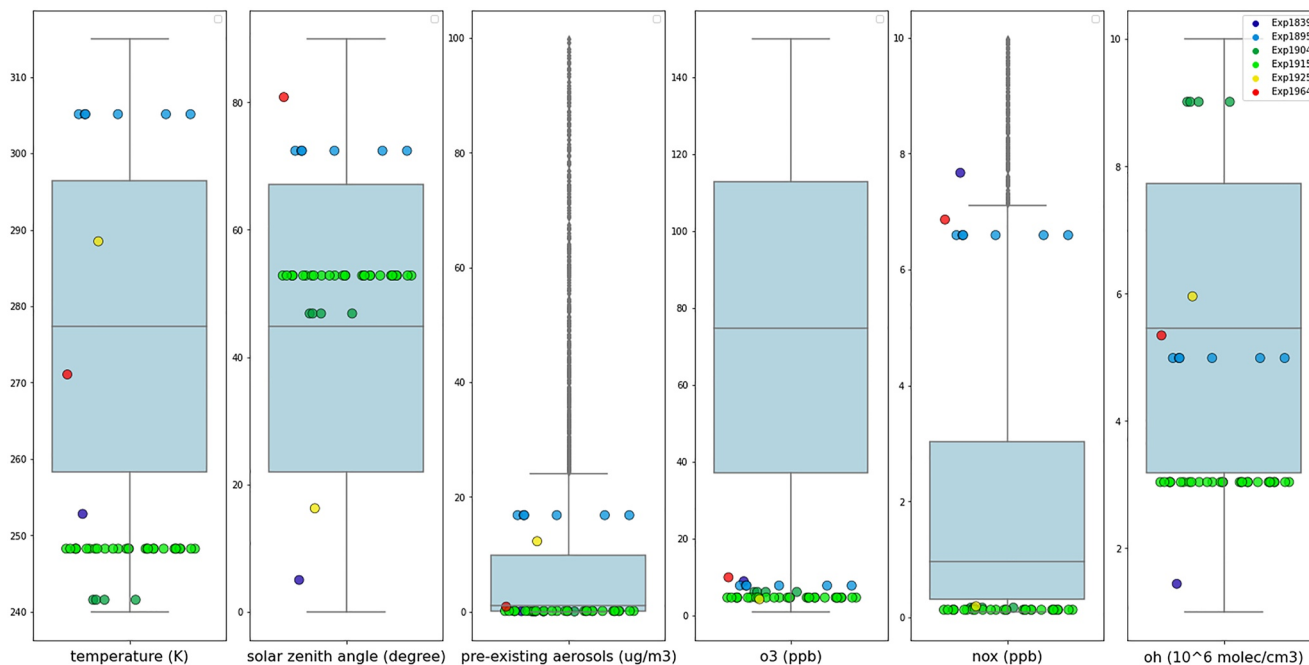
### Appendix C: Additional Results



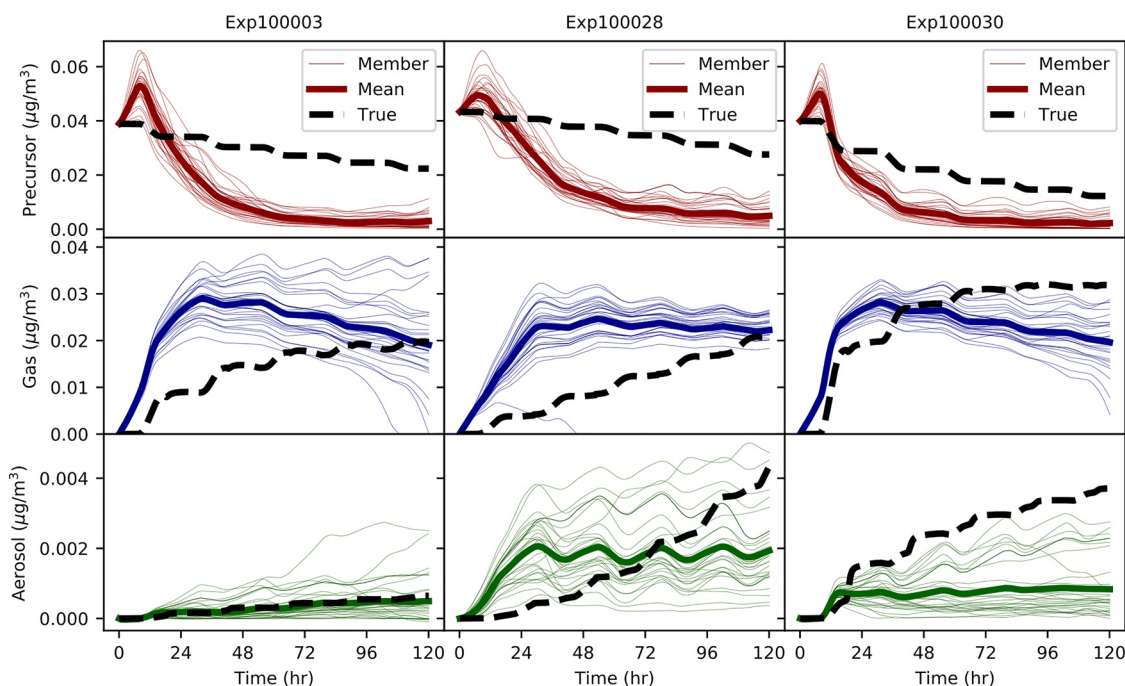
**Figure C1.** Three toluene sample experiment trajectories for the multilayer perceptron model. Solid (thick) lines show the mean Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A) trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.



**Figure C2.** Three toluene sample experiment trajectories for the gated recurrent unit model. Solid (thick) lines show the mean Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A) trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.



**Figure C3.** Unstable  $\alpha$ -pinene simulations. The blue box (IQR) and whisker ( $Q3/Q1 \pm IQR \times 1.5$ ) plots represent the training data initial conditions distribution for each input parameter. The colored dots represent where in the input space (vertical) the unstable experiment is. Different experiments are coded by color and different ensemble members for each experiment spread out horizontally.



**Figure C4.** Examples of multilayer perceptron box simulations where the environmental variables not including temperature and solar zenith angle were allowed to vary with time, compared with the Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere simulations.

Three example experiment trajectories are shown in Figure C4 for the 3-task MLP model, for the full 5-day box simulations which all began at midnight (the same examples for the GRU model are shown in Figure 10). The simulations performed at other starting times covered a longer 4-day window compared to the 1-day simulations shown in the main text in Figure 11.

### Data Availability Statement

We would like to acknowledge high-performance computing support from Cheyenne and Casper (Computational and Information Systems Laboratory, CISL, 2020) provided by NCAR's Computational and Information Systems Laboratory, sponsored by the National Science Foundation. The emulators described here and simulation code used to train and test the models are archived at <https://github.com/NCAR/gecko-ml>. All GECKO-A data sets created for this study are available at <https://doi.org/10.5281/zenodo.5790042>.

### Acknowledgments

This material is based upon work supported by the National Center for Atmospheric Research, which is a major facility sponsored by the National Science Foundation under Cooperative Agreement No. 1852977. Jinkyul Cho acknowledges funding from NASA 80NSSC20K0214.

