

Neural network emulation of the formation of organic aerosols based on the explicit GECKO-A chemistry model

John Schreck¹, Charles Becker¹, David John Gagne¹, Keely Lawrence¹, Siyuan Wang², Camille Mouchel-Vallon³, Jinkyul Choi⁴, and Alma Hodzic⁵

¹National Center for Atmospheric Research

²CIRES/NOAA/CSL

³LISA, CNRS

⁴Unknown

⁵National Center for Atmospheric Research (UCAR)

November 23, 2022

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. Recent applications of atmospheric chemistry emulation with machine learning (ML) applied to the simpler chemical mechanisms of tropospheric ozone have shown its ability to produce realistic predictions and significantly reduce the computational cost. This study proves that ML can accurately emulate SOA formation from an explicit chemistry model for several precursors with 100 to 100,000 times speedup over GECKO-A, making it computationally usable in a chemistry-climate model. To train the ML emulator, we generated thousands of GECKO-A box simulations sampled from a broad range of initial environmental conditions, and focused on the chemistry of three representative SOA precursors: the oxidation by OH of two anthropogenic (toluene, dodecane), and one biogenic VOC (alpha-pinene). We compare fully-connected and recurrent neural network methods and use an ensemble approach to quantify their underlying uncertainty and robustness. The SOA predictions generally remain stable over a simulation period of 5 days with an approximate error of 2-8%.

1 **Neural network emulation of the formation of organic**
2 **aerosols based on the explicit GECKO-A chemistry**
3 **model**

4 **John S. Schreck¹, Charles Becker¹, David John Gagne¹, Keely Lawrence¹,**
5 **Siyuan Wang^{2,3}, Camille Mouchel-Vallon⁴, Jinkyul Choi⁵, Alma Hodzic¹**

6 ¹National Center for Atmospheric Research (NCAR), Boulder, CO, USA

7 ²Cooperative Institute for Research in Environmental Sciences (CIRES), University of Colorado, Boulder,
8 CO, USA

9 ³National Oceanic and Atmospheric Administration (NOAA), Chemical Sciences Laboratory (CSL),
10 Boulder, CO, USA

11 ⁴Laboratoire d'Aérologie, Université de Toulouse, CNRS, UPS, Toulouse, France

12 ⁵Environmental Engineering Program, University of Colorado, Boulder, CO, USA

13 **Key Points:**

- 14 • Incorporation of explicit organic chemistry into 3D chemistry-simulations requires
15 emulation.
- 16 • We developed two types of neural network emulators for the GECKO-A chemistry
17 model.
- 18 • The emulators produced accurate and stable simulations for three precursor species.

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. Recent applications of atmospheric chemistry emulation with machine learning (ML) applied to the simpler chemical mechanisms of tropospheric ozone have shown its ability to produce realistic predictions and significantly reduce the computational cost. This study proves that ML can accurately emulate SOA formation from an explicit chemistry model for several precursors with 100 to 100,000 times speedup over GECKO-A, making it computationally usable in a chemistry-climate model. To train the ML emulator, we generated thousands of GECKO-A box simulations sampled from a broad range of initial environmental conditions, and focused on the chemistry of three representative SOA precursors: the oxidation by OH of two anthropogenic (toluene, dodecane), and one biogenic VOC (alpha-pinene). We compare fully-connected and recurrent neural network methods and use an ensemble approach to quantify their underlying uncertainty and robustness. The SOA predictions generally remain stable over a simulation period of 5 days with an approximate error of 2-8%.

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

51 the solution of an explicit representation of organic chemistry which would be compu-
52 tationally feasible to link with current air quality and climate models.

53 **1 Introduction**

54 Secondary organic aerosols (SOA) have been an active area of research in the past
55 decade with the goal of improving their representation in air quality and climate mod-
56 els (Tsigaridis et al., 2014; Hodzic et al., 2016), which is essential for predicting their ef-
57 fect on human health (Mauderly & Chow, 2008) and their contribution to radiative forc-
58 ing in the climate system (Boucher et al., 2013). The misrepresentation of SOA forma-
59 tion pathways in 3D models has led to a long-standing discrepancy between observed and
60 modeled organic aerosol concentrations that has been reported from urban to remote re-
61 gions (de Gouw, 2005; Hodzic et al., 2020). Unlike sulfate and other inorganic aerosols,
62 which are made from a few dominant chemical pathways, SOAs result from the conden-
63 sation of a very large number of partly oxidized gases. These gases are generated from
64 the multi-generational oxidation of volatile organic compounds (VOCs) emitted from an-
65 thropogenic and natural sources. This complexity is not included in current 3D mod-
66 els that rely on simplified SOA parameterizations that have been developed and opti-
67 mized based on laboratory measurements or ambient aircraft data (Ng et al., 2007; Hodzic
68 & Jimenez, 2011). This empirical approach does not include the mechanistic understand-
69 ing of processes leading to SOA formation, and the adequate sensitivity to environmen-
70 tal variables modulating SOA concentration.

71 Detailed chemistry models such as the widely used near-explicit Master Chemical
72 Mechanism (MCM) (Jenkin et al., 2003) or the "fully-explicit" Generator of Explicit Chem-
73 istry and Kinetics of Organics in the Atmosphere (GECKO-A) (Aumont et al., 2005; Cam-
74 redon et al., 2007) provide a mechanistic representation of the organic aerosol chemistry
75 and relevant process, and lead to an improved agreement with ambient SOA measure-
76 ments (Lee-Taylor et al., 2011; Mouchel-Vallon et al., 2020). Chemical mechanisms gener-
77 ated by GECKO-A typically include millions to tens of millions of reactions, and hun-
78 dreds of thousands intermediate species (Aumont et al., 2005). MCM mechanisms are
79 handwritten and much smaller than GECKO-A ones as they represent only 2-3 first gen-
80 erations of chemistry. Due to the remarkable computational cost, no air quality mod-
81 els or chemistry-climate models can afford to run with GECKO-A included in the fore-
82 seeable future. To our knowledge, only the study by (Li et al., 2015) attempted to in-

83 include the MCM organic chemistry mechanism into a regional 3D model and faced com-
84 putational challenges.

85 Recent years have seen quite a few inspiring applications in developing machine learn-
86 ing emulators using explicit/process-level models and implementing the trained emula-
87 tors into large-scale models (Brenowitz & Bretherton, 2018; Beucler et al., 2020; Get-
88 telman et al., 2021). Replacing complex processes with ML emulators have potential ad-
89 vantages by learning non-linear relationships that can represent underlying physical or
90 chemical processes not captured in simple empirical characterizations, as well as mul-
91 tiple orders of magnitude speedups in computation when compared to fully coupled process-
92 based models. However, maintaining both an acceptable level of accuracy and a system
93 that remains numerically stable through an adequate amount of time with emulators re-
94 mains challenging.

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

107 Our primary aim in this work is to extend atmospheric chemistry emulation to or-
108 ganic aerosols for which current climate models do not currently account for. Addition-
109 ally, this proof of concept study evaluates two different types of neural network archi-
110 tectures: (1) a feed-forward, fully-connected network and (2) a recurrent neural network
111 (RNN). Both models were designed to be feasibly integrated into current 3D transport
112 models, that can provide fast and accurate predictions for organic aerosol concentrations.
113 The RNN was chosen to determine if it could help to overcome numerical stability prob-
114 lems, as observed by others using fully-connected model architectures (Brenowitz & Brether-

115 ton, 2018; Kelp et al., 2020), as they come equipped with feedback connections that can
116 store information about previous events in the form of a latent vector, e.g. RNNs pos-
117 sess memory (Hochreiter, 1991). As these architectures were developed to learn repre-
118 sentations of sequential data to solve temporal problems, they offer the ability to use multi-
119 length inputs. However, incorporating multi-length inputs into 3D models would require
120 us to dissociate chemistry production from other processes (e.g. transport, removal). We
121 address this with the development of a novel "1-step" RNN that only requires a single
122 time step of input, but relies on a separate simple neural network to initialize the hid-
123 den state vector for the very first time-step.

124 The paper is organized as follows: Section 2 outlines the data generation, train-
125 ing, hyperparameter tuning, and evaluation procedures for the reference model and both
126 neural network types. In Section 3, we characterize the two models performance rela-
127 tive to the GECKO-A data sets for different precursor species, and compared to each
128 other to assess the overall strengths and weaknesses of each model architecture. Both
129 model types are also tested on data sets that help assess the ability of the models to gen-
130 eralize into new domains. Sections 4 and 5 provide a brief discussion about the pros and
131 cons of these different model architectures and the ongoing challenges regarding numer-
132 ical stability, computational cost, and interpretability for emulating SOA production with
133 ML.

134 2 Methods

135 2.1 Description of the reference model

136 To provide reference chemical mechanisms, we used the GECKO-A chemical gen-
137 erator (Aumont et al., 2005; Camredon et al., 2007), which describes in great details the
138 chemical oxidation of organic compounds in the atmosphere. The resulting chemical mech-
139 anisms for each SOA precursor species are complete (down to the ultimate products CO_2
140 and H_2O), and explicit by preserving knowledge of the molecular structures of all the
141 intermediate compounds. Compared to other widely used semi-explicit chemical mod-
142 els such as MCM, GECKO-A can consider many generations of oxidative chemistry i.e.,
143 20 generations are considered here (Lee-Taylor et al., 2011). This has important impli-
144 cations for the formation of organic aerosols as SOA formation arises from a multitude
145 of partly oxidized compounds, rather than from a few dominant molecules. In addition,

146 the SOA formation timescale varies greatly for different precursors. For example, at am-
147 bient conditions, it takes only a few hours to form SOA from dodecane vs. several days
148 to form SOA from toluene (Hodzic et al., 2014). For the gas-particle partitioning of or-
149 ganic molecules, dynamic partitioning is used. It is reasonable to consider GECKO-A
150 simulations as a benchmark for building an emulator for SOA chemistry given its rea-
151 sonable agreement with observations shown for both comparisons with chamber mea-
152 surements e.g. for alkanes and alkenes compounds (La et al., 2016) and ambient mea-
153 surements e.g. during MIRAGE, BEACHON and Go-AMAZON, (Lee-Taylor et al., 2011,
154 2015; Mouchel-Vallon et al., 2020).

155 **2.2 GECKO-A generated training datasets**

156 We ran the GECKO-A model with systematically varied input parameters to gener-
157 ate the dataset used to train the machine learning emulator. At this stage we focus
158 on three representative SOA precursors: toluene, dodecane, and α -pinene. Toluene and
159 dodecane are emitted from a wide range of anthropogenic sources, while α -pinene is one
160 of the major SOA precursors emitted by vegetation. These compounds together contribute
161 substantially to the global SOA burden. We generate chemical mechanisms and corre-
162 sponding datasets for these precursors including the OH oxidation mechanisms of toluene,
163 dodecane and α -pinene. For each oxidation mechanism, the reactions of the precursor
164 with oxidants other than OH were not considered. Thus, the considered chemistry is mostly
165 representative of daytime conditions.

166 Based on our current understanding of atmospheric chemistry and the common chemistry-
167 climate modeling frameworks, we identified the following six variables that are key to
168 predicting SOA formation from VOC oxidation under tropospheric conditions: (1) tem-
169 perature, (2) solar zenith angle, (3) pre-existing aerosol mass, (4) ozone concentrations,
170 (5) nitrogen oxides (NO_x) concentrations, and (6) OH radical concentrations. The range
171 of variability considered for these parameters and the associated sampling scheme is sum-
172 marized in Table 1 and illustrated in Figure 1. Additionally, a diurnal variation in tem-
173 perature of an average 5 degrees amplitude is used. In each training data set, we use the
174 Latin Hypercube sampling approach to obtain two-thousand environmental input com-
175 binations. Temperature, solar zenith angle, ozone, and OH were all sampled uniformly,
176 whereas pre-existing aerosol concentrations and NO_x were sampled as a logarithmic dis-
177 tribution. The combination of these ranges is relevant for a wide range of tropospheric

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

Table 1. Environmental parameter ranges used in GECKO-A simulations.

178 conditions, from remote to moderately polluted environments. These environmental vari-
 179 ables, with the exception of temperature, are held constant in a given 5-day GECKO-
 180 A box model simulation in an attempt to coerce the ML models to generalize better and
 181 be more robust to over fitting.

182 In each dataset, the initial concentrations of a given SOA precursor were set to an
 183 arbitrary low value of 10 ppt similar to previous studies (Lannuque et al., 2018; Hodzic
 184 et al., 2014, 2015). Although the tropospheric concentrations of the precursor can be higher
 185 in polluted regions, this low value is representative of the remote atmosphere, and was
 186 chosen so that the amount of aerosol produced from the given precursor is small com-
 187 pared to preexisting OA and will not impact the gas/particle partitioning, nor the over-
 188 all photochemical reactivity. Precursor spans several orders of magnitude in our 5-day
 189 GECKO-A simulations as it decays exponentially and is effectively consumed before the
 190 end of each simulation leading to the production of thousands of intermediate organic
 191 gases. As shown on Figure 1, complexities of the organic chemistry are illustrated by the
 192 wide variation of produced SOA mass with respect to each environmental variable. For
 193 example, significantly higher SOA mass concentrations are produced at colder temper-
 194 atures.

195 As the precursor mass is exponentially distributed, before using the precursor data
 196 as input to the ML model, we transform the precursor values by taking the base-10 log-
 197 arithm to avoid any stiffness in the system. Next, each input variable X_j in the train-
 198 ing data, including chemical concentrations and environmental variables, were standard-
 199 ized independently into z-scores according to the formula $X_j = (X_j - u)(X_j - s)^{-1}$,
 200 where u and s are the mean and standard deviation of X_j . Standardization was chosen

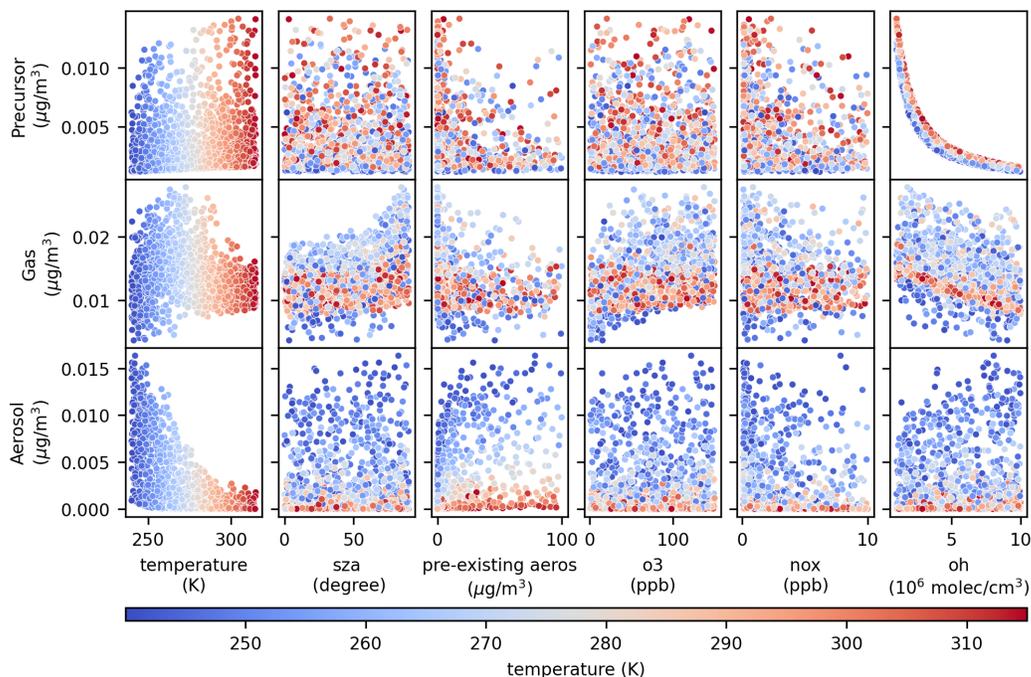


Figure 1. Training distributions (mean value through 5-days) of targets (rows) vs. environmental variables (columns).

201 over other common transformations because many features are not normally distributed
 202 in the data sets for the three species. Hence the transformation recasts the values of the
 203 input and output channels into a format where the values of each variable are centered
 204 and have similar spread. This is especially important when computing the error on the
 205 model predictions against the training data values when the weights in the model are
 206 being updated, for example, to prevent the model from over-fitting on the quantities hav-
 207 ing the largest spread.

208 A total of 2,000 experiments were run in GECKO-A for each precursor species, and
 209 output every five minutes over the course of five days. Thus, a total of 2,880,000 sam-
 210 ples were generated per species. However, as the target variables are a subset of the in-
 211 put feature variables at the previous time step, we removed the first (final) time step of
 212 the output (input) variables from each experiment leaving a total of 2,878,000 samples.
 213 These simulations were then split into three subsets: training (80%), validation (10%),
 214 and testing (10%). The training set is used to optimize the weights of an ML model, while
 215 the validation set is used to select the hyperparameters, or meta-settings describing the

216 ML model architecture, such as the number of neurons in the hidden layer, that result
 217 in the best-performing model across the space of possible models (see below). The final
 218 testing set is not used to train or adjust the machine learning models and is only used
 219 in the final evaluations described herein.

220 2.3 Neural network models

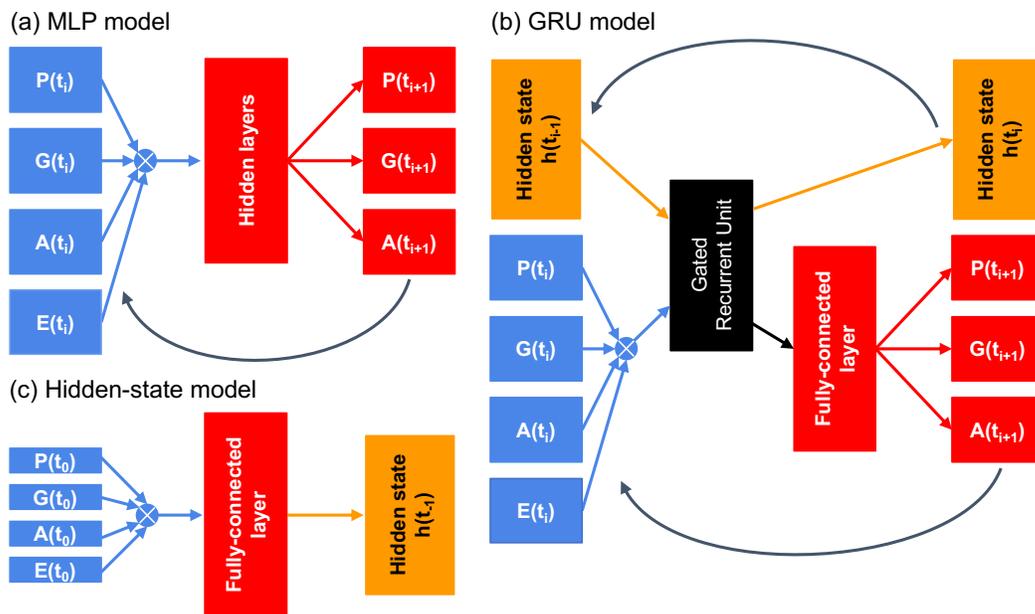


Figure 2. Model architectures showing the (a) MLP (multi-layer perceptron) and (b) GRU (Gated Recurrent Unit) models for predicting concentrations, and (c) the model for the initialization of the recurrent hidden state given an initial condition as input. In all three models, inputs at time t_i are colored 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.

221 Figure 2 shows several neural network models that we consider here as emulators
 222 for GECKO-A. In Figure 2(a), a multi-layer perceptron (MLP) architecture is shown (re-
 223 ferred to as the “MLP model”). The MLP model accepts as input the scaled values of
 224 the precursor, gas, aerosol, and environmental variables at time t_i (blue boxes in the fig-

225 ure). These quantities are concatenated into a vector, and denoted $X(t_i) \equiv (P(t_i), G(t_i), A(t_i), E(t_i))$.
226 The model outputs the precursor, gas, and aerosol values at time $t_{i+1} = t_i + \delta t$ (red boxes
227 in the figure), denoted $Y(t_{i+1}) \equiv (P(t_i + \delta t), G(t_i + \delta t), A(t_i + \delta t))$, where δt here is
228 300 seconds. The MLP is an artificial neural network that contains, in addition to in-
229 put and output layers, at least one hidden layer (horizontal red block in Figure 2(a)),
230 and is the simplest neural network architecture available. Each hidden layer contains a
231 set of perceptrons, which mathematically are linear regressions on the outputs of the pre-
232 vious layer followed by a non-linear transformation called the activation function. For
233 our MLP, we use the Rectified Linear Unit (ReLU) activation function, which sets neg-
234 ative values to 0 and allows positive values to pass through unchanged. The final hid-
235 den layer is connected to the output layer, which is a linear regression on the hidden layer
236 outputs. We tested using multiple hidden layers but found that a single hidden layer pro-
237 duced the lowest validation set error. The MLP model for each of the precursor species
238 make future predictions for the values of the chemical quantities of interest, using only
239 the current chemical state of the atmosphere. As such, any MLP model satisfies the Marko-
240 vian condition, and possesses no memory about atmospheric chemical states visited in
241 the past beyond the previous timestep.

242 As GECKO-A generates sequential trajectories describing the evolution of species'
243 concentrations over time, we have also applied a recurrent neural network (RNN) model
244 to investigate whether it may have an advantage over the MLP model due to its abil-
245 ity to utilize its memory about the past at $\{t_{i-n}, \dots, t_{i-1}\}$ to make the next prediction
246 at t_i . As such, a temporal model may be better equipped compared to the state-less MLP
247 model to describe the changes occurring to the quantities over time, in addition to lim-
248 iting and/or preventing runaway error propagation. The input to an RNN is a sequence
249 and a hidden state. The sequence elements could be scalars, vectors, or other higher-dimensional
250 tensors. The first (or last) element in the sequence is ingested by the RNN along with
251 a starting hidden state, producing an encoding of the element and a hidden state for the
252 current encoding. The next element in the sequence is then used as input along with the
253 current hidden state, which produces an encoding of the second element and another hid-
254 den state. This encoding represents not just the second element, but the elements that
255 came before it, due to the fact that the model leverages its feedback connections to pro-
256 duce the encoding. This process continues until all of the elements in the sequence have

257 been seen by the model, which produces a final encoding of the entire sequence, as well
 258 as a hidden state.

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

Similar to the MLP, the GRU model shown in Figure 2(b) 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 prob-
 lem by computing these quantities according to

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

268 where the sigmoid function $\sigma(x) = (1 + \exp^{-x})^{-1}$ projects input values to be within
 269 $[0, 1]$. The softmax function for a K -component z is $\exp(z_k) / \sum_{j=1}^K \exp(z_j)$. The ten-
 270 sors W_* and U_* , and bias terms b_* , contain the fit parameters that are updated through
 271 back-propagation during training. The \odot symbol refers to element-wise multiplication.

272 The quantity $z(t_i)$ is the update gate, and $r(t_i)$ is the reset gate, which determines
 273 how much information over past time steps to forget. The quantity $c(t_i)$ represents the
 274 current memory content in the layer and utilizes the reset gate to store information from
 275 the past. The equation for the current hidden state $h(t_i)$ uses the update gate to deter-
 276 mine how much information from the current time step to collect and how much to col-

277 lect from past time steps. The encoded information at the current time step, $Y^*(t_i)$, is
278 computed using the current hidden state. Then, as Figure 2(b) illustrates, it is passed
279 through a ReLU activation function (black arrow), and then through a fully-connected
280 layer (tall red box), which outputs the scalar precursor, gas, and aerosol values at the
281 time t_{i+1} .

282 It is common that the input sequence to an RNN is longer than length one, in which
283 case a hidden state can be immediately informed by the sequence (trajectory) to make
284 the next prediction. Applied to GECKO-A trajectories, we are free to choose the length
285 of the input sequence, which does not have to be fixed. Indeed, in our GECKO-A box
286 model simulations concentrations of organic species are only undergoing chemistry pro-
287 cessing, whereas in 3D models other processes (e.g., transport, dry and wet removal) are
288 included. Thus, we must also consider the added complexity of incorporating RNNs into
289 3D transport models. For a RNN that uses a sequence for input, each chemical variable
290 needs to be transported and stored for the number of time steps needed as input, which
291 could be memory intensive and programmatically challenging.

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

301 When there is no initial hidden state available, the initial condition $X(t_0)$ is passed
302 through a separate MLP, referred to as the hidden-state model, to obtain $h(t_{-1})$. Fig-
303 ure 2(c) illustrates that this model’s architecture accepts as input values at some t_0 for
304 precursor, gas, aerosol, and the environmental variables ($P(t_0), G(t_0), A(t_0), E(t_0)$), and
305 outputs a hidden state $h(t_{-1})$. The hidden-state model contains one fully-connected layer
306 with a linear activation. The input size is equal to the length of $X(t_0)$ while the output
307 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

340 rate value to update the adjustable parameters of both the GRU and the hidden state
341 models in tandem. This procedure is repeated until the second-to-last time step in an
342 experiment trajectory is used as input to the model.

343 During training of the GRU model, we select random experiments to create batches
344 of data when computing the total loss at each time step. One epoch is defined as all train-
345 ing experiment trajectories passing through the model once, so the same data as with
346 training the MLP model, except that the data is ordered by time (and randomized by
347 experiment). At the end of every epoch, the model is put into evaluation mode and used
348 to predict the trajectories for the validation experiments. The MAE is then computed
349 between the model predictions and the validation experiments. After each epoch the val-
350 idation MAE is used to measure improvement of the model predictions in two ways: (1)
351 to anneal the learning rate if the models performance does not improve after some num-
352 ber of epochs, e.g. it “over-fits” on the validation experiments, and (2) to stop the train-
353 ing entirely once the model does not improve on the validation experiments after some
354 number of epochs. We chose 3 and 7 epochs in (1) and (2) respectively.

