



**Barcelona  
Supercomputing  
Center**  
*Centro Nacional de Supercomputación*

# Bringing the complexity of organic chemistry to climate models with machine learning techniques

Camille Mouchel-Vallon, Alma Hodzic, John Schreck, Charlie Becker, Keeley Lawrence, Siyuan Wang, Jinkyul Choi, Daven Henze, David John Gagne



#PlanDeRecuperación

**TR** Plan de Recuperación,  
Transformación  
y Resiliencia



PASC24, Zürich, 2024/06/05

# Context

---



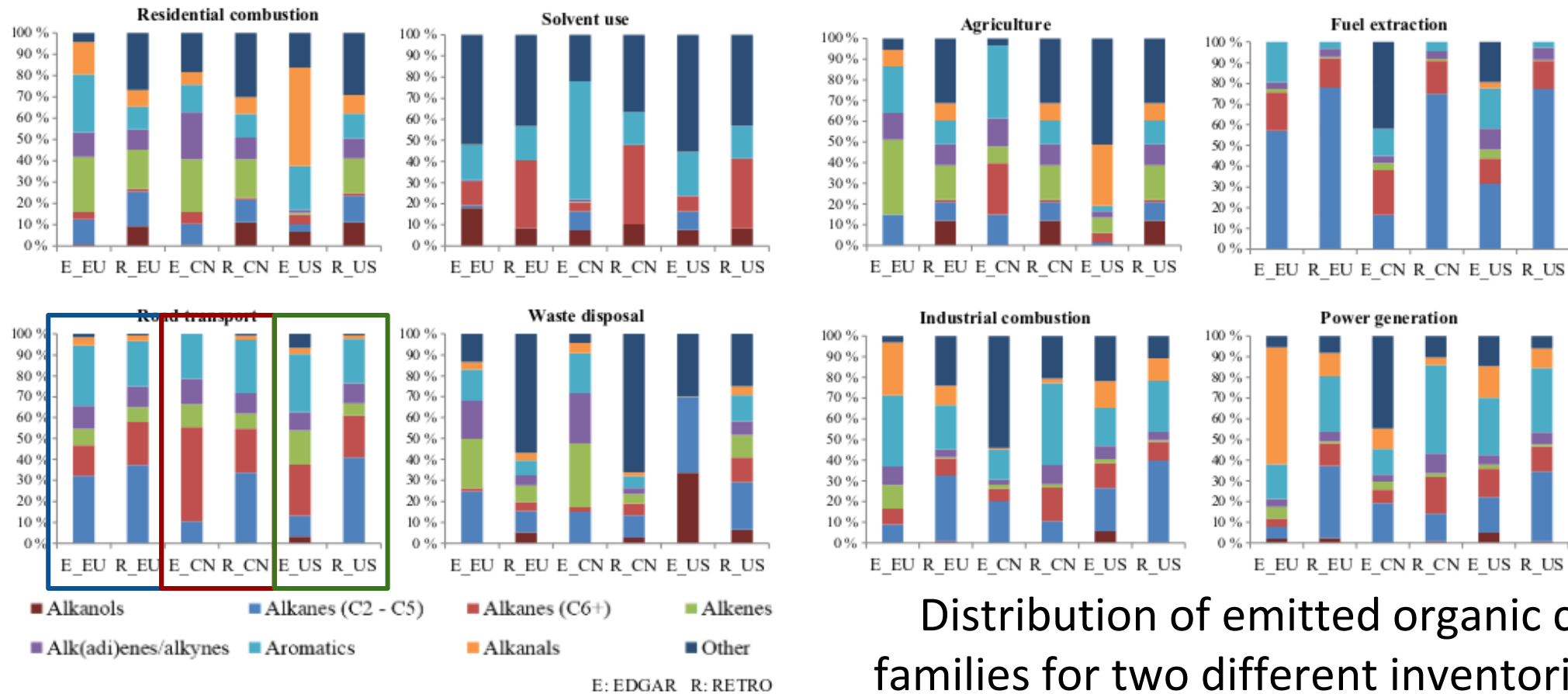
Plan de Recuperación,  
Transformación  
y Resiliencia

**1<sup>st</sup> objective:** A high resolution spatial (up to 2.5x2.5 km) and temporal (hourly) emissions system to compute emissions of **primary atmospheric pollutants** for Spain to be used in the national **air quality prediction** system.

Gaseous primary atmospheric pollutants  
NO<sub>x</sub>, CO, SO<sub>x</sub>, NH<sub>3</sub>, **NMVOCs**

Non-Methane Volatile Organic Compounds are emitted by human and natural activities and their chemical evolution in the atmosphere has an impact on air quality, health and climate

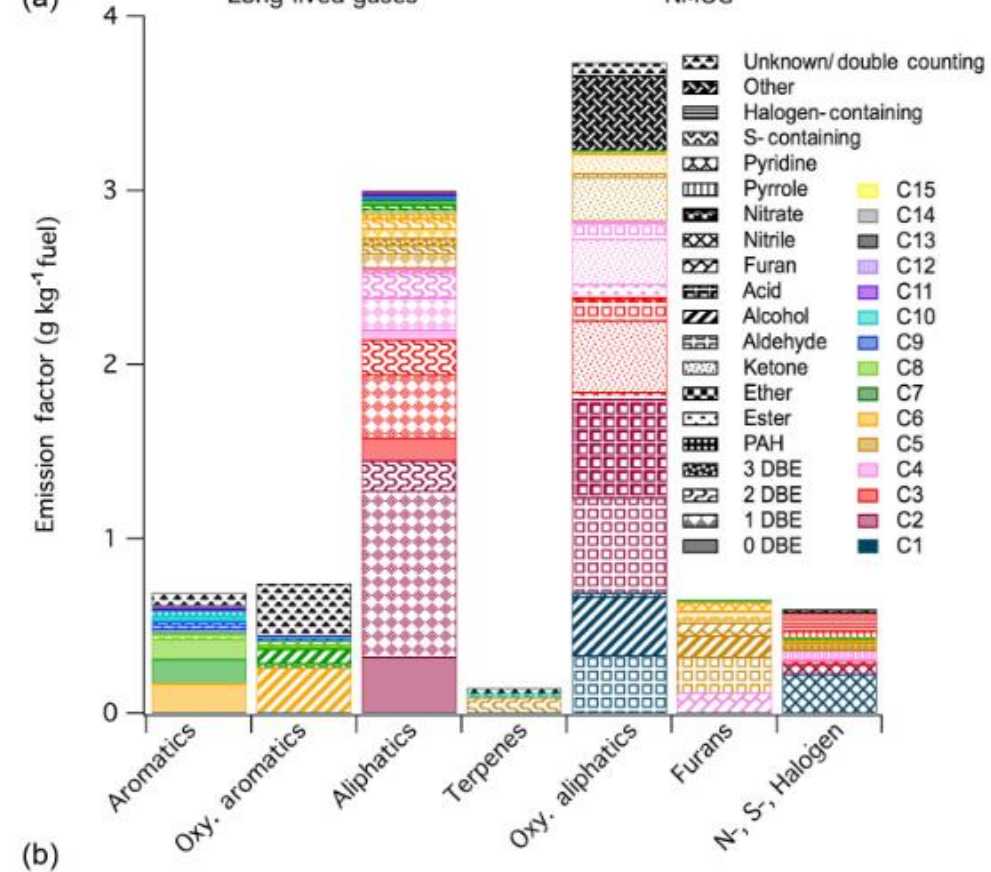
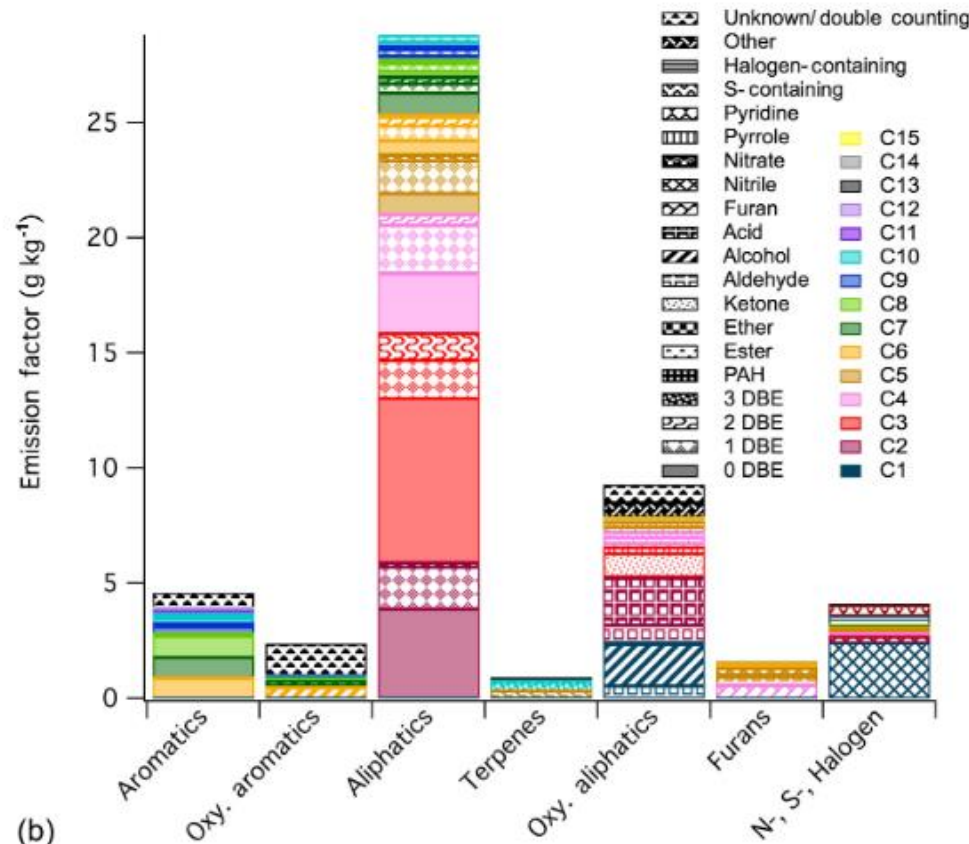
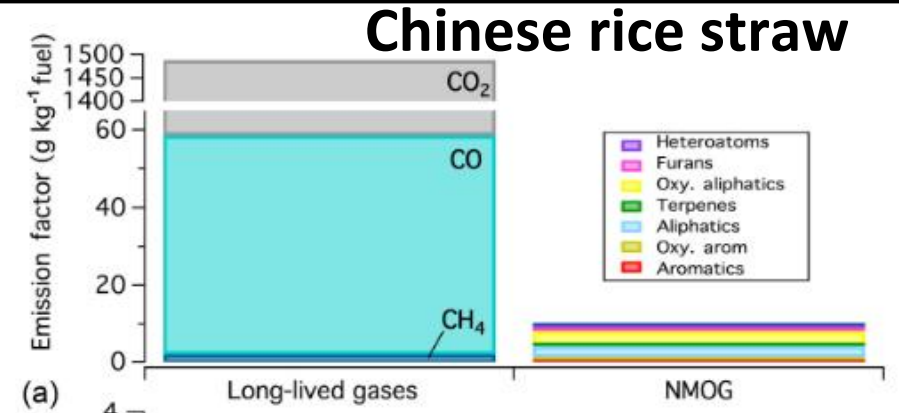
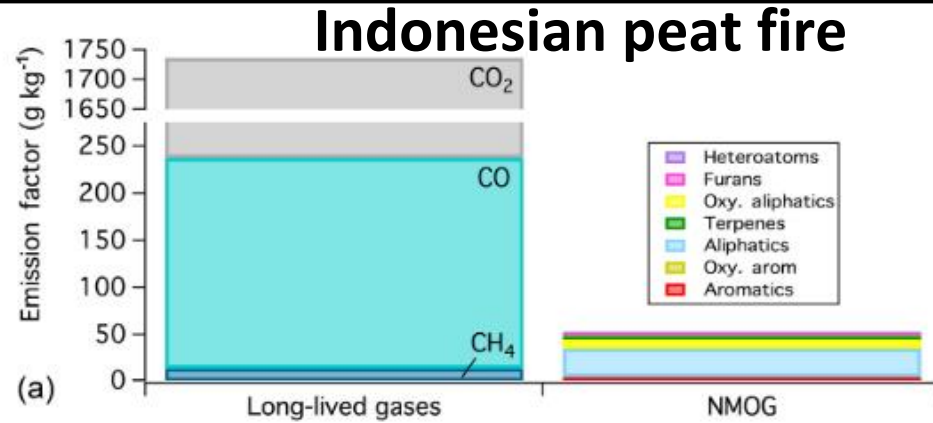
# Primary organic compounds: many sources, high diversity, high uncertainty



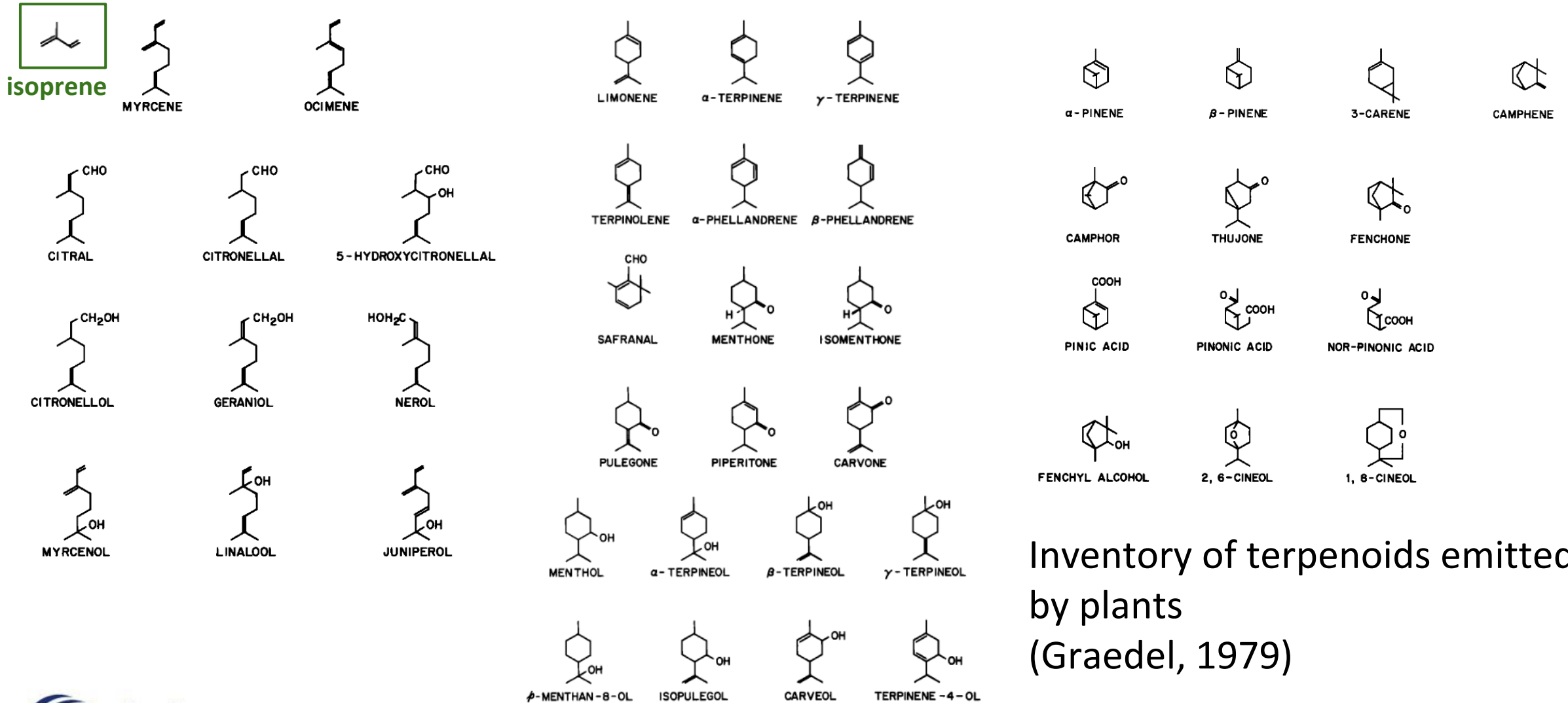
Distribution of emitted organic compounds families for two different inventories in **Europe**, **China** and **USA**  
(Huang et al., 2017)

