

# Atmospheric Chemistry Modeling using Machine Learning

**Christoph A. Keller**

NASA Global Modeling and Assimilation Office (GMAO)  
Universities Space Research Association (USRA)

**Mat J. Evans**

Wolfson Atmospheric Chemistry Laboratories, University of York  
National Centre for Atmospheric Sciences, University of York

EGU Annual Meeting  
11 April 2019

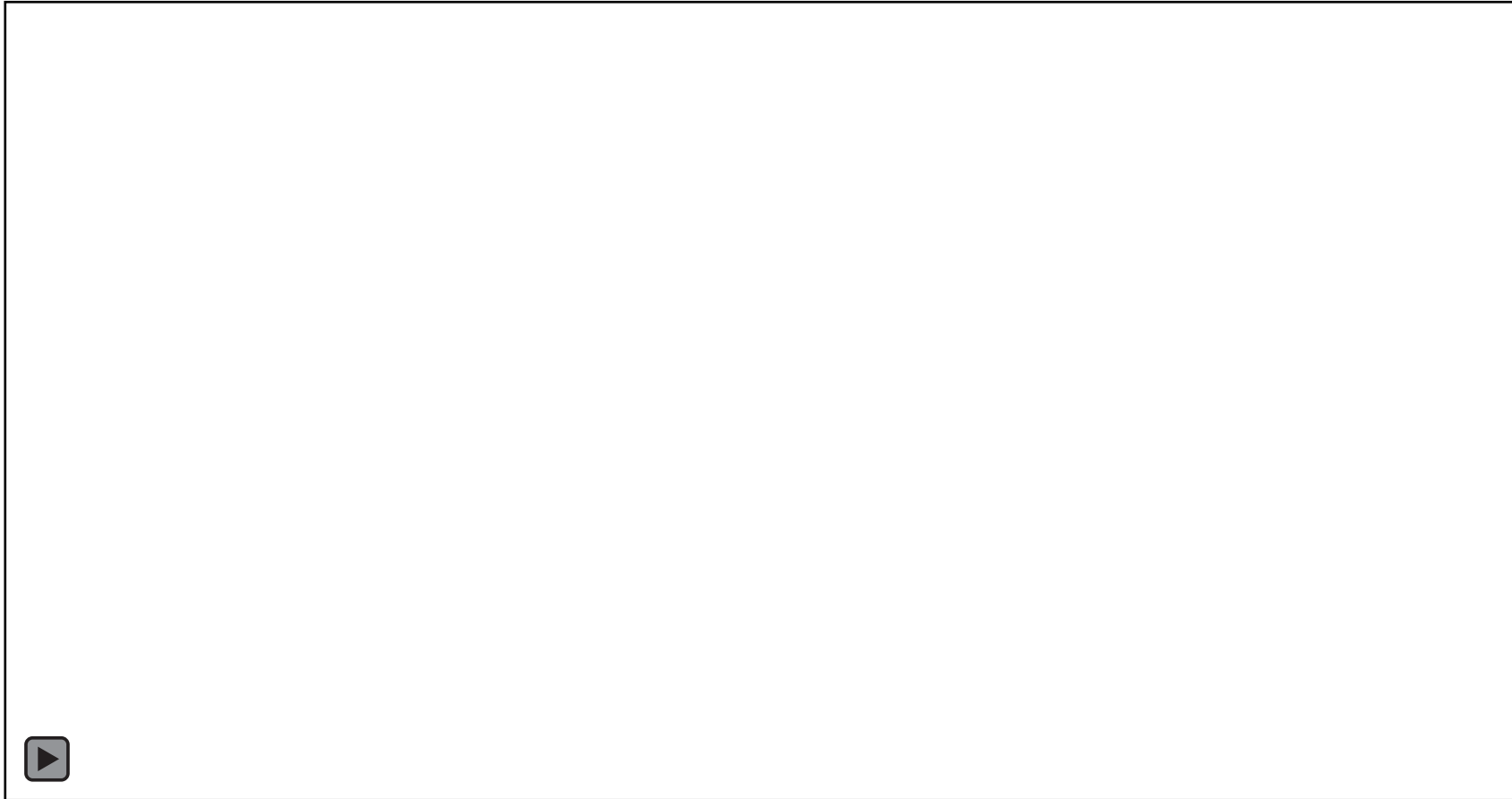


UNIVERSITY  
*of York*



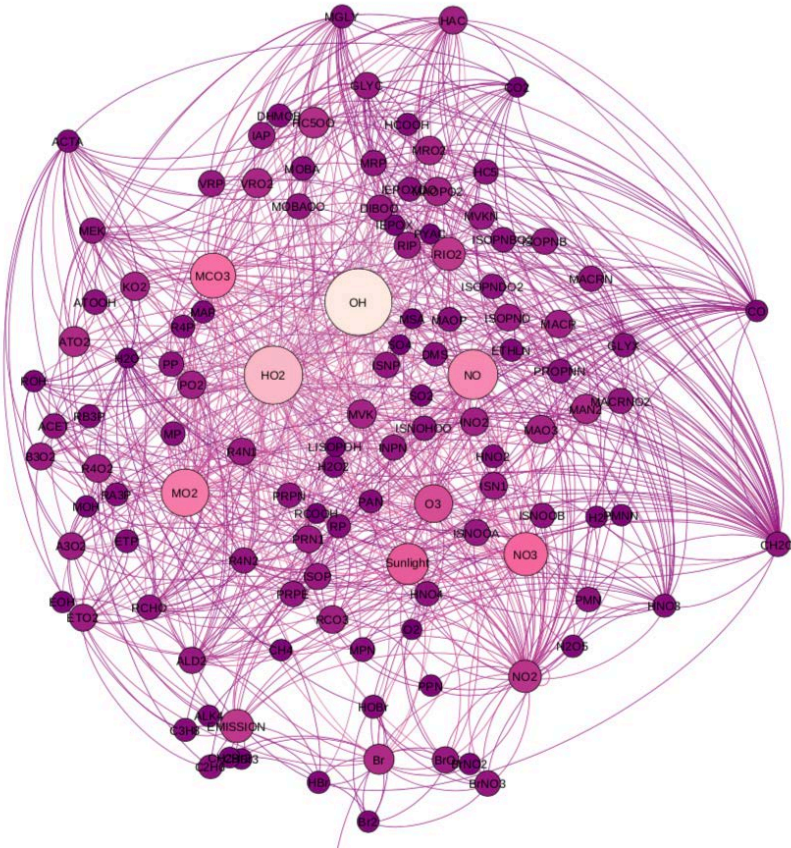
**National Centre for  
Atmospheric Science**  
NATURAL ENVIRONMENT RESEARCH COUNCIL