355 2.5 Evaluation Procedure

356 The ability of MLP and GRU-based models to predict the time evolution of pre-
357 cursor, gas, and aerosol mass concentrations is evaluated by comparing the box model
358 predictions against the benchmark values as produced by the GECKO-A model. Here
359 each model is placed into evaluation mode, which disables any stochastic components
360 such as the recurrent dropout used when training the GRU, and is used to make pre-
361 dictions on the hold-out validation set of data that was not used during the training to
362 influence the weight updates in each model. For each validation experiment, a starting
363 amount of precursor, gas, aerosol, and environmental variables at time t_0 is passed through
364 the MLP network to obtain the predicted quantities at the first time step t_1 , where $t_1 =$
365 $t_0 + \delta t$. These predictions are then used along with the environmental variables at the
366 next time step as the next input to the model, and so forth for the length of the exper-
367 iment (see Fig. 2(b)).

368 **2.5.1 Ensembles**

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

381 **2.6 Hyperparameter optimization**

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

391 To estimate such a dependency, we use the package Earth Computing Hyperpa-
 392 rameter Optimization (ECHO) developed by the authors at NCAR (Schreck & Gagne,
 393 2021), and perform hyper parameter optimization given an objective metric for the three
 394 species for both MLP and GRU models. The objective metric for both the MLP and GRU
 395 models is the box MAE on the validation holdout set. With the MLP model, a box sim-
 396 ulation begins with the initial precursor concentration at t_0 , while for the GRU model
 397 the MAE for box simulations is computed for a set of starting times $\{t_0, t_i, \dots, t_j\}$ and
 398 added together, to also test the hidden-state model on different initial precursor amounts.

399 MLP and GRU models were optimized with ECHO for the three species, with the
400 outcomes described in appendix Appendix B. Tables C1 and C2 list the best hyperpa-
401 rameters found in each optimization study for the two models. Using the best hyperpa-
402 rameter set, an ensemble of 30 models were trained, where each model had a different
403 random weight initialization. See appendix Appendix C for more details.

404 Although it is rather trivial to train a model to output realistic predictions one time
405 step ahead from the truth (primarily due to high auto-correlation), limiting cumbersome
406 error accumulation when propagated through time is highly sensitive to model param-
407 eters in complex problems such as this. We found that efficient hyperparameter searches
408 were crucial for finding models that could successfully stabilize and limit error accumu-
409 lation through the length of the simulation. See appendix Appendix B for further de-
410 tails.

411 **3 Results**

412 **3.1 Performance of trained MLP and GRU models**

413 Table 2 lists the bulk validation performance metrics for the MLP and GRU mod-
414 els for toluene, dodecane and α -pinene. The metrics are the Pearson coefficient and the
415 Hellinger distance computed for each prediction task, and the number of unstable or run-
416 away experiments observed. Here, an experiment is considered as unstable when predicted
417 values exceed $1 \mu\text{g}/\text{m}^3$ which corresponds to an unrealistic formation yield from 10 ppt
418 of initial precursor. Unstable experiments were not used in the computed metrics reported
419 below.

420 Figures 3 and 4 illustrate the MLP and GRU models' performance on reproduc-
421 ing the experimental data for three experiments selected from the test set of toluene ex-
422 periments. Overall, they show that both the MLP and GRU models can predict exper-
423 iment trajectories that resemble the GECKO-A ones within a factor of two. The pre-
424 dicted ensemble mean matched closely with the GECKO-A trajectories for the three pre-
425 diction tasks, for both models. For toluene, Table 2 shows that all of the model predic-
426 tion tasks led to Pearson coefficients greater than 0.97. Additionally, the Hellinger dis-
427 tances for each task are all low for both MLP and GRU models, indicating that the tem-
428 poral distributions predicted for different initial conditions matched closely with those
429 generated by GECKO-A.

Task / Metric	Toluene			Dodecane			α -pinene		
	Pearson	Hellinger	Unstable	Pearson	Hellinger	Unstable	Pearson	Hellinger	Unstable
MLP precursor	0.989	0.0008	0	0.950	0.0002	0.01	0.556	0.0015	0.0005
MLP gas	0.977	0.0045	0	0.952	0.0033	0.01	0.971	0.0028	0.0005
MLP aerosol	0.927	0.0150	0	0.894	0.0161	0.01	0.939	0.0153	0.0005
GRU precursor	0.993	0.0018	0	0.950	0.0023	0	0.867	0.0025	0
GRU gas	0.990	0.0012	0	0.984	0.0020	0	0.991	0.0007	0
GRU aerosol	0.975	0.0105	0	0.961	0.0083	0	0.976	0.0091	0

Table 2. Table of computed metrics for MLP and GRU models, for each of toluene, dodecane, and α -pinene. The average Pearson coefficient and Hellinger distance are listed for the precursor, gas, and aerosol prediction tasks. The fraction of experiments that went unstable is listed for each model and task. All reported metrics for both models were computed using the testing hold-out set of experiments.

430 For the other precursor species, each model degraded in performance in different
431 ways. The GRU model mainly did not perform as well at predicting precursor. For ex-
432 ample, the Pearson coefficient was 0.886 for α -pinene compared with 0.992 for toluene.
433 The Hellinger distance for the other two species also modestly increased compared to
434 that for toluene. On gas and aerosol predictions, Table 2 indicates that the GRU per-
435 formed about the same for all three species according to these two metrics, and that none
436 of the predicted numerical values became unstable during the box simulations.

437 Table 2 shows the MLP model performance on precursor prediction was the best
438 for toluene and then α -pinene, which had Pearson values of 0.989 and 0.950, respectively.
439 However, the MLP struggled by comparison for dodecane, where the Pearson value was
440 0.556. On the gas and aerosol predictions, the MLP model was mostly consistent for the
441 three species, with gas prediction performing better by comparison to aerosol prediction.
442 Furthermore, out of 200 experiments that went unstable during a box simulation there
443 were 13 for dodecane and 2 for α -pinene, despite the fact that the Pearson score for do-
444 decane remained high across the prediction tasks.

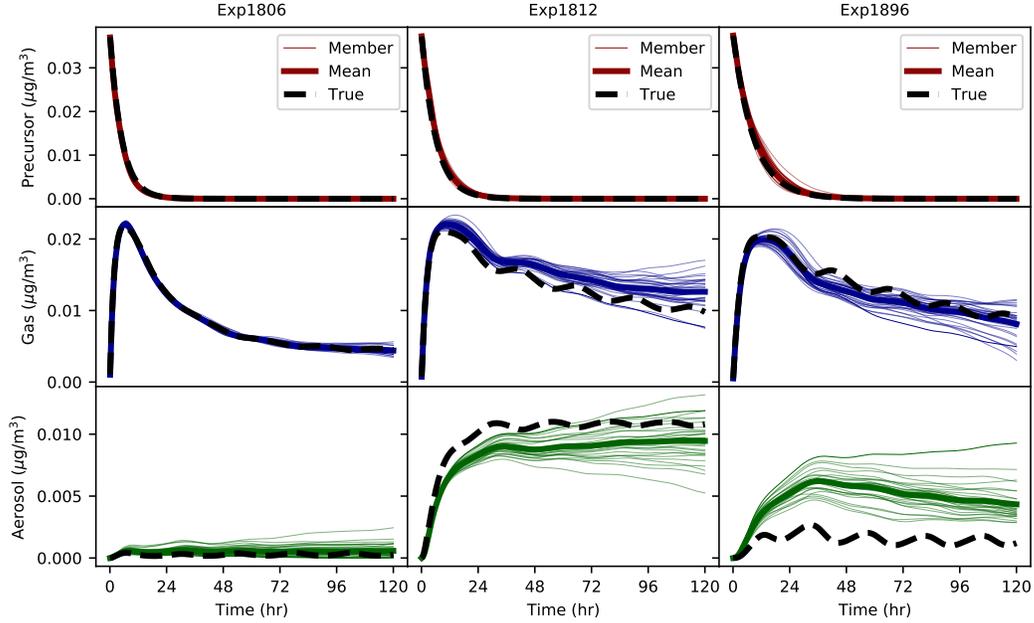


Figure 3. Three toluene sample experiment trajectories for the MLP model. Solid (thick) lines show the mean GECKO-A trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.

3.2 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 5 and the mean standard deviation (Figure 6) between the different species across all 30 ensemble members. Figure 5 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 5 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.

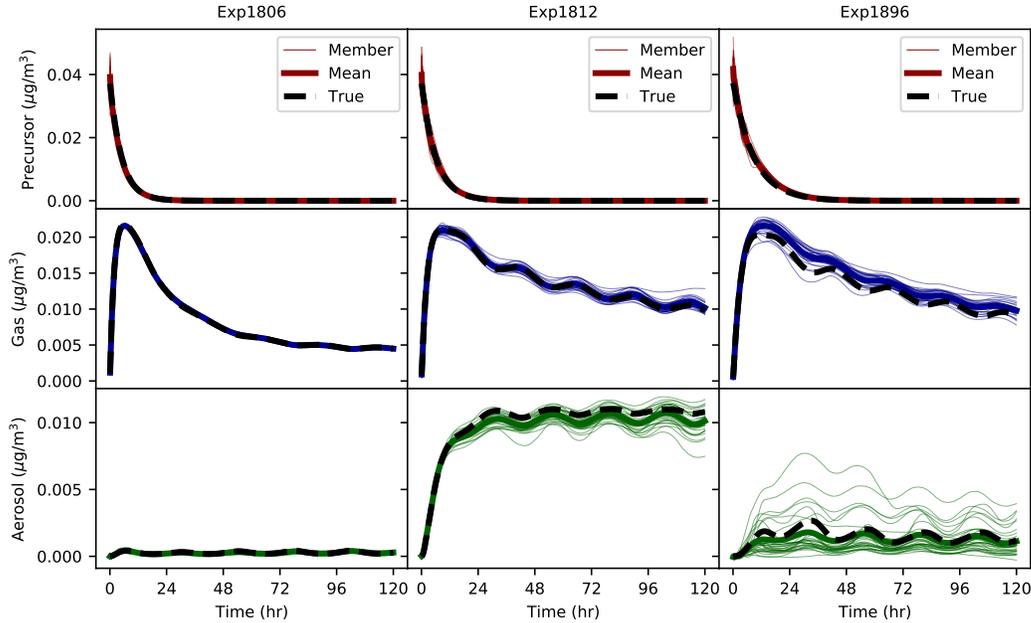


Figure 4. Three toluene sample experiment trajectories for the GRU model. Solid (thick) lines show the mean GECKO-A trajectories from 30 ensemble members, dashed lines show the reference GECKO-A trajectory, and the thin lines show each of the 30 ensemble member predictions.

460 Figure 6 shows a notable difference in the variation among ensemble members be-
 461 tween the MLP and GRU, specifically the slope of the gas and aerosol trajectories. The
 462 steep positive slope of the MLP demonstrates the inherent growing uncertainty in the
 463 model itself as it progresses further from the starting condition, which is also seen in the
 464 ensemble spread on figure 3. The GRU has a much flatter trajectory, especially in the
 465 later time steps due to the short-term memory of the trajectory being encoded into the
 466 hidden state, which could potentially make it much more suitable for maintaining sta-
 467 bility in much longer running simulations. Additionally, if it is computationally feasi-
 468 ble, one could choose to run the ensemble suite and take the mean to increase accuracy.
 469 With this approach, the mean absolute percentage errors for the GRU are less than 2%
 470 for the gas and aerosol partitions, and approximately 3-8% for the MLP.

471 3.3 Performance dependency on initial conditions

472 Next the models' performance dependency on different initial conditions was probed
 473 by selecting initial values $X(t_N)$ for some t_N in a GECKO-A experiment trajectory as

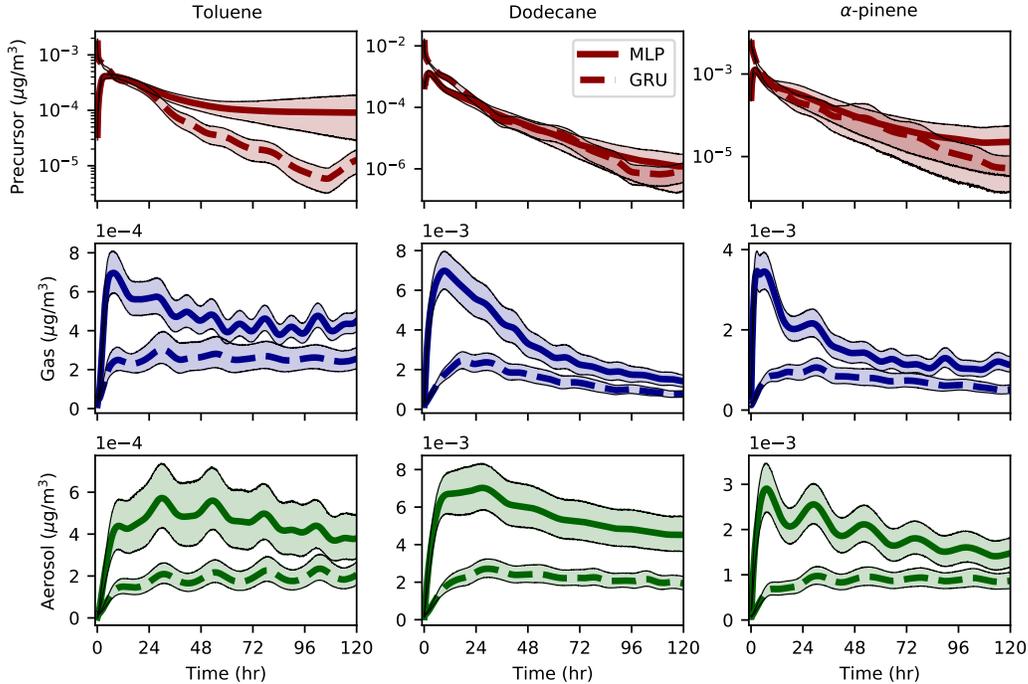


Figure 5. Ensemble bootstrapped CRPS through time. MLP (solid) and GRU (dashed) with the shaded regions representing the 95% confidence interval.

474 the initial starting point in a box simulation. For the GRU model, $X(t_N)$ is initially passed
 475 through the hidden-state model to obtain a starting hidden state. As the experiments
 476 contain 1440 total time steps, box simulations were left to run for $1440 - N$ time steps
 477 once the initial time was selected. As the instabilities observed in the MLP model dis-
 478 cussed above occurred at a range of different time steps after the box simulation was first
 479 started, we wanted to check each model's stability over as many possible time steps as
 480 there was data available. This means that box simulations started at earlier times in ex-
 481 periments will run for more time steps compared to those which started at later times
 482 in the experiments. Figure 7 shows the average ensemble Pearson coefficient for MLP
 483 and GRU models versus the initial box simulation start time. A similar plot for the Hellinger
 484 distance is shown in Figure D1.

485 Figure 7 illustrates the broad variation in GRU and MLP model performance at
 486 predicting precursor versus initial simulation start time. The variation is more signif-
 487 icant in dodecane and α -pinene, compared to toluene. For example, the lowest Pearson
 488 values are observed when the MLP box simulation is started about 2-3 days into the GECKO-

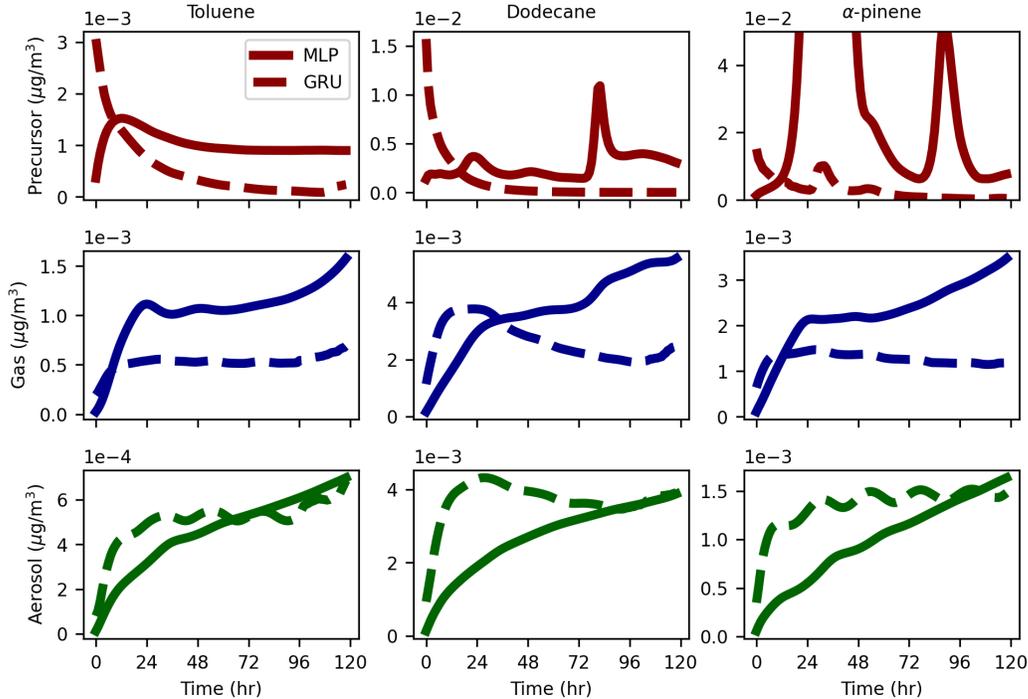


Figure 6. Mean ensemble standard deviation across all validation experiments as a function of simulation hour.

489 A experiment and runs for a similar amount of time. Then, some recovery is observed
 490 for the MLP model as observed by increasing Pearson coefficients at later start times for
 491 the prediction tasks, in particular on days 4 and 5, for the shorter box simulations. For
 492 α -pinene in particular, the GRU was also observed to go unstable at earlier start times.
 493 However, by comparison to the MLP model for all initial starting times, the total num-
 494 ber of experiments having gone unstable was significantly less.

495 On the gas prediction task, the GRU model typically performed better than the
 496 MLP at earlier start times (longer box simulations), while at the later start times (shorter
 497 box simulations) the MLP had higher Pearson scores for precursor and gas prediction.
 498 We observed in some experiments the GRU struggling at later start times to reproduce
 499 the GECKO-A predicted gas values as accurately as the MLP, in particular, when those
 500 concentration values were very small compared to the initial experiment precursor amount,
 501 but the model predictions remained stable at these times.

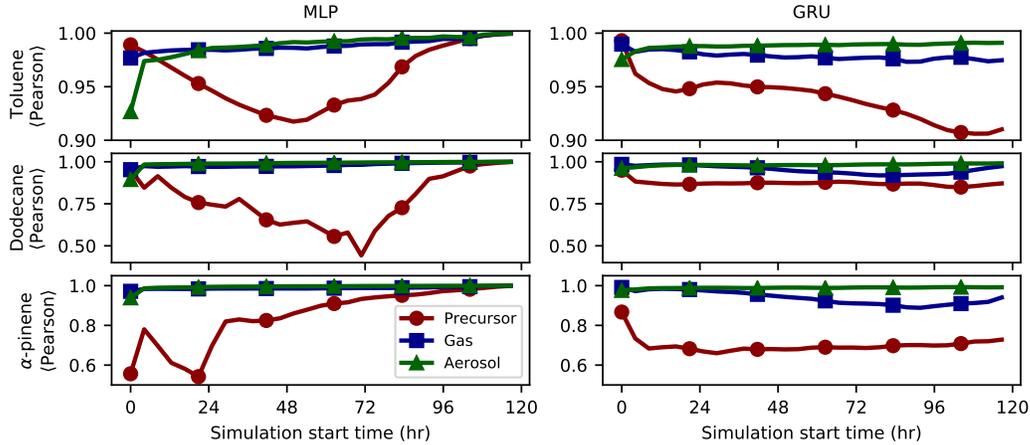


Figure 7. The mean Pearson coefficient versus the initial box simulation start time for all validation experiments.

3.4 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.

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. Thus, we consider MLP and GRU models that only perform gas and aerosol prediction, and not precursor, to probe whether reducing the total number of prediction targets will improve model performance on gas and aerosol predictions. 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. The Pearson coefficients and

	Toluene			Dodecane			α -pinene		
Task / Metric	Pearson	Hellinger	Unstable	Pearson	Hellinger	Unstable	Pearson	Hellinger	Unstable
MLP gas	0.980	0.0034	0	0.911	0.0093	0	0.957	0.0050	0
MLP aerosol	0.957	0.0095	0	0.908	0.0248	0	0.918	0.0321	0
GRU gas	0.989	0.0012	0	0.991	0.0014	0	0.990	0.0010	0
GRU aerosol	0.985	0.0148	0	0.990	0.0097	0	0.986	0.0105	0

Table 3. Table of computed metrics for MLP and GRU models which are tasked with prediction of gas and aerosol for each of toluene, dodecane, and α -pinene. The average Pearson coefficient and average Hellinger distance are listed for the two prediction tasks. The fraction of experiments that went unstable is listed for each model and task. All reported metrics for both models were computed using the testing set of experiments.

521 Hellinger distances are comparable for both model types and architectures. An overall
 522 improvement in scores is seen when predicting the precursor concentrations was not a
 523 target, but the most notable difference is that all model runs remained stable.

524 3.5 GECKO-A emulator evaluation with external datasets

525 The performance abilities of both MLP and GRU models were tested by expand-
 526 ing the data sets to include additional simulations, (1) for 10 times (X10, = 100 ppt) and
 527 100 times (X100, = 1 ppb) higher initial concentrations of the precursor, which are more
 528 representative of somewhat polluted atmospheric conditions, and (2) for simulating the
 529 diurnal variation in the precursor levels, that was not present in the original data sets
 530 which did not include the daily variability on the emissions, and removal of the precur-
 531 sor.

532 3.5.1 Model performance on increased precursor concentrations

533 The simulations performed to create X10 and X100 data sets were carried out iden-
 534 tically compared with the reference simulations starting at 10 ppt precursor concentra-
 535 tions, except that the initial precursor concentrations were increases by a factor of 10

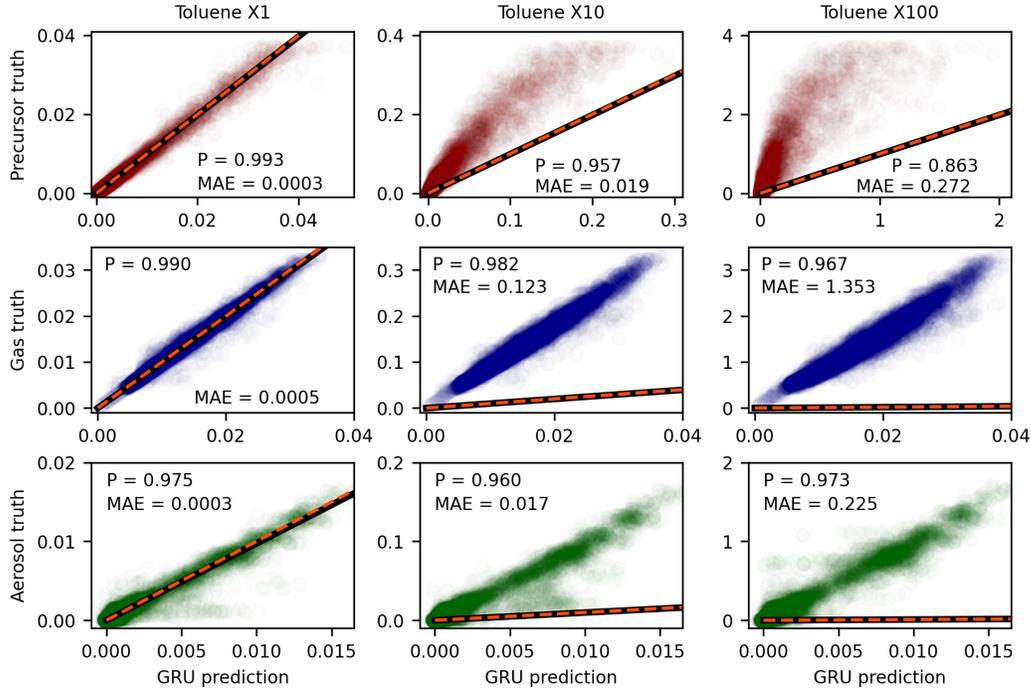


Figure 8. Scatter plots comparing the GECKO-A value (y-axis) against the GRU predicted value (x-axis) for models trained on X1 data and applied to test sets with 10X and 100X higher initial concentrations. The dashed orange line shows $y = x$, while the solid black line shows the linear relationship for models trained on the respective X^* data set.

536 (X10) and 100 (X100). The new data sets were then split into train, validation, and test
 537 data sets just as before, then transformed using the fitted scaling transformations on the
 538 original data sets. Then, the X10 and X100 test data sets were passed through MLP and
 539 GRU models that were trained on the original data set containing the smaller initial value
 540 of the precursor.

541 Figure 8 shows the predictions of the GRU model on the reference test set of ex-
 542 periments (left column), and the expanded X10 and X100 test data sets (middle and right
 543 columns, respectively). The Pearson coefficient and MAE for each prediction task are
 544 listed in the sub-panels. The figure shows that the GRU trained on the smaller initial
 545 precursor concentrations made predictions on the X10 and X100 data sets that corre-
 546 lated strongly with the true values for precursor, gas, and aerosol, as is seen by high val-
 547 ues of the Pearson coefficients for the different prediction tasks, but the MAE for each
 548 task increased by orders of magnitude with larger starting precursor concentrations.