# Primary organic compounds: many sources, high diversity, high uncertainty

Measured fire emissions composition (Hatch et al., 2017)

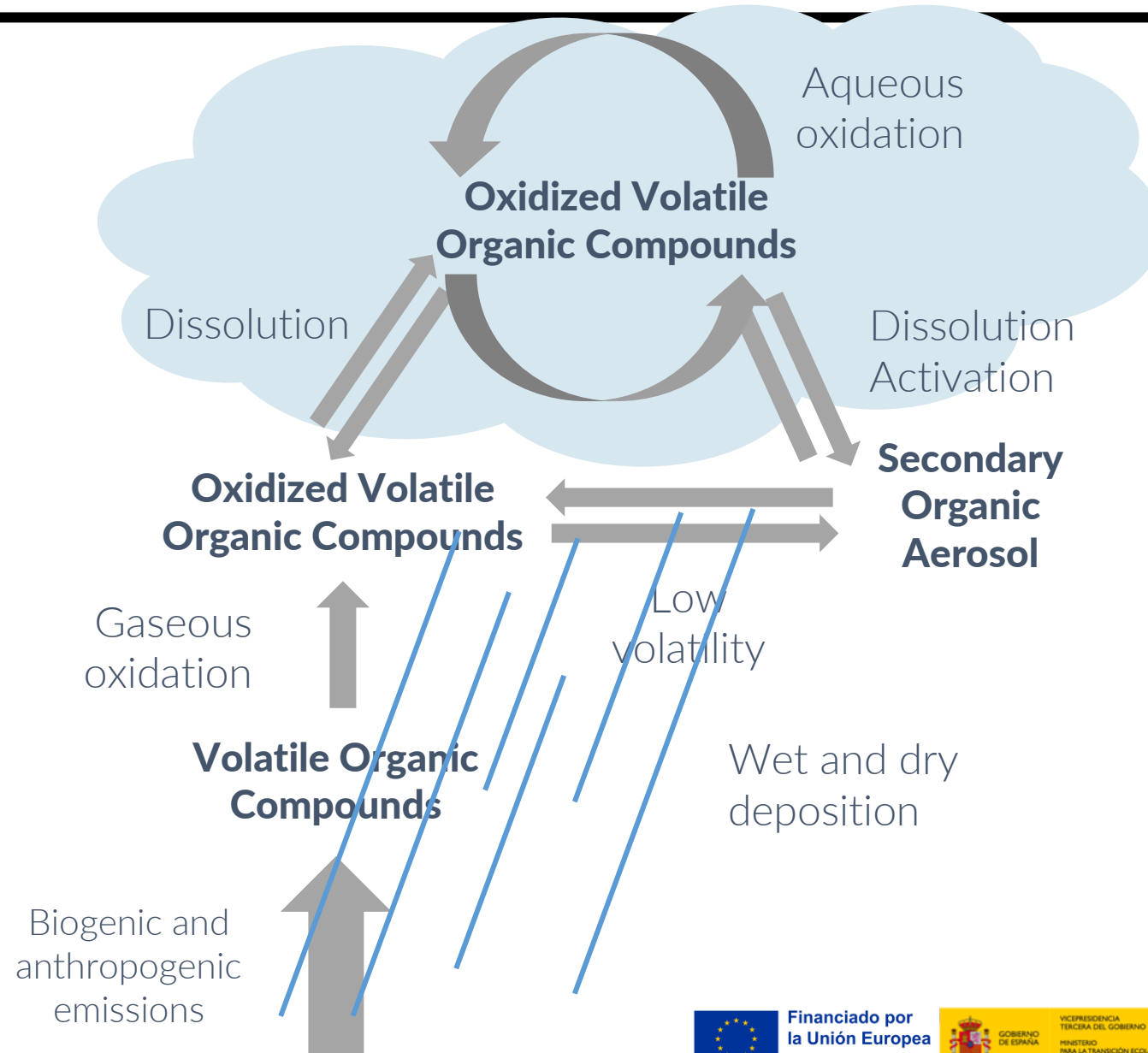
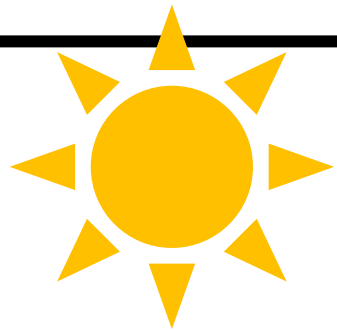


# Primary organic compounds: many sources, high diversity, high uncertainty

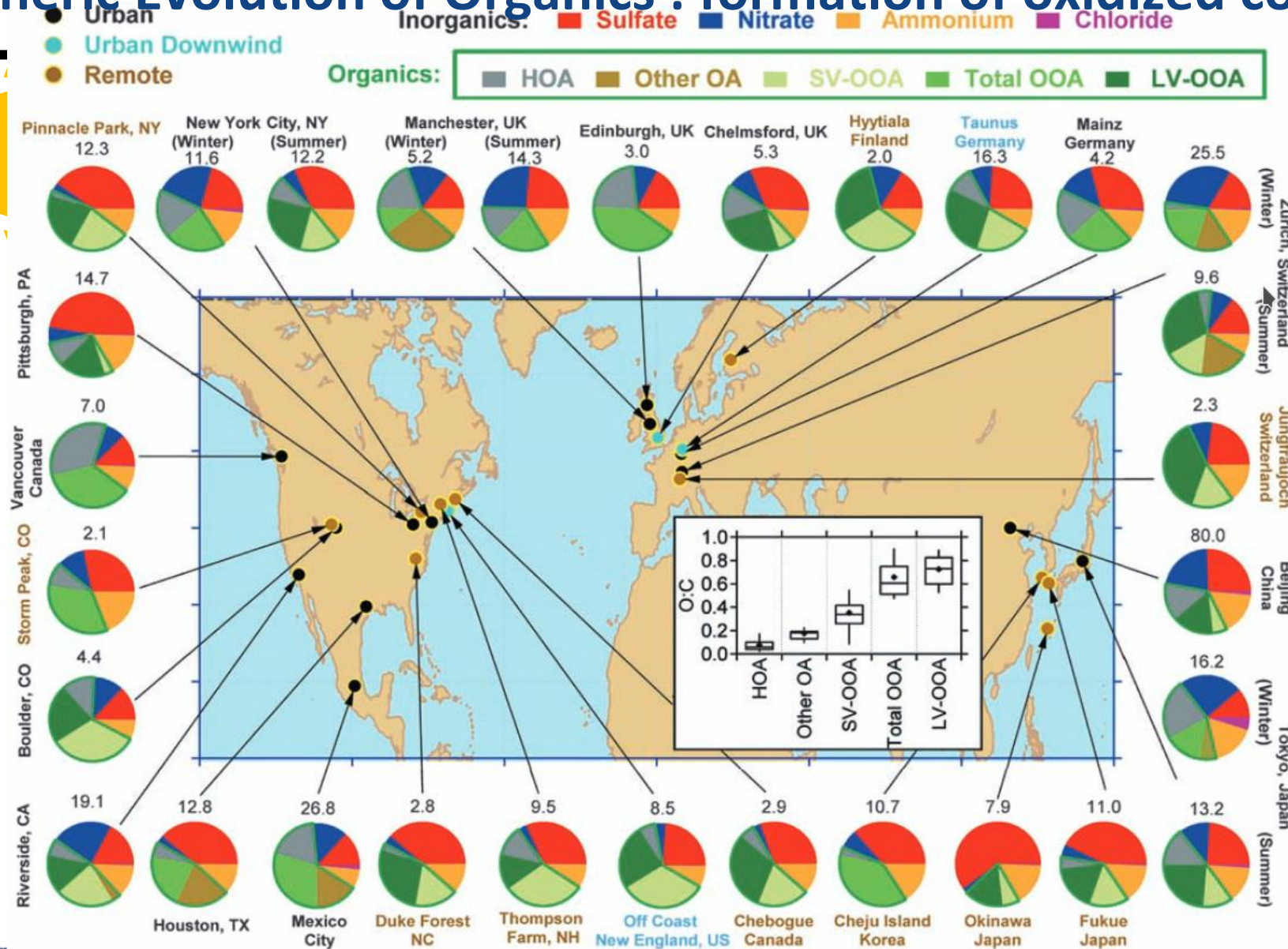


Inventory of terpenoids emitted by plants (Graedel, 1979)

# Atmospheric Evolution of Organics : formation of oxidized compounds

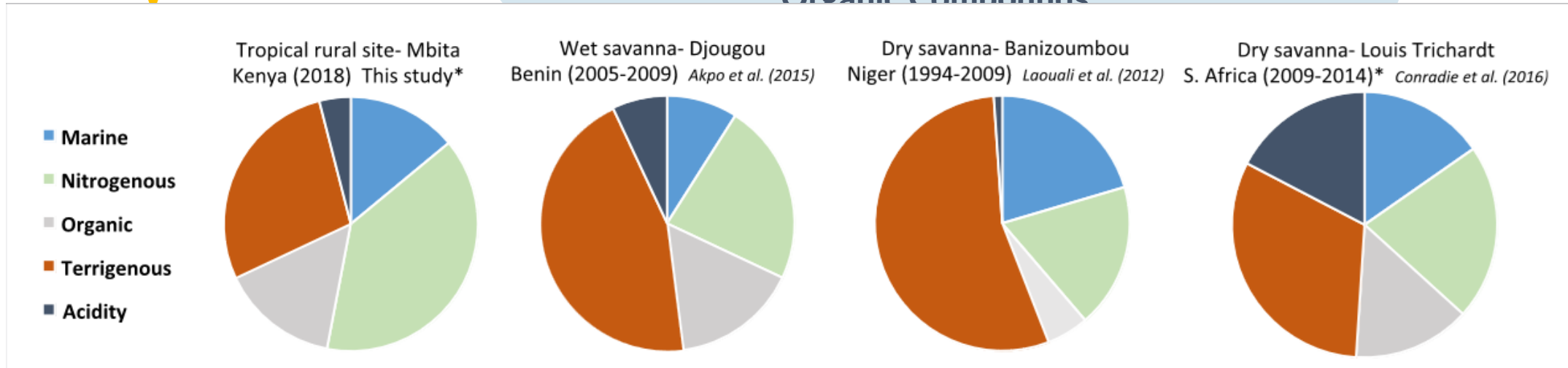
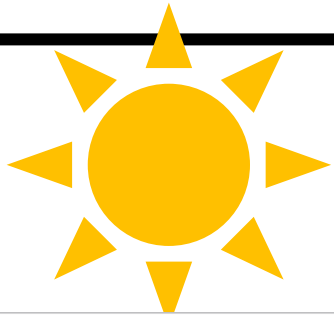


# Atmospheric Evolution of Organics : formation of oxidized compounds



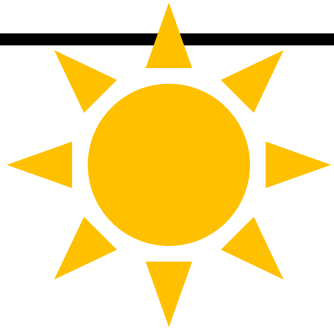
Measured PM1 composition  
Jimenez et al. (2009)