# Numerical simulation of atmospheric chemistry



- 0.25° resolution (~ 25km), 72 levels, 250 chemical species

# Numerical simulation of atmospheric chemistry



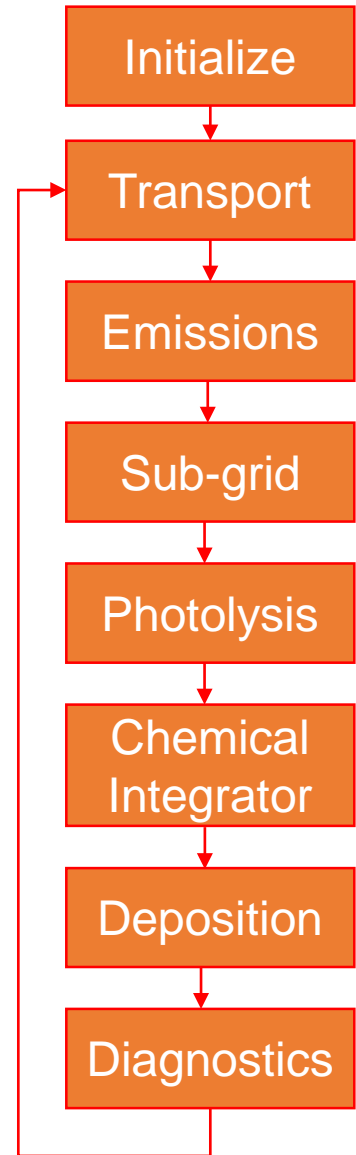
**Transport process:** Move chemicals across grid boxes

**Chemistry process:** In each grid box, solve chemical reactions, i.e. solve stiff ordinary differential equations (ODEs)

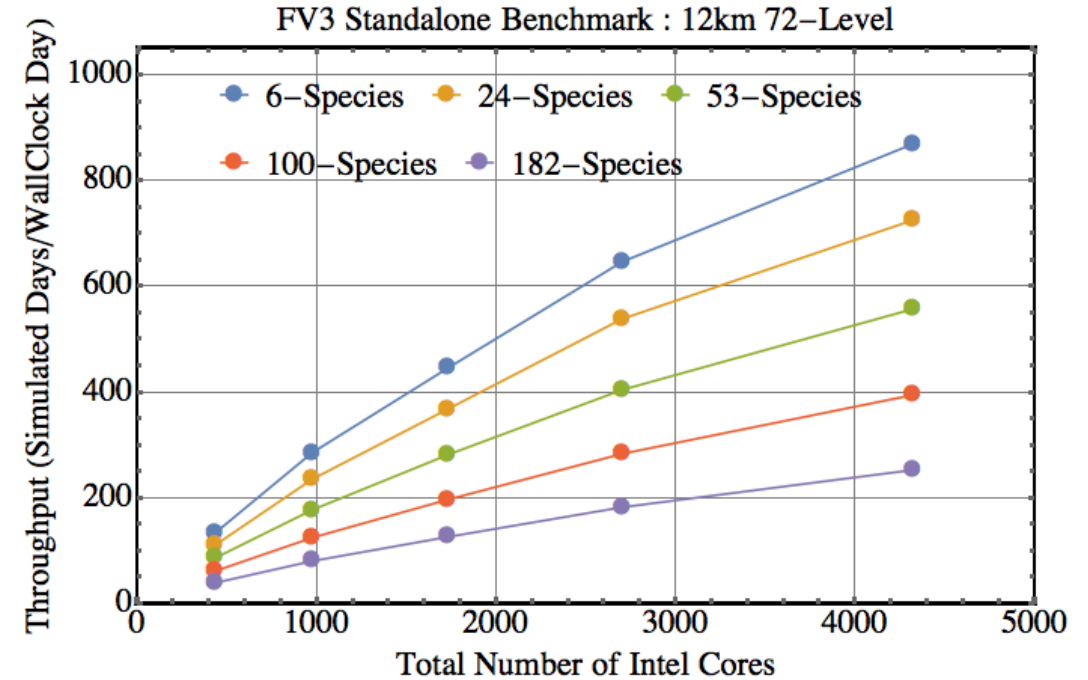
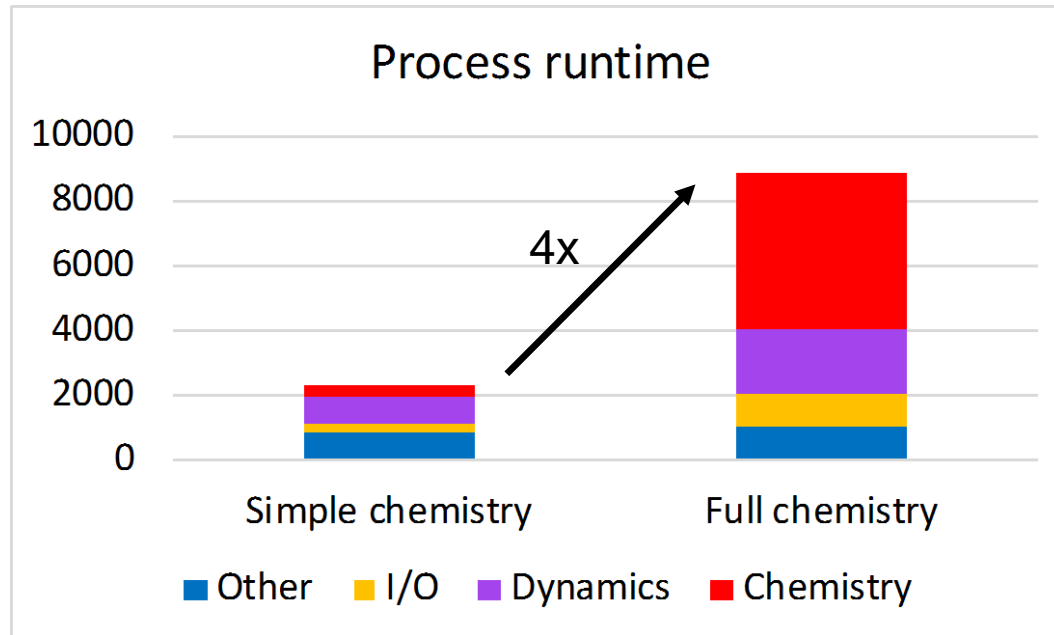


its rate is calculated as

$$-\frac{d}{dt}[A] = -\frac{d}{dt}[B] = \frac{d}{dt}[C] = \frac{d}{dt}[D] = k[A][B]$$