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

560 **3.5.2 Evaluation with varying environmental conditions**

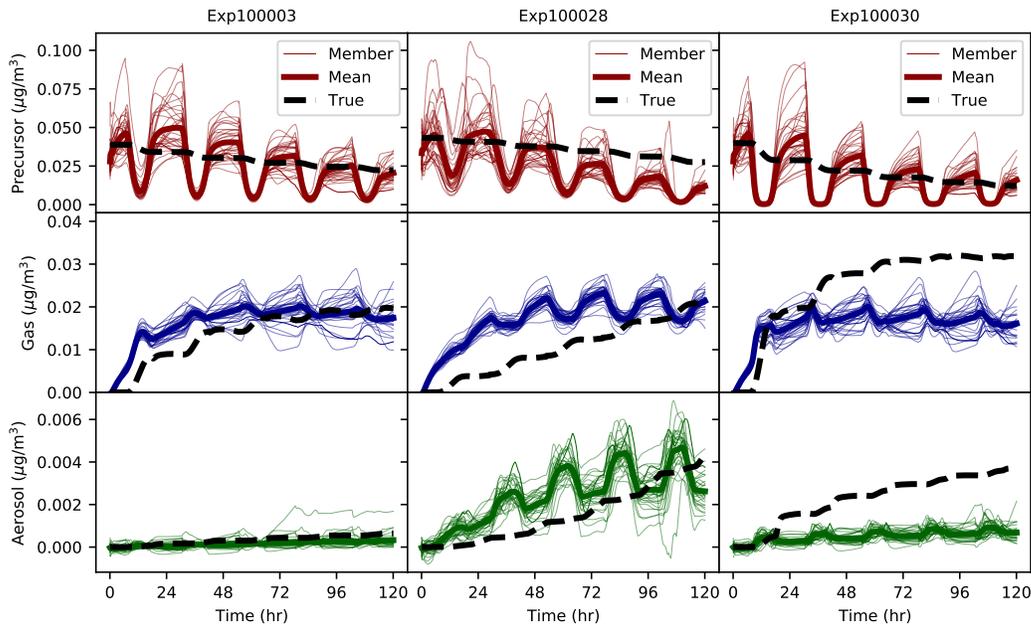


Figure 9. Examples of GRU box simulations where select environmental variables were allowed to vary with time.

561 Lastly, the models’ performance was tested on 36 experiments run for toluene that
 562 simulated daily varying conditions for five days. Like for the training data set, the pre-
 563 cursor’s initial concentration was set to 10 ppt. Initial temperature, pre-existing aerosol

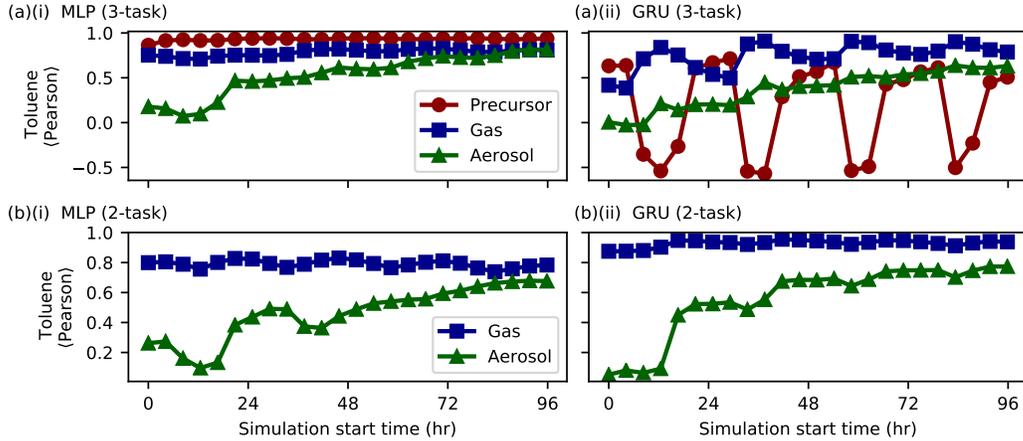


Figure 10. The average Pearson coefficient versus the initial box simulation start time computed from 36 experiments for toluene. (a) The results for the three-task MLP and GRU models are shown in (i) and (ii), respectively. (b) The same quantities as in (a) except for the two-task MLP and GRU models. All box simulations ran for 24 hours.

564 seed, ozone and NO_x were randomly selected in the same ranges as the training data set
 565 (Tab. 1). CO mixing ratio was initialized to 100 ppb. Relative humidity was held con-
 566 stant to a random value picked in the 50-80 % range. The latitude was also randomly
 567 selected in the 80S-80N range. Contrary to the training data set, after initialization, all
 568 chemical concentrations were free to evolve with the diurnal cycle to simulate a realis-
 569 tic atmospheric degradation of toluene and the subsequent organic aerosol formation.

570 Because the experiments simulated a diurnal cycle and started at midnight, box
 571 simulations were performed with MLP and GRU models at different starting times in
 572 the experiments to assess the impact of training the models on daytime oxidation only.
 573 Three example experiment trajectories are shown in Figure 9 for the 3-task GRU model,
 574 for the full 5-day box simulations which all began at midnight (the same examples for
 575 the MLP model are shown in Figure D4). The simulations performed at other starting
 576 times covered a shorter 1-day window. Figure 10 shows the average Pearson coefficient
 577 for these shorter simulations for both the 3- and 2-task MLP and GRU models.

578 Figure 9 shows the predicted curves for simulations starting at midnight are no-
 579 tably different compared with those in Figure 4. In particular, the GRU model seems
 580 to have captured some of the diurnal changes, where oxidation appears to proceed dur-
 581 ing the day, but not at night. Drastic changes in the predicted precursor amounts are

582 observed during day-night transition periods. However, the predicted concentration val-
583 ues are clearly not in agreement with the true values. Although it is less obvious, close
584 inspection of the MLP predicted precursor values in Figure D4 shows that it too responded
585 to the diurnal variation in the extended experiments, but with poor numerical accuracy.
586 Both MLP and GRU models predicted that all experiments remained stable at all sim-
587 ulation times.

588 Figure 10(a) shows that over a 1-day simulation window, the performance still dropped
589 relative to the experiments where the environmental variables were held constant for toluene.
590 The MLP model had comparably high Pearson score for precursor prediction across the
591 start times, but gas and aerosol performance was lower by comparison. The periodic re-
592 sponse of the GRU to the diurnal signal in Figure 10(a)(ii) is indicated by the sign-change
593 in the average Pearson value for predicted precursor, which goes negative when box sim-
594 ulations were started during day-time hours, while those started overnight stayed posi-
595 tive. The GRUs performance on gas and aerosol prediction also peaked for simulations
596 that started during the middle of the day-time, and was poorest by comparison for those
597 started late at night. Figure 10(b) shows that MLP and GRU models, which were only
598 tasked with predicting gas and aerosol, performed mostly similar to the 3-task models,
599 with notable gas performance improvement for the 2-task GRU.

600 **3.6 Computational performance of emulator models and GECKO-A**

601 In addition to being able to reproduce reasonably well the evolution of concentra-
602 tions of organic compounds on the test data sets for the three species, the MLP and GRU
603 emulators also led to significant computational gains. Table A1 lists estimates for the
604 time required by GECKO-A, MLP, and GRU models to advance one time step, e.g. five
605 minutes of simulation time, for the three precursor species. For toluene, GECKO-A re-
606 quires 0.9 seconds and is about 78 and 244 times faster than dodecane and α -pinene, re-
607 spectively. By comparison, both MLP and GRU models require about the same time for
608 the three precursor species on the CPU, typically a few microseconds, with the MLP faster
609 than the GRU by up to a factor of five. Thus, for toluene both neural network models
610 could be expected to perform hundreds of times faster, while for α -pinene the expected
611 speed-up could be up to 4-6 orders of magnitude faster than the explicit model.

612 4 Discussion

613 In general, the comparative differences seen between the MLP and novel 1-step GRU
 614 applications show the advantages of a recurrent neural network emulator in a few key
 615 areas. First, its overall accuracy, especially at integrated time steps further away from
 616 its starting conditions, is notably higher than the MLP model. Furthermore, the encod-
 617 ing of a hidden state which can represent the trajectory of each input, the key feature
 618 of a recurrent network architecture, appears to help constrain model uncertainty and ul-
 619 timately, numerical stability. Most neural network applications for atmospheric chem-
 620 istry have not yet begun to examine such model uncertainties. Our use of training a suite
 621 of ensemble members using the exact same architecture for each model, and only initial-
 622 izing the weights differently prior to training, provides some evidence that model uncer-
 623 tainty can be sensitive to, and better constrained by, certain model types. Related, we
 624 also note that our recurrent model remains numerically stable for all species and for most
 625 starting initial conditions, which is not true for a small percentage of MLP member /
 626 experiment combinations. This insight may not have been detected without the inspec-
 627 tion of an ensemble suite, as most of the MLP models remained stable.

628 Maintaining numerical stability with the use of emulators for atmospheric chem-
 629 istry and other atmospheric parameterizations is a known issue and initial steps have been
 630 taken to address it (Brenowitz & Bretherton, 2018; Kelp et al., 2020, 2021). These re-
 631 cent studies found some performance improvements by using a “recurrent training” scheme,
 632 where a model was rolled out in time for n time steps during training, and a loss was
 633 calculated on the sum of n time steps, instead of a single time step. However, the mod-
 634 els used in these studies were not recurrent neural networks, as the network architectures
 635 were that of an MLP (Brenowitz & Bretherton, 2018) and an encoder/decoder frame-
 636 work (Kelp et al., 2020, 2021), which only utilized feed-forward connections. Rather, train-
 637 ing these models relied on the multi-time step loss function as a means to update the
 638 model weights using a sequence instead of a single length input. Our GRU model pro-
 639 vides an alternative approach, by rolling out the model to the end of the training exper-
 640 iment and calculating the loss at successive single time steps, the feedback connections’
 641 memory of the trajectory through $t-1$ is simply used as input at t along side the cur-
 642 rent values of the precursor, gas and aerosol.

643 Recurrent neural networks have not yet been thoroughly explored in 3D atmospheric
644 modeling, although there have been applications in other earth systems areas including
645 hydrology (Kratzert et al., 2018; Ardabili et al., 2019), earthquake magnitude predic-
646 tion (Mousavi & Beroza, 2020), rain-runoff (Boulmaiz et al., 2020) and wind velocity fore-
647 casting (Irrgang et al., 2020), as well as vegetation growth estimation (Reddy & Prasad,
648 2018). One reason for this might be the lower dimensional nature of many of these mod-
649 els, which would be computationally less burdensome to put into production as opposed
650 to integrating a model that requires multiple time steps of input into a full 3D climate
651 or weather model. For this reason, we have developed a method that still only requires
652 one time step of input but maintains the advantage of having an encoded memory of past
653 time steps. A small disadvantage of our framework is that it does require an additional
654 model to predict the initial hidden state prior to running the GRU. However, if the com-
655 munity ultimately finds that it is computationally and programmatically feasible to cou-
656 ple large recurrent networks into full 3D transport models, investigation of training re-
657 current models with multiple time steps of input would be a recommended pathway.

658 Although there are clear benefits demonstrated from use of a recurrent network,
659 there are computational limitations. The hidden state increases the input needed for each
660 prediction from 9 for the MLP model to 1000 for the GRU. This is not problematic for
661 1D validation efforts, but would become too memory intensive if this model were inte-
662 grated into 3D simulations. A smaller GRU hidden state is possible but may result in
663 drops in performance. If directly shrinking the vector is not feasible, lossy compression
664 of the vector with principal component analysis or an autoencoder may balance a smaller
665 performance loss with slightly more computation. For this reason, despite the lower per-
666 formance metrics, we still find value in simplified neural networks such as the MLP if
667 they can still approximate a solution within the given tolerance. Additionally, some per-
668 formance could be sacrificed for a smaller GRU model (see Figure B3).

669 Our results demonstrate some ML generalization challenges involving the selection
670 and training of neural networks on experiments with both small initial precursor con-
671 centrations and select static environmental variables. For example, low precursor con-
672 centrations of 10 ppt were chosen primarily to limit the influence of a single precursor
673 on the photochemical reactivity, and gas/particle partitioning in GECKO-A. However,
674 this had a large impact on the generalizability outside the training range (Fig 8). If we
675 were to implement our models into a 3D climate model, they would need to be trained

676 on a larger range of precursor values that are also representative of more polluted at-
677 mospheric conditions. Additionally, environmental variables outside of temperature were
678 held constant in an effort to help the models generalize better, but this was not success-
679 ful. While there did not appear to be any direct evidence of over-fitting to the training
680 data, an open question remains of how to properly configure the reference box models
681 to provide data to best capture the physical relationships in a complex chemical system.
682 One speculation is that the GRU could generalize more effectively than observed here
683 by having varying environmental fields within an experiment, to better parameterize the
684 models feedback connections. Many other generalization questions also remain, such as
685 the inclusion of night chemistry (oxidation with O₃ and NO₃), as well as reactions be-
686 tween species originated from various precursors.

687 To our knowledge, this is the first neural network emulation of organic atmospheric
688 chemistry. As a result, there are many areas that warrant further exploration: (1) cou-
689 pling both the MLP and GRU models to a 3D chemistry-climate model, such as WRF-
690 or GEOS-Chem, to better understand their successes and shortcomings, (2) further quan-
691 tification of the underlying uncertainties in model predictions to determine whether the
692 error sources originate from the data or the model architecture choices, or both, (3) test-
693 ing of different data sets, training regimes, and model architectures to better general-
694 ize across different chemical regimes (such as daytime vs. nighttime chemistry), (4) ap-
695 plication of transfer learning for domain adaption (Kouw & Loog, 2018), and for poten-
696 tially managing the cumbersome production of data sets, (5) incorporating physical con-
697 straints into the model architecture or training procedure as a means for constraining
698 model outputs, for example the total mass or the number of C atoms needs to be con-
699 served, and (6) utilizing explainable and interpretable methodologies to better under-
700 stand what the model has learned, and what it is using to drive its predictions.

701 **5 Conclusions**

702 In summary, the neural network emulators proposed here, especially the GRU model,
703 appear to provide fast and accurate representations of complex chemical processes. As
704 such they may be incorporated into 3D models to potentially provide insight into im-
705 portant chemical processes currently absent climate models. The recurrent neural net-
706 work considered contained feedback connections, and was generally more stable over longer
707 box simulations and maintained higher numerical accuracy with the GECKO-A data sets,

708 as compared to the MLP architectures, which did not possess any memory capabilities
 709 by way of feedback connections. Additionally, models with only two output tasks for pre-
 710 dicting gas and aerosol quantities and not precursor led to further performance and sta-
 711 bility improvements in both models. Furthermore, extensive hyper-parameter search was
 712 a crucial step in finding the best models in each case. The novelty of our recurrent neu-
 713 ral network that only requires one time step of input data allows for a similar ease of trans-
 714 fer compared to those already explored such as random forests and MLPs. This approach
 715 does not depend on the specific data set used for training and validation, and was de-
 716 signed so that a recurrent model can be integrated into current 3D models without adding
 717 additional transport complexity. Thus, this “1-step” approach could be applied in other
 718 areas where emulators are being used for the prediction of time-ordered quantities.

719 Appendix A Average time step comparison

Model	Toluene		Dodecane		α -pinene	
GECKO-A	0.9 s	1	71 s	1	220 s	1
MLP CPU	2.1 μ s	430	0.8 μ s	8.88×10^4	1.6 μ s	1.38×10^5
MLP GPU	0.08 μ s	11250	0.07 μ s	1.01×10^6	0.08 μ s	2.75×10^6
GRU CPU	3.1 μ s	290	3.2 μ s	2.22×10^4	3.3 μ s	6.67×10^4
GRU GPU	0.38 μ s	2368	0.38 μ s	1.87×10^5	0.38 μ s	5.79×10^5

Table A1. The average time each model required to advance 300 seconds of simulation time. For each species, the first column shows the time step in seconds while the second column shows the ratio of the GECKO-A time step to each model time step.

720 Table A1 compares the average speed in which GECKO-A and MLP and GRU mod-
 721 els take to advance 300 seconds of simulation time. The neural network models were eval-
 722 uated on NCAR’s casper supercomputer, on a node that contained an 18-core 2.3-GHz
 723 Intel Xeon Gold 6140 processors (CPU) and a NVIDIA Tesla V100 32GB graphics cards
 724 (GPU). The GECKO-A simulations were performed on NCARs cheyenne supercomputer,
 725 on a node that contained a 2.3-GHz Intel Xeon E5-2697V4 (Broadwell) processor (CPU).

726 Appendix B Hyperparameter optimization

727 The ECHO package is based on Optuna (Akiba et al., 2019), and facilitates opti-
 728 mization using the high-performance computing clusters available at NCAR. The opti-
 729 mization procedure begins by initiating a “study” and performing the first “trial” where
 730 values are selected for a set of hyper parameters within specified ranges, the model is trained,
 731 and its performance measured in box simulations. The outcome of the trial is saved to
 732 the study along with other metadata. For the current objective, any trial is independent
 733 from any other trial. Trials are ran until the optimization converges or the number of
 734 trials saved to a study reaches a predetermined number. Upon the completion of a study
 735 the relative importance of each hyper parameter on model performance may be estimated.
 736 To sample hyper-parameters, we selected the Tree-structured Parzen Estimator (TPE)
 737 (Bergstra et al., 2011) for this task. For each hyper parameter, the TPE method fits a
 738 Gaussian Mixture model (GMM) $l(x)$ to the set of parameter values associated with the
 739 best MAE, for all of the trials that have been carried out to completion. TPE addition-
 740 ally fits a second GMM to the leftover parameter values, and then chooses the next pa-
 741 rameter value that maximizes $l(x) / g(x)$. We initially delay using the TPE sampler and
 742 use random sampling instead. We have observed for the present models that this initial
 743 step helps to inform the TPE sampler by initially supplying observations that the GMM
 744 models may leverage to make better informed parameter selections. For additional de-
 745 tails about hyper parameter optimization, see the supporting information.

746 Once a study is complete, the parameters sampled in each trial along with the box-
 747 MAEs are used to compute the relative parameter importance. There are a variety of
 748 approaches for such estimation, including mean decrease impurity evaluation (MDI) (Louppe
 749 et al., 2013) and functional analysis of variance (functional ANOVA, or fANOVA hence-
 750 forth) (Hutter et al., 2014). Both approaches utilize tree-based ensemble methods to es-
 751 timate the relationship between the values of a set of hyperparameters used to train a
 752 model, and the optimization objective value that resulted. The MDI estimation of a hy-
 753 perparameter is zero when it depends only on the relevant variables, hence it is irrele-
 754 vant to making a prediction. The most relevant hyperparameter has the largest estima-
 755 tion value. Similarly, in fANOVA when the estimated variance between input x_i and out-
 756 put y_j is low or zero, it is not an important input feature, and vice versa.

757 We also compute the partial dependence plot for each parameter, which estimates
 758 the marginal effect one (or more) features have on the predicted outcome of a machine
 759 learning model (Friedman, 2001). A partial dependence plot can show whether the re-
 760 lationship between the box-MAE and an input feature is linear, monotonic or more com-
 761 plex. For example, when applied to a linear regression model, partial dependence plots
 762 always show a linear relationship.

763 MLP and GRU model optimization for the three species was performed and the
 764 best model parameterization was selected from each optimization study. Figure B1(a)
 765 plots the optimization objective for GRU model to the toluene data set versus number
 766 of optimization trials. Figure B1(b) illustrates partial dependence curves for the GRU
 767 layer size for the three species. Figures B2 and B3 illustrate partial dependence curves
 768 for each hyper parameter varied by ECHO for MLP and GRU models respectively. Ta-
 769 bles B1 and B2 list the hyperparameter importance as measured by MDI and fANOVA.

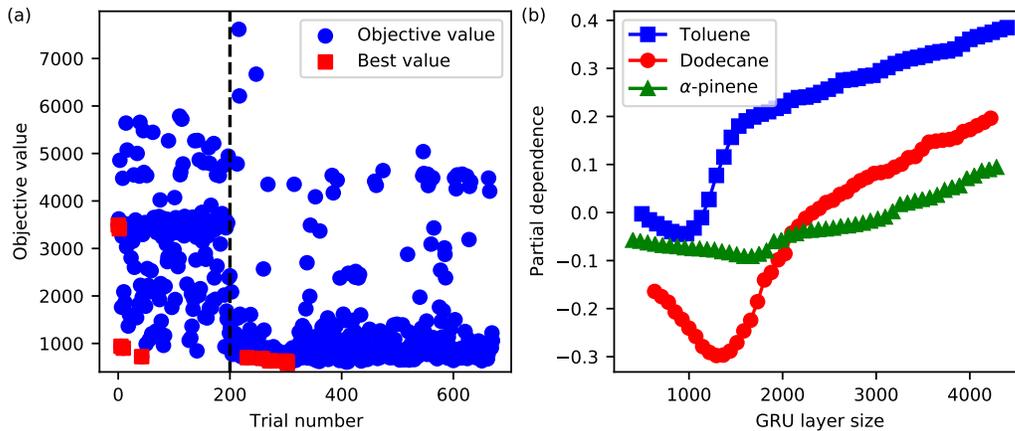


Figure B1. (a) The optimization history for the GRU model trained. on the α -pinene data set. (b) Partial dependence versus GRU layer size.

770 The blue dots in Figure B1(a) show outcomes of trails. A red dot indicates when
 771 a set of hyperparameters is the best performing one in a study. The horizontal line in-
 772 dicates the two stages of the algorithm. To the left, random sampling was used to se-
 773 lect a set of hyperparameters. To the right, TPE sampling was used. As the figure shows,
 774 the optimization procedure mostly converges to better performing models in the second
 775 stage, but no improvement in the study was observed after about 300 trial attempts in
 776 this example.

777 Figure B1(b) shows how the optimization metric depends on changes to values of
 778 hyperparameters, referred to as the partial dependence. In the figure, the partial depen-
 779 dence shows how the GRU model MAE summed over 1439 time steps depends on the
 780 GRU model layer size, and a random seed in python 3.7 of 1000. The best layer size for
 781 each species is the one that leads to the most negative value of the partial dependence,
 782 and is at a curve’s global minimum. In this example, toluene and dodecane are more sen-
 783 sitive to the value of the GRU layer size relative to α -pinene, especially for layer sizes
 784 that are near the best value. Table B2 additionally shows that the computed MDI and
 785 fANOVA values for the hidden size of the GRU are larger than that for α -pinene. The
 786 curve for α -pinene is also more flat in appearance and encompasses a smaller range of
 787 values of the partial dependence compared to toluene and dodecane. Overall, the most
 788 important hyperparameters in the optimization studies for the three species were the loss
 789 weight for hidden state model, the initial learning rate, and the GRU layer size. For do-
 790 decane, the aerosol loss weight and the batch size were also estimated to be important
 791 training parameters. For the MLP models, Table B1 shows that both MDI and fANOVA
 792 score the learning rate as the most important training parameter for all three species,
 793 with the batch size the second most important.

794 Figures B2 and B3 illustrate partial dependence curves for hyper parameters used
 795 in each optimization study, for MLP and GRU models, respectively. Tables B1 and B2
 796 list the hyperparameter importance value estimations using the MDI and fANOVA meth-
 797 ods, for MLP and GRU models, respectively.

798 **Appendix C Model and training parameters**

799 Tables C1 and C2 list the best hyperparameter values in optimization studies for
 800 the MLP and GRU models respectively. The parameters listed in these tables were used
 801 to train ensembles of models that were then used to produce the results shown in the
 802 figures in the main text. Figure C1 shows the CRPS for MLP and GRU models which
 803 are tasked with gas and aerosol prediction only.

804 **Appendix D Additional results**

805 Figure D1 shows the average Hellinger distance for MLP and GRU models tasked
 806 with predicting precursor, gas, and aerosol, versus the start time of the box simulation.
 807 Figures D2 and D3 show the average Pearson coefficient and the average Hellinger dis-

MLP	Toluene		Dodecane		α -pinene	
Parameter	fANOVA	MDI	fANOVA	MDI	fANOVA	MDI
Learning rate	0.969	0.945	0.759	0.540	0.759	0.618
Batch size	-	-	0.107	0.179	0.169	0.134
Hidden layer size	0.022	0.034	0.063	0.076	0.017	0.068
Epochs	0.006	0.013	0.055	0.124	0.037	0.120
L2 penalty	0.002	0.004	0.003	0.041	0.018	0.044
L1 penalty	0.001	0.001	0.013	0.039	0.001	0.016

Table B1. Hyperparameter importance values for optimization studies of a MLP model trained on toluene, dodecane, and α -pinene GECKO-A experiment trajectories. The maximum number of trees used and the maximum depth was set to 1000 in all estimations. The batch size was fixed at 8192 for toluene.

GRU	Toluene		Dodecane		α -pinene	
Parameter	fANOVA	MDI	fANOVA	MDI	fANOVA	MDI
Hidden loss weight	0.310	0.386	0.040	0.020	0.393	0.518
GRU hidden size	0.200	0.140	0.198	0.112	0.186	0.046
Learning rate	0.118	0.205	0.152	0.297	0.315	0.194
Precursor loss weight	0.101	0.052	0.054	0.024	0.026	0.081
Batch size	0.079	0.052	0.165	0.053	0.025	0.032
Aerosol loss weight	0.072	0.042	0.184	0.363	0.017	0.022
Gas loss weight	0.060	0.043	0.057	0.019	0.017	0.039
GRU dropout	0.045	0.042	0.112	0.023	0.013	0.032
L2 penalty	0.013	0.038	0.039	0.088	0.008	0.034

Table B2. Hyperparameter importance values for optimization studies of a GRU model trained on toluene, dodecane, and α -pinene GECKO-A experiment trajectories. The maximum number of trees used and the maximum depth was set to 1000 in all estimations.

Parameter	Toluene	Dodecane	α -pinene
Learning rate	1.39×10^{-5}	4.41×10^{-6}	6.39×10^{-6}
Batch size	8192	2538	6907
Hidden layer size	4902	2655	4049
Epochs	841	907	1450
L2 penalty	3.49×10^{-4}	2.60×10^{-3}	6.61×10^{-5}
L1 penalty	1.39×10^{-5}	1.22×10^{-11}	1.76×10^{-5}

Table C1. The values of the best hyperparameters in the optimization studies for the MLP models for the three species. The batch size for toluene was fixed at 8192. The leaky ReLU activation function was used after the hidden layer.

Parameter	Toluene	Dodecane	α -pinene
Hidden loss weight	0.161	0.980	1.896
GRU hidden size	1215	1253	1850
Learning rate	6.926×10^{-5}	5.275×10^{-5}	2.474×10^{-5}
Precursor loss weight	0.812	0.537	0.805
Batch size	1426	980	767
Aerosol loss weight	0.421	0.911	0.894
Gas loss weight	0.151	0.962	0.621
GRU dropout	0.122	0.137	0.415
L2 penalty	2.269×10^{-8}	4.138×10^{-8}	1.171×10^{-8}

Table C2. The values of the best hyperparameters in the optimization studies for the GRU models for the three species. Other fixed parameters used were an early stopping patience of 6 and the learning rate annealing patience of the 2.