# Atmospheric Evolution of Organics : formation of oxidized compounds

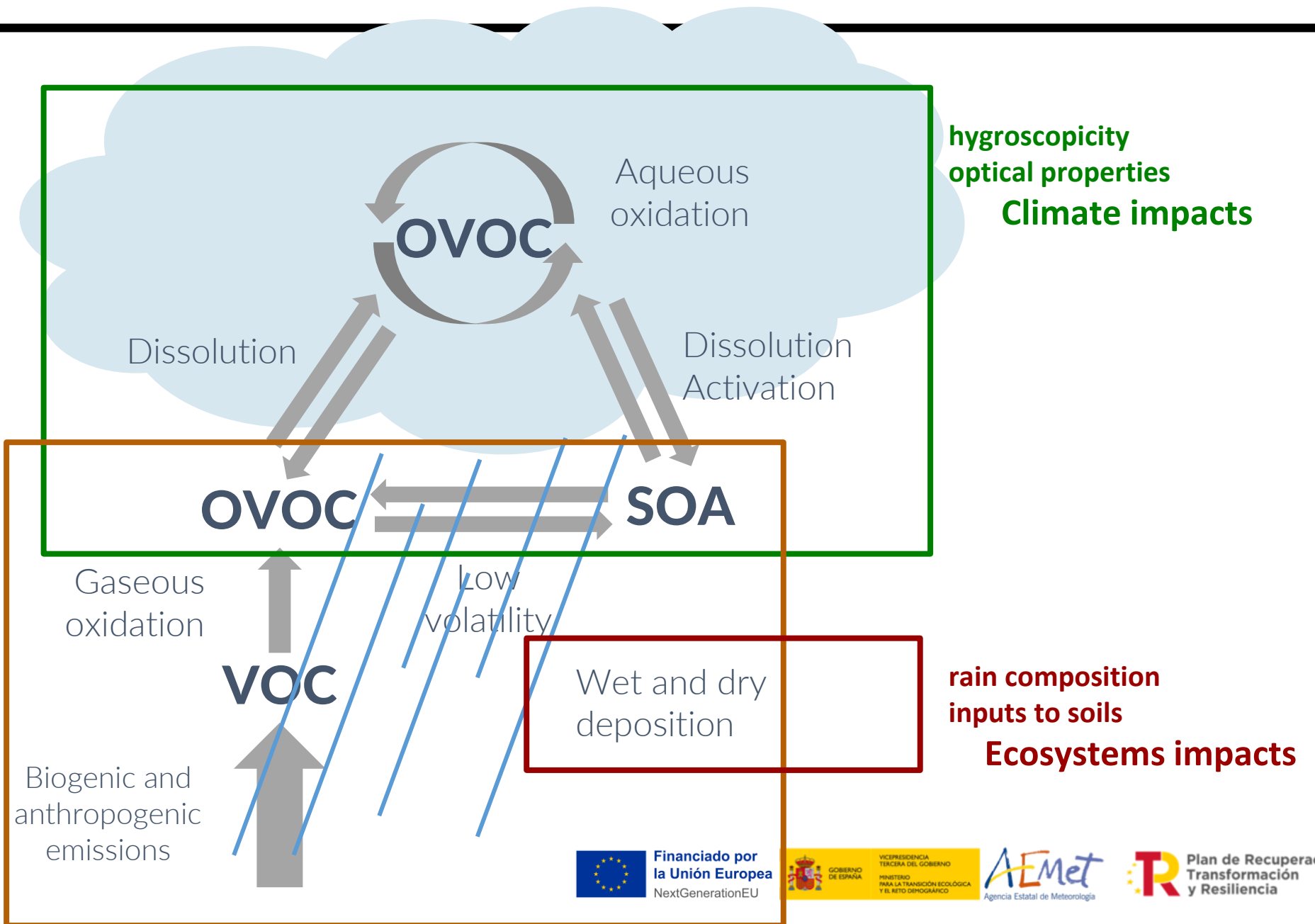




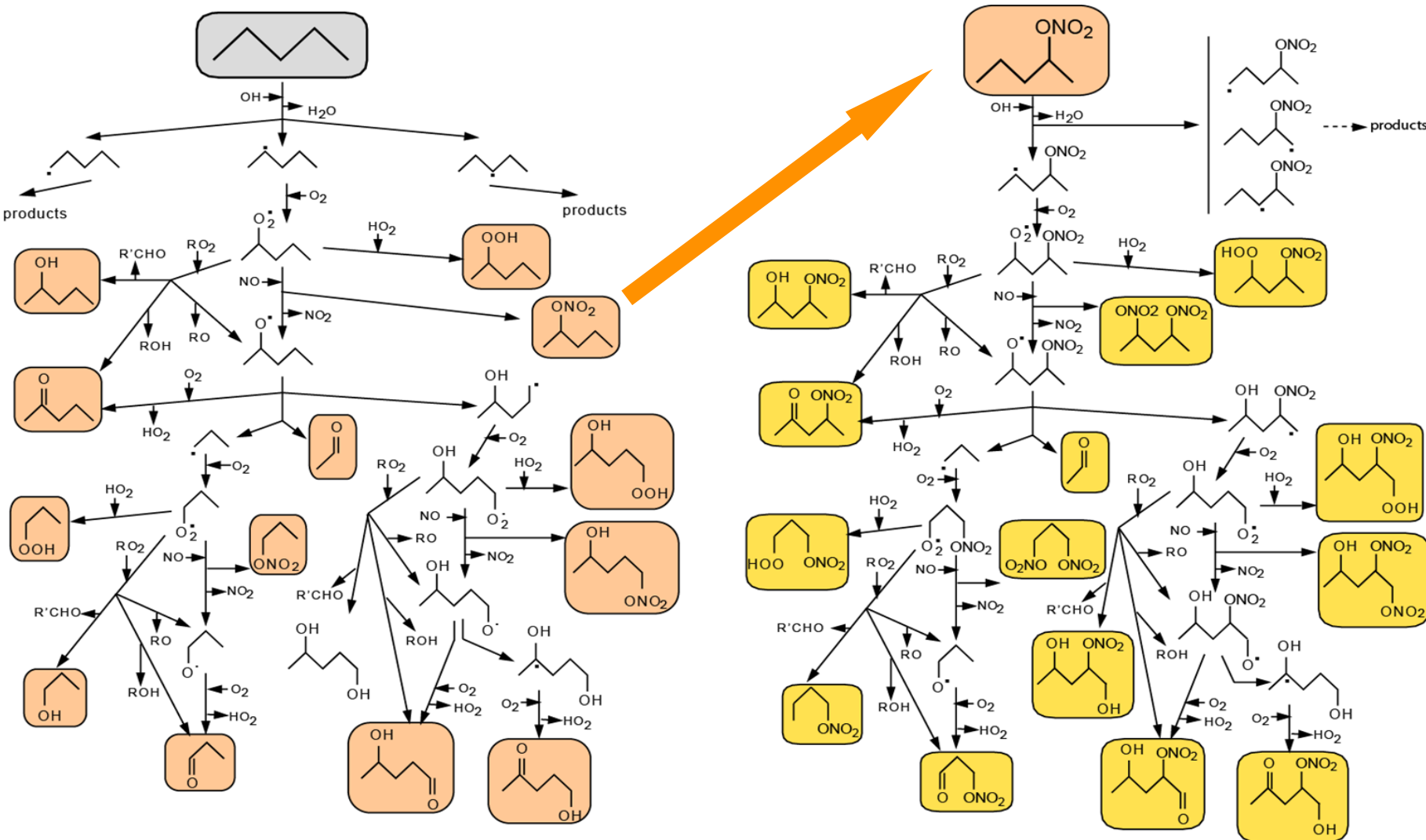
# Impacts



ozone chemistry  
particulate matter composition  
**Air quality and health impacts**



# Progressive and complex oxidation of organics



Oxidized Volatile Organic Compounds

Gaseous oxidation

↑

Volatile Organic Compounds

How to bring this complexity to air quality models?

CO<sub>2</sub>

Parent Hydrocarbon

1<sup>st</sup> Generation Species

2<sup>nd</sup> Generation Species

# How I Learned to Stop Worrying and Love the Complexity

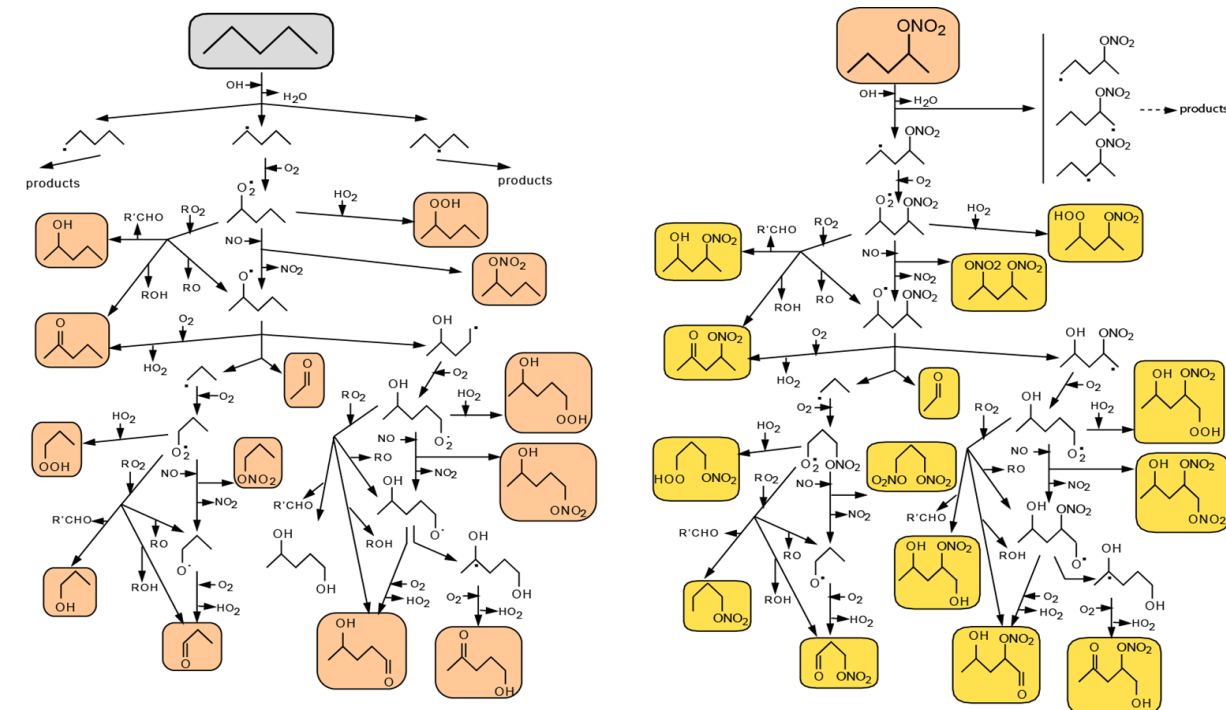
**Master Chemical Mechanism** (Jenkin et al., 1997;  
15yrs 2-3people ?)

- 143 precursors from  $C_1$  to  $C_{12}$
- 17000 reactions of 6700 species
- Simplified!

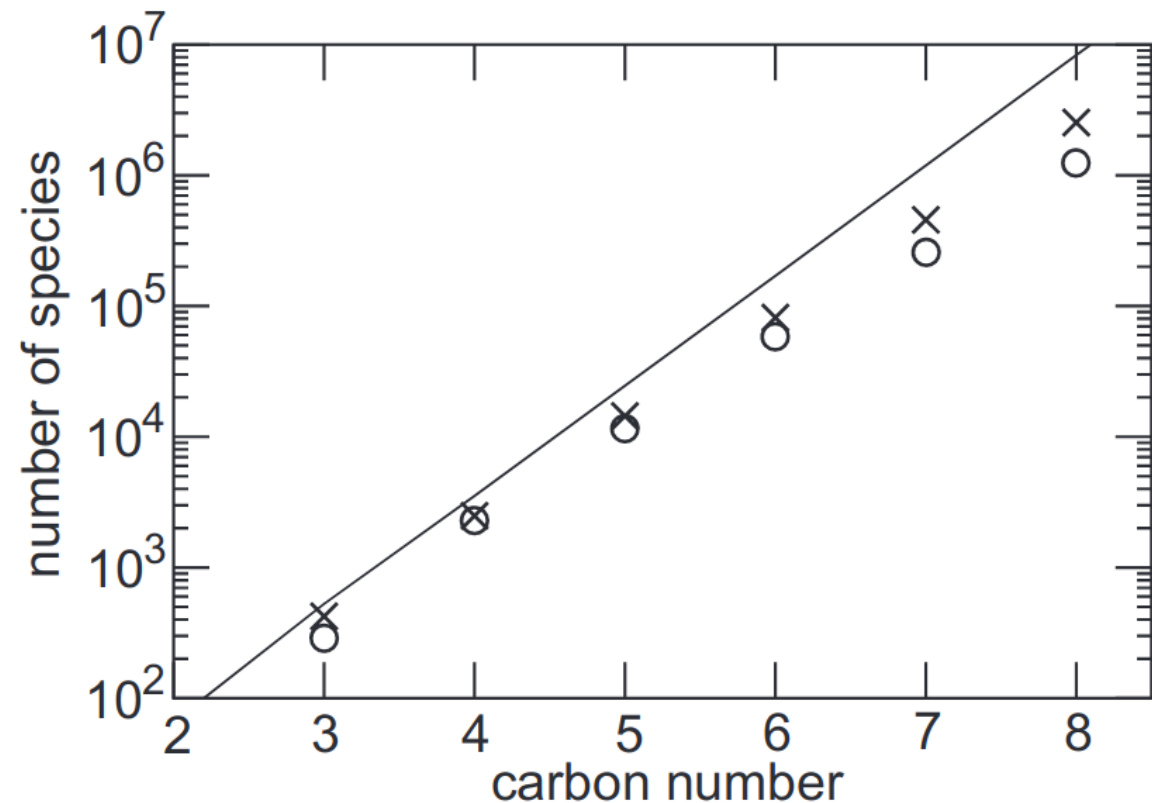
**CLEPS** (Mouchel-Vallon et al., 2017; 2yrs 2people)

- Cloud oxidation of  $C_1$ - $C_4$  products from isoprene oxidation
- 1315 reactions involving 717 chemical species
- Simplified!

Handwritten detailed mechanisms: high potential for errors and time consuming (creation and maintenance)



# How I Learned to Stop Worrying and Love the Complexity



Potential number of distinct species in a full oxidation mechanism as a function of the size of the precursor

Aumont et al. (2005)

**Master Chemical Mechanism** (Jenkin et al., 1997; 15yrs 2-3people ?)

- 143 precursors from  $C_1$  to  $C_{12}$
- 17000 reactions of 6700 species
- Simplified!

**CLEPS** (Mouchel-Vallon et al., 2017; 2yrs 2people)

- Cloud oxidation of  $C_1$ - $C_4$  products from isoprene oxidation
- 1315 reactions involving 717 chemical species
- Simplified!