# The current solution: wait, wait, wait



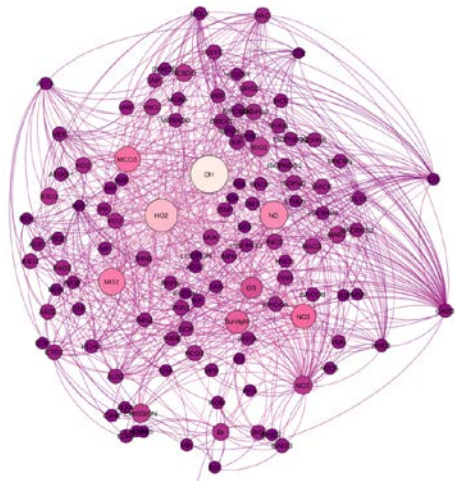
Courtesy of W. Putman, NASA GMAO

- High-resolution chemistry simulation requires 3416 CPU's
- Can simulate approx. 20 days in 24 hours
- Outputting the full chemical state is 1.5 TB / simulation day

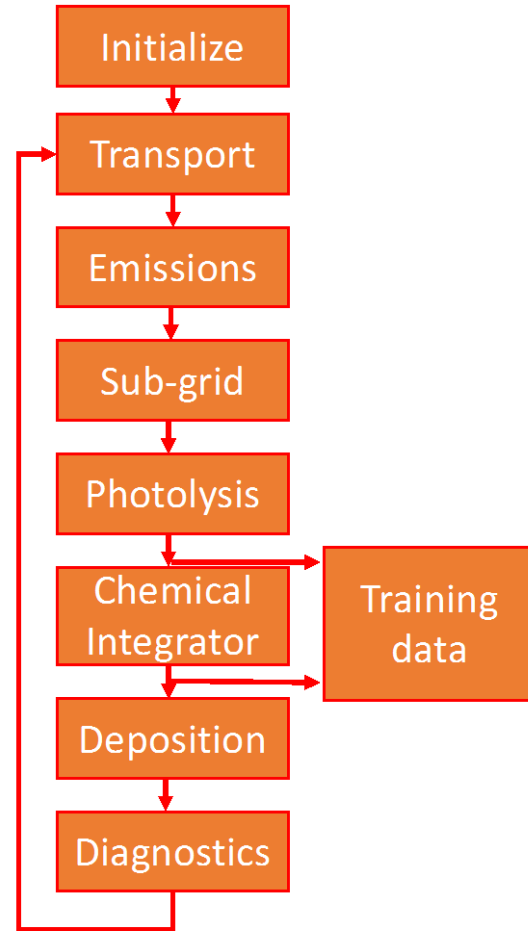




# Replace chemical integrator with machine learning model

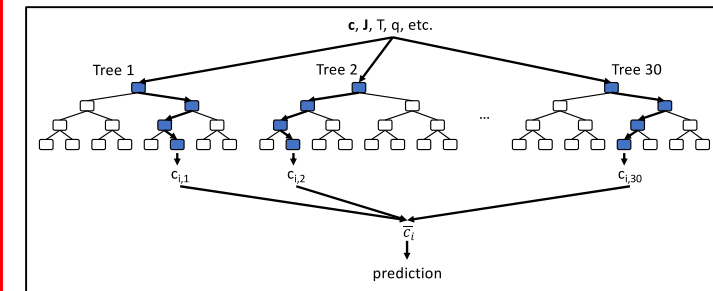
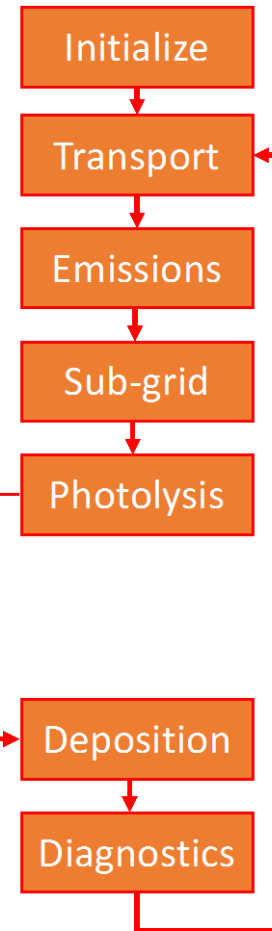


## Numerical model



Random forest

## Emulator



# Machine learning for atmospheric chemistry modeling

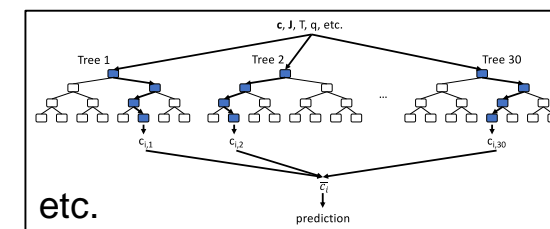
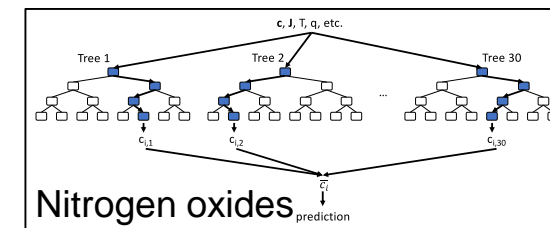
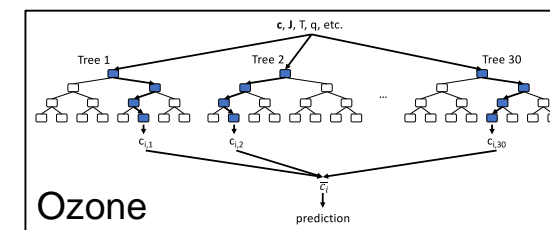
- 143 chemical species
- 91 photolysis rates
- Temperature
- Pressure
- Rel. humidity
- Solar zenith angle