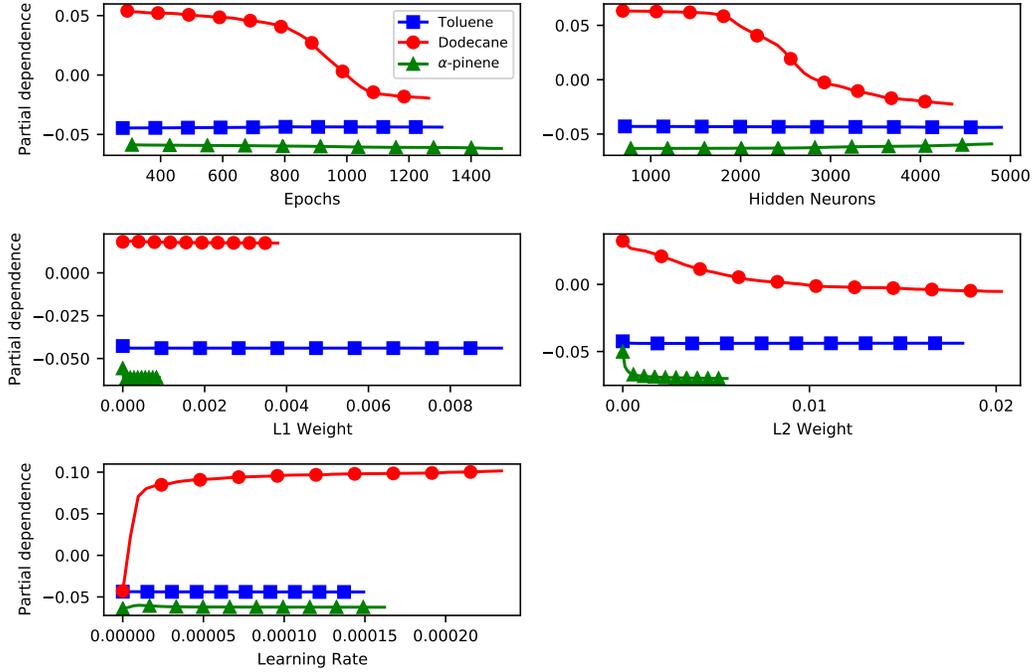


Figure B2. A partial dependence plot for each parameter varied by ECHO for the MLP model for toluene.

808 tance for MLP and GRU models tasked with predicting gas and aerosol, versus the start
 809 time of the box simulation. In these three figures, a box simulation began at the start
 810 time and continued until no more time steps were available to compare with the GECKO-
 811 A trajectories.

812 Three example experiment trajectories are shown in Figure D4 for the 3-task MLP
 813 model, for the full 5-day box simulations which all began at midnight (the same exam-
 814 ples for the GRU model are shown in Figure 9). The simulations performed at other start-
 815 ing times covered a longer 4-day window compared to the 1-day simulations shown in
 816 the main text in Figure 10. Figure D5 shows the average Pearson coefficient for these
 817 4-day simulations for both the 3- and 2-task MLP and GRU models.

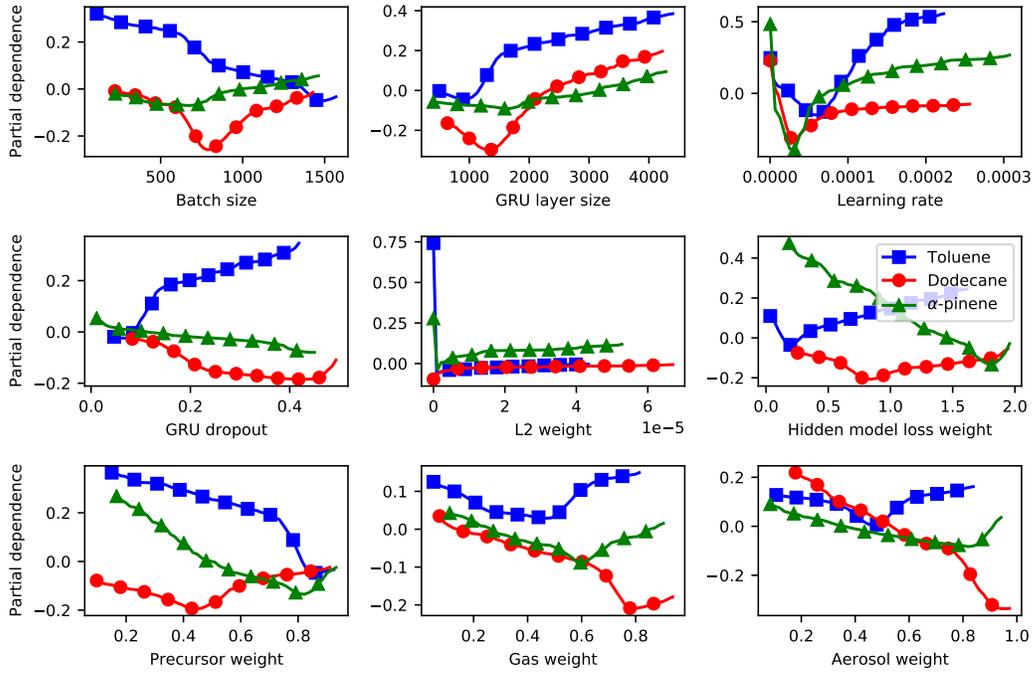


Figure B3. A partial dependence plot for each parameter varied by ECHO for the GRU model for toluene.

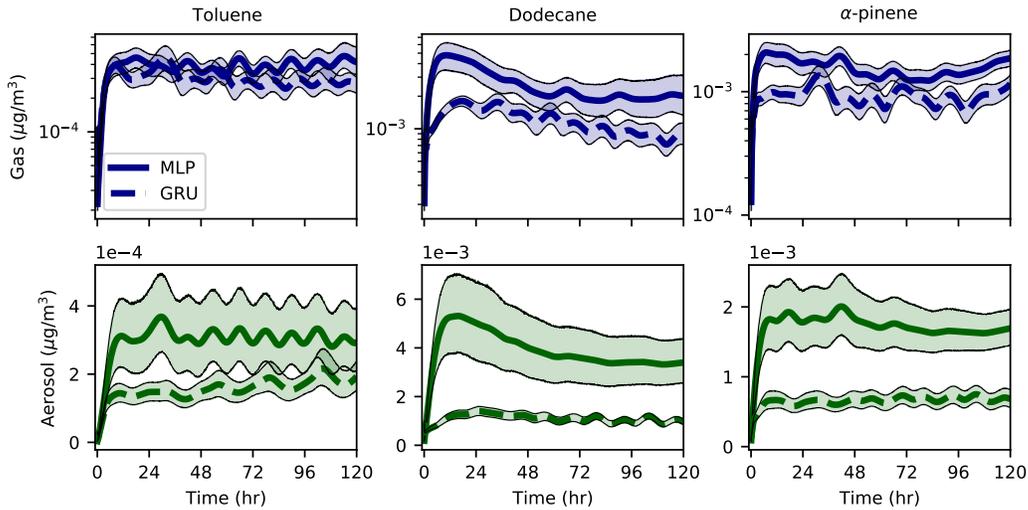


Figure C1. Computed CRPS versus simulation time for toluene, dodecane and α -pinene, for MLP and GRU models (solid and dashed lines, respectively) predicting gas and aerosol but not precursor. The precursor value at some time is still used as input to both models. The shaded regions show the 95% confidence interval.

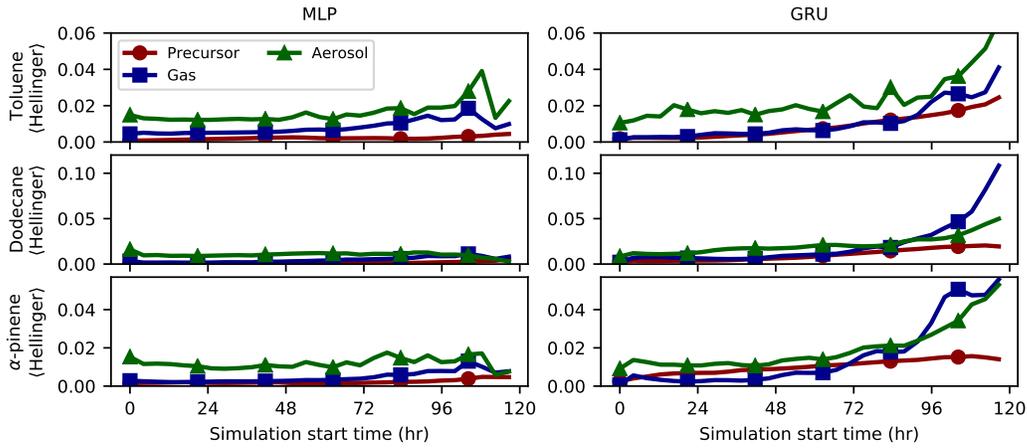


Figure D1. The average Hellinger distance versus the initial box simulation start time, computed from the 200 test experiments for the three species considered. The MLP and GRU results are shown in panels on the left and right, respectively.

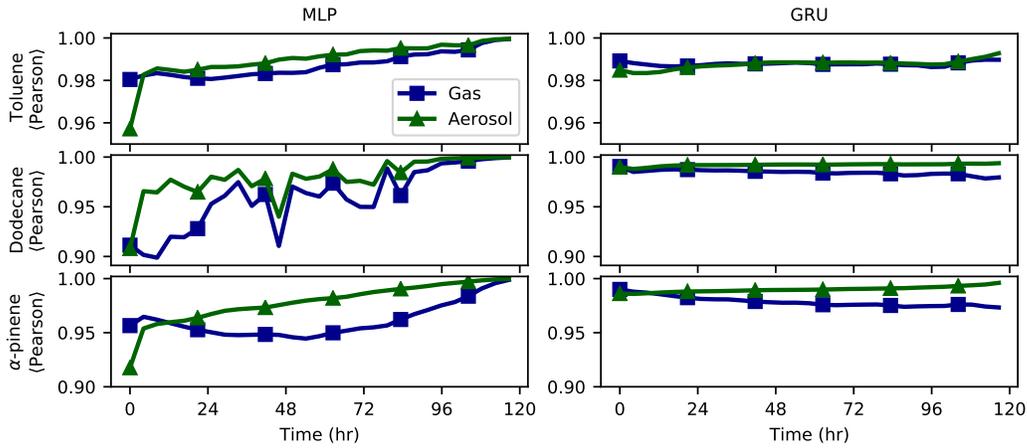


Figure D2. The average Pearson coefficient versus the initial box simulation start time, computed from the 200 test experiments for the three species considered. Both model types were tasked with gas and aerosol prediction only, but otherwise were the same as the versions which predicted precursor. The MLP and GRU results are shown in panels on the left and right, respectively.

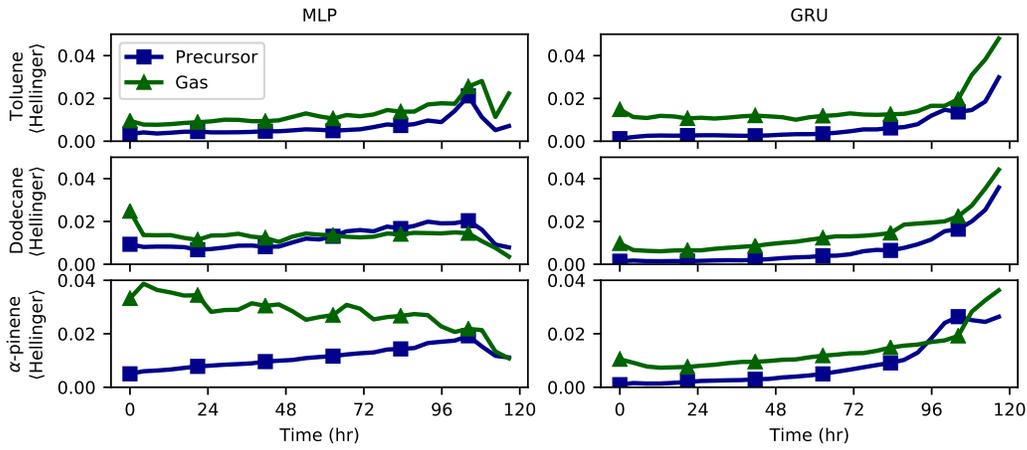


Figure D3. The average Hellinger distance versus the initial box simulation start time, computed from the 200 test experiments for the three species considered. Both model types were tasked with gas and aerosol prediction only, but otherwise were the same as the versions which predicted precursor. The MLP and GRU results are shown in panels on the left and right, respectively.

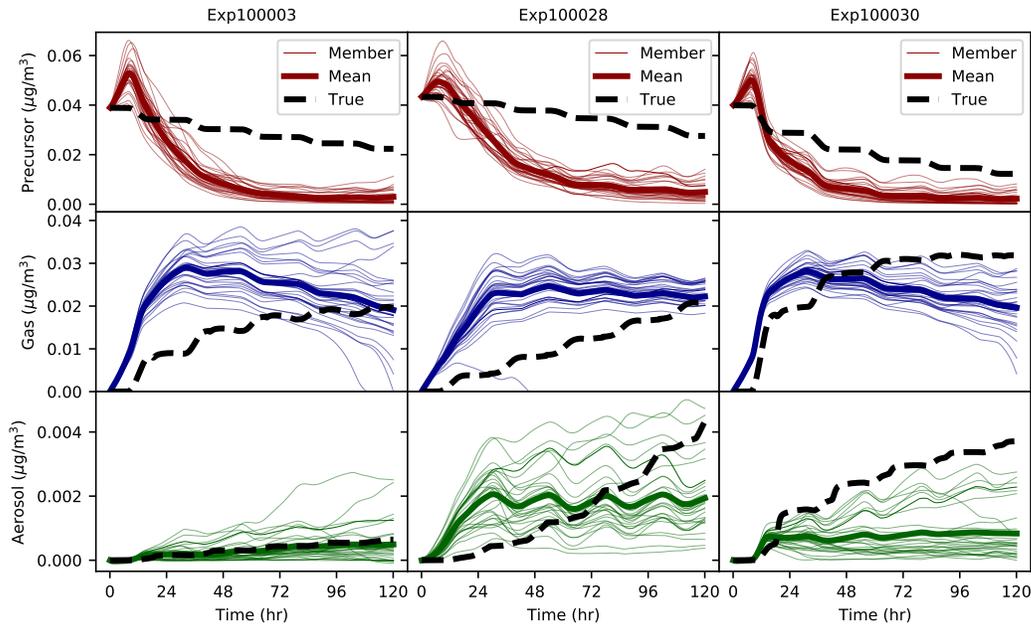


Figure D4. Examples of MLP box simulations where the environmental variables not including temperature and solar zenith angle were allowed to vary with time, compared with the GECKO-A simulations.

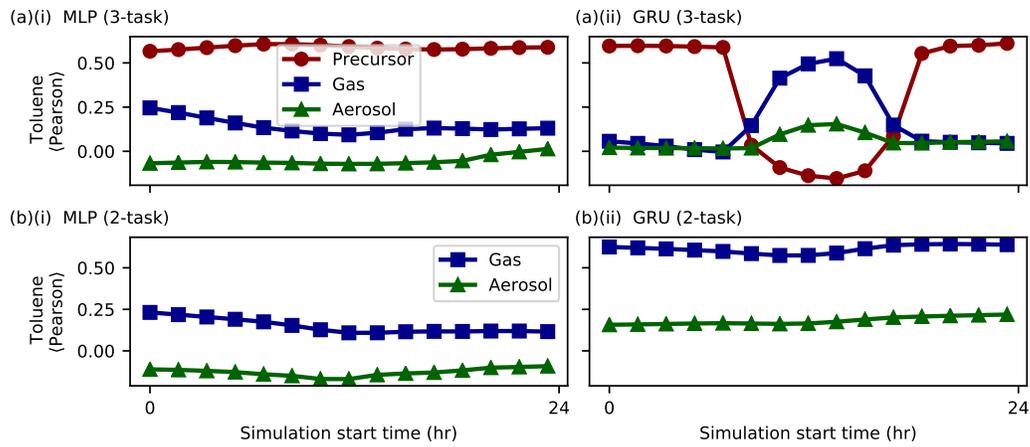


Figure D5. The Pearson coefficient versus the initial box simulation start time computed from 36 experiments for toluene. In (a), the results for the three-task MLP and GRU models are shown in (i) and (ii), respectively, while (b) shows the two-task MLP and GRU models. All box simulations ran for 96 hours.

818 **Acknowledgments**

819 This material is based upon work supported by the National Center for Atmospheric Re-
 820 search, which is a major facility sponsored by the National Science Foundation under
 821 Cooperative Agreement No. 1852977. We would like to acknowledge high-performance
 822 computing support from Cheyenne and Casper (Computational and Information Sys-
 823 tems Laboratory, CISL, 2020) provided by NCAR’s Computational and Information Sys-
 824 tems Laboratory, sponsored by the National Science Foundation. The emulators described
 825 here and simulation code used to train and test the models are archived at [https://github](https://github.com/NCAR/gecko-ml)
 826 [.com/NCAR/gecko-ml](https://github.com/NCAR/gecko-ml). All GECKO-A data sets created for this study are available at
 827 <https://doi.org/10.5281/zenodo.5790042>

828 **References**

829 Akiba, T., Sano, S., Yanase, T., Ohta, T., & Koyama, M. (2019). Optuna: A next-
 830 generation hyperparameter optimization framework. In *Proceedings of the 25th*
 831 *acm sigkdd international conference on knowledge discovery and data mining*.

832 Ardabili, S., Mosavi, A., Dehghani, M., & Várkonyi-Kóczy, A. R. (2019). Deep
 833 learning and machine learning in hydrological processes climate change and
 834 earth systems a systematic review. In *Int. res. conf. high educ.* (pp. 52–62).

835 Aumont, B., Szopa, S., & Madronich, S. (2005). Modelling the evolution of or-
 836 ganic carbon during its gas-phase tropospheric oxidation: development of an
 837 explicit model based on a self generating approach. *Atmospheric Chemistry*
 838 *and Physics*, 5(9), 2497–2517. Retrieved from [https://acp.copernicus.org/](https://acp.copernicus.org/articles/5/2497/2005/)
 839 [articles/5/2497/2005/](https://acp.copernicus.org/articles/5/2497/2005/) doi: 10.5194/acp-5-2497-2005

840 Bergstra, J., Bardenet, R., Bengio, Y., & Kégl, B. (2011). Algorithms for hyper-
 841 parameter optimization. *Advances in neural information processing systems*,
 842 24.

843 Beucler, T., Pritchard, M., Gentine, P., & Rasp, S. (2020). Towards Physically-
 844 Consistent, Data-Driven Models of Convection. *International Geoscience*
 845 *and Remote Sensing Symposium (IGARSS)(1)*, 3987–3990. doi: 10.1109/
 846 [IGARSS39084.2020.9324569](https://doi.org/10.1109/IGARSS39084.2020.9324569)

847 Boucher, O., Randall, D., Artaxo, P., Bretherton, C., Feingold, G., Forster, P., ...
 848 Zhang, X.-Y. (2013). Clouds and aerosols. In T. [Stocker et al. (Eds.), *Climate*
 849 *change 2013: The physical science basis. contribution of working group i to*

- 850 *the fifth assessment report of the intergovernmental panel on climate change.*
851 Cambridge, United Kingdom: Cambridge University Press.
- 852 Boulmaiz, T., Guermoui, M., & Boutaghane, H. (2020). Impact of training data size
853 on the lstm performances for rainfall–runoff modeling. *Modeling Earth Systems*
854 *and Environment*, 6, 2153–2164.
- 855 Brenowitz, N. D., & Bretherton, C. S. (2018). Prognostic Validation of a Neu-
856 ral Network Unified Physics Parameterization. *Geophysical Research Letters*,
857 45(12), 6289–6298. doi: 10.1029/2018GL078510
- 858 Camredon, M., Aumont, B., Lee-Taylor, J., & Madronich, S. (2007). The
859 soa/voc/no_x system: an explicit model of secondary organic aerosol for-
860 mation. *Atmospheric Chemistry and Physics*, 7(21), 5599–5610. Re-
861 trieved from <https://acp.copernicus.org/articles/7/5599/2007/> doi:
862 10.5194/acp-7-5599-2007
- 863 Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014a). Empirical evaluation
864 of gated recurrent neural networks on sequence modeling. *arXiv preprint*
865 *arXiv:1412.3555*.
- 866 Chung, J., Gulcehre, C., Cho, K., & Bengio, Y. (2014b). Empirical evaluation of
867 gated recurrent neural networks on sequence modeling. arxiv 2014. *arXiv*
868 *preprint arXiv:1412.3555*.
- 869 Computational and Information Systems Laboratory, CISL. (2020). *Cheyenne:*
870 *HPE/SGI ICE XA System (NCAR Community Computing)* (Tech. Rep.). Na-
871 tional Center for Atmospheric Research. Retrieved from [https://doi.org/10](https://doi.org/10.5065/D6RX99HX)
872 [.5065/D6RX99HX](https://doi.org/10.5065/D6RX99HX)
- 873 de Gouw, J. A. (2005). Budget of organic carbon in a polluted atmosphere: Results
874 from the new england air quality study in 2002. *J. Geophys. Res.*, 110(D16).
875 doi: 10.1029/2004JD005623
- 876 Friedman, J. H. (2001). Greedy function approximation: a gradient boosting ma-
877 chine. *Ann. Stat.*, 1189–1232.
- 878 Gettelman, A., Gagne, D. J., Chen, C. C., Christensen, M. W., Lebo, Z. J.,
879 Morrison, H., & Gantos, G. (2021). Machine Learning the Warm Rain
880 Process. *Journal of Advances in Modeling Earth Systems*, 13(2). doi:
881 10.1029/2020MS002268
- 882 Hochreiter, S. (1991). Untersuchungen zu dynamischen neuronalen netzen. *Diploma*,

- 883 *Technische Universität München*, 91(1).
- 884 Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. *Neural Compu-*
885 *tation*, 9(8), 1735–1780.
- 886 Hodzic, A., Aumont, B., Knote, C., Madronich, S., & Tyndall, G. (2014). Volatility
887 dependence of Henry’s law constants of condensable organics: Application to
888 estimate depositional loss of secondary organic aerosols. *Geophysical Research*
889 *Letters*, 4795–4804. doi: 10.1002/2014GL060649. Received
- 890 Hodzic, A., Campuzano-Jost, P., Bian, H., Chin, M., Colarco, P. R., Day, D. A.,
891 ... Jimenez, J. L. (2020). Characterization of organic aerosol across
892 the global remote troposphere: a comparison of ATom measurements and
893 global chemistry models. *Atmos. Chem. Phys.*, 20(8), 4607–4635. doi:
894 10.5194/acp-20-4607-2020
- 895 Hodzic, A., & Jimenez, J. L. (2011). Modeling anthropogenically controlled sec-
896 ondary organic aerosols in a megacity: a simplified framework for global
897 and climate models. *Geoscientific Model Development*, 4(4), 901–917. Re-
898 trieved from <https://gmd.copernicus.org/articles/4/901/2011/> doi:
899 10.5194/gmd-4-901-2011
- 900 Hodzic, A., Kasibhatla, P. S., Jo, D. S., Cappa, C. D., Jimenez, J. L., Madronich,
901 S., & Park, R. J. (2016). Rethinking the global secondary organic aerosol (soa)
902 budget: stronger production, faster removal, shorter lifetime. *Atmospheric*
903 *Chemistry and Physics*(16), 7917–7941. doi: 10.5194/acp-16-7917-2016
- 904 Hodzic, A., Madronich, S., Kasibhatla, P. S., Tyndall, G., Aumont, B., Jimenez,
905 J. L., ... Orlando, J. (2015). Organic photolysis reactions in tropo-
906 spheric aerosols: effect on secondary organic aerosol formation and life-
907 time. *Atmospheric Chemistry and Physics*, 15(16), 9253–9269. Retrieved
908 from <https://acp.copernicus.org/articles/15/9253/2015/> doi:
909 10.5194/acp-15-9253-2015
- 910 Hutter, F., Hoos, H., & Leyton-Brown, K. (2014). An efficient approach for assess-
911 ing hyperparameter importance. In *International conference on machine learn-*
912 *ing* (pp. 754–762).
- 913 Irrgang, C., Saynisch-Wagner, J., & Thomas, M. (2020). Machine learning-based
914 prediction of spatiotemporal uncertainties in global wind velocity reanalyses.
915 *Journal of Advances in Modeling Earth Systems*, 12(5), e2019MS001876.