Handwritten detailed mechanisms: high potential for errors and time consuming (creation and maintenance)

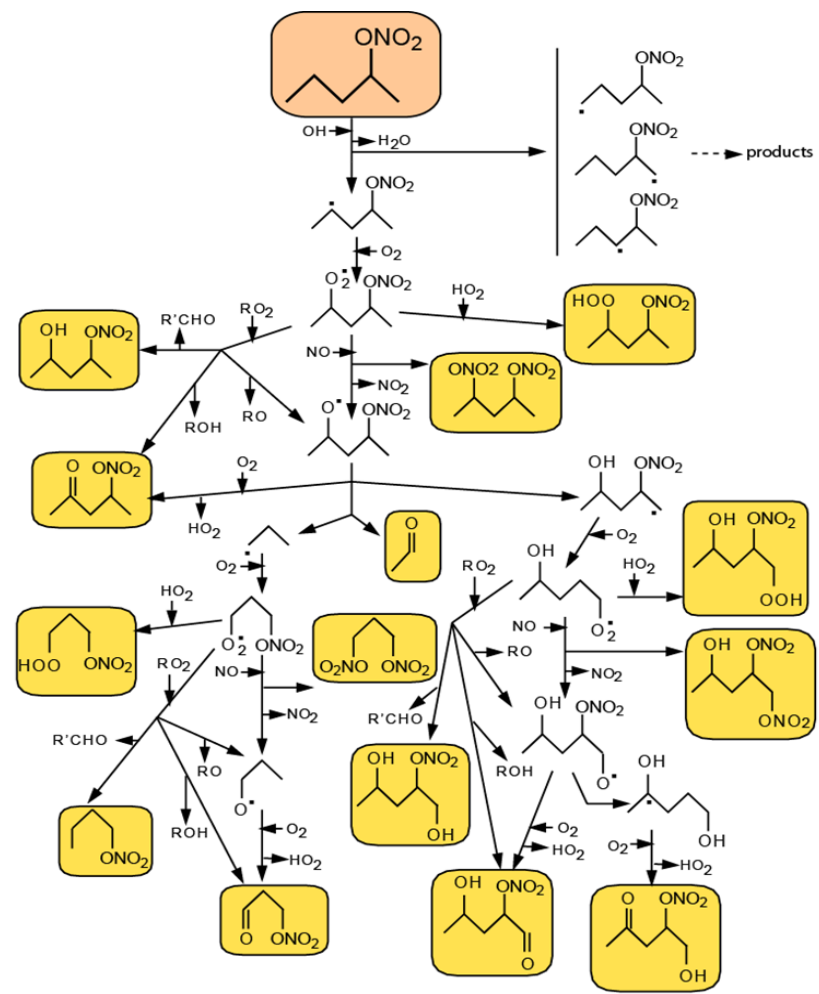
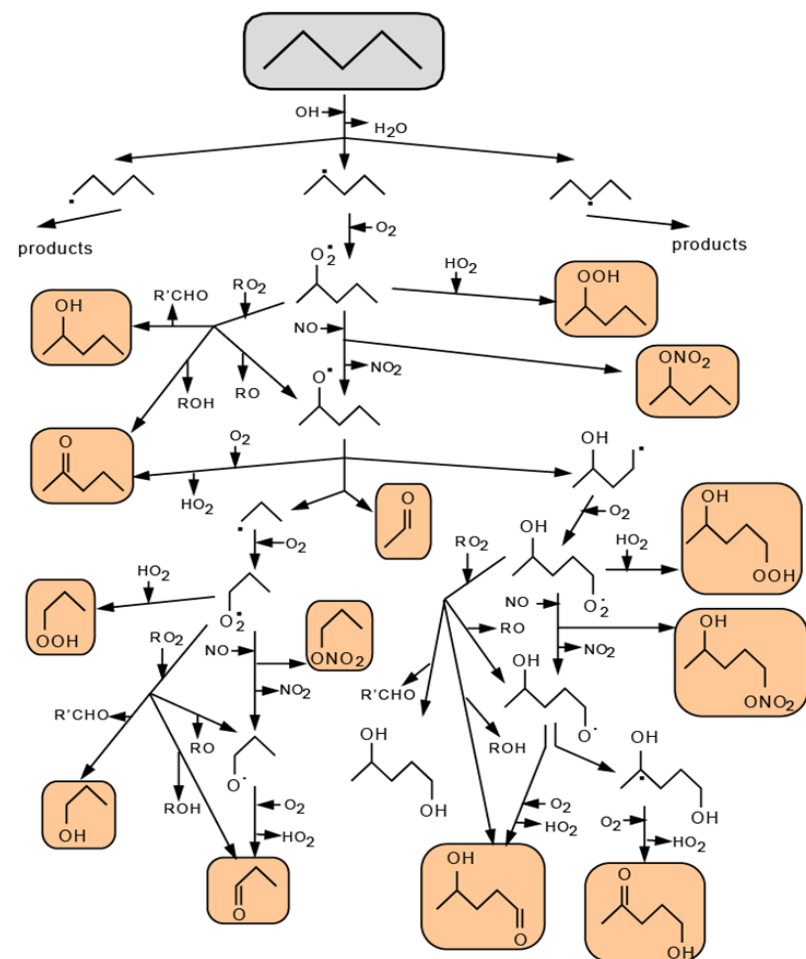
# We need to go faster

---

Writing the chemical mechanism

Solving the ODEs in a (3D) model

# Automating boring things

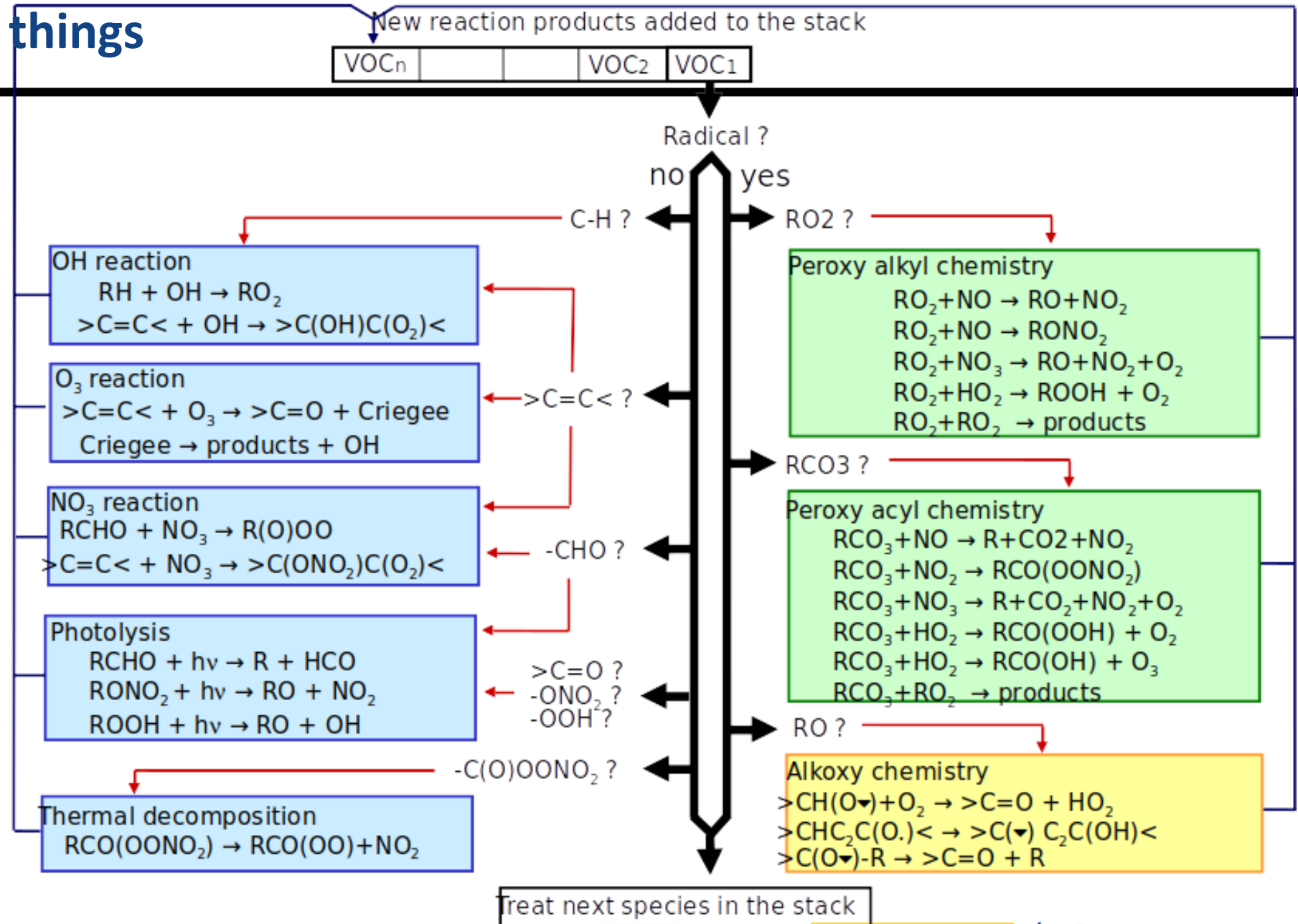


Systematic and repeated oxidation steps

Generator for Explicit Chemistry and Kinetics of Organics in the Atmosphere (GECKO-A, Aumont et al., 2005)

# Automating boring things

Aumont et al., 2005



# Accelerating things

Writing the chemical mechanism

**GECKO-A**

12 biogenic, 53 anthropogenic precursors



23 million reactions involving 4.4 million species

(Mouchel-Vallon et al., 2020)



Solving the ODEs in a (3D) model

Simulating 2 days in 2 grid cells

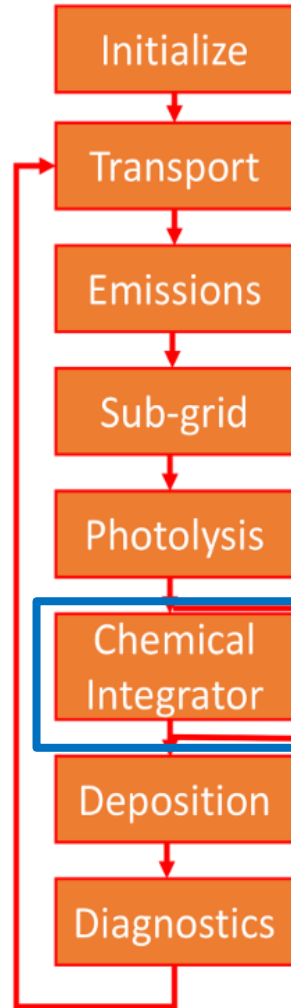
≈ 36 hours on 16 cores

Machine learning?



# Emulating atmospheric chemistry

## Numerical model



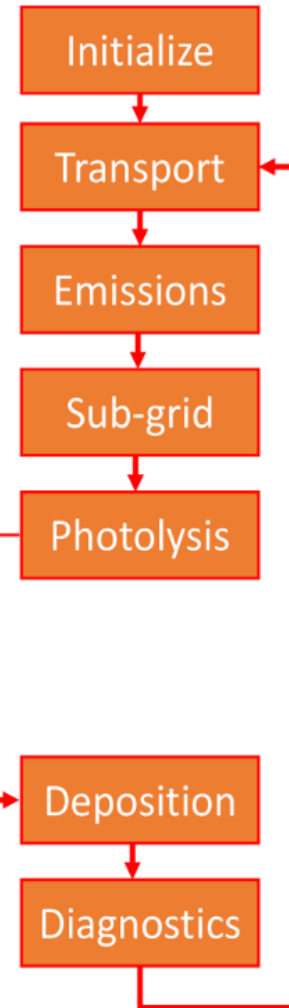
Computationally expensive

Training data



Random forest

## Emulator



Use GECKO-A 0D explicit simulations as training dataset

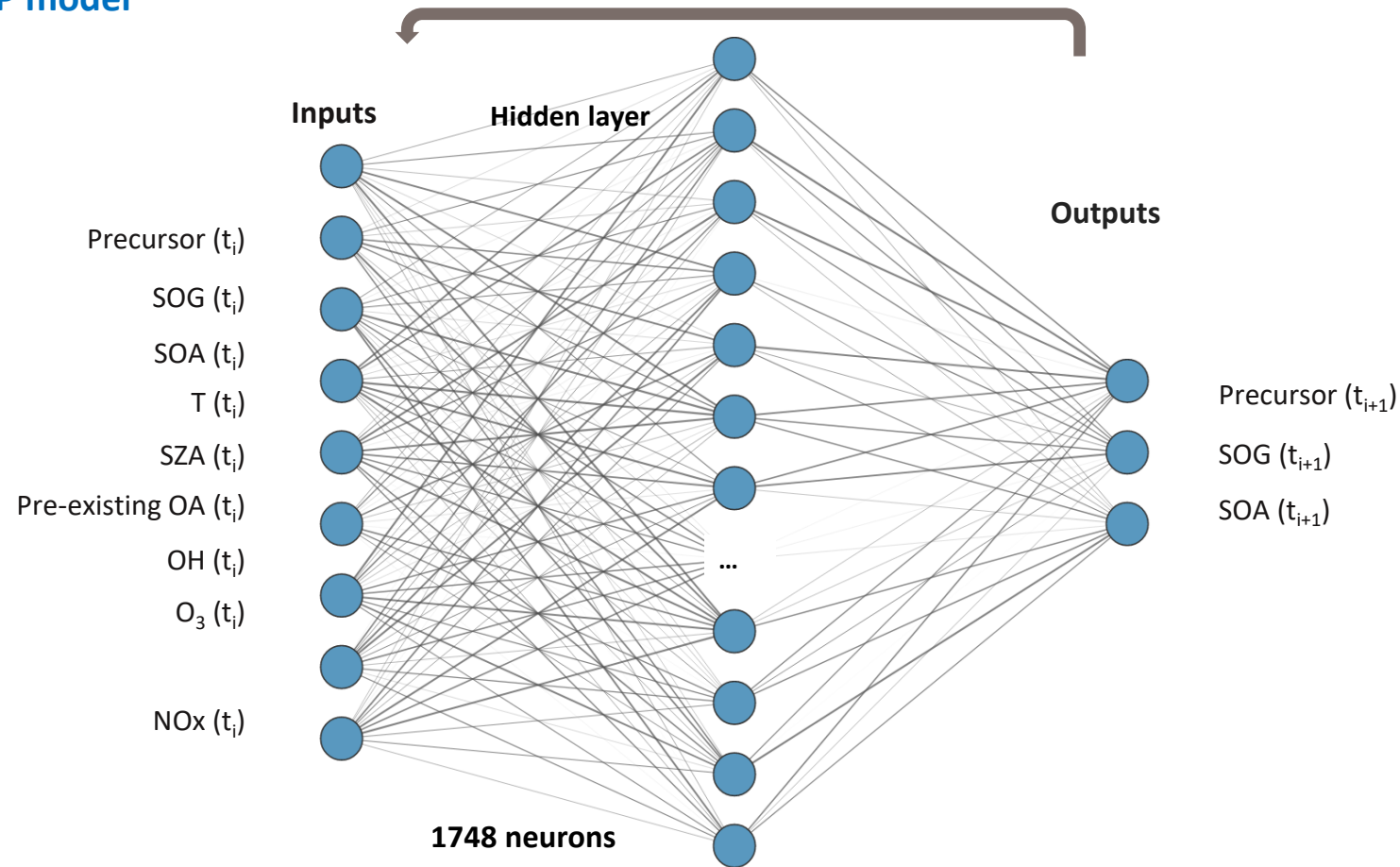


Emulate the behavior of the detailed model with machine learning

- Neural Networks (NN)
- Random Forest (RF)