### References

- Ardabili, S., Mosavi, A., Dehghani, M., & Várkonyi-Kóczy, A. R. (2019). Deep learning and machine learning in hydrological processes climate change and Earth systems a systematic review. In *International conference on global research and education* (pp. 52–62).
- Aumont, B., Szopa, S., & Madronich, S. (2005). Modelling the evolution of organic carbon during its gas-phase tropospheric oxidation: Development of an explicit model based on a self generating approach. *Atmospheric Chemistry and Physics*, 5(9), 2497–2517. <https://doi.org/10.5194/acp-5-2497-2005>
- Beucler, T., Pritchard, M., Gentine, P., & Rasp, S. (2020). Towards physically-consistent, data-driven models of convection. In *International geoscience and remote sensing symposium (IGARSS) 2020-2020 IEEE international geoscience and remote sensing symposium* (pp. 3987–3990). <https://doi.org/10.1109/IGARSS39084.2020.9324569>
- Beucler, T., Pritchard, M., Rasp, S., Ott, J., Baldi, P., & Gentine, P. (2021). Enforcing analytic constraints in neural networks emulating physical systems. *Physical Review Letters*, 126(9), 098302. <https://doi.org/10.1103/physrevlett.126.098302>
- Boucher, O., Randall, D., Artaxo, P., Bretherton, C., Feingold, G., Forster, P., et al. (2013). Clouds and aerosols. In T. Stocker (Eds.) *Climate change 2013: The physical science basis. Contribution of working group I to the fifth assessment report of the intergovernmental panel on climate change*. Cambridge University Press.
- Boulmaiz, T., Guermoui, M., & Boutaghane, H. (2020). Impact of training data size on the 1stm performances for rainfall–runoff modeling. *Modeling Earth Systems and Environment*, 6(4), 2153–2164. <https://doi.org/10.1007/s40808-020-00830-w>
- Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic validation of a neural network unified physics parameterization. *Geophysical Research Letters*, 45(12), 6289–6298. <https://doi.org/10.1029/2018GL078510>
- Camredon, M., Aumont, B., Lee-Taylor, J., & Madronich, S. (2007). The soa/voc/no<sub>x</sub> system: An explicit model of secondary organic aerosol formation. *Atmospheric Chemistry and Physics*, 7(21), 5599–5610. <https://doi.org/10.5194/acp-7-5599-2007>

- Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014a). Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv preprint arXiv:1412.3555.
- Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014b). Empirical evaluation of gated recurrent neural networks on sequence modeling. arXiv 2014. arXiv preprint arXiv:1412.3555.
- Computational and Information Systems Laboratory, CISL. (2020). *Cheyenne: HPE/SGI ICE XA system (NCAR community computing) (Tech. Rep.)*. National Center for Atmospheric Research. <https://doi.org/10.5065/D6RX99HX>
- de Gouw, J. A. (2005). Budget of organic carbon in a polluted atmosphere: Results from the new England air quality study in 2002. *Journal of Geophysical Research*, *110*(D16), D16305. <https://doi.org/10.1029/2004JD005623>
- Fry, J. L., Draper, D. C., Barsanti, K. C., Smith, J. N., Ortega, J., Winkler, P. M., et al. (2014). Secondary organic aerosol formation and organic nitrate yield from NO<sub>3</sub> oxidation of biogenic hydrocarbons. *Environmental Science and Technology*, *48*(20), 11944–11953. <https://doi.org/10.1021/es502204x>
- Gagne, D. J., Christensen, H. M., Subramanian, A. C., & Monahan, A. H. (2020). Machine learning for stochastic parameterization: Generative adversarial networks in the Lorenz '96 model. *Journal of Advances in Modeling Earth Systems*, *12*(3), e2019MS001896. <https://doi.org/10.1029/2019ms001896>
- Gettelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J., Morrison, H., & Gantos, G. (2021). Machine learning the warm rain process. *Journal of Advances in Modeling Earth Systems*, *13*(2). <https://doi.org/10.1029/2020MS002268>
- Hochreiter, S. (1991). Untersuchungen zu dynamischen neuronalen netzen. Diploma. *Technische Universität München*, *91*(1).
- Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, *9*(8), 1735–1780. <https://doi.org/10.1162/neco.1997.9.8.1735>
- Hodzic, A., Aumont, B., Knote, C., Madronich, S., & Tyndall, G. (2014). Volatility dependence of Henry's law constants of condensable organics: Application to estimate depositional loss of secondary organic aerosols. *Geophysical Research Letters*, *41*(13), 4795–4804. <https://doi.org/10.1002/2014GL060649>
- Hodzic, A., Campuzano-Jost, P., Bian, H., Chin, M., Colarco, P. R., Day, D. A., et al. (2020). Characterization of organic aerosol across the global remote troposphere: A comparison of ATom measurements and global chemistry models. *Atmospheric Chemistry and Physics*, *20*(8), 4607–4635. <https://doi.org/10.5194/acp-20-4607-2020>
- Hodzic, A., & Jimenez, J. L. (2011). Modeling anthropogenically controlled secondary organic aerosols in a megacity: A simplified framework for global and climate models. *Geoscientific Model Development*, *4*(4), 901–917. <https://doi.org/10.5194/gmd-4-901-2011>
- Hodzic, A., Kasibhatla, P. S., Jo, D. S., Cappa, C. D., Jimenez, J. L., Madronich, S., & Park, R. J. (2016). Rethinking the global secondary organic aerosol (soa) budget: Stronger production, faster removal, shorter lifetime. *Atmospheric Chemistry and Physics*, *16*(12), 7917–7941. <https://doi.org/10.5194/acp-16-7917-2016>
- Hodzic, A., Madronich, S., Kasibhatla, P. S., Tyndall, G., Aumont, B., Jimenez, J. L., et al. (2015). Organic photolysis reactions in tropospheric aerosols: Effect on secondary organic aerosol formation and lifetime. *Atmospheric Chemistry and Physics*, *15*(16), 9253–9269. <https://doi.org/10.5194/acp-15-9253-2015>
- Irrgang, C., Saynisch-Wagner, J., & Thomas, M. (2020). Machine learning-based prediction of spatiotemporal uncertainties in global wind velocity reanalyses. *Journal of Advances in Modeling Earth Systems*, *12*(5), e2019MS001876. <https://doi.org/10.1029/2019ms001876>
- Jenkin, M. E., Saunders, S. M., Wagner, V., & Pilling, M. J. (2003). Protocol for the development of the Master Chemical Mechanism, MCM v3 (Part B): Tropospheric degradation of aromatic volatile organic compounds. *Atmospheric Chemistry and Physics*, *3*(1), 181–193. <https://doi.org/10.5194/acp-3-181-2003>