- 916 Jenkin, M. E., Saunders, S. M., Wagner, V., & Pilling, M. J. (2003, February).
 917 Protocol for the development of the Master Chemical Mechanism, MCM v3
 918 (Part B): Tropospheric degradation of aromatic volatile organic compounds.
 919 *Atmospheric Chem. Phys.*, *3*(1), 181–193. doi: 10.5194/acp-3-181-2003
- 920 Keller, C. A., & Evans, M. J. (2019). Application of random forest regression
 921 to the calculation of gas-phase chemistry within the geos-chem chemistry
 922 model v10. *Geoscientific Model Development*, *12*(3), 1209–1225. Re-
 923 trieved from <https://gmd.copernicus.org/articles/12/1209/2019/> doi:
 924 10.5194/gmd-12-1209-2019
- 925 Kelp, M. M., Jacob, D. J., Kutz, J. N., Marshall, J. D., & Tessum, C. W. (2020).
 926 Toward stable, general machine-learned models of the atmospheric chem-
 927 ical system. *Journal of Geophysical Research: Atmospheres*, *125*(23),
 928 e2020JD032759. Retrieved from [https://agupubs.onlinelibrary.wiley](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020JD032759)
 929 [.com/doi/abs/10.1029/2020JD032759](https://agupubs.onlinelibrary.wiley.com/doi/abs/10.1029/2020JD032759) (e2020JD032759 2020JD032759) doi:
 930 <https://doi.org/10.1029/2020JD032759>
- 931 Kelp, M. M., Jacob, D. J., Lin, H., & Sulprizio, M. P. (2021). An online-learned
 932 neural network chemical solver for stable long-term global simulations of at-
 933 mospheric chemistry. *Journal of Advances in Modeling Earth Systems*, Under
 934 Review. Retrieved from <https://eartharxiv.org/repository/view/2886/>
 935 doi: <https://doi.org/10.31223/X52K7J>
- 936 Kelp, M. M., Tessum, C. W., & Marshall, J. D. (2018). Orders-of-magnitude
 937 speedup in atmospheric chemistry modeling through neural network-based
 938 emulation. *arXiv*(206), 1–23.
- 939 Kouw, W. M., & Loog, M. (2018). An introduction to domain adaptation and trans-
 940 fer learning. *arXiv preprint arXiv:1812.11806*.
- 941 Kratzert, F., Klotz, D., Brenner, C., Schulz, K., & Herrnegger, M. (2018, Novem-
 942 ber). Rainfall–runoff modelling using long Short-Term memory (LSTM)
 943 networks. *Hydrol. Earth Syst. Sci.*, *22*(11), 6005–6022.
- 944 La, Y. S., Camredon, M., Ziemann, P. J., Valorso, R., Matsunaga, A., Lannuque,
 945 V., . . . Aumont, B. (2016). Impact of chamber wall loss of gaseous organic
 946 compounds on secondary organic aerosol formation: explicit modeling of soa
 947 formation from alkane and alkene oxidation. *Atmospheric Chemistry and*
 948 *Physics*, *16*(3), 1417–1431. Retrieved from <https://acp.copernicus.org/>

- 949 [articles/16/1417/2016/](https://doi.org/10.5194/acp-16-1417-2016) doi: 10.5194/acp-16-1417-2016
- 950 Lannuque, V., Camredon, M., Couvidat, F., Hodzic, A., Valorso, R., Madronich,
 951 S., ... Aumont, B. (2018). Exploration of the influence of environmen-
 952 tal conditions on secondary organic aerosol formation and organic species
 953 properties using explicit simulations: development of the vbs-gecko parame-
 954 terization. *Atmospheric Chemistry and Physics*, 18(18), 13411–13428. Re-
 955 trieved from <https://acp.copernicus.org/articles/18/13411/2018/> doi:
 956 10.5194/acp-18-13411-2018
- 957 Lee-Taylor, J., Hodzic, A., Madronich, S., Aumont, B., Camredon, M., & Valorso,
 958 R. (2015). Multiday production of condensing organic aerosol mass in urban
 959 and forest outflow. *Atmospheric Chemistry and Physics*, 15(2), 595–615. Re-
 960 trieved from <https://acp.copernicus.org/articles/15/595/2015/> doi:
 961 10.5194/acp-15-595-2015
- 962 Lee-Taylor, J., Madronich, S., Aumont, B., Baker, A., Camredon, M., Hodzic, A.,
 963 ... Zaveri, R. A. (2011). Explicit modeling of organic chemistry and sec-
 964 ondary organic aerosol partitioning for mexico city and its outflow plume.
 965 *Atmospheric Chemistry and Physics*, 11(24), 13219–13241. Retrieved
 966 from <https://acp.copernicus.org/articles/11/13219/2011/> doi:
 967 10.5194/acp-11-13219-2011
- 968 Li, J., Cleveland, M., Ziemba, L. D., Griffin, R. J., Barsanti, K. C., Pankow, J. F.,
 969 & Ying, Q. (2015). Modeling regional secondary organic aerosol using the
 970 Master Chemical Mechanism. *Atmos. Environ.*, 102(February), 52–61. doi:
 971 10.1016/j.atmosenv.2014.11.054
- 972 Louppe, G., Wehenkel, L., Sutera, A., & Geurts, P. (2013). Understanding vari-
 973 able importances in forests of randomized trees. *Advances in neural informa-*
 974 *tion processing systems*, 26, 431–439.
- 975 Mauderly, J. L., & Chow, J. C. (2008, February). Health effects of organic aerosols.
 976 *Inhal. Toxicol.*, 20(3), 257–288. doi: 10.1080/08958370701866008
- 977 Mouchel-Vallon, C., Lee-Taylor, J., Hodzic, A., Artaxo, P., Aumont, B., Cam-
 978 redon, M., ... Madronich, S. (2020). Exploration of oxidative chem-
 979 istry and secondary organic aerosol formation in the amazon during the
 980 wet season: explicit modeling of the manaus urban plume with gecko-
 981 a. *Atmospheric Chemistry and Physics*, 20(10), 5995–6014. Retrieved

- 982 from <https://acp.copernicus.org/articles/20/5995/2020/> doi:
983 10.5194/acp-20-5995-2020
- 984 Mousavi, S. M., & Beroza, G. C. (2020). A machine-learning approach for
985 earthquake magnitude estimation. *Geophysical Research Letters*, *47*(1),
986 e2019GL085976.
- 987 Ng, N. L., Kroll, J. H., Chan, A. W. H., Chhabra, P. S., Flagan, R. C., & Seinfeld,
988 J. H. (2007). Secondary organic aerosol formation from m-xylene, toluene,
989 and benzene. *Atmospheric Chemistry and Physics*, *7*(14), 3909–3922. Re-
990 trieved from <https://acp.copernicus.org/articles/7/3909/2007/> doi:
991 10.5194/acp-7-3909-2007
- 992 Reddy, D. S., & Prasad, P. R. C. (2018). Prediction of vegetation dynamics us-
993 ing ndvi time series data and lstm. *Modeling Earth Systems and Environment*,
994 *4*(1), 409–419.
- 995 Rumelhart, D. E., Hinton, G. E., & Williams, R. J. (1986). Learning representations
996 by back-propagating errors. *Nature*, *323*(6088), 533–536.
- 997 Schreck, J. S., & Gagne, D. J. (2021). *Earth computing hyperparameter optimization*.
998 GitHub. Retrieved from <https://github.com/NCAR/echo-opt>
- 999 Tsigaridis, K., Daskalakis, N., Kanakidou, M., Adams, P. J., Artaxo, P., Bahadur,
1000 R., ... Zhang, X. (2014). The aerocom evaluation and intercomparison of
1001 organic aerosol in global models. *Atmospheric Chemistry and Physics*, *14*(19),
1002 10845–10895. Retrieved from [https://acp.copernicus.org/articles/14/](https://acp.copernicus.org/articles/14/10845/2014/)
1003 [10845/2014/](https://acp.copernicus.org/articles/14/10845/2014/) doi: 10.5194/acp-14-10845-2014

Figure 1.

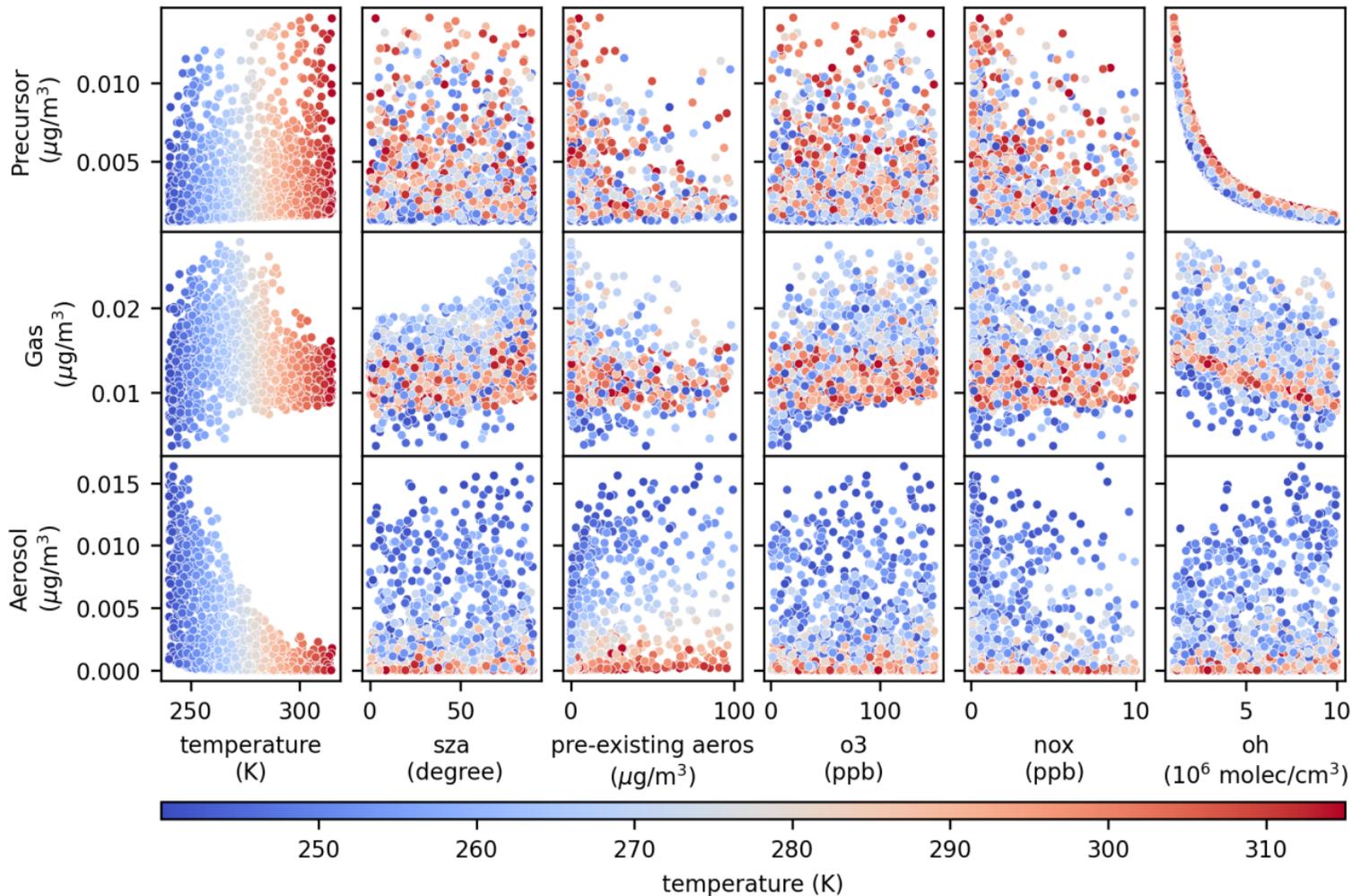
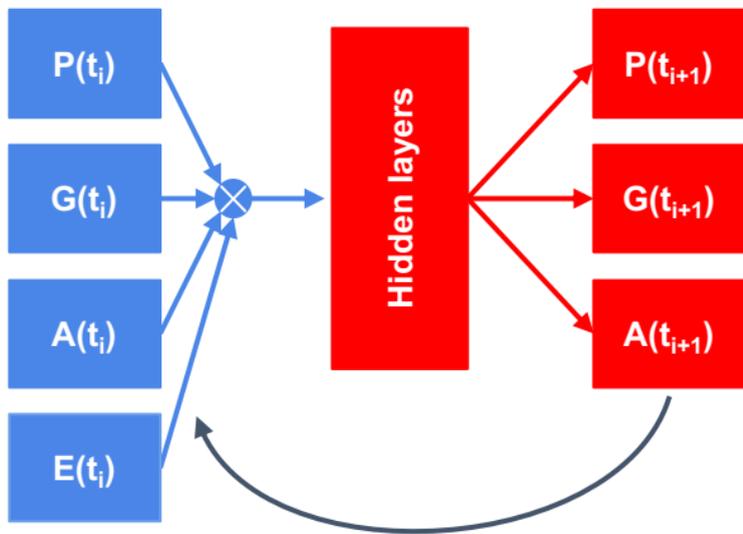
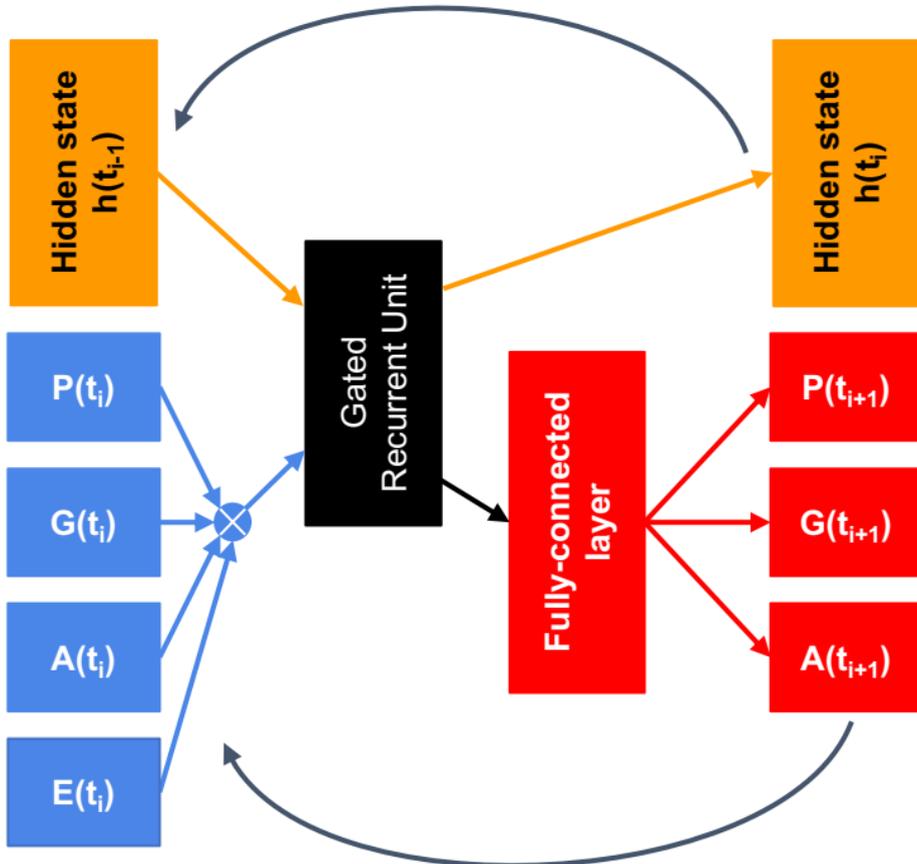


Figure 2.

(a) MLP model



(b) GRU model



(c) Hidden-state model

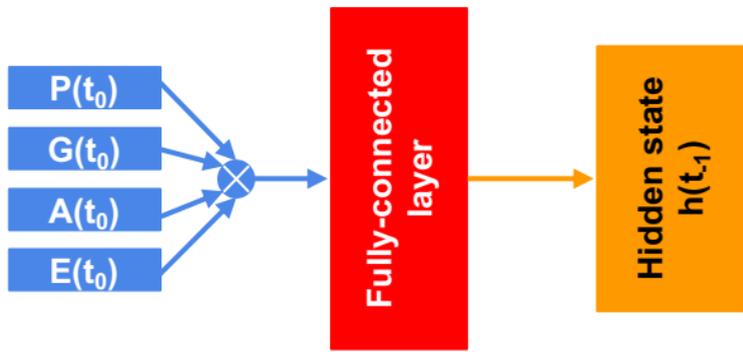
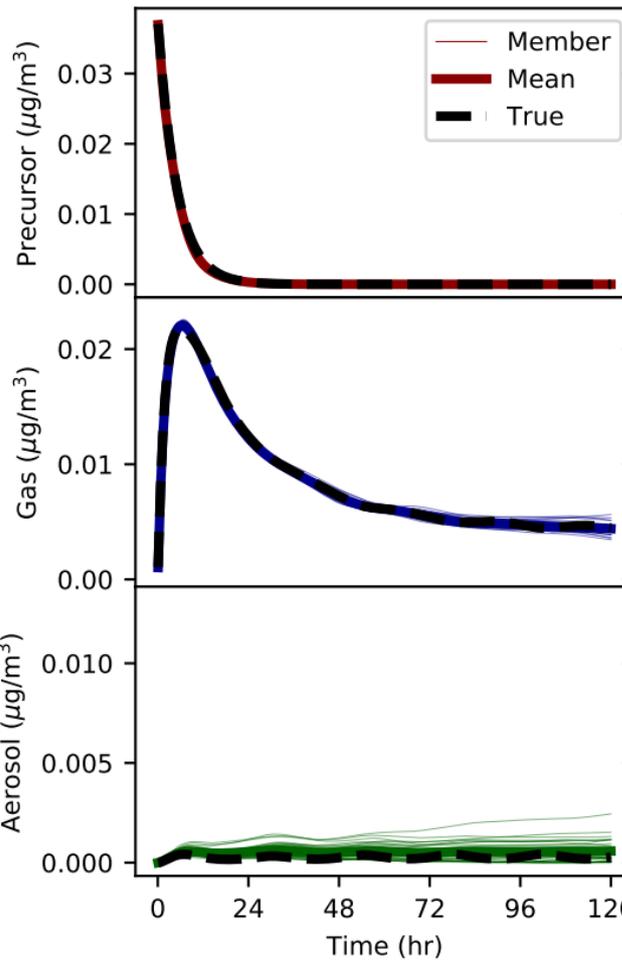
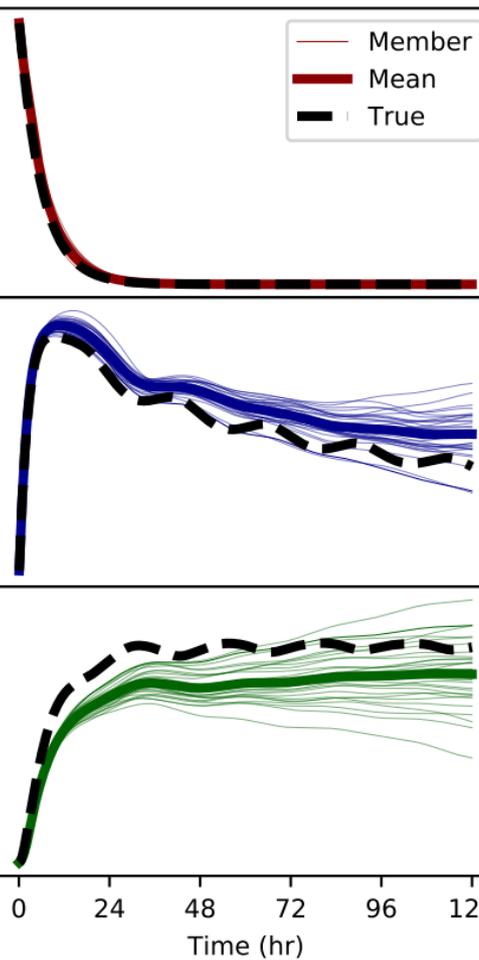


Figure 3.

Exp1806



Exp1812



Exp1896

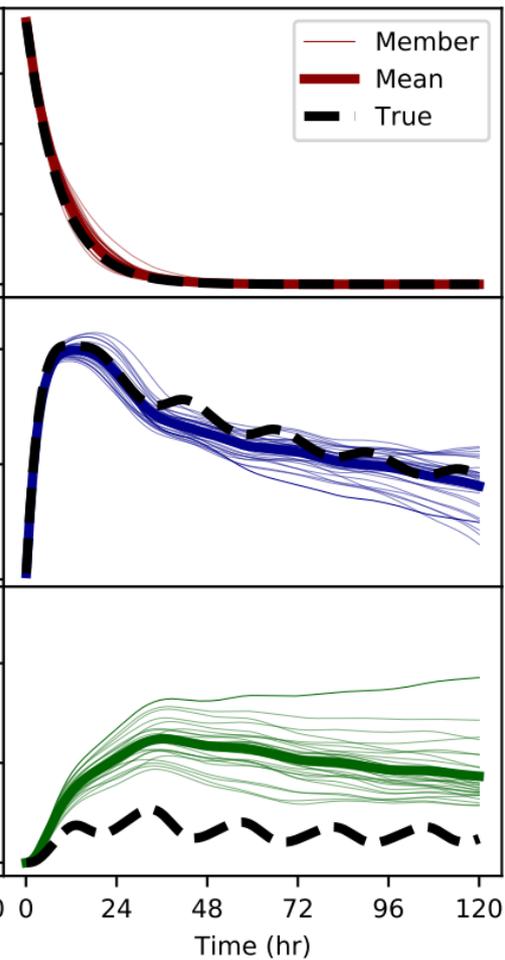
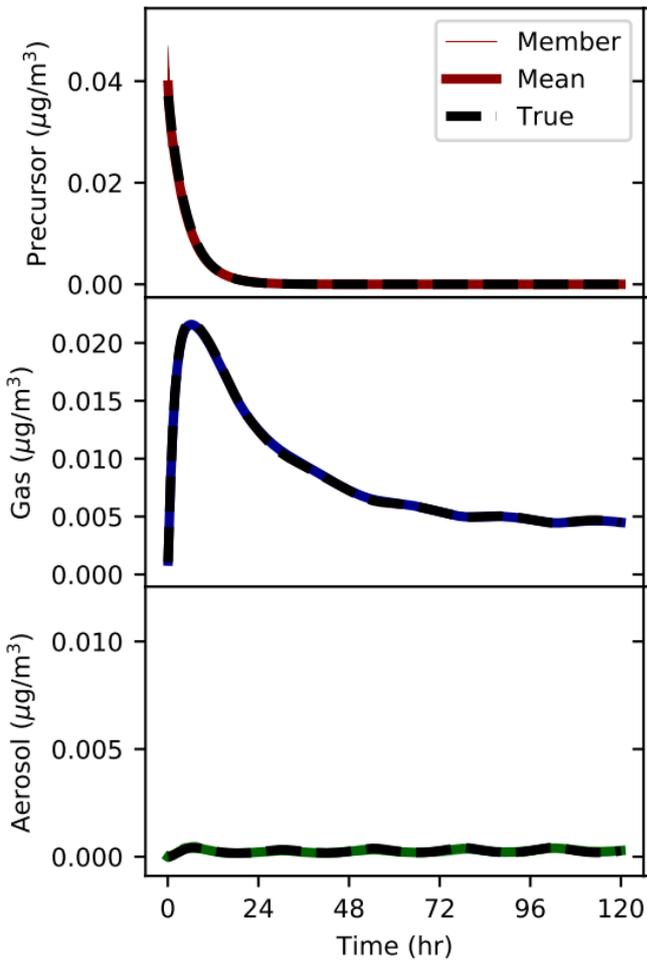
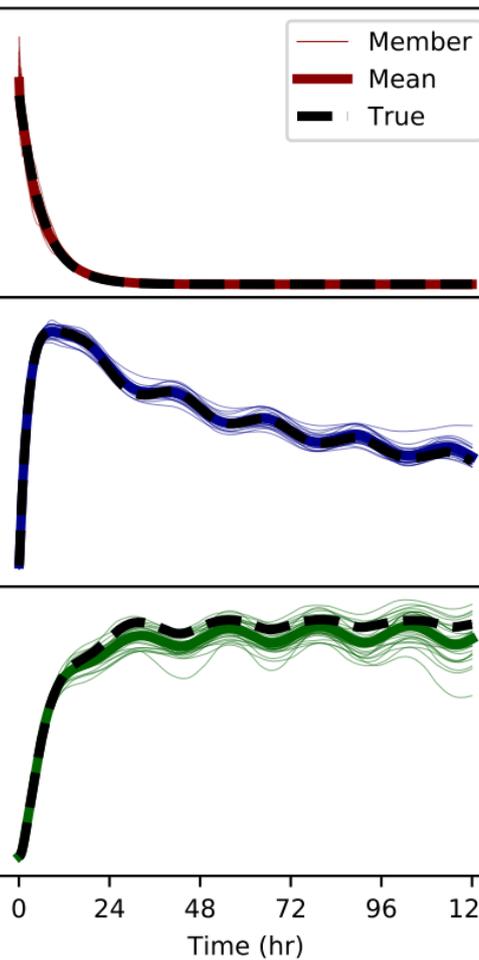


Figure 4.

Exp1806



Exp1812



Exp1896

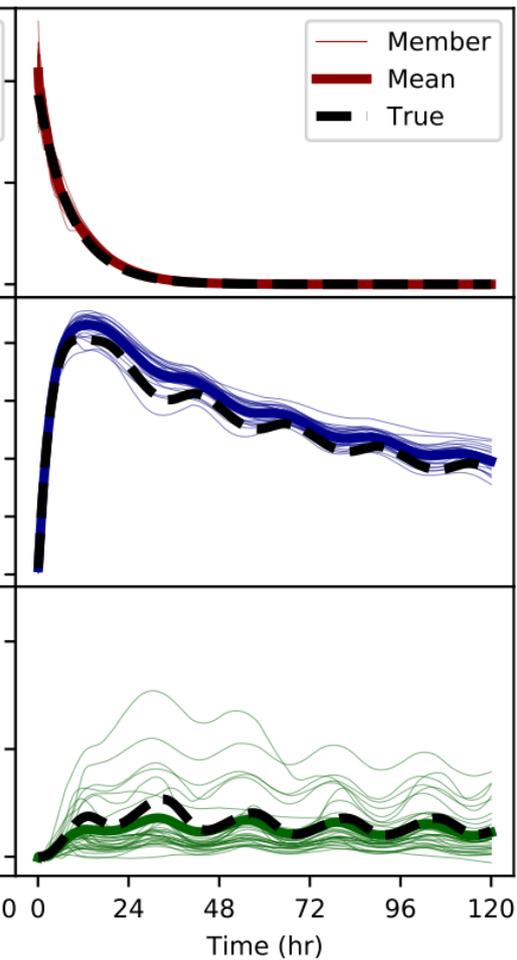
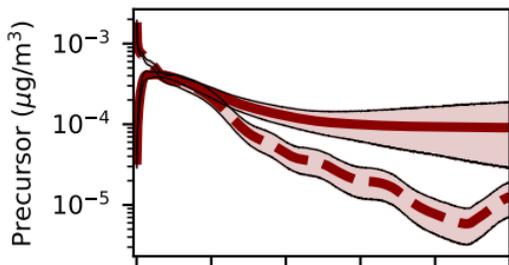


Figure 5.