# NN approach for predicting time-series of concentrations

## 1. Multi-layer perceptron MLP model

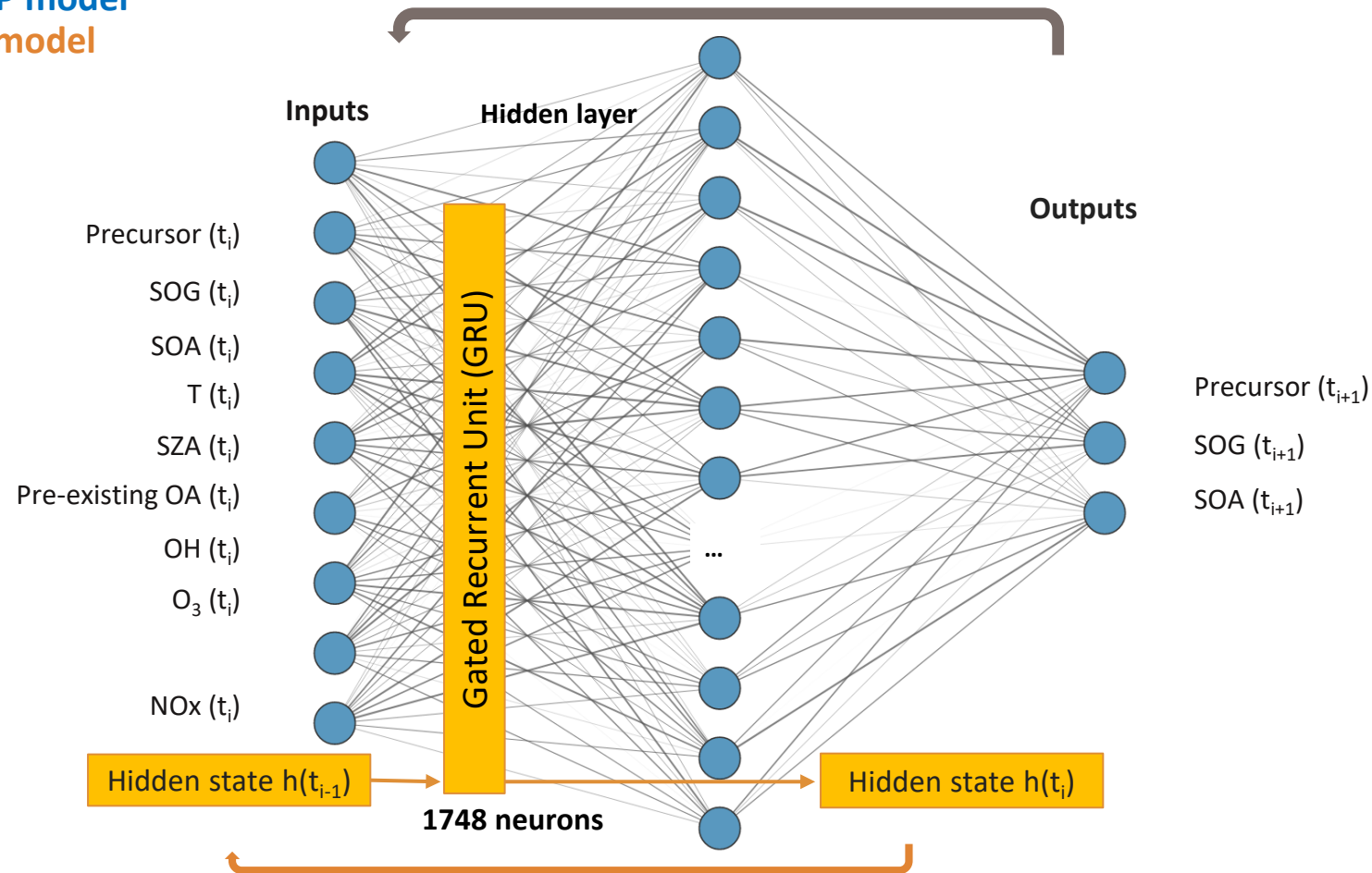


Schreck et al. JAMES 2022: Neural network emulation of the formation of OA based on the explicit GECKO-A chemistry model

# NN approach for predicting time-series of concentrations

## 1. Multi-layer perceptron MLP model

## 2. Gated-recurrent unit GRU model

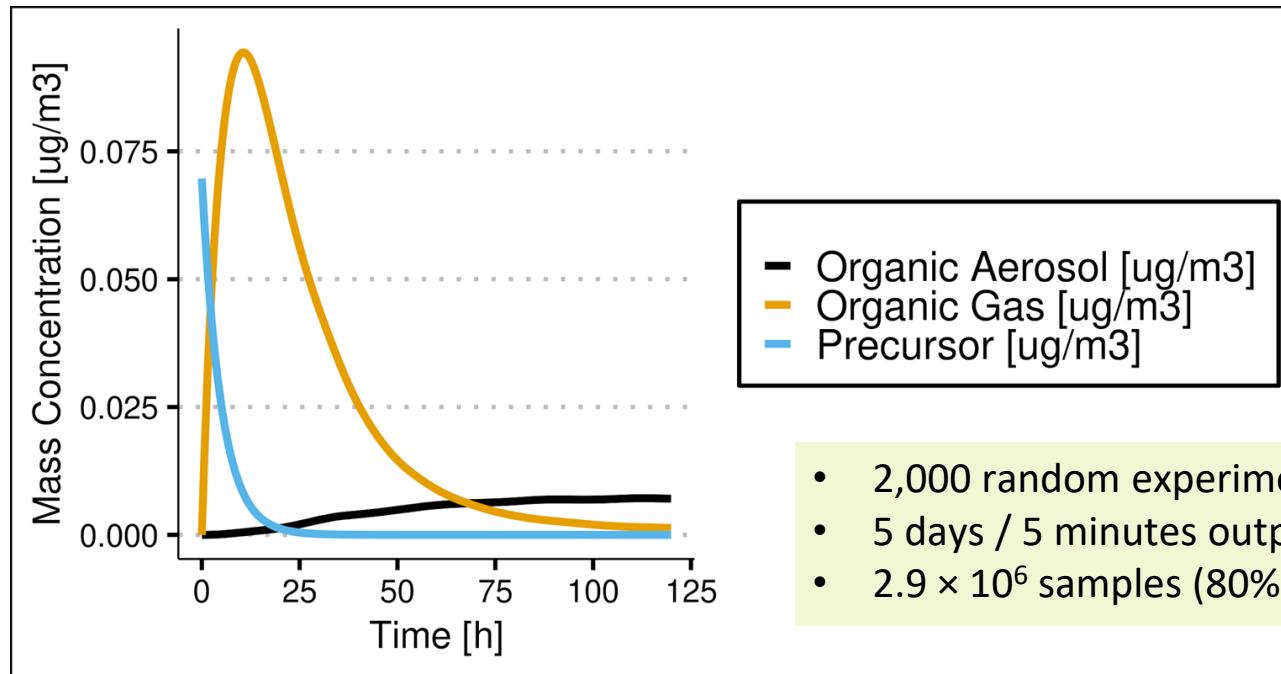


Schreck et al. JAMES 2022: Neural network emulation of the formation of OA based on the explicit GECKO-A chemistry model

# Generate Neural Network training dataset with GECKO-A

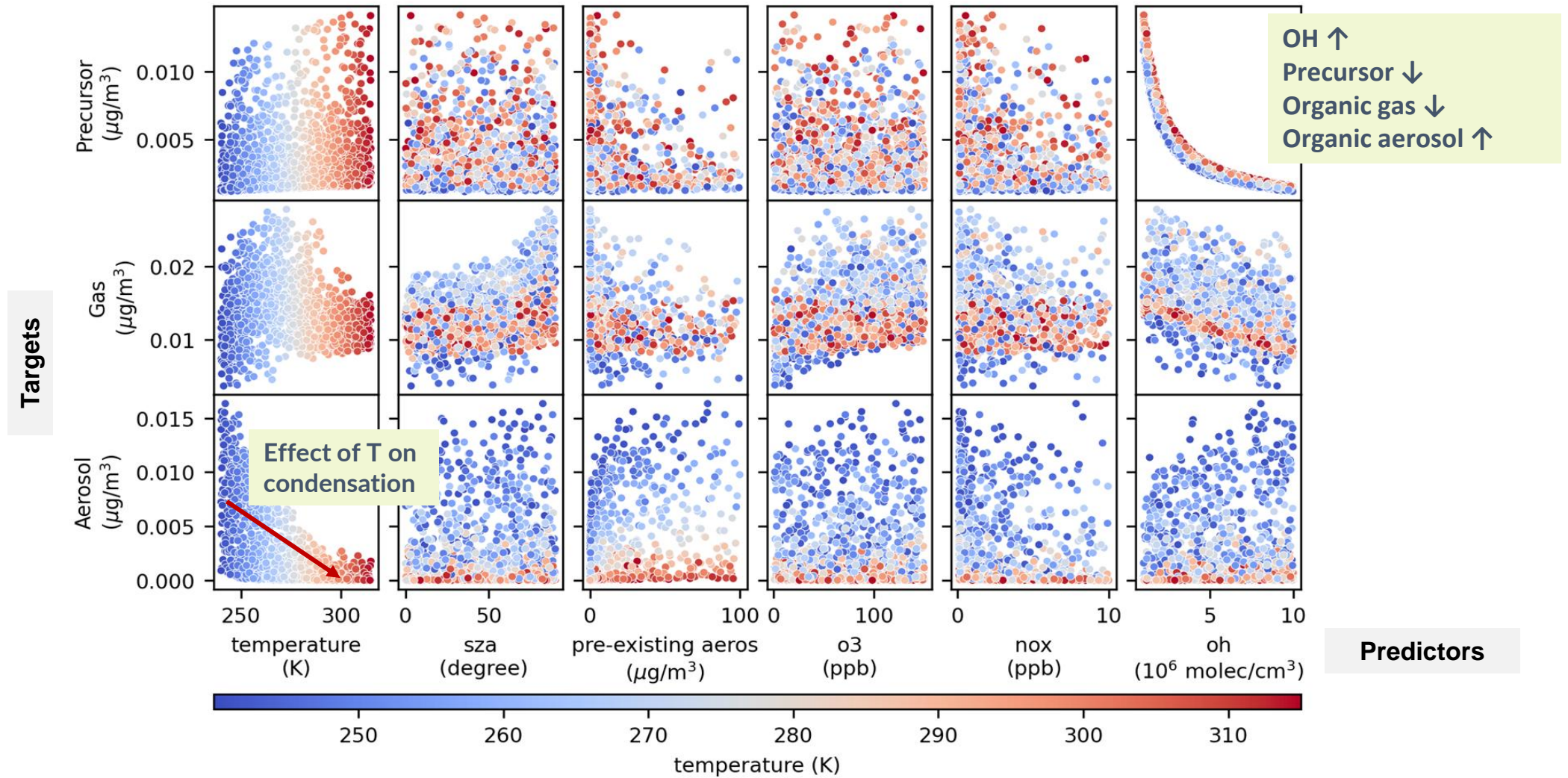
- 3 precursors: **toluene, dodecane, a-pinene**
- Random environmental conditions
- No diurnal variations
- Initial precursor amounts of 10 ppt, 0.1 and 1 ppb

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
NOx	0.01-10 ppb	Logarithmic
OH	$10^1$ - $10^6$ molecules/ $\text{cm}^3$	Uniform



- 2,000 random experiments for each precursor
- 5 days / 5 minutes output
- $2.9 \times 10^6$  samples (80% training, 20% validation)

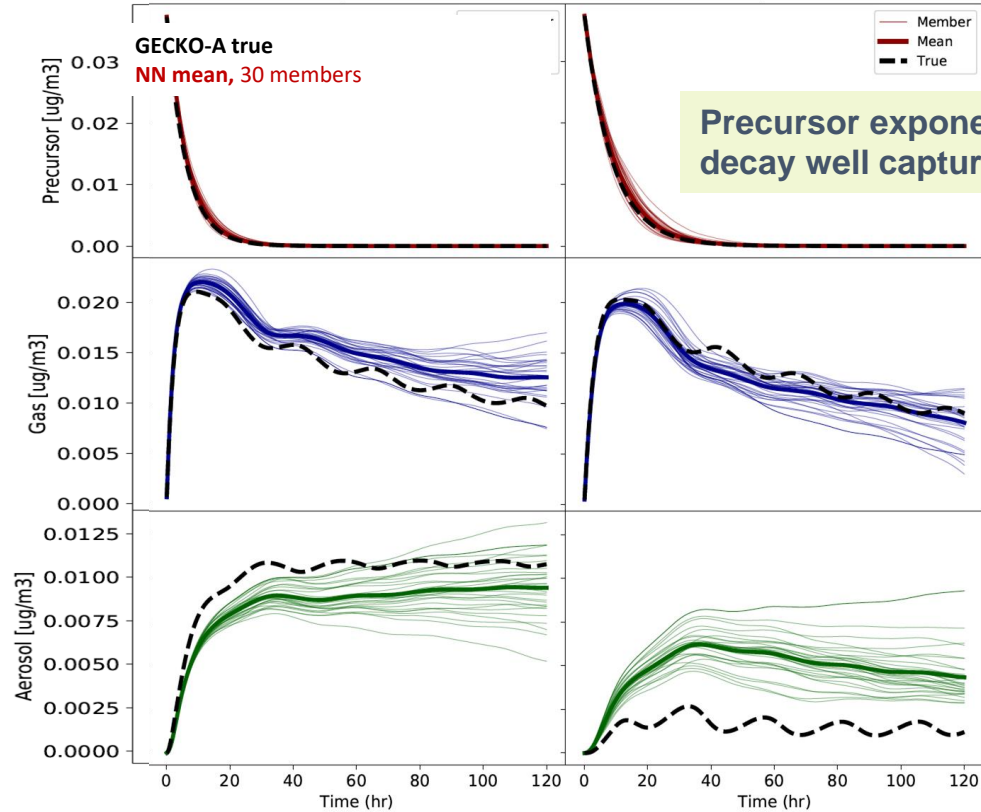
# Neural Network training dataset



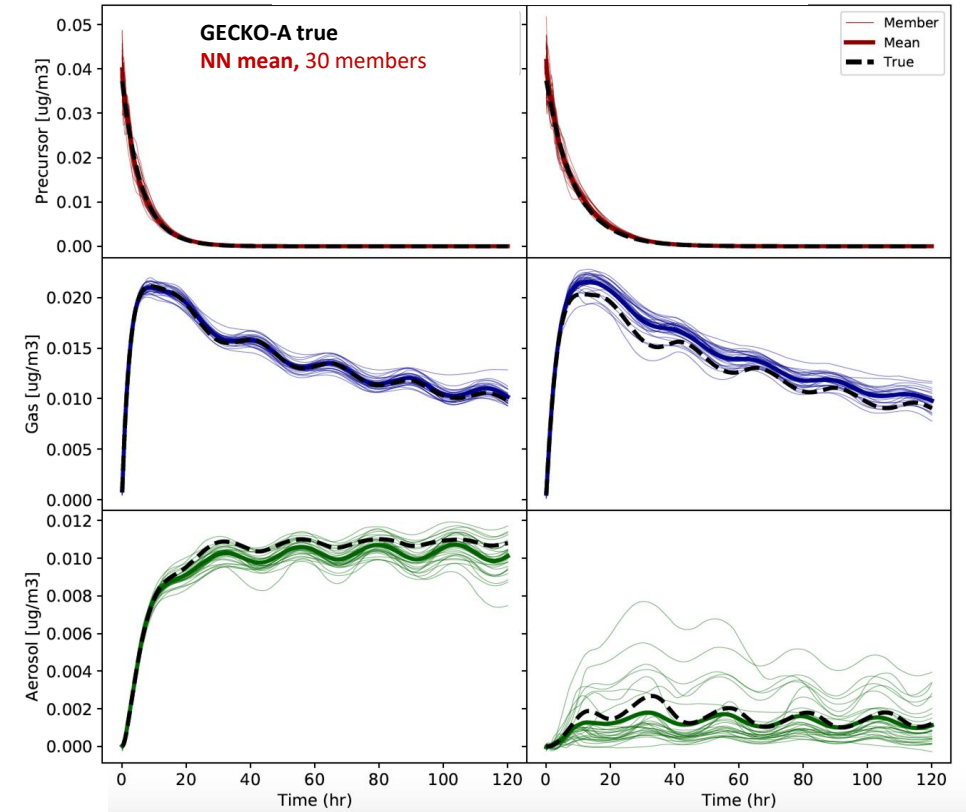
# Multi-Layer Perceptron vs. Gated-Recurrent Unit network