Concentrations  
after chemistry

- Training data set has 2.7 billion data points (44 GB)
- Tested: (neural network), random forest and XGBoost

Separate model  
for each species



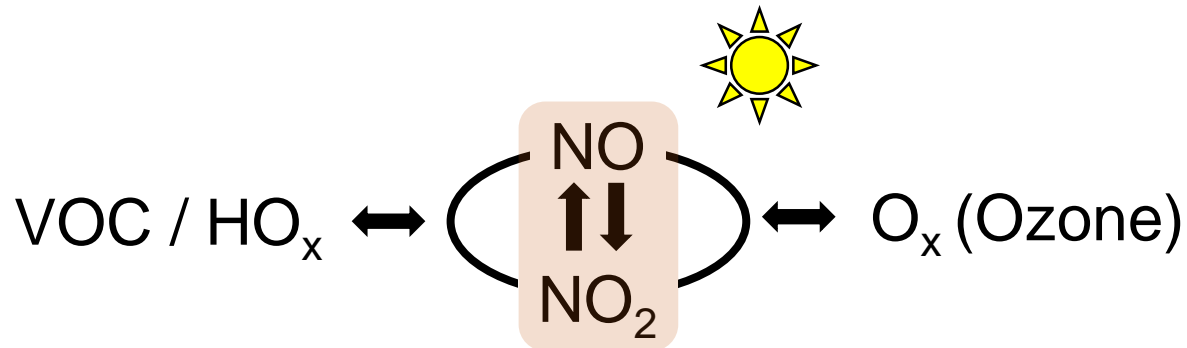
# Impose chemical constraints on ML model to improve (long-term) accuracy

## 1. Distinguish between short-term vs. long-term species

Long-lived (tendencies):  $[X_i]_{T+\Delta T} = [X_i]_T + f(\mathbf{k}, \mathbf{J}, [\mathbf{X}])$

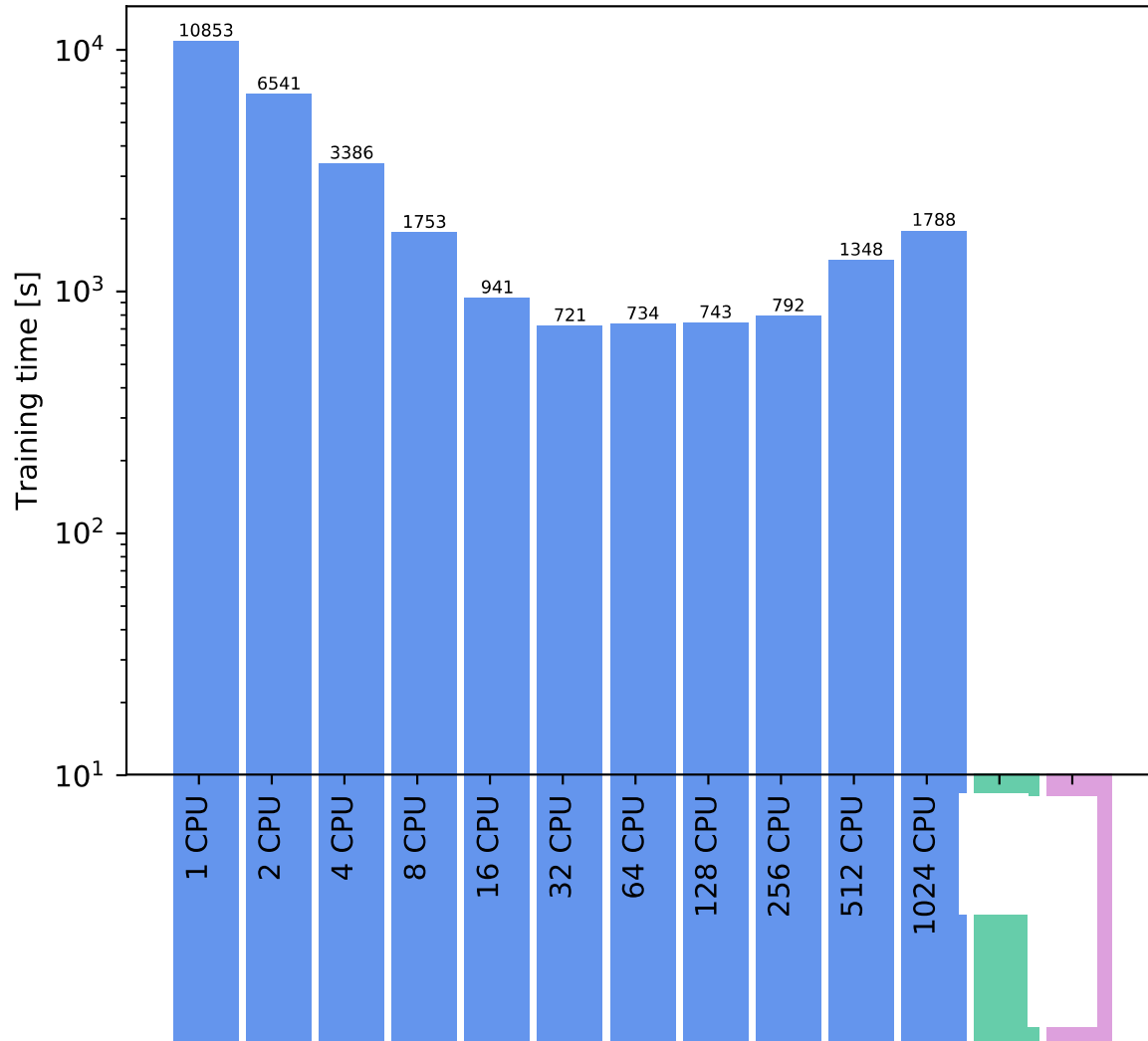
Short-lived (steady state):  $[X_i]_{T+\Delta T} = f(\mathbf{k}, \mathbf{J}, [\mathbf{X}])$