Toluene



Dodecane

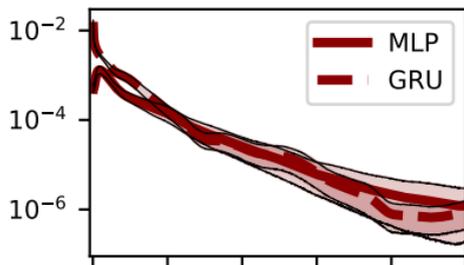
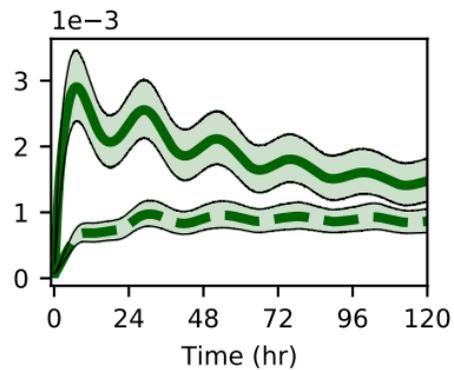
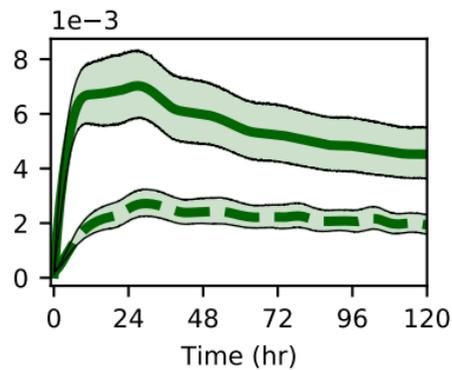
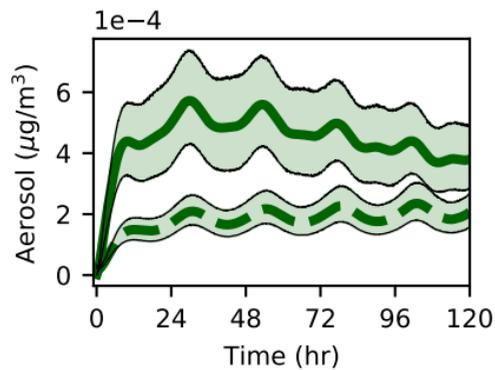
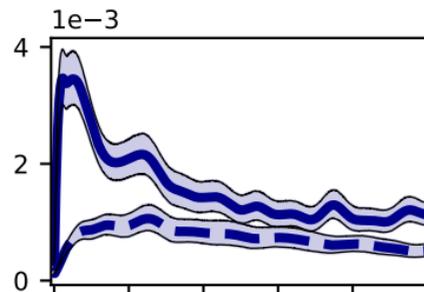
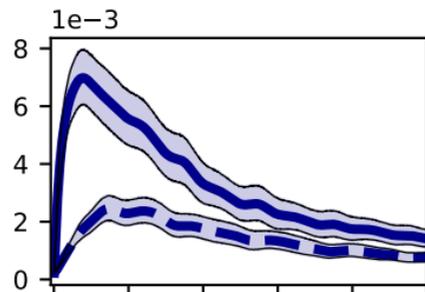
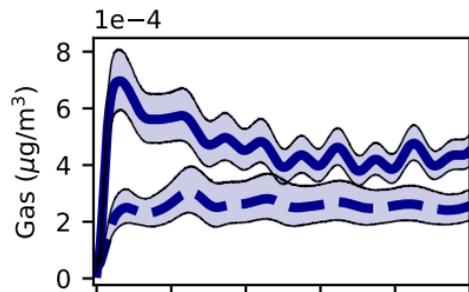
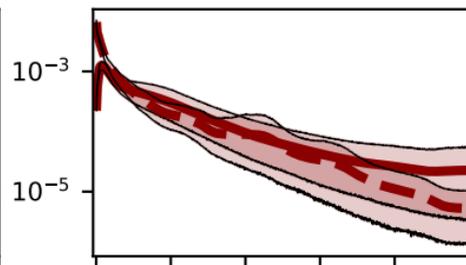
 α -pinene

Figure 6.

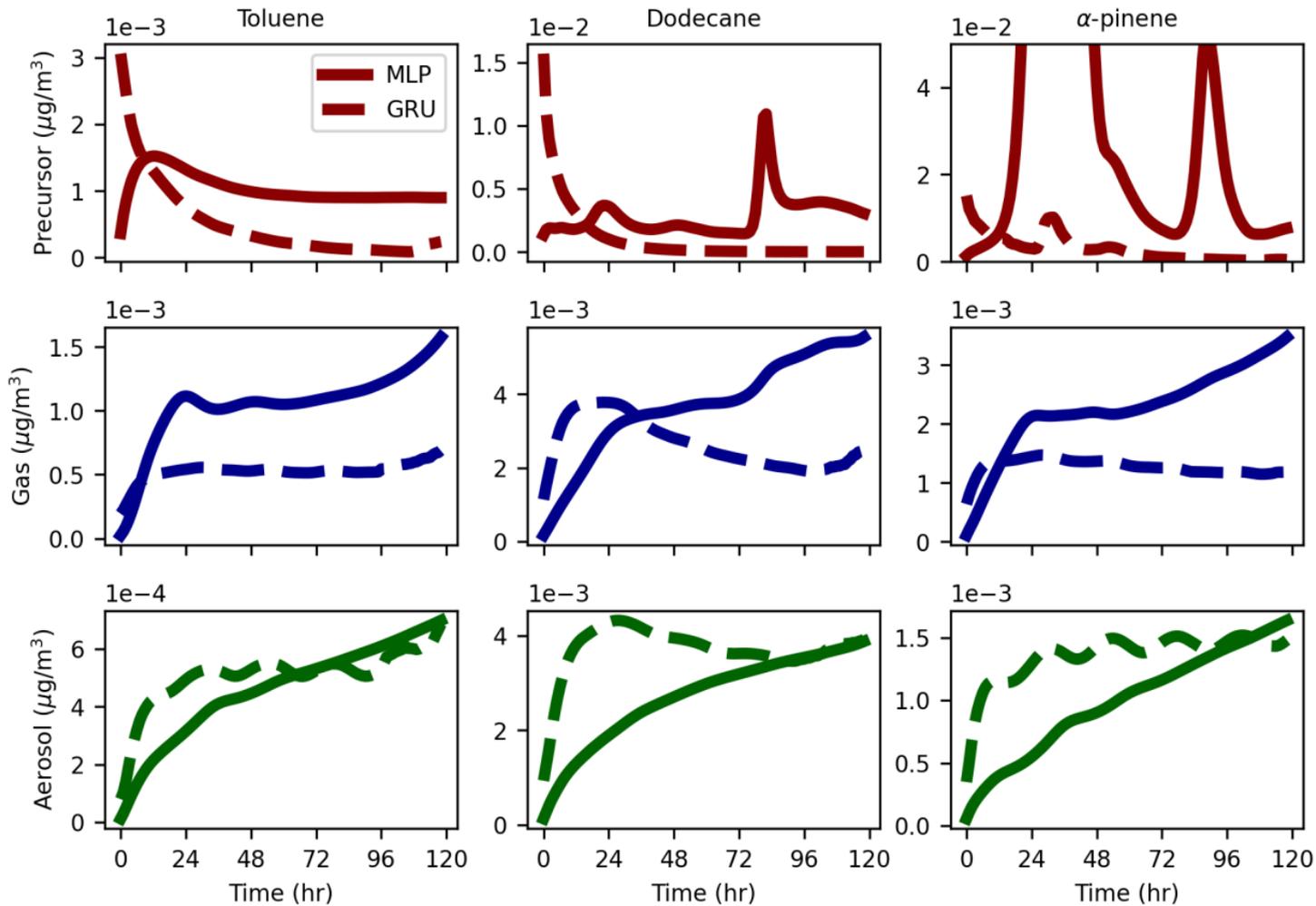


Figure 7.

MLP

GRU

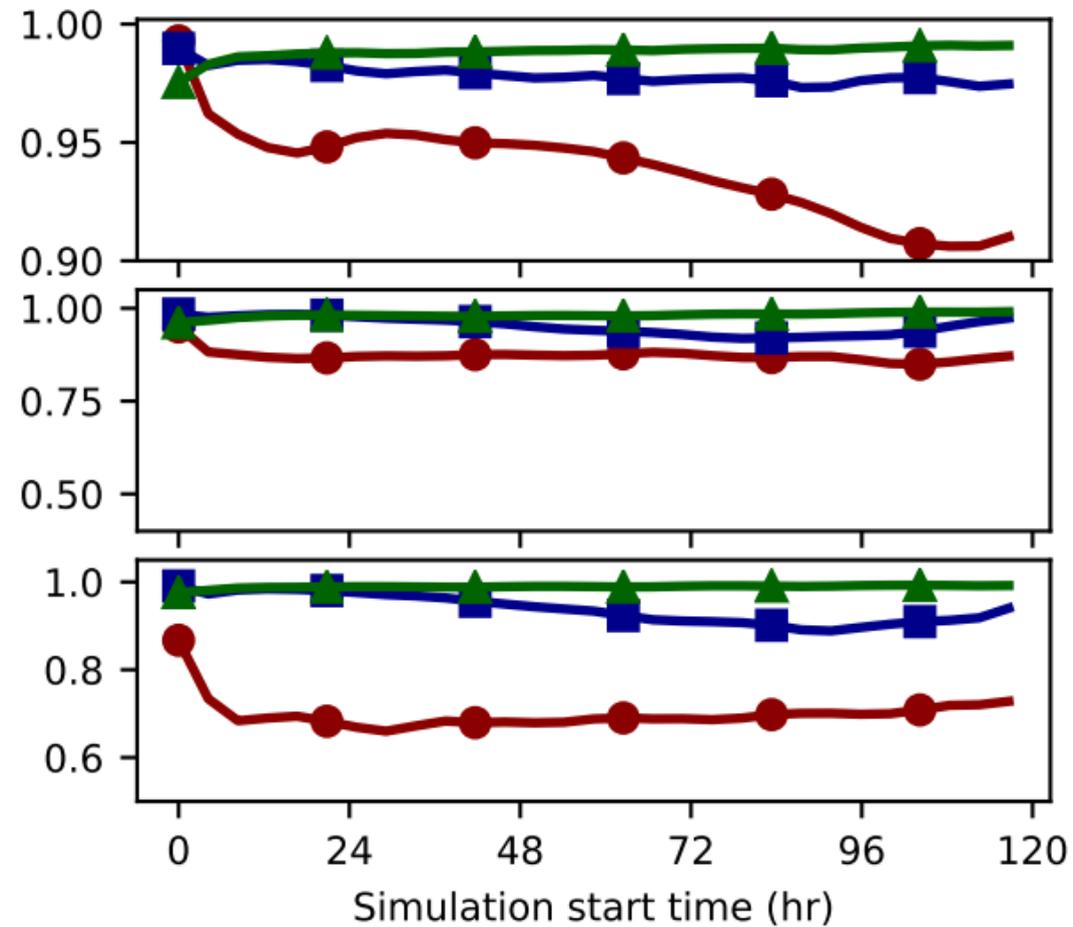
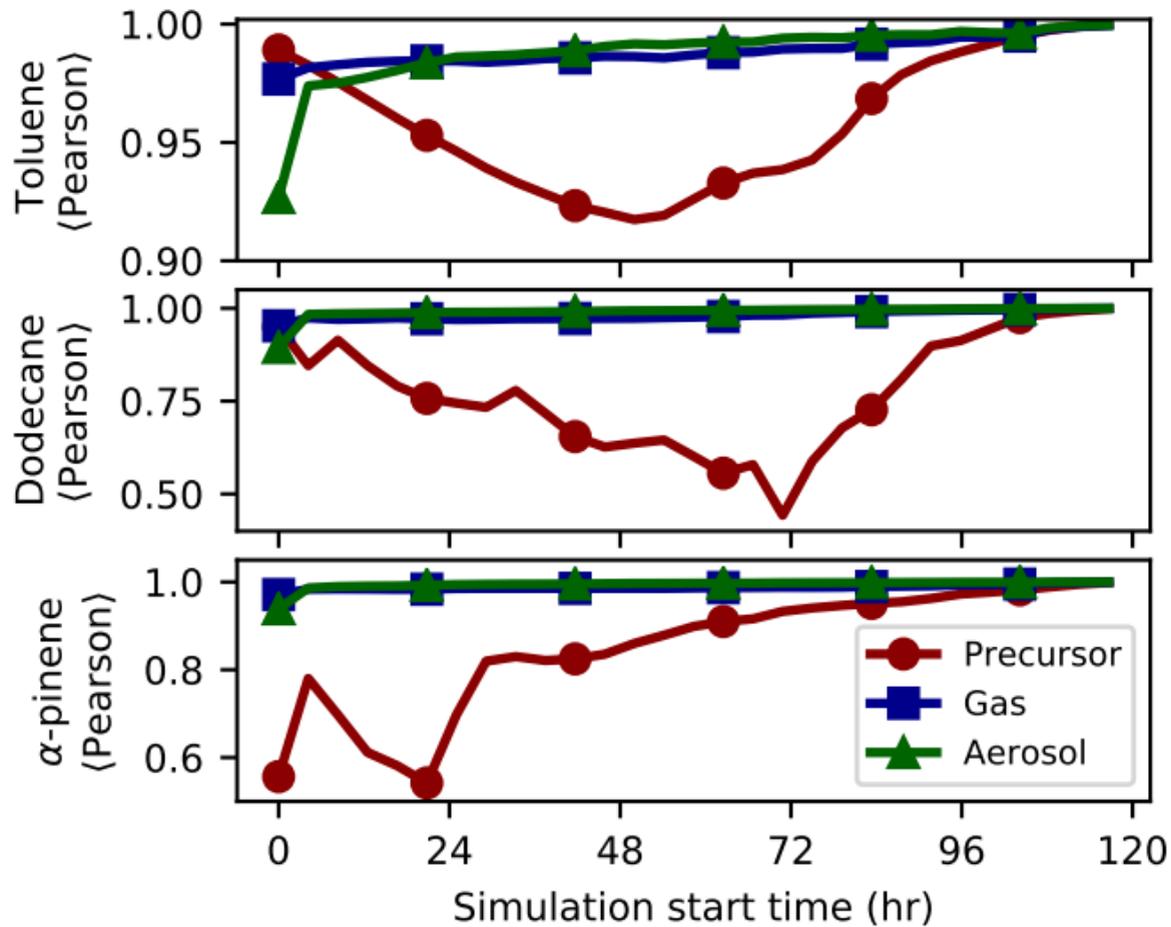
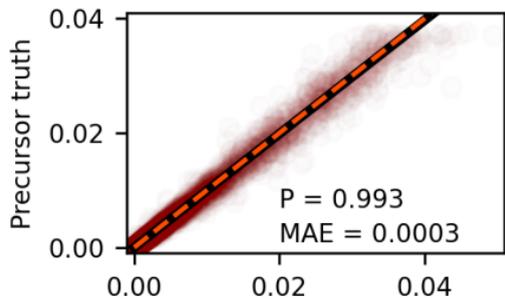
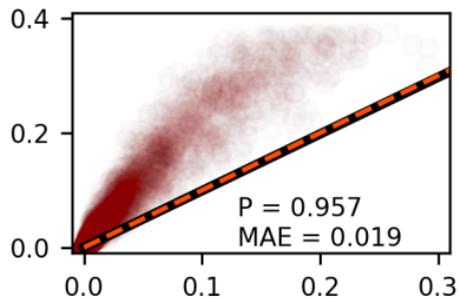


Figure 8.

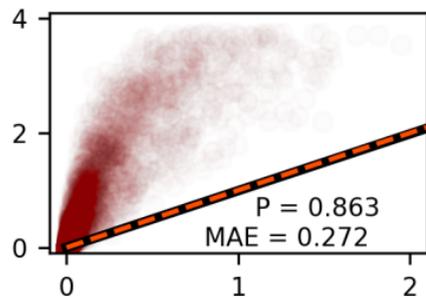
Toluene X1



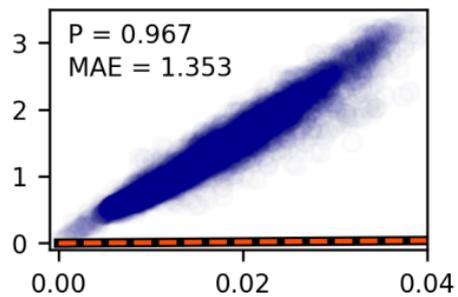
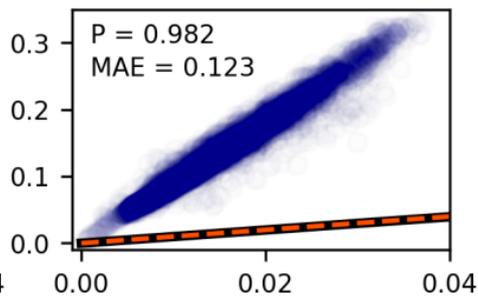
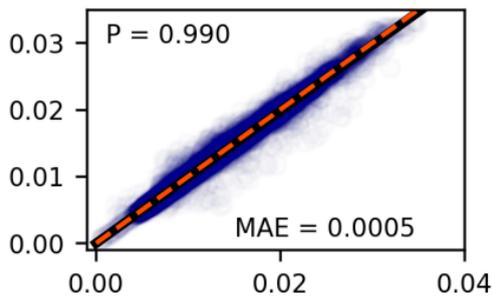
Toluene X10



Toluene X100



Gas truth



Aerosol truth

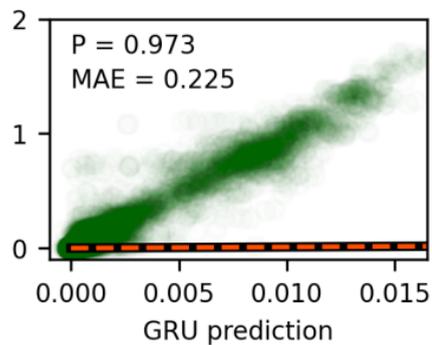
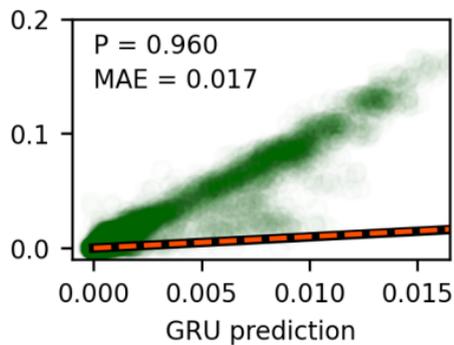
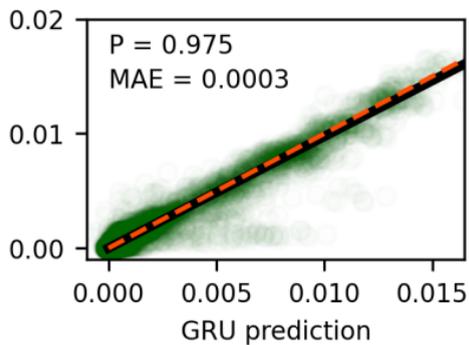


Figure 9.

Exp100003

Exp100028

Exp100030

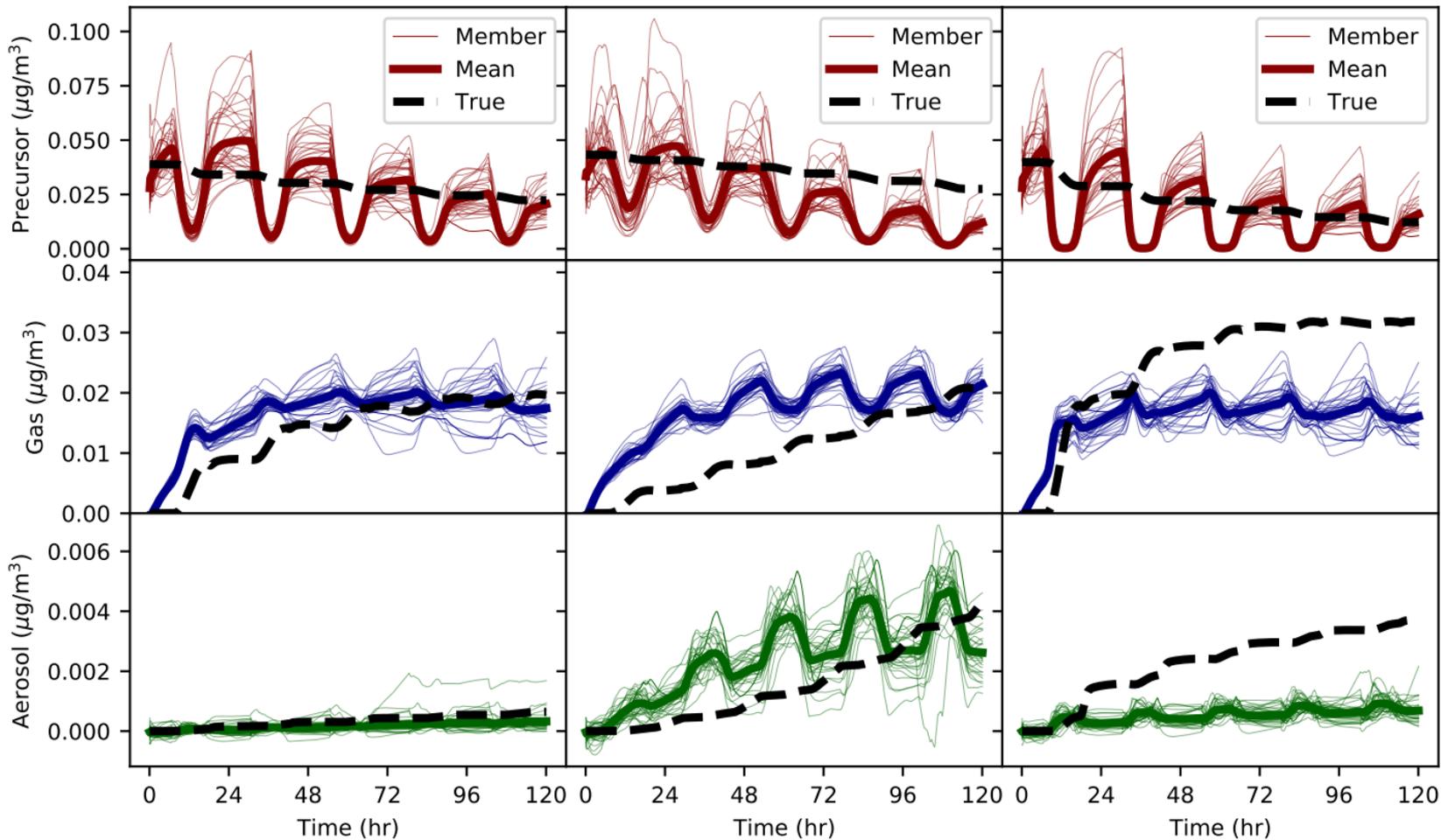
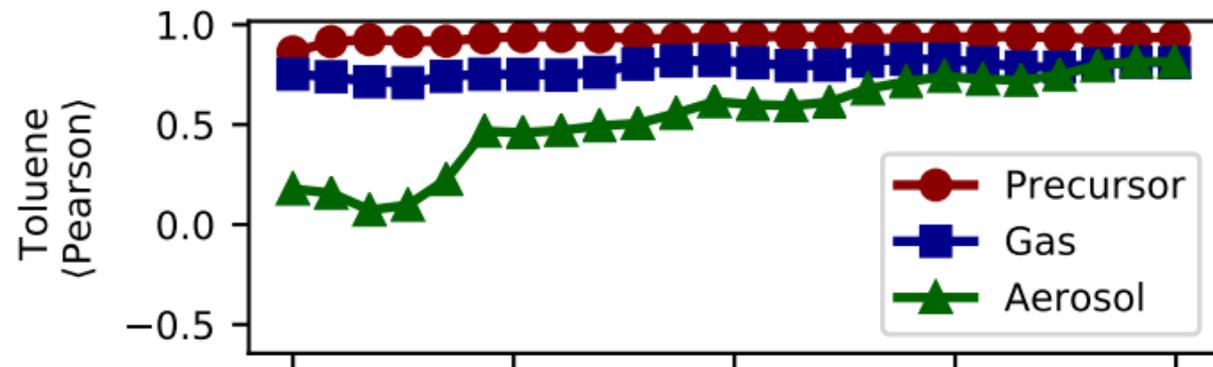
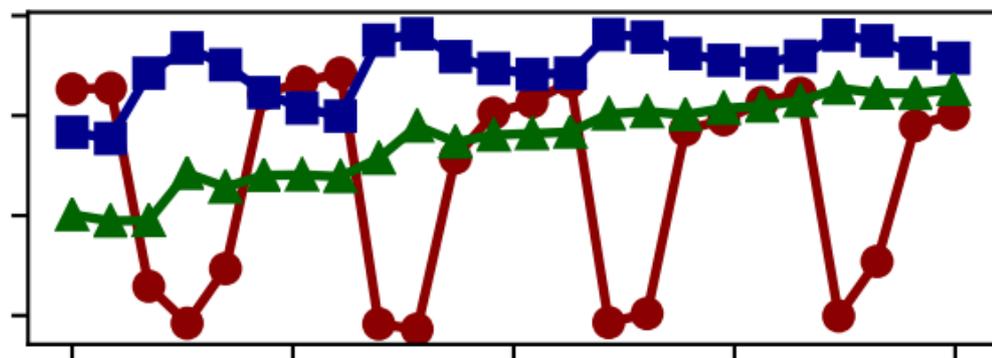


Figure 10.