30-ensemble member predictions for Toluene

Multi-layer perceptron



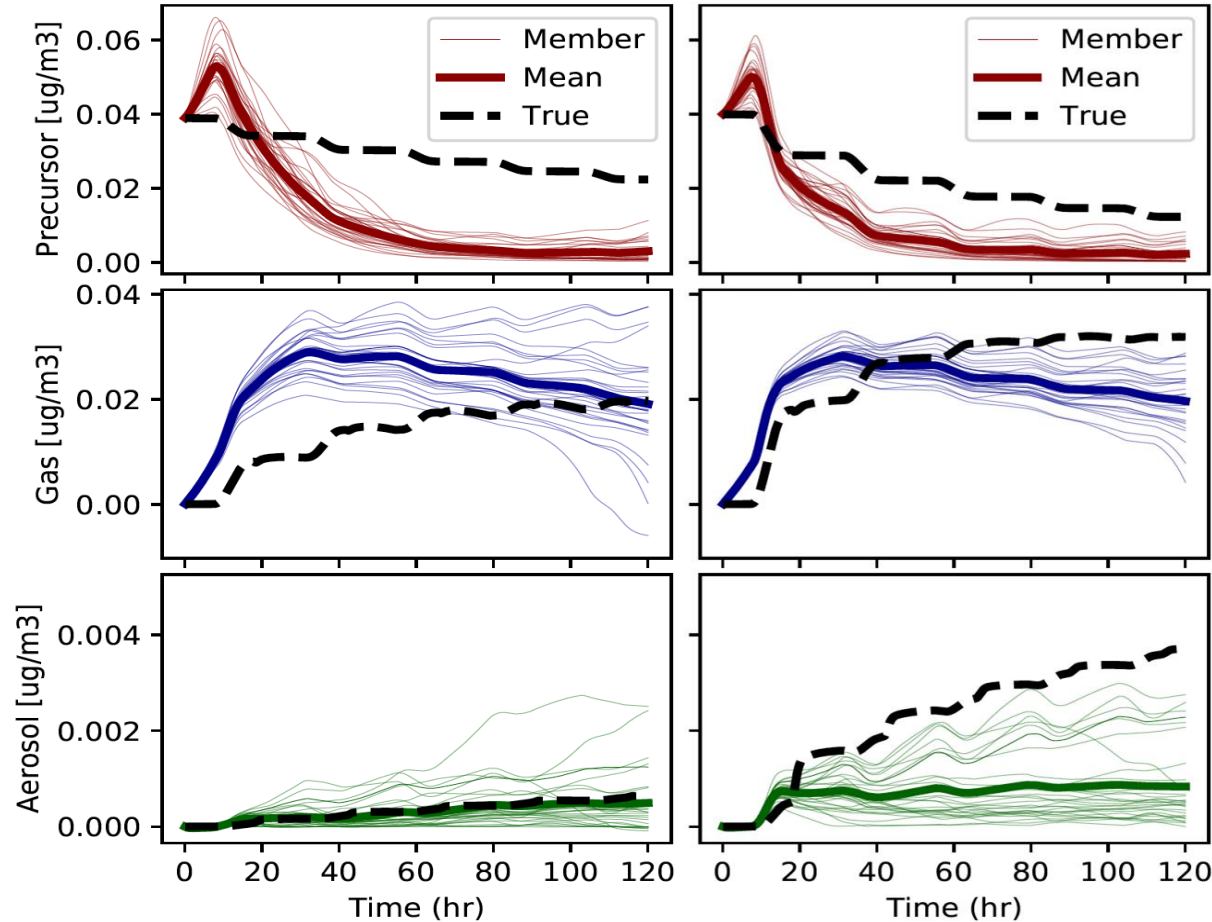
Gated-recurrent unit



⇒ GRU performs better but is challenging to implement in a 3D model

# Application in a box model for diurnally varying conditions

GECKO-A vs. GECKO-MLP



⇒ Neural Networks trained on datasets built under constant conditions cannot reproduce realistic diurnal cycles

⇒ Raises the question of the representativeness of the training data sets

# Computational gain GECKO-A vs. GECKO-NN

- For Toluene, GECKO-NN is  $4 \times 10^2$  times faster than GECKO-A on CPU, and  $10^4$  times faster on GPU

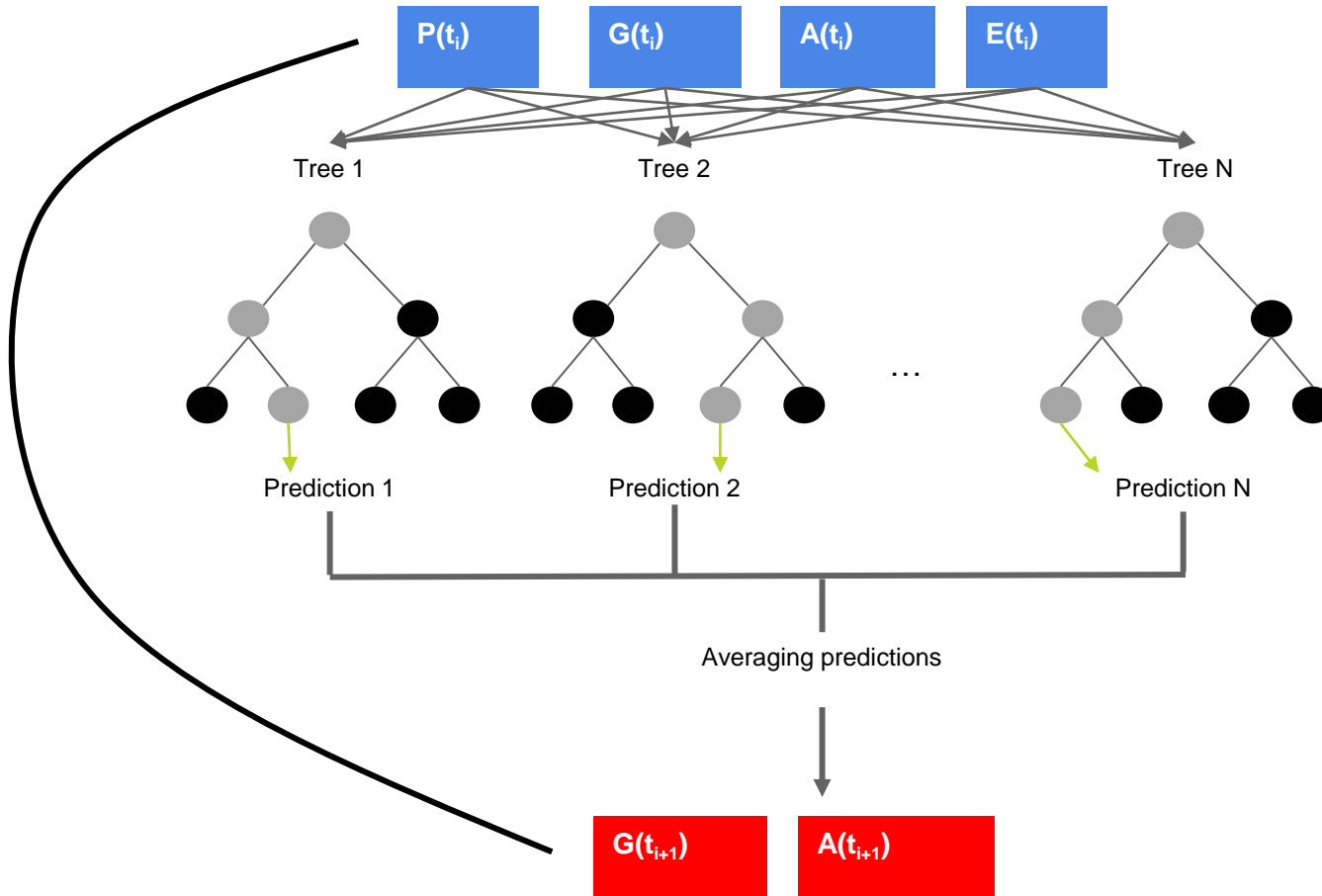
Model	Toluene		Dodecane		$\alpha$ -pinene	
	GECKO-A	GECKO-NN	GECKO-A	GECKO-NN	GECKO-A	GECKO-NN
GECKO-A	0.9 s	1	71 s	1	220 s	1
MLP CPU	2.1 $\mu$ s	430	0.8 $\mu$ s	$8.88 \times 10^4$	1.6 $\mu$ s	$1.38 \times 10^5$
MLP GPU	0.08 $\mu$ s	11250	0.07 $\mu$ s	$1.01 \times 10^6$	0.08 $\mu$ s	$2.75 \times 10^6$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

Timing per  
5min timestep

Ratio of GECKO-A / GECKO-NN



# Random Forest approach



## Inputs :

- Precursor concentration  $P(t_i)$
- Gas organic concentration  $G(t_i)$
- Aerosol organic concentration  $A(t_i)$
- Temperature
- Solar zenith angle
- Pre-existing aerosol
- $O_3$
- $NO_x$
- $OH$

$E(t_i)$

## Output :

- Gas organic concentration  $G(t_{i+1})$
- Aerosol organic concentration  $A(t_{i+1})$

*Mouchel-Vallon & Hodzic, JGR, 2023:*

*Toward emulating an explicit organic chemistry mechanism with random forest models.*

# Random Forest: Realistic training dataset

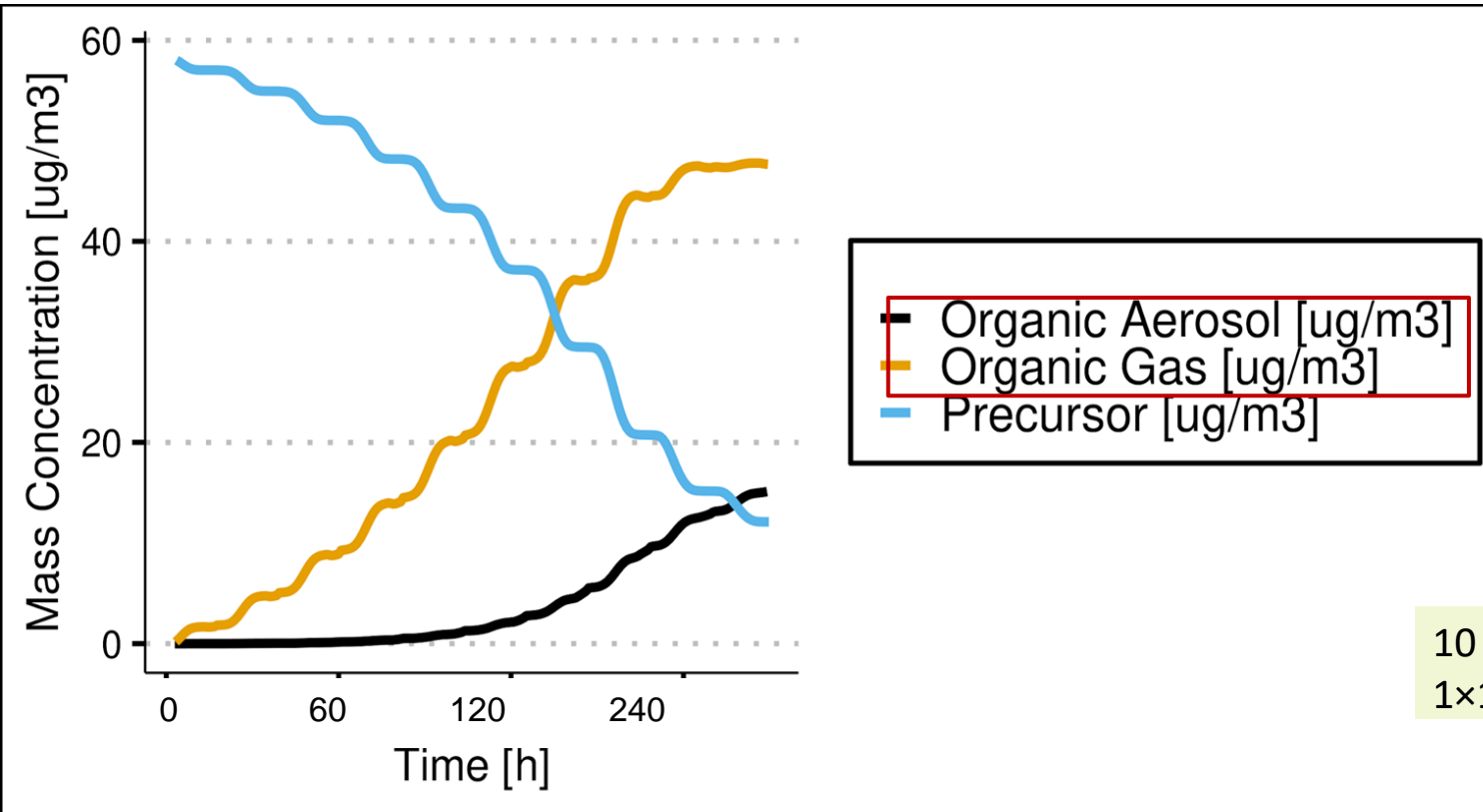
- 2 precursors: **toluene, dodecane**
- Random environmental conditions
- With diurnal variations
- Random initial precursor concentration