## 2. Predict NO + NO<sub>2</sub> combined (NO<sub>x</sub> family approach)



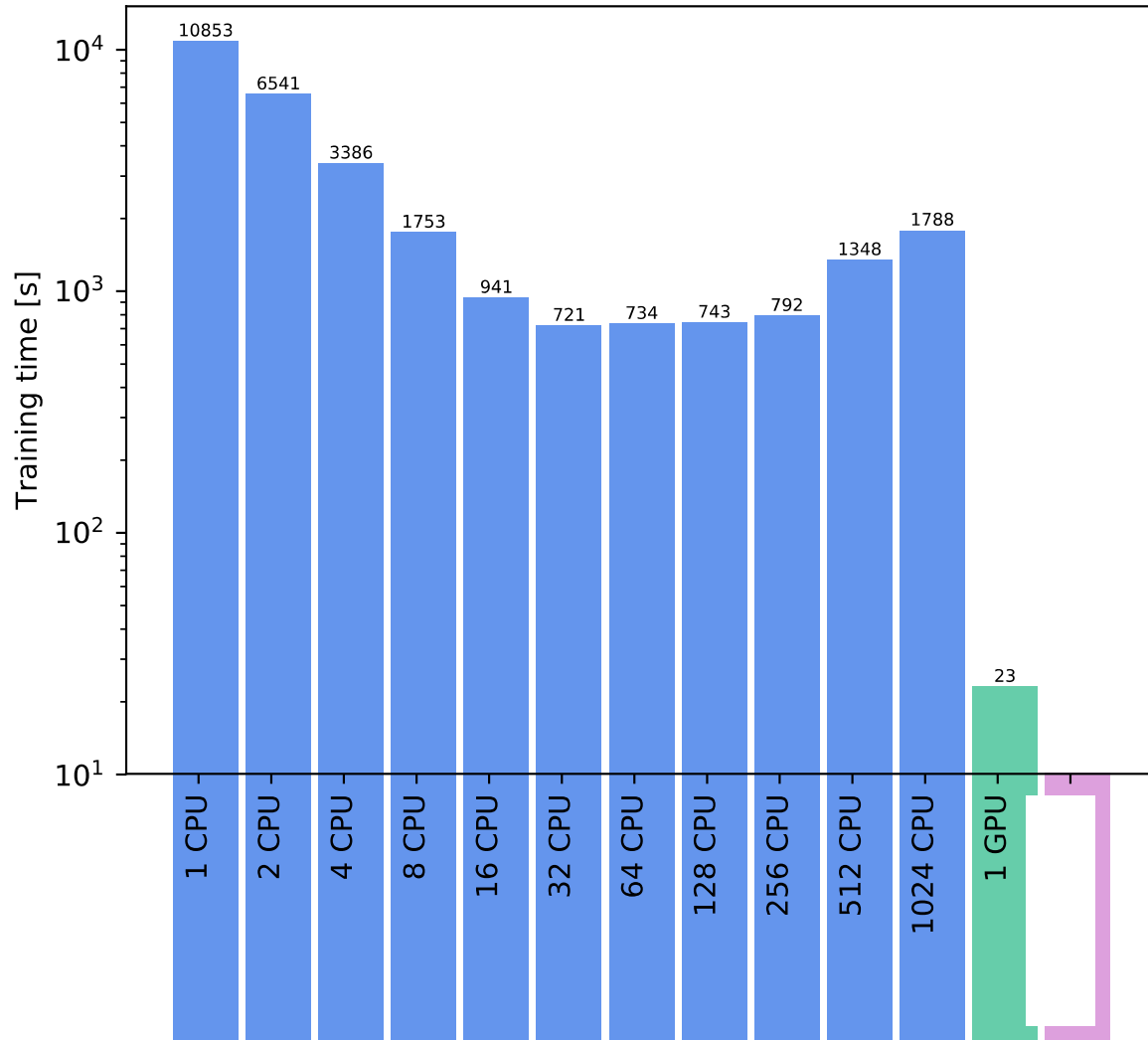
# Random forest / XGBoost training benchmarks

Comparison of XGBoost training time (data set = 44 GB)



# Random forest / XGBoost training benchmarks

Comparison of XGBoost training time (data set = 44 GB)