- Keller, C. A., & Evans, M. J. (2019). Application of random forest regression to the calculation of gas-phase chemistry within the geos-chem chemistry model v10. *Geoscientific Model Development*, *12*(3), 1209–1225. <https://doi.org/10.5194/gmd-12-1209-2019>
- Kelp, M. M., Jacob, D. J., Kutz, J. N., Marshall, J. D., & Tessum, C. W. (2020). Toward stable, general machine-learned models of the atmospheric chemical system. *Journal of Geophysical Research: Atmospheres*, *125*(23), e2020JD032759. <https://doi.org/10.1029/2020JD032759>
- Kelp, M. M., Jacob, D. J., Lin, H., & Sulprizio, M. P. (2022). An online-learned neural network chemical solver for stable long-term global simulations of atmospheric chemistry. *Journal of Advances in Modeling Earth Systems*, *14*(6), e2021MS002926. <https://doi.org/10.1029/2021ms002926>
- Kelp, M. M., Tessum, C. W., & Marshall, J. D. (2018). Orders-of-magnitude speedup in atmospheric chemistry modeling through neural network-based emulation. arXiv(206), 1–23.
- Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018). Rainfall–runoff modelling using long Short-Term memory (LSTM) networks. *Hydrology and Earth System Sciences*, *22*(11), 6005–6022. <https://doi.org/10.5194/hess-22-6005-2018>
- La, Y. S., Camredon, M., Ziemann, P. J., Valorso, R., Matsunaga, A., Lannuque, V., et al. (2016). Impact of chamber wall loss of gaseous organic compounds on secondary organic aerosol formation: Explicit modeling of SOA formation from alkane and alkene oxidation. *Atmospheric Chemistry and Physics*, *16*(3), 1417–1431. <https://doi.org/10.5194/acp-16-1417-2016>
- Lannuque, V., Camredon, M., Couvidat, F., Hodzic, A., Valorso, R., Madronich, S., et al. (2018). Exploration of the influence of environmental conditions on secondary organic aerosol formation and organic species properties using explicit simulations: Development of the vbs-gecko parameterization. *Atmospheric Chemistry and Physics*, *18*(18), 13411–13428. <https://doi.org/10.5194/acp-18-13411-2018>
- Lee-Taylor, J., Hodzic, A., Madronich, S., Aumont, B., Camredon, M., & Valorso, R. (2015). Multiday production of condensing organic aerosol mass in urban and forest outflow. *Atmospheric Chemistry and Physics*, *15*(2), 595–615. <https://doi.org/10.5194/acp-15-595-2015>
- Lee-Taylor, J., Madronich, S., Aumont, B., Baker, A., Camredon, M., Hodzic, A., et al. (2011). Explicit modeling of organic chemistry and secondary organic aerosol partitioning for Mexico City and its outflow plume. *Atmospheric Chemistry and Physics*, *11*(24), 13219–13241. <https://doi.org/10.5194/acp-11-13219-2011>
- Li, J., Cleveland, M., Ziemba, L. D., Griffin, R. J., Barsanti, K. C., Pankow, J. F., & Ying, Q. (2015). Modeling regional secondary organic aerosol using the Master Chemical Mechanism. *Atmospheric Environment*, *102*, 52–61. <https://doi.org/10.1016/j.atmosenv.2014.11.054>
- Mauderly, J. L., & Chow, J. C. (2008). Health effects of organic aerosols. *Inhalation Toxicology*, *20*(3), 257–288. <https://doi.org/10.1080/08958370701866008>
- Mouchel-Vallon, C., Lee-Taylor, J., Hodzic, A., Artaxo, P., Aumont, B., Camredon, M., et al. (2020). Exploration of oxidative chemistry and secondary organic aerosol formation in the amazon during the wet season: Explicit modeling of the Manaus urban plume with gecko-a. *Atmospheric Chemistry and Physics*, *20*(10), 5995–6014. <https://doi.org/10.5194/acp-20-5995-2020>
- Mousavi, S. M., & Beroza, G. C. (2020). A machine-learning approach for earthquake magnitude estimation. *Geophysical Research Letters*, *47*(1), e2019GL085976. <https://doi.org/10.1029/2019gl085976>
- Ng, N. L., Kroll, J. H., Chan, A. W. H., Chhabra, P. S., Flagan, R. C., & Seinfeld, J. H. (2007). Secondary organic aerosol formation from m-xylene, toluene, and benzene. *Atmospheric Chemistry and Physics*, *7*(14), 3909–3922. <https://doi.org/10.5194/acp-7-3909-2007>
- Reddy, D. S., & Prasad, P. R. C. (2018). Prediction of vegetation dynamics using ndvi time series data and lstm. *Modeling Earth Systems and Environment*, *4*(1), 409–419. <https://doi.org/10.1007/s40808-018-0431-3>



- Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations by back-propagating errors. *Nature*, *323*(6088), 533–536. <https://doi.org/10.1038/323533a0>
- Schreck, J. S., & Gagne, D. J. (2021). Earth computing hyperparameter optimization. GitHub Retrieved from <https://github.com/NCAR/echo-opt>
- Tsigaridis, K., Daskalakis, N., Kanakidou, M., Adams, P. J., Artaxo, P., Bahadur, R., et al. (2014). The aerocom evaluation and intercomparison of organic aerosol in global models. *Atmospheric Chemistry and Physics*, *14*(19), 10845–10895. <https://doi.org/10.5194/acp-14-10845-2014>
- Valorso, R., Aumont, B., Camredon, M., Raventos-Duran, T., Mouchel-Vallon, C., Ng, N. L., et al. (2011). Explicit modelling of soa formation from  $\alpha$ -pinene photooxidation: Sensitivity to vapour pressure estimation. *Atmospheric Chemistry and Physics*, *11*(14), 6895–6910. <https://doi.org/10.5194/acp-11-6895-2011>
- Whitehouse, L., Tomlin, A., & Pilling, M. (2004). Systematic reduction of complex tropospheric chemical mechanisms, part i: Sensitivity and time-scale analyses. *Atmospheric Chemistry and Physics*, *4*(7), 2025–2056. <https://doi.org/10.5194/acp-4-2025-2004>
- Ying, Q., & Li, J. (2011). Implementation and initial application of the near-explicit master chemical mechanism in the 3D community multiscale air quality (cmaq) model. *Atmospheric Environment*, *45*(19), 3244–3256. <https://doi.org/10.1016/j.atmosenv.2011.03.043>
- Yuval, J., & O’Gorman, P. A. (2020). Stable machine-learning parameterization of subgrid processes for climate modeling at a range of resolutions. *Nature Communications*, *11*(1), 1–10. <https://doi.org/10.1038/s41467-020-17142-3>