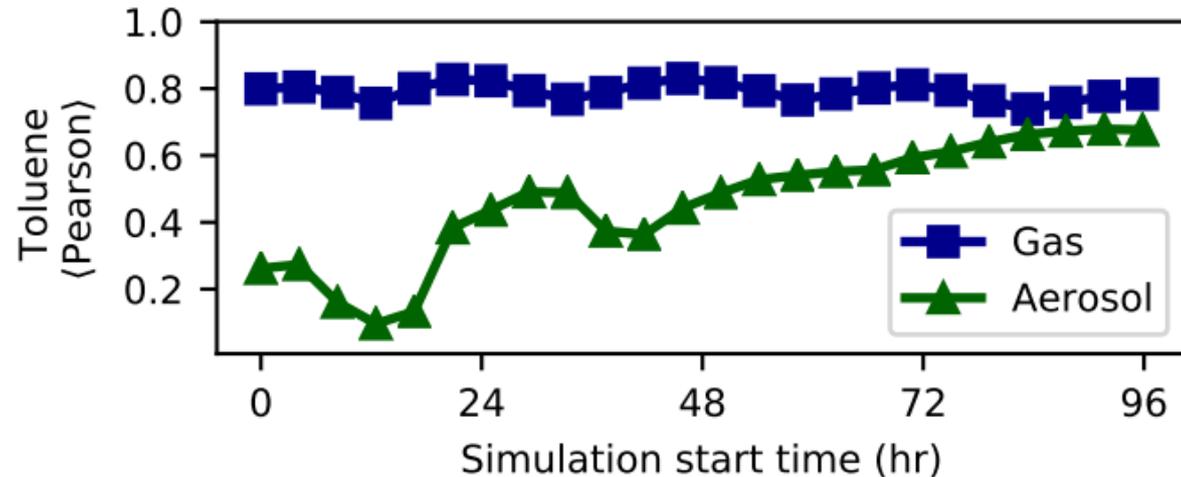
(a)(i) MLP (3-task)



(a)(ii) GRU (3-task)



(b)(i) MLP (2-task)



(b)(ii) GRU (2-task)

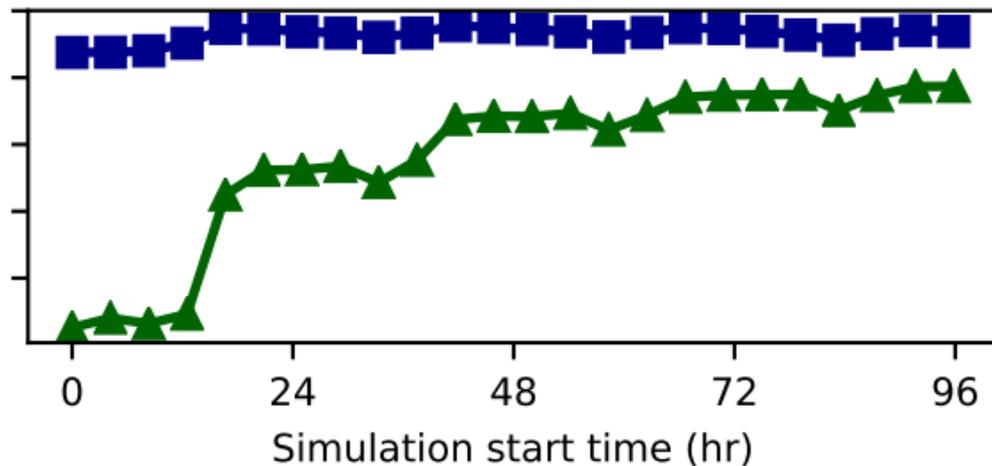


Figure 11.

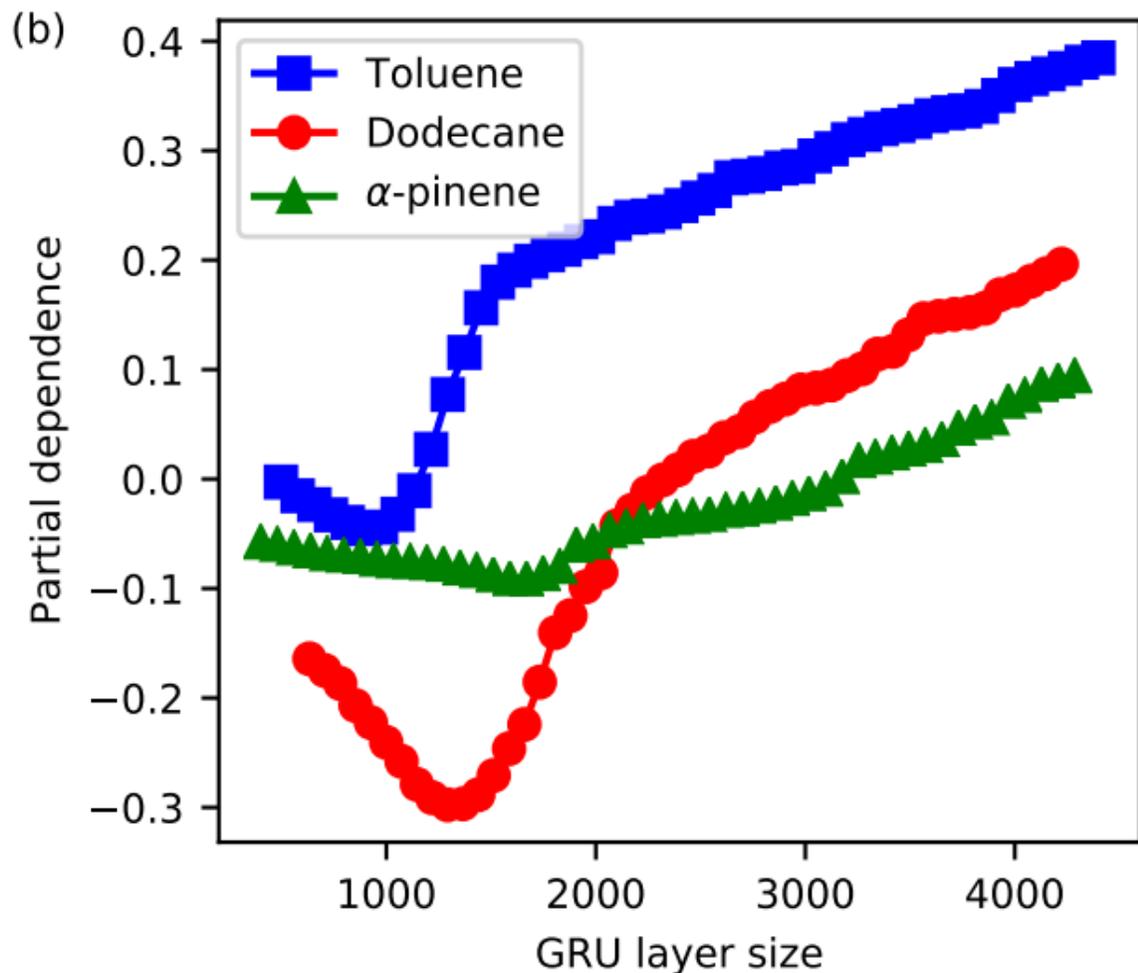
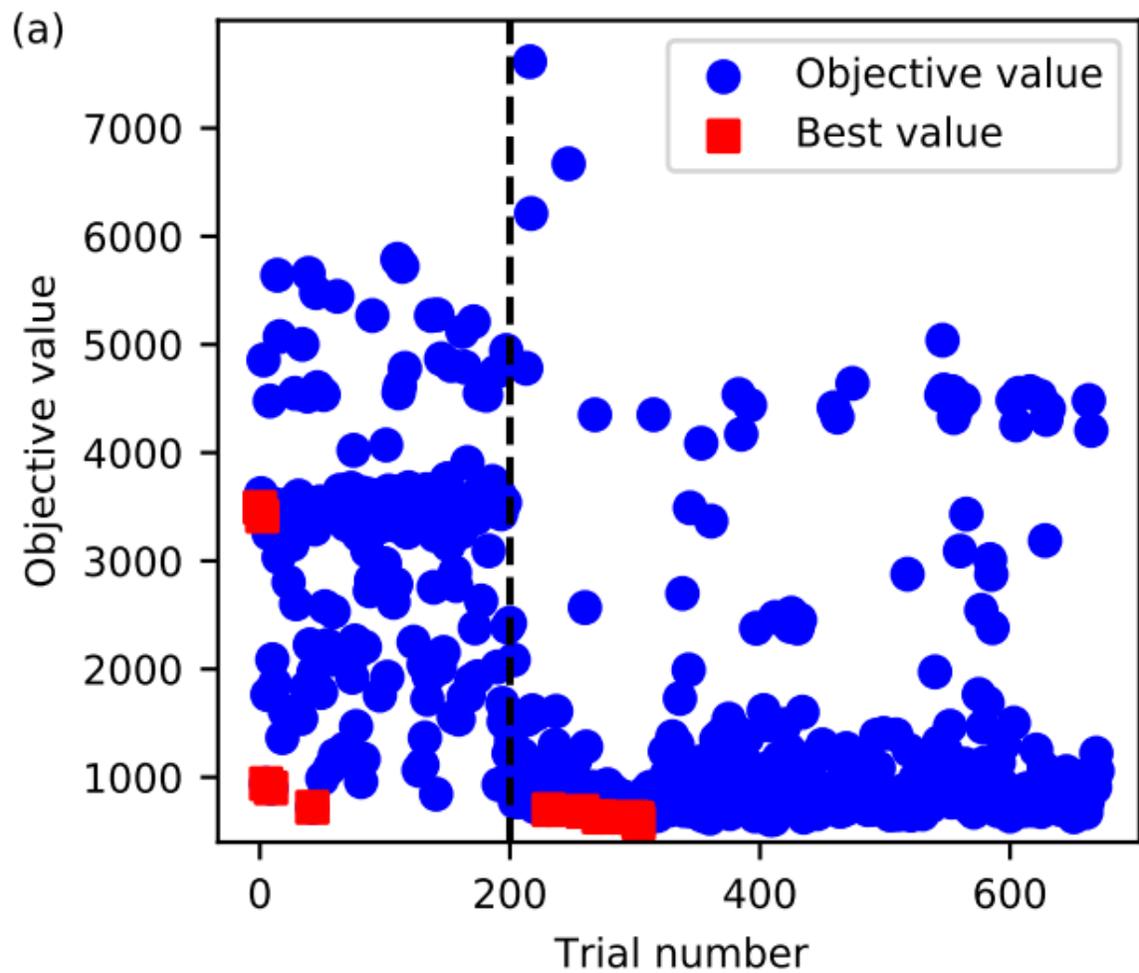


Figure12.

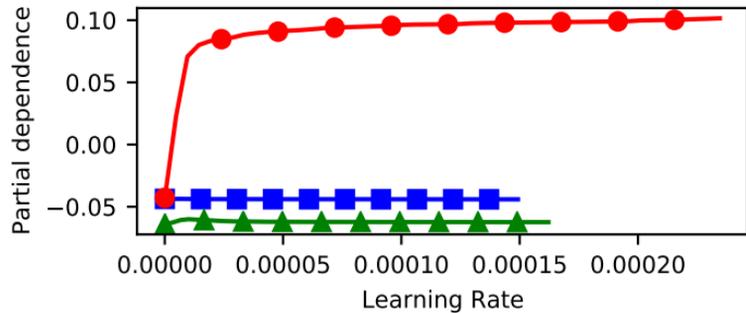
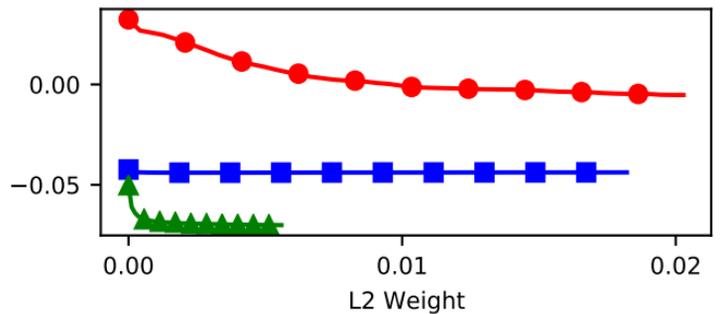
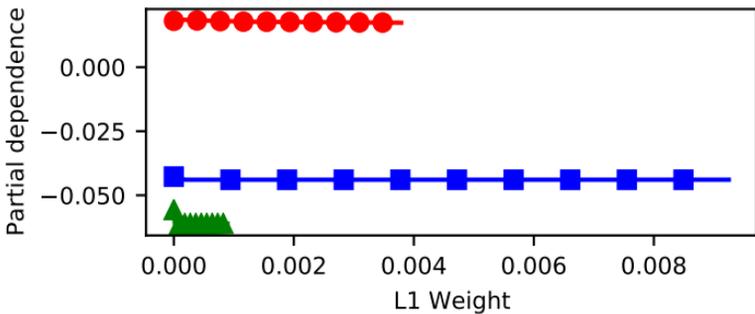
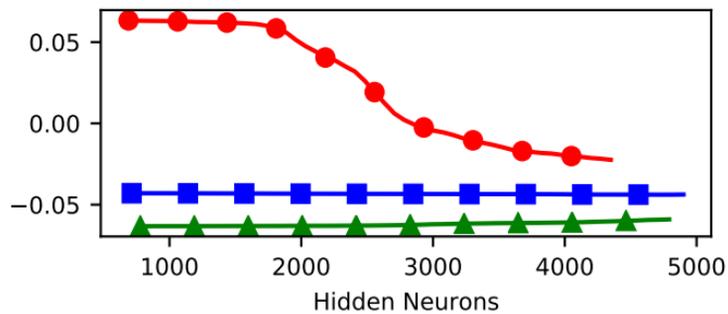
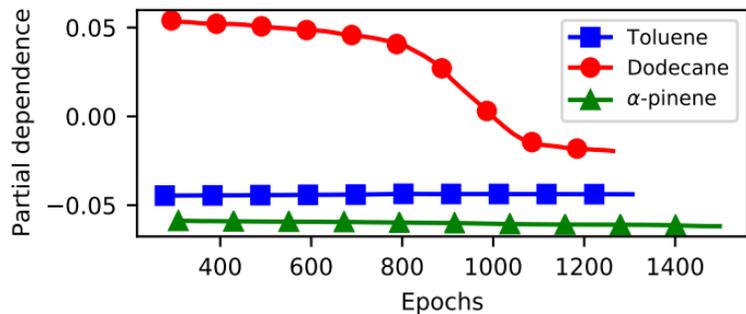


Figure 13.

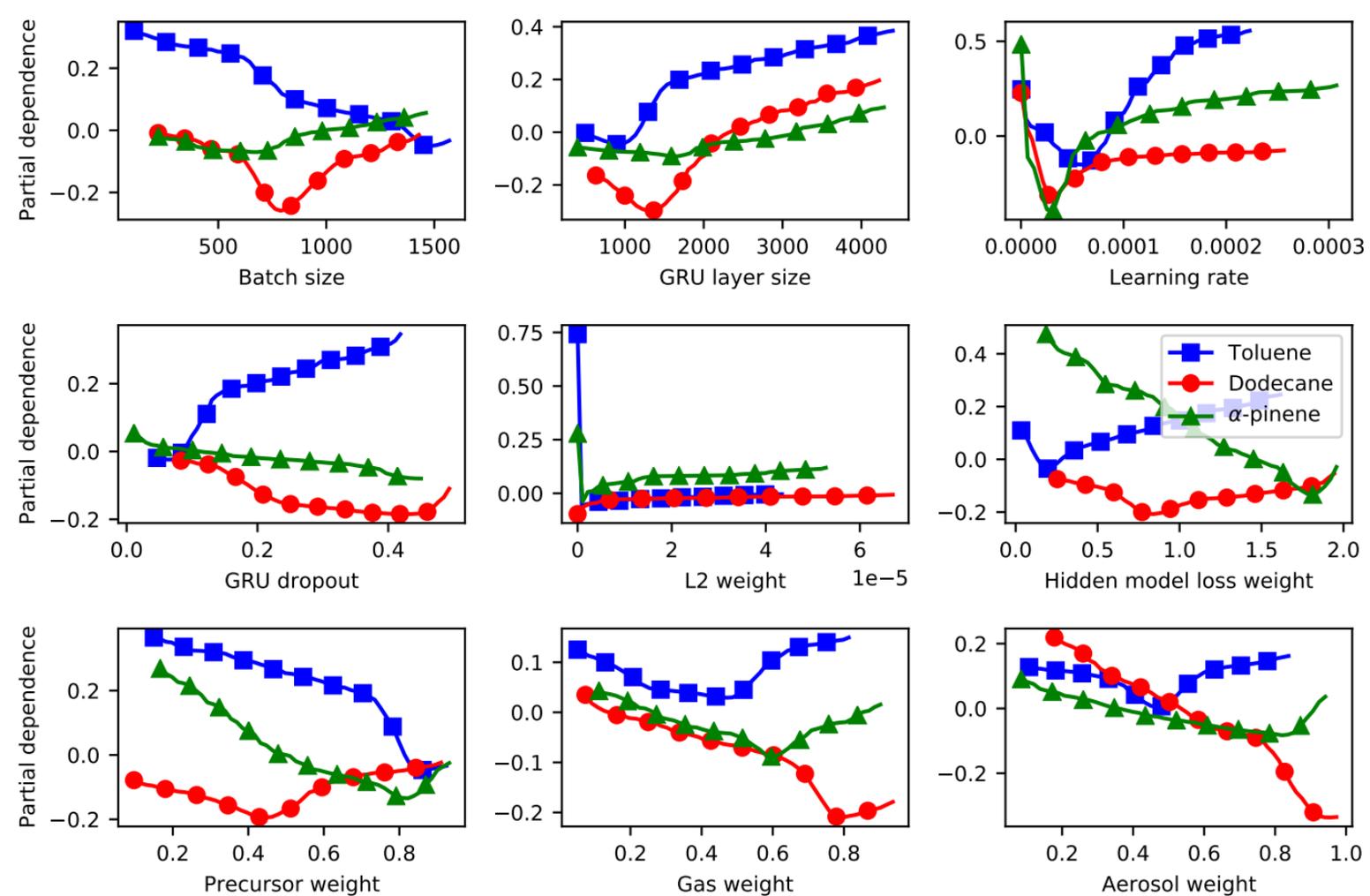
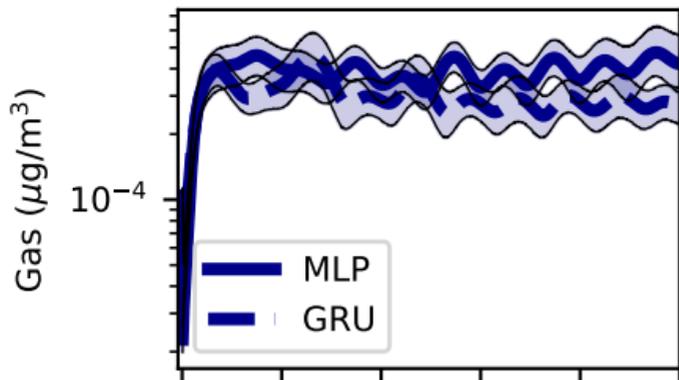


Figure 14.

Toluene



Dodecane

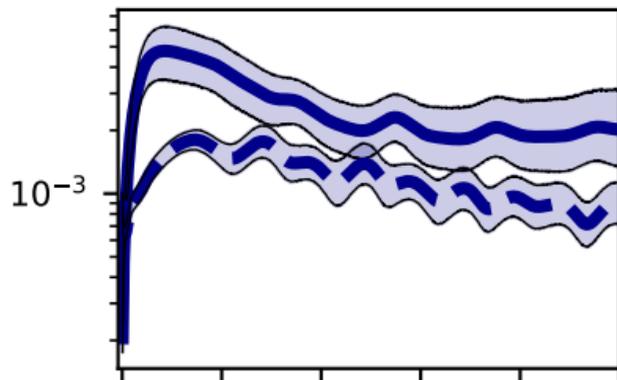
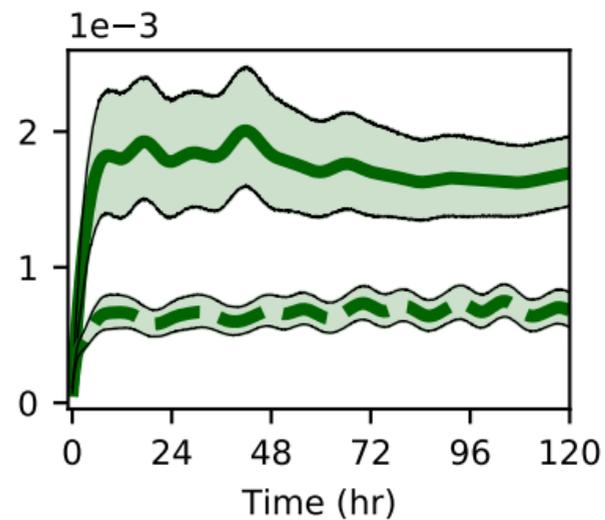
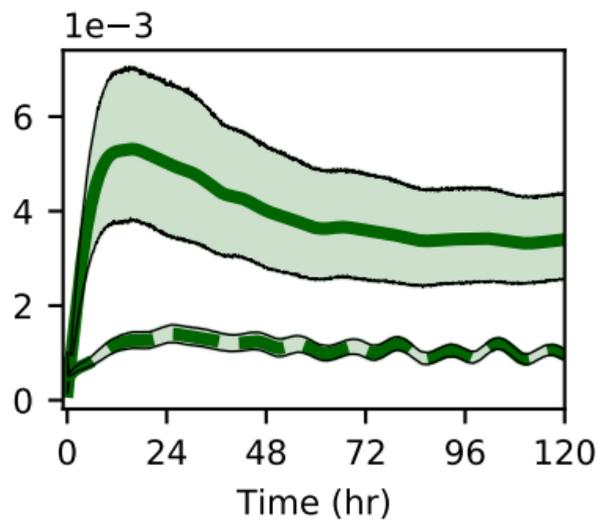
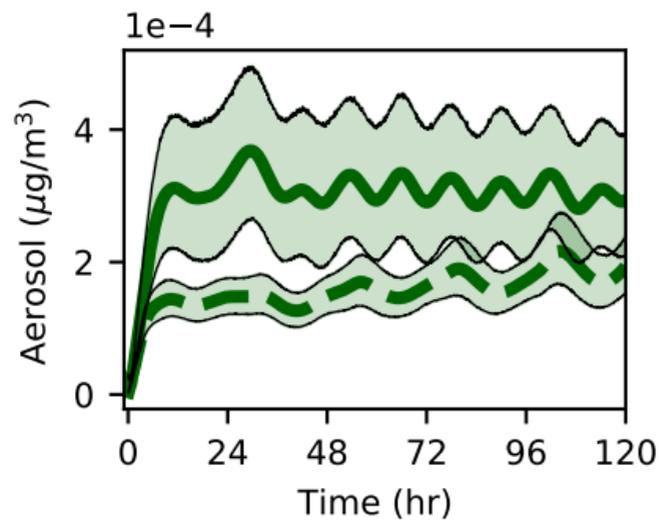
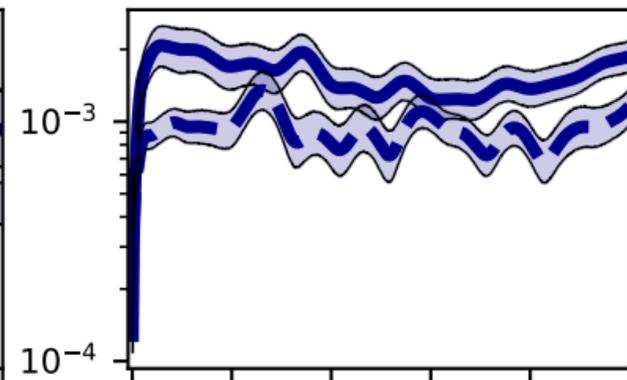
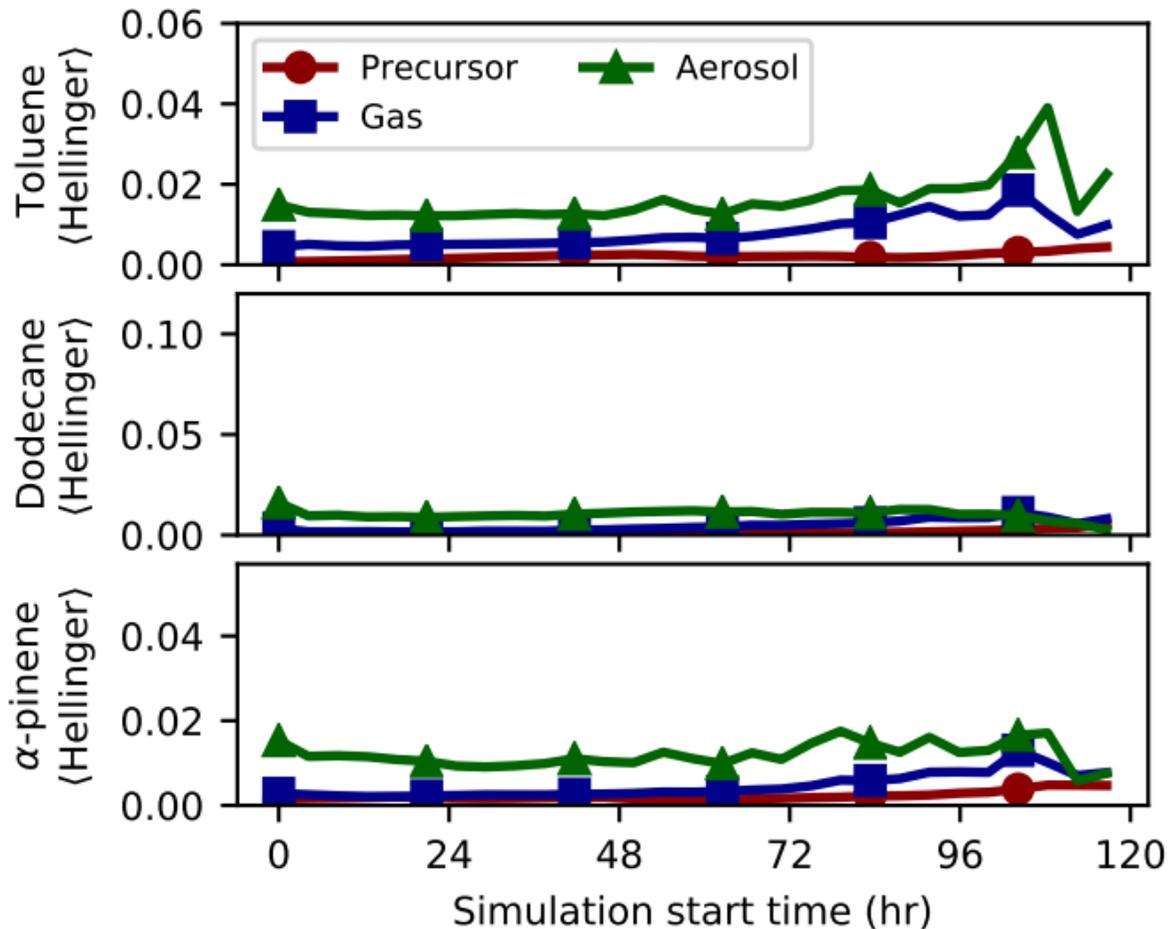
 α -pinene

Figure 15.