Parameter	Range	Parameter	Range
Latitude [°]	80S–80N	Relative humidity [%]	3–102
Temperature [K]	216–313	Atmospheric pressure [atm]	0.5–1.02
Preexisting aerosol seed [ $\mu\text{g m}^{-3}$ ]	0.03–340	Initial $\text{NO}_x$ [ppb]	$10^{-4}$ –42
Initial precursor [ppb]	0–16	Initial CO	33–1012
Initial $\text{O}_3$ [ppb]	1–100	NO Emission [ $\text{molec cm}^{-2} \text{s}^{-1}$ ]	$10^7$ – $10^9$

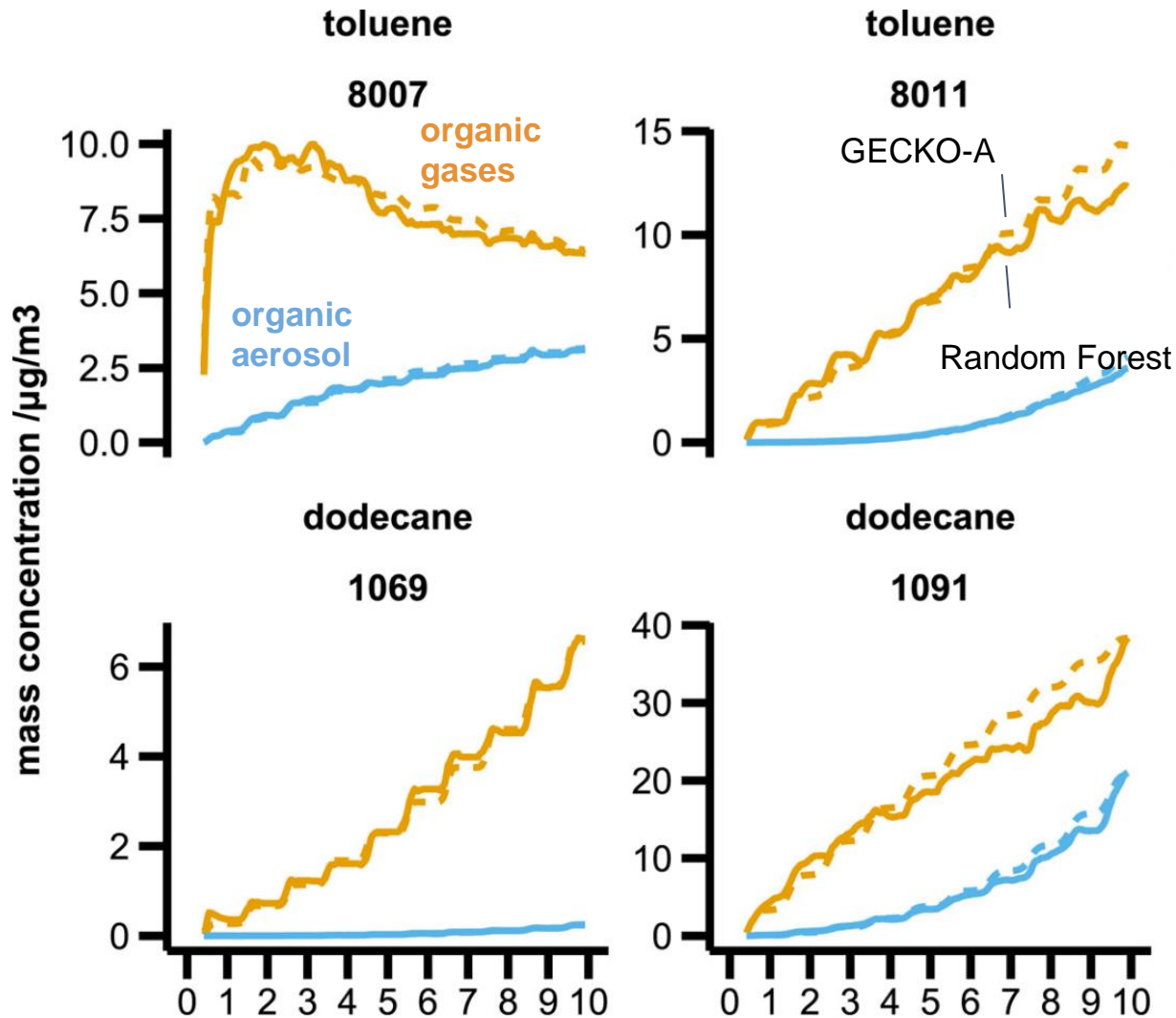
From 3D model  
tropospheric ranges

Only 2 targets

10 days / 5 min output  
 $1 \times 10^6$  samples/precursor

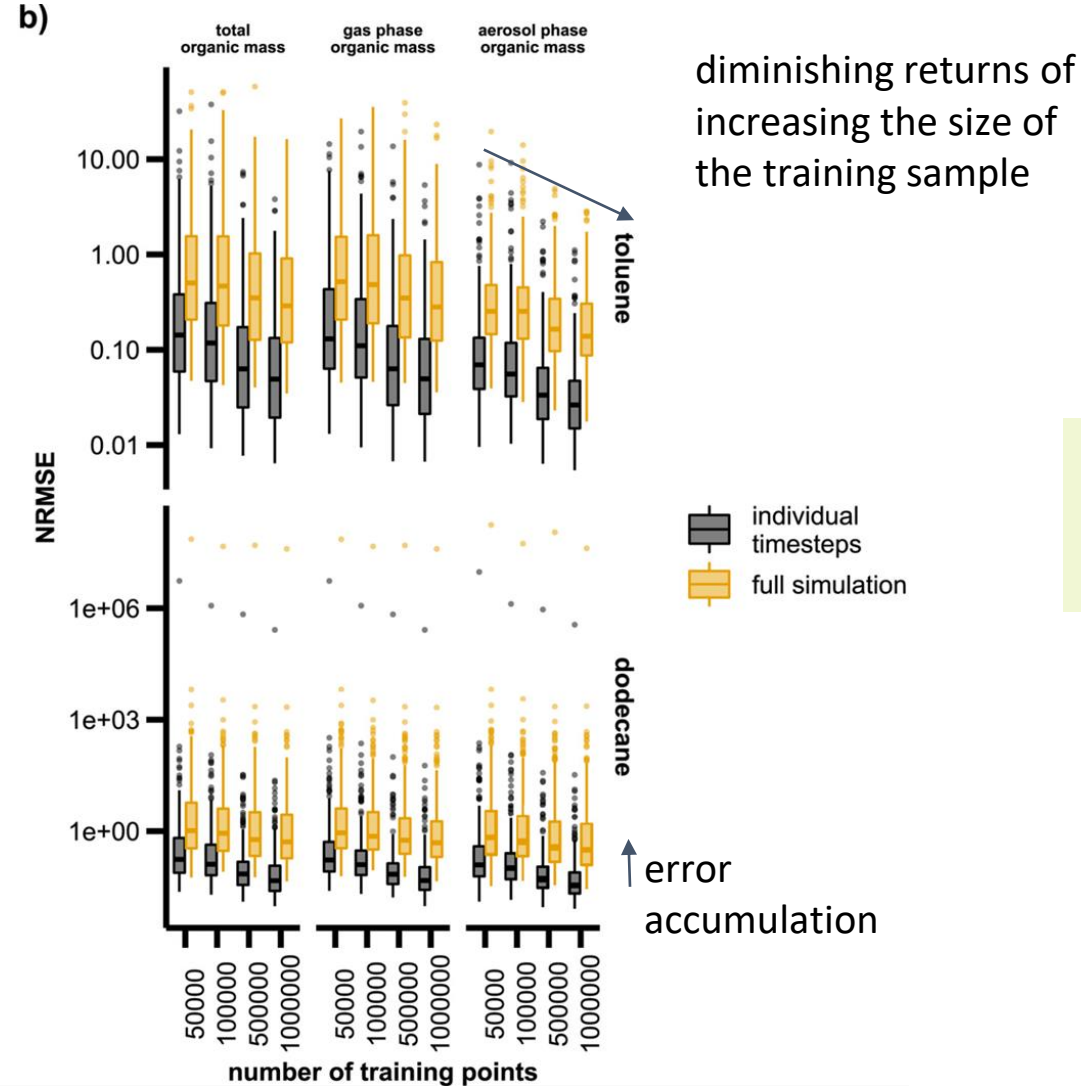
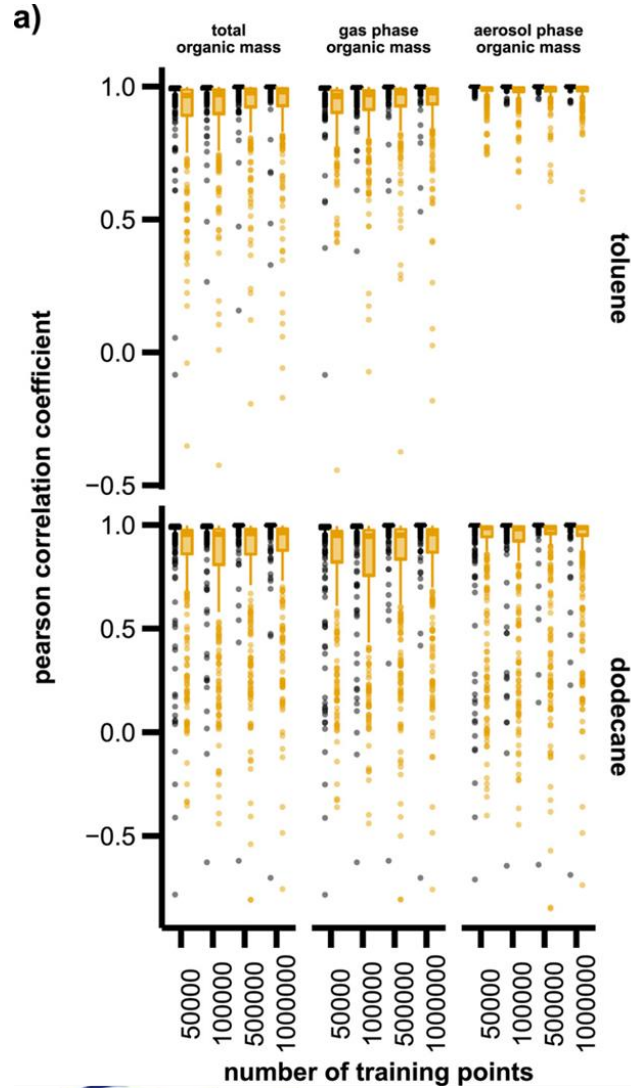


# Random Forest results: examples



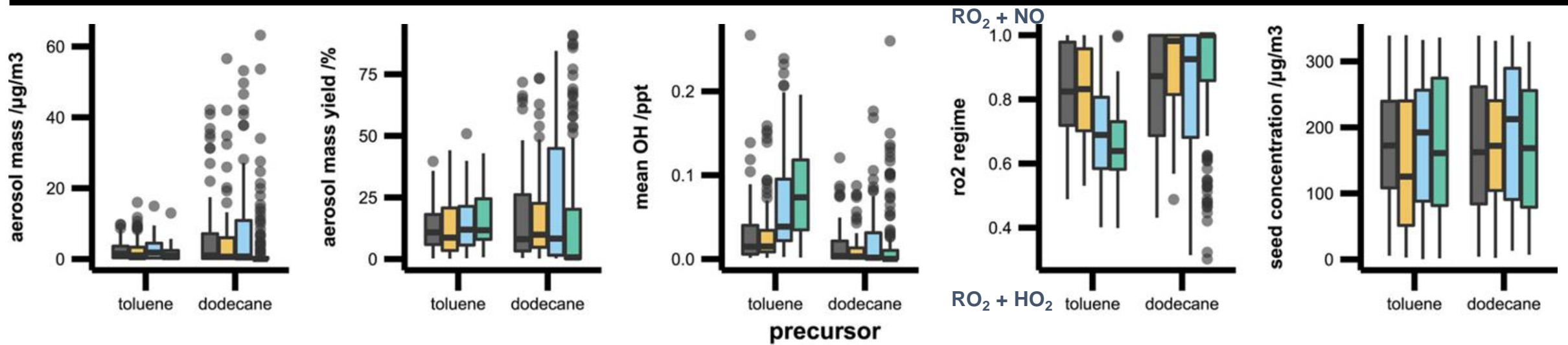
RF perform similarly to NN  
The training set allows  
reproducing the diurnal cycles