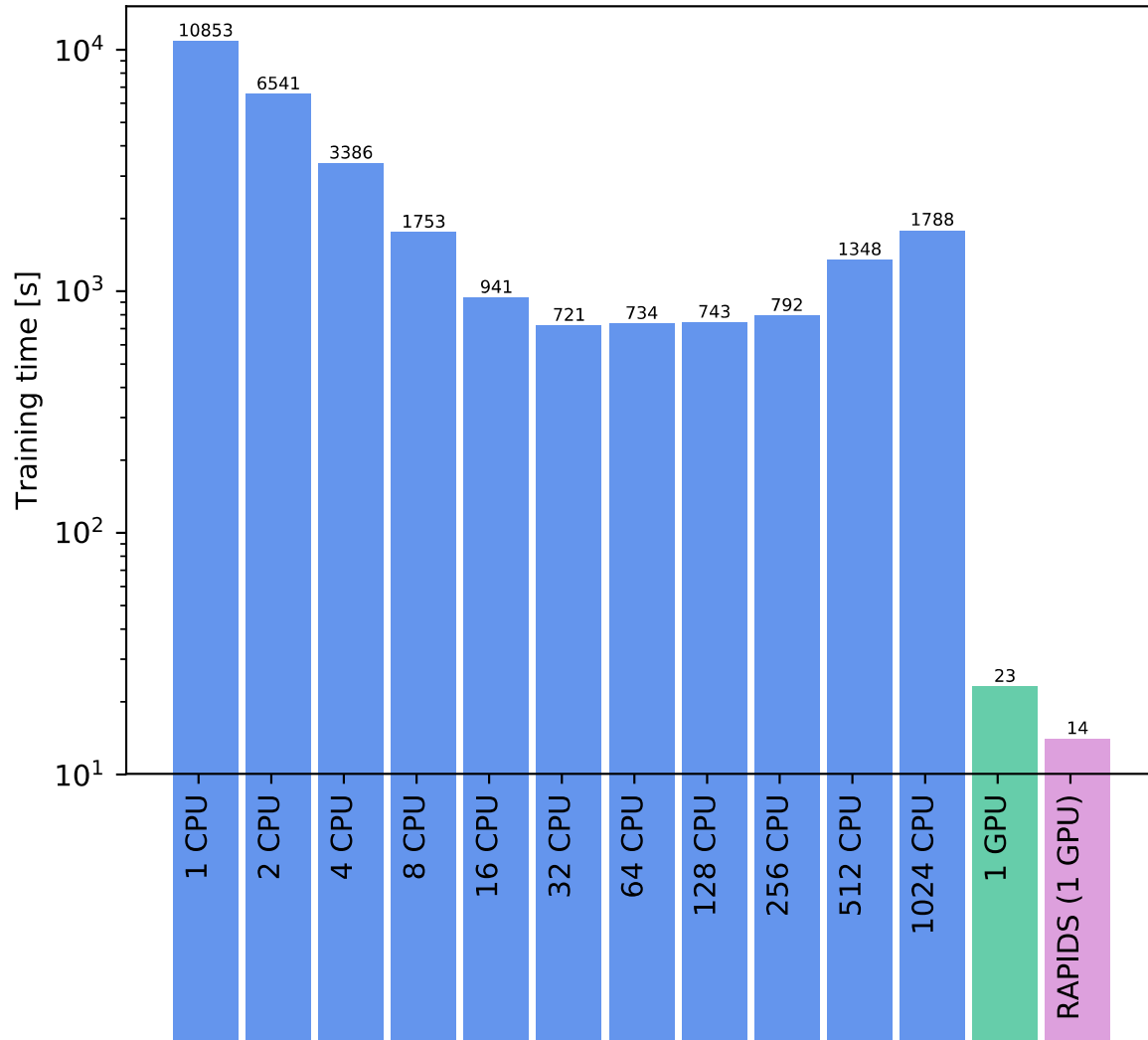






# Random forest / XGBoost training benchmarks

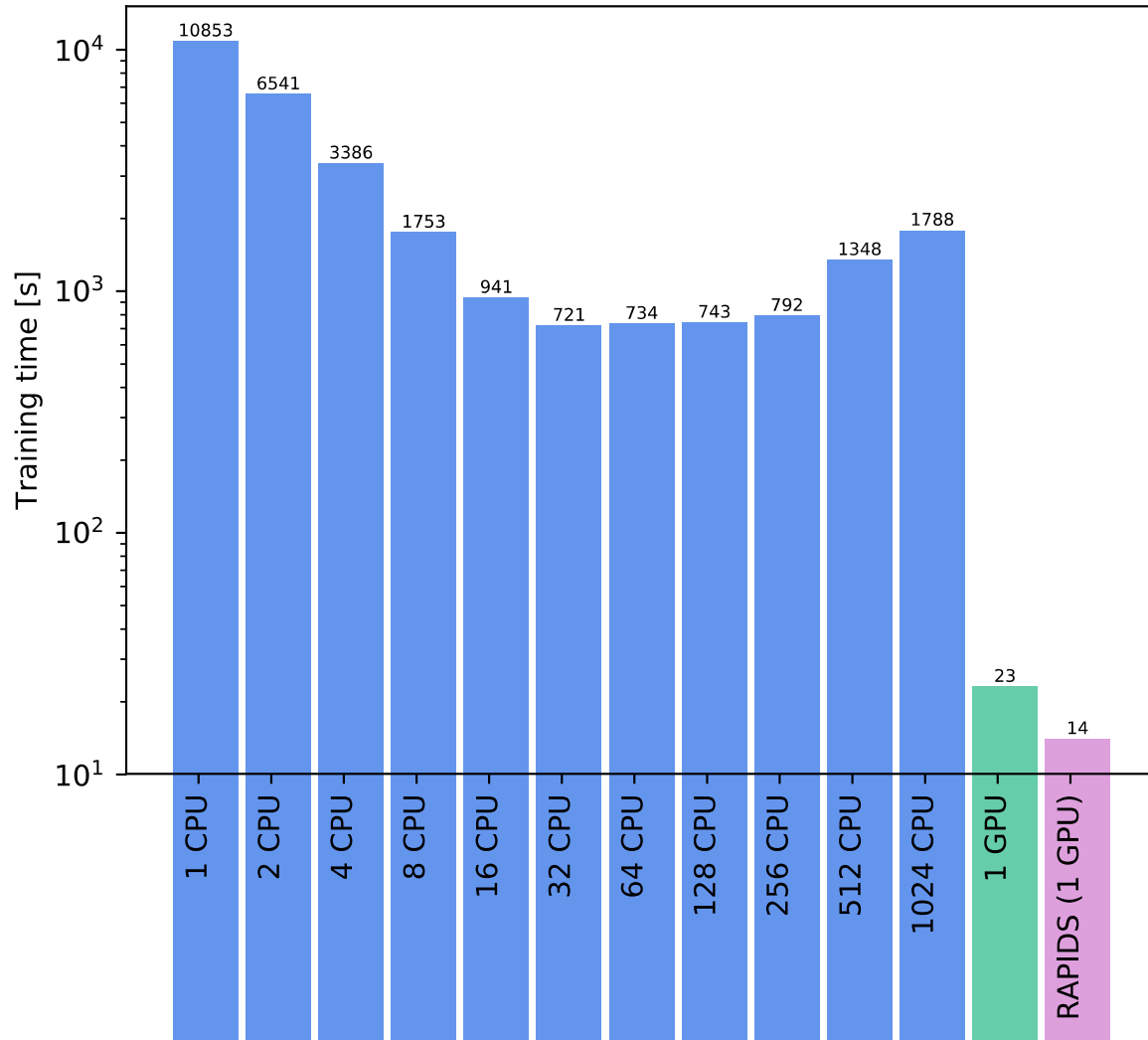
Comparison of XGBoost training time (data set = 44 GB)