MLP



GRU

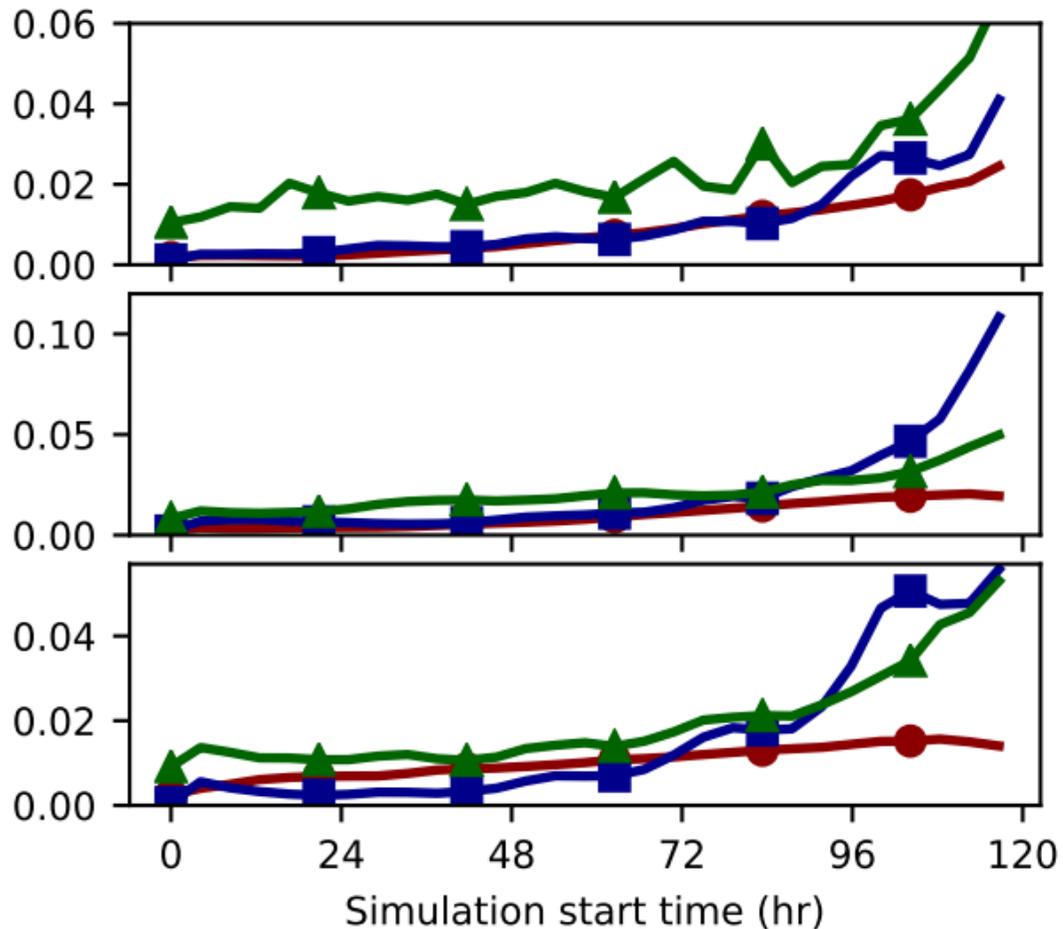
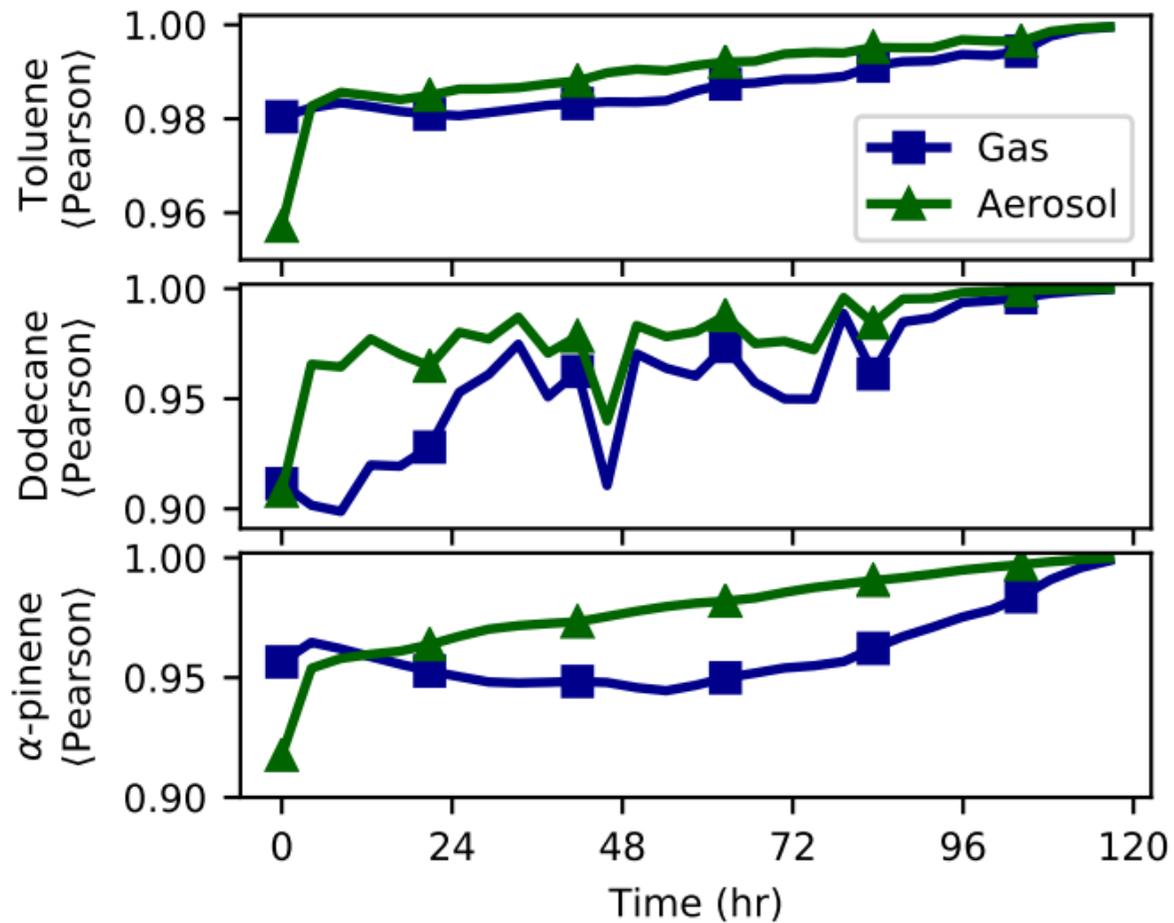


Figure 16.

MLP



GRU

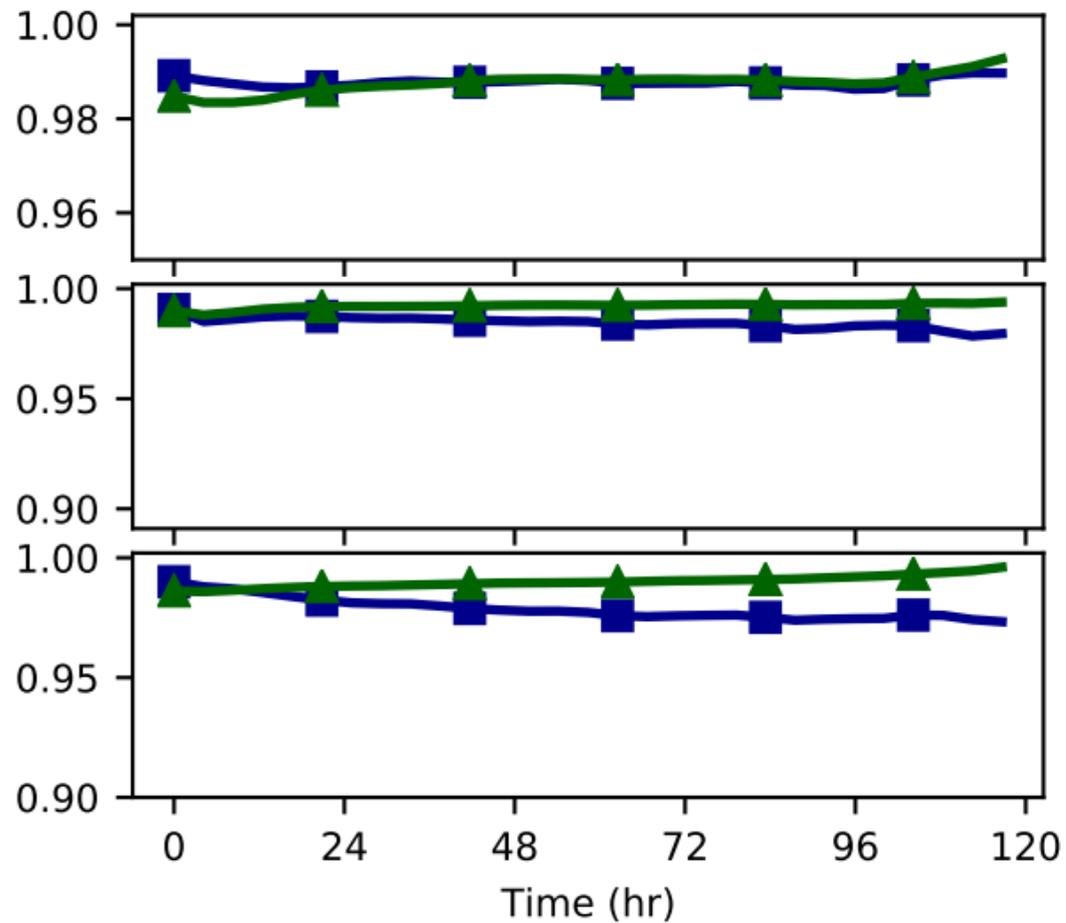
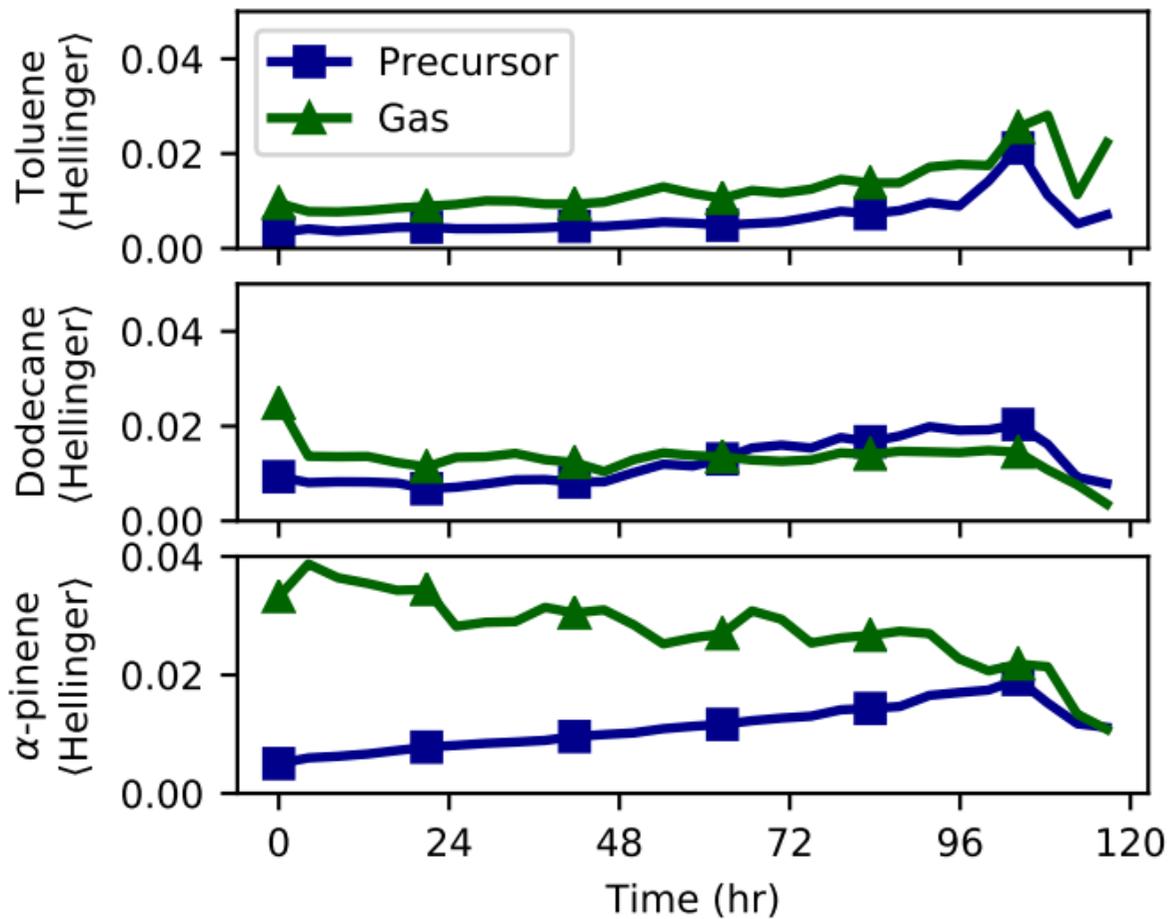


Figure 17.

MLP



GRU

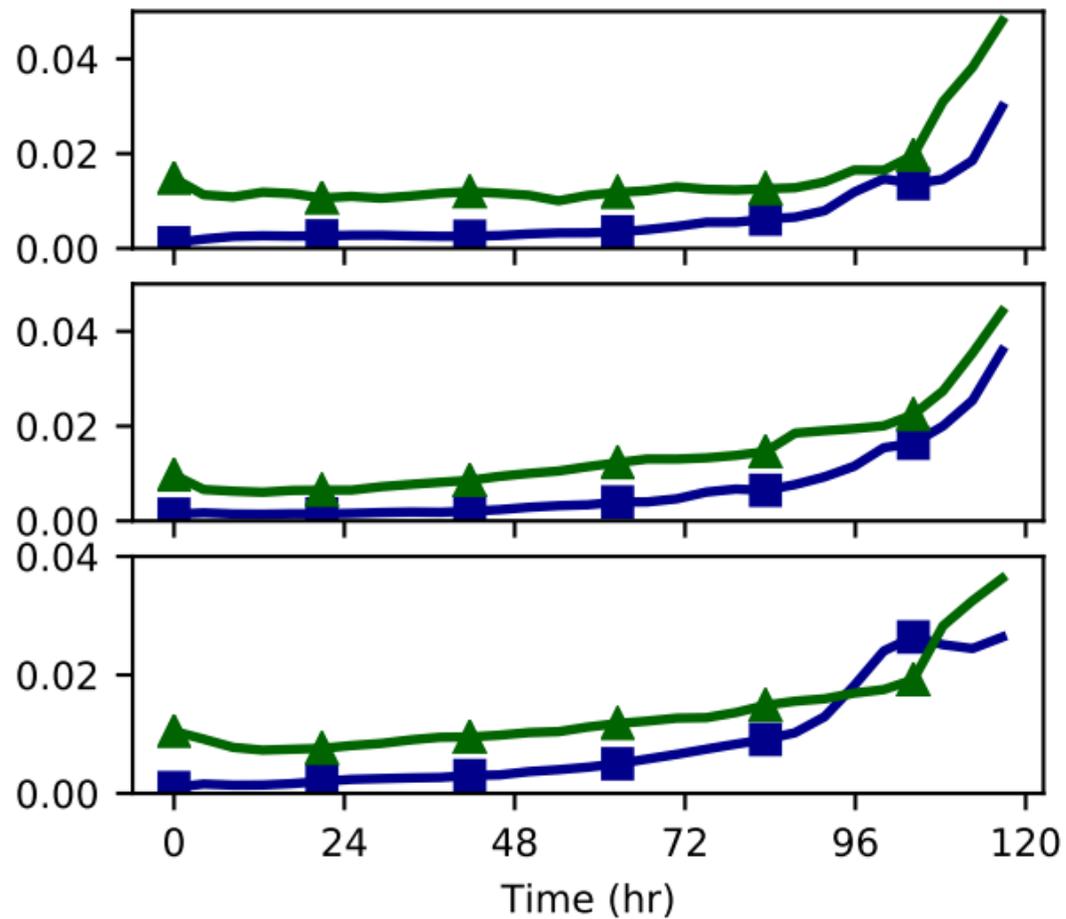
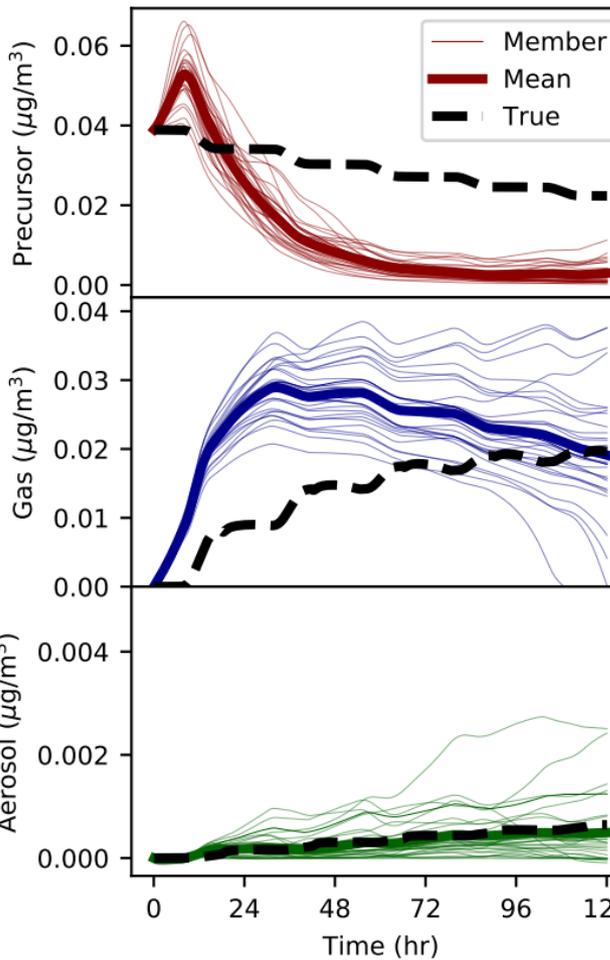
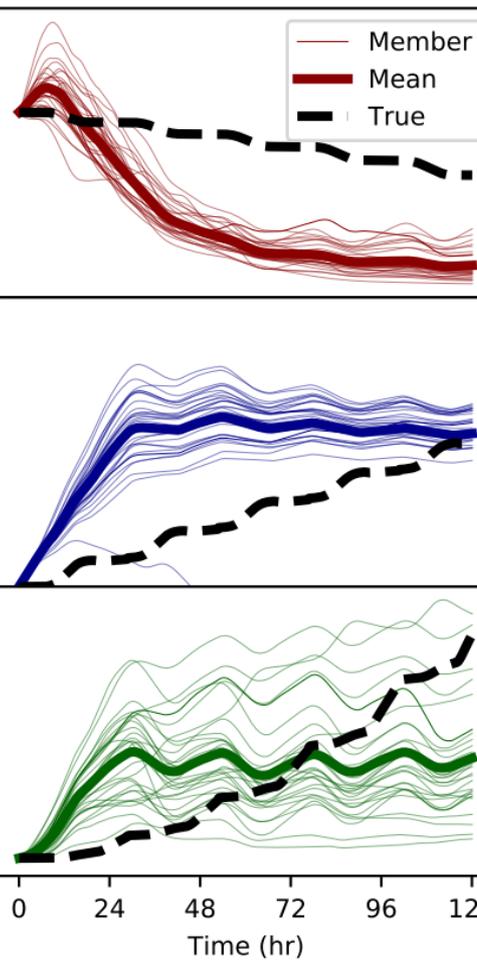


Figure 18.

Exp100003



Exp100028



Exp100030

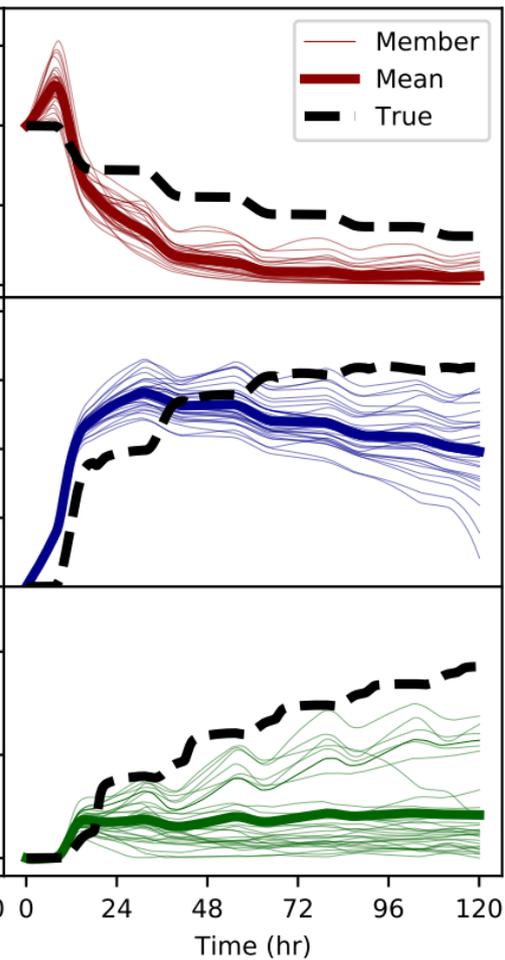
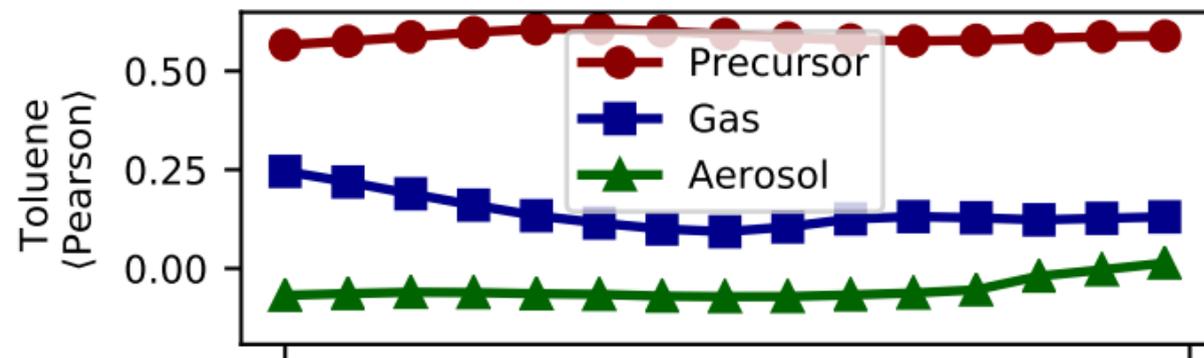
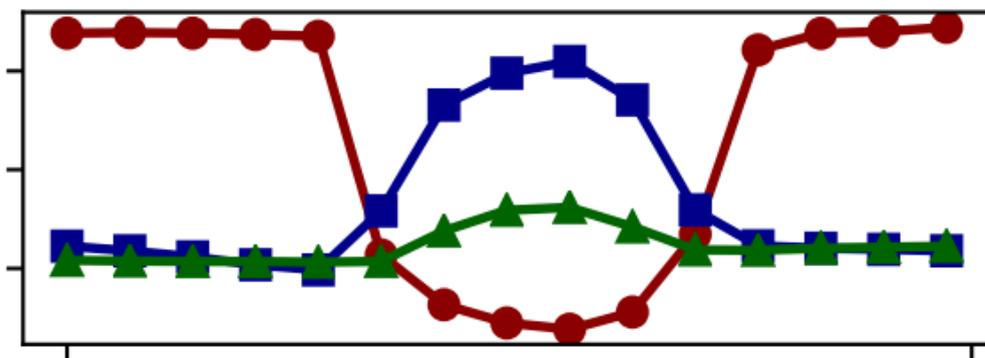


Figure 19.

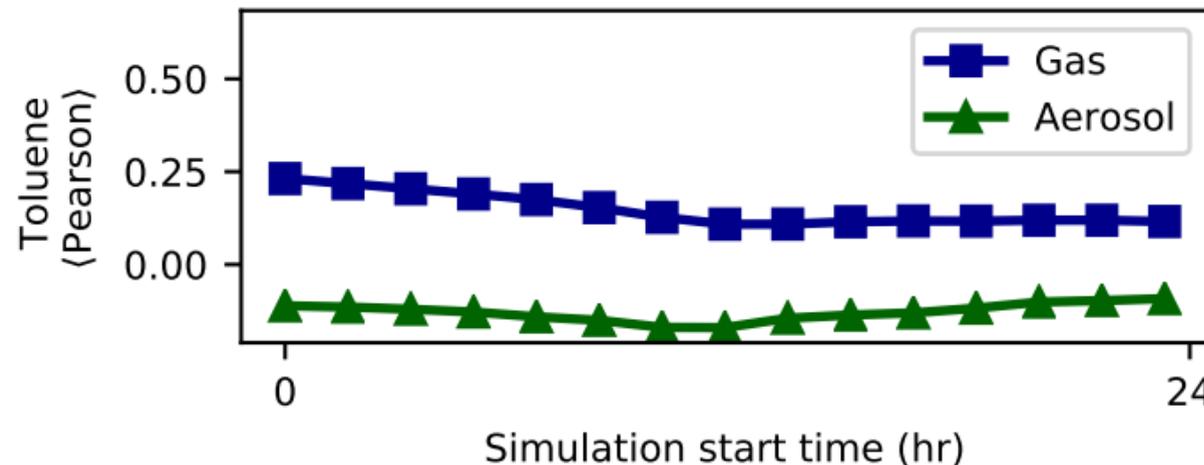
(a)(i) MLP (3-task)



(a)(ii) GRU (3-task)



(b)(i) MLP (2-task)



(b)(ii) GRU (2-task)