# Random Forest results: performance



Increasing the number of training samples improves performance

# Random Forest results: errors distribution

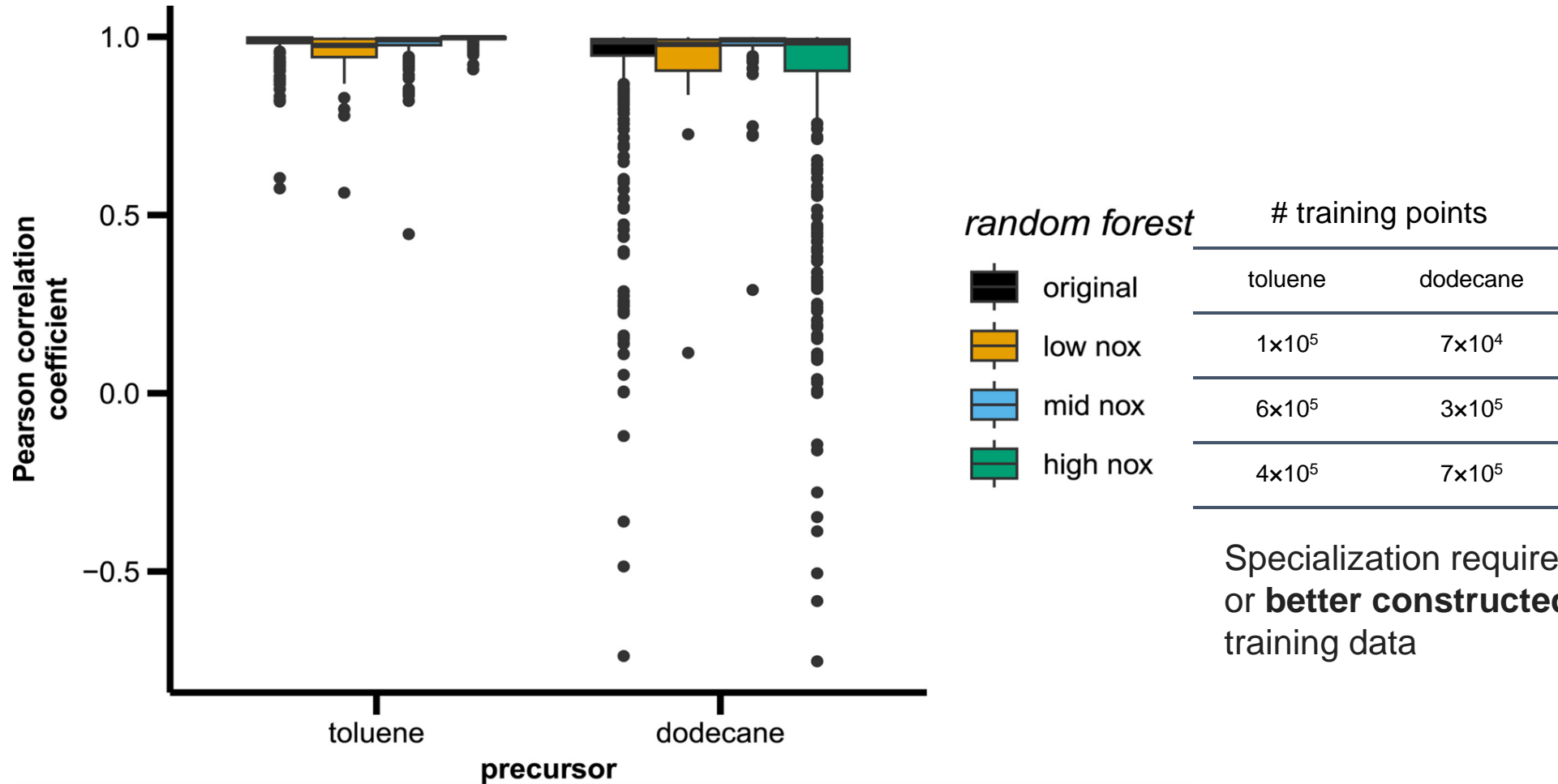


NRMSE quartile 1st 2nd 3rd 4th  
 increasing error

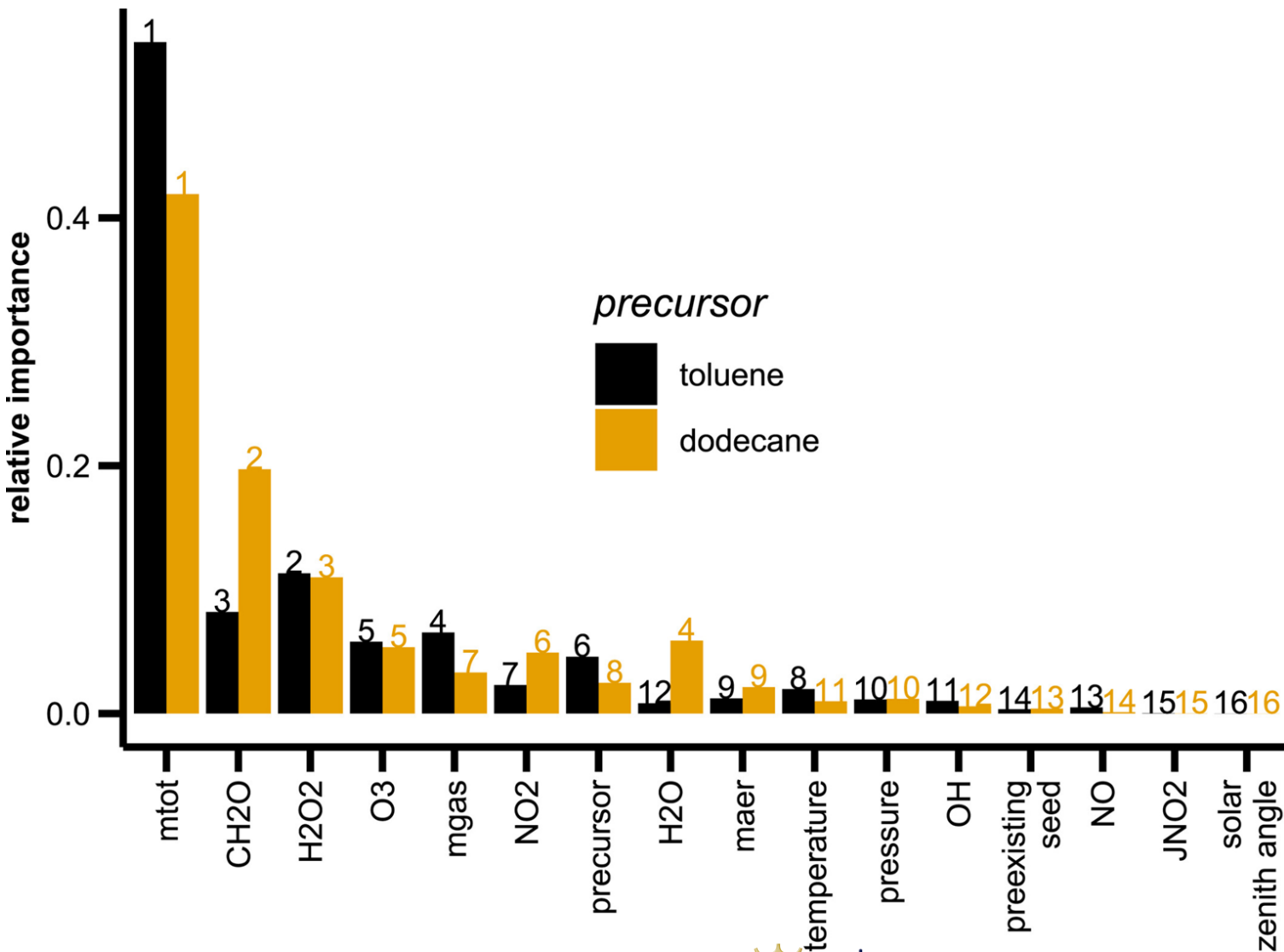
Little sensitivity of errors to the aerosol quantity to predict (seed, yield, mass)  
 Errors sensitive to NO<sub>x</sub> regime and OH mixing ratios: underrepresented regimes in training set

# Random Forest results: specializing random forests

Do performances improve if the range of chemical regimes to predict is reduced?



# Random Forest results: predictors relative importance

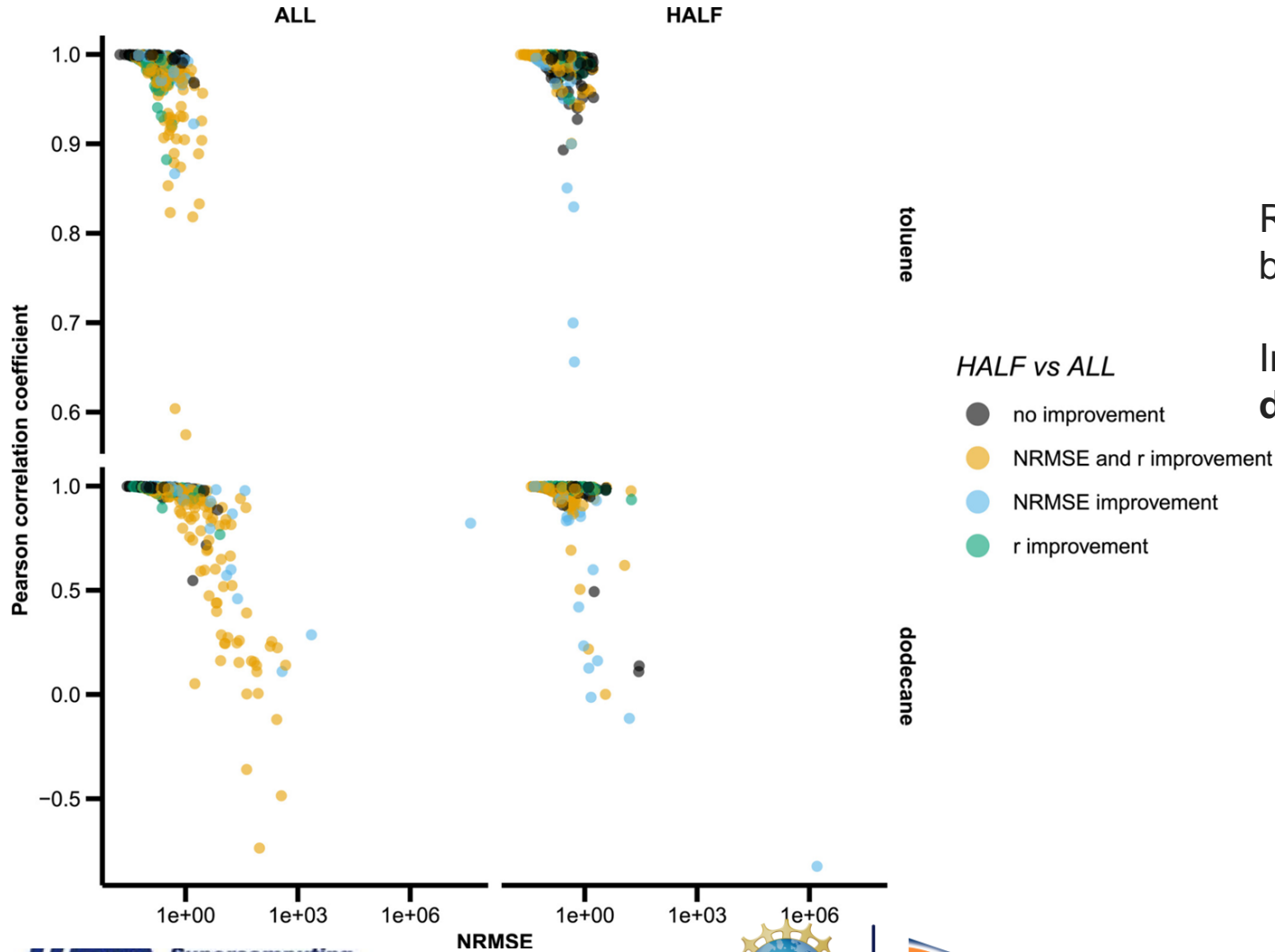


The importance of predictors (mean decrease impurity) is similar for **both** species

Distribution of importance and predictive ability among **correlated variables**

# Random Forest results: reducing the number of predictors

Do performances improve if the number of predictor is reduced?



Reducing the number of (**correlated?**) predictors is beneficial for the worst performing simulations

Importance of **selecting predictors** and **dimensionality reduction**

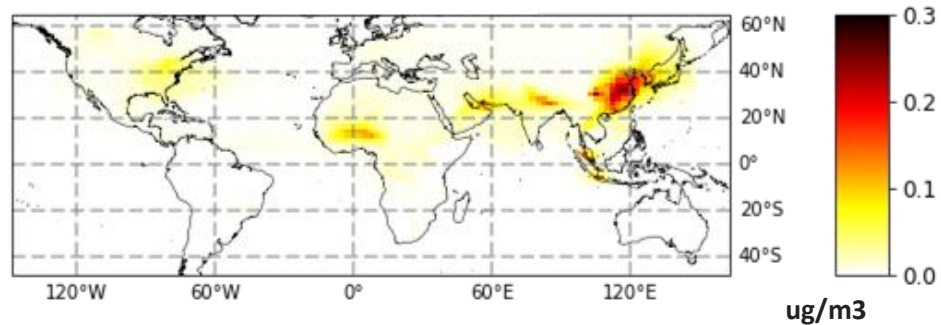


# First tests in a global model: VBS vs. GECKO-NN

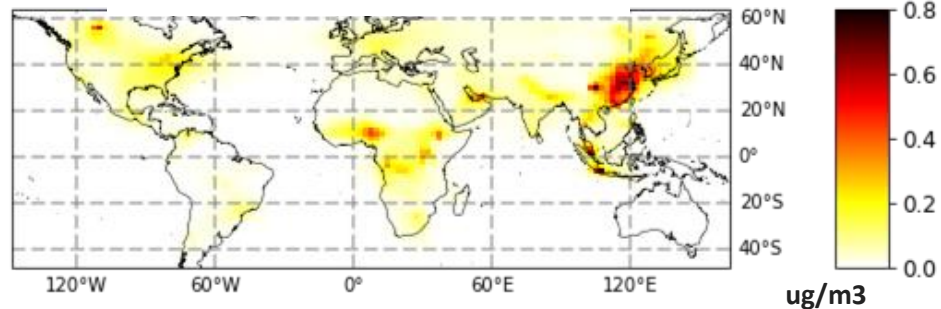
Implementing GECKO-NN in Geos-Chem

## GEOS-Chem Monthly average Toluene SOA (May 2016)

GECKO-NN-MLP



VBS Hodzic et al., 2016



GECKO-NN simulations are

- stable over several months
- within a factor of 2 of the VBS parametrization for Toluene-SOA.

# Conclusions and Outlook

---

- It is possible to emulate the behavior of detailed atmospheric chemistry models with machine learning
- Long term stability can be achieved for recurrent neural networks with GRU
- The training dataset must be carefully constructed to cover all environmental conditions
- Random forests can perform similarly to NN+GRU
- Predictors selection is crucial

## Current and future works at BSC

- Bring the complexity of organic chemistry to air quality models, built on the development of detailed Spanish emissions
- Explore ML use, with lessons learned from this work: start again from the basics and systematic exploration by (i) progressively increasing chemical complexity (from toy mechanisms to GECKO-A complexity) and (ii) testing multiple families of ML techniques (RF, NN, GraphNet ...)
- Implementation in the MONARCH air quality model



**Barcelona  
Supercomputing  
Center**  
*Centro Nacional de Supercomputación*

# Thank you for your attention

[camille.mouchel@bsc.es](mailto:camille.mouchel@bsc.es)



Financiado por  
la Unión Europea  
NextGenerationEU



VICEPRESIDENCIA  
TERCERA DEL GOBIERNO  
MINISTERIO  
PARA LA TRANSICIÓN ECOLÓGICA  
Y EL RETO DEMOGRÁFICO



Agencia Estatal de Meteorología



Plan de Recuperación,  
Transformación  
y Resiliencia