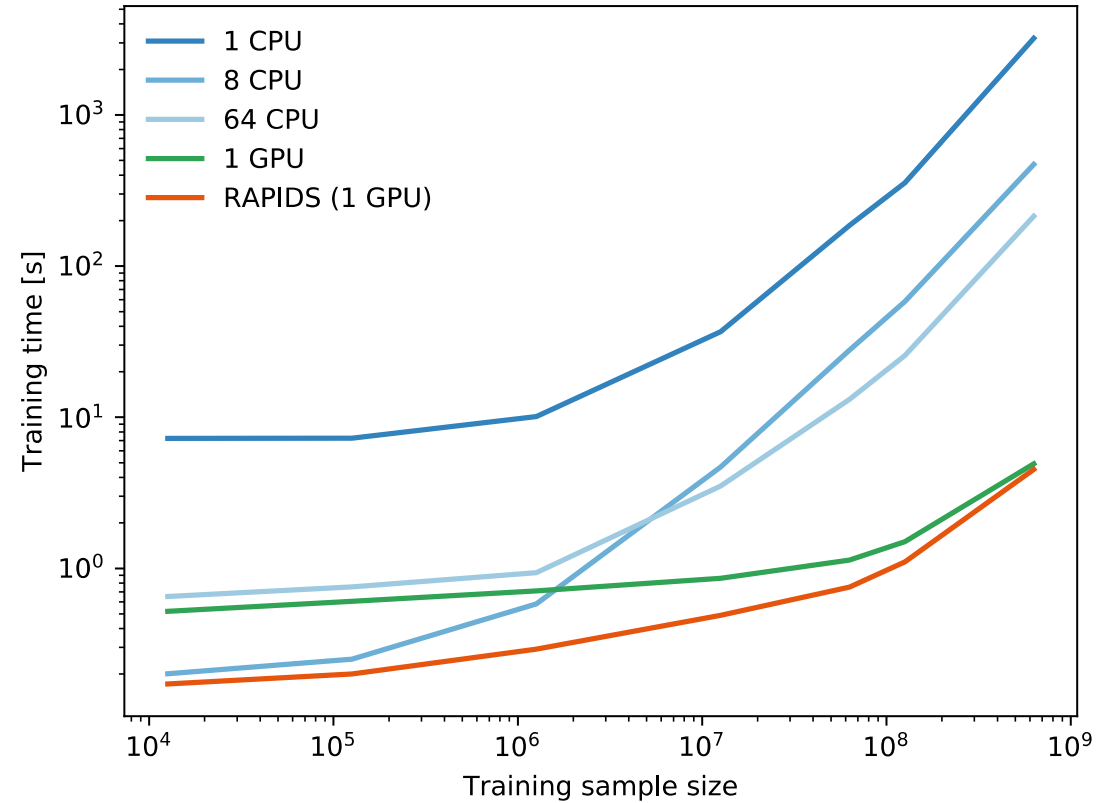


# Random forest / XGBoost training benchmarks

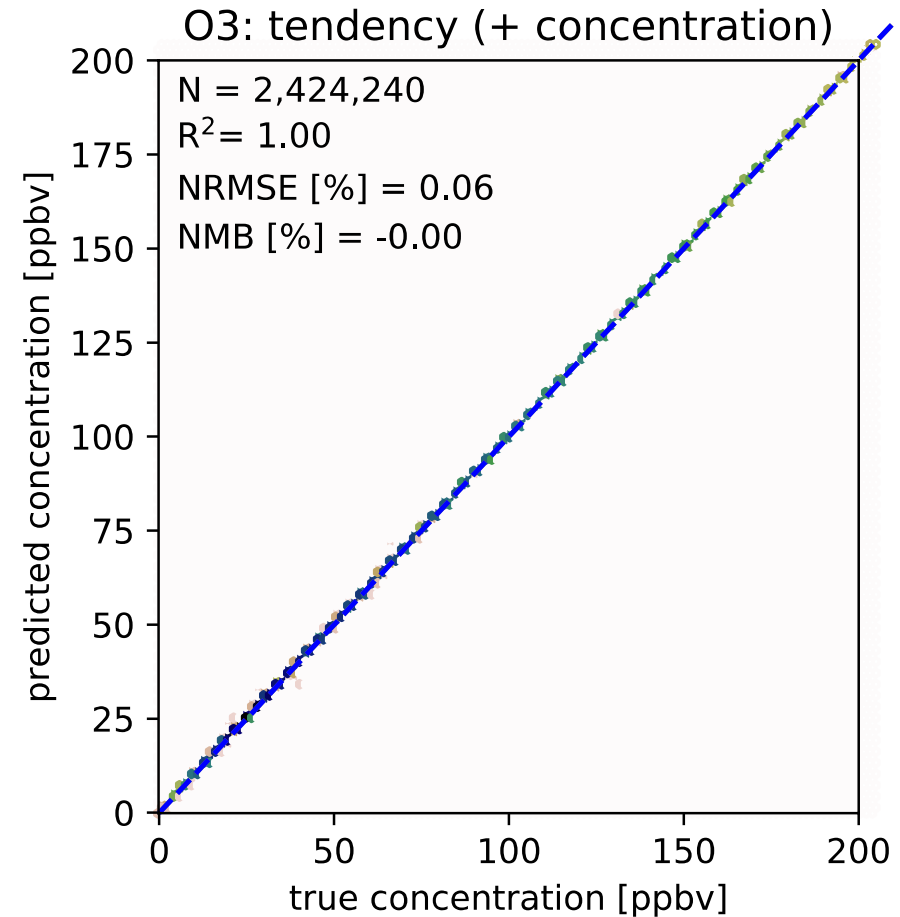
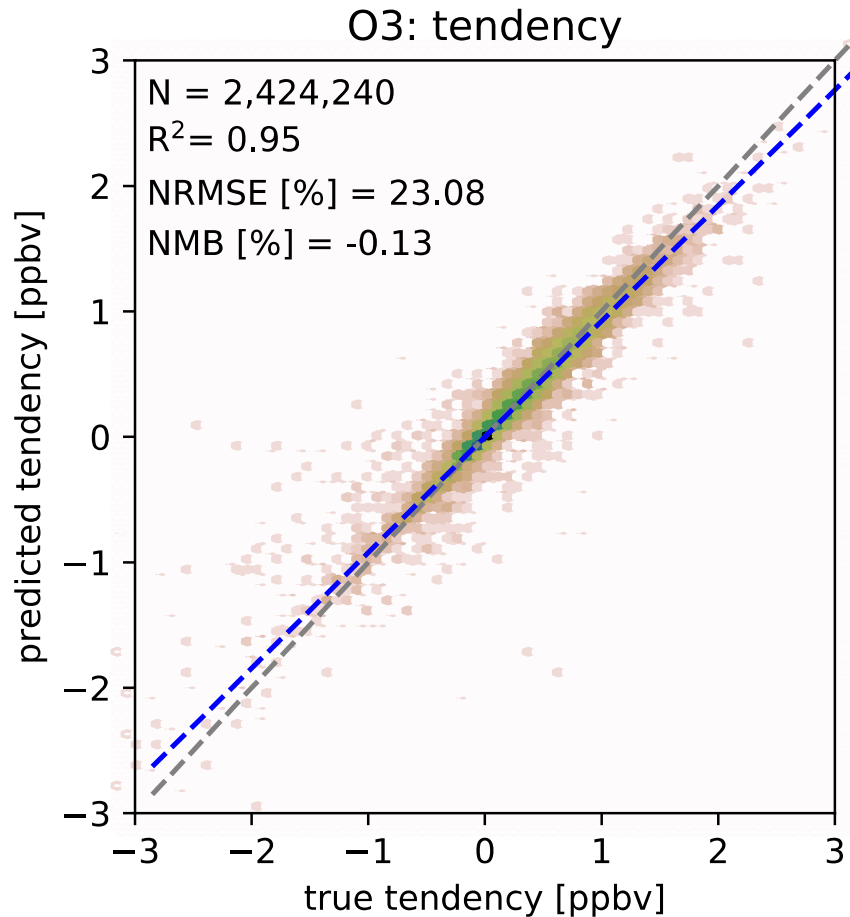
Comparison of XGBoost training time (data set = 44 GB)



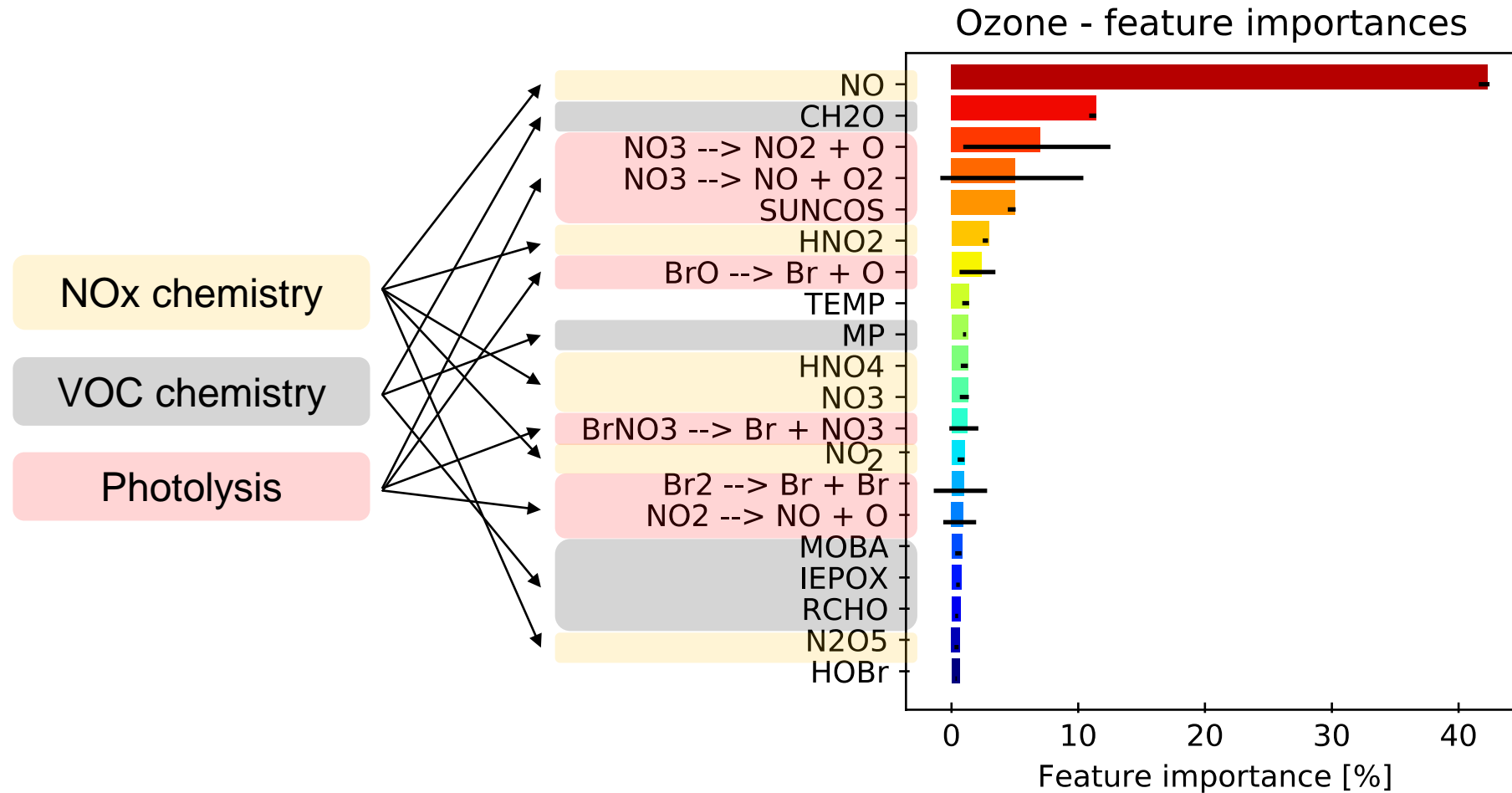
XGBoost training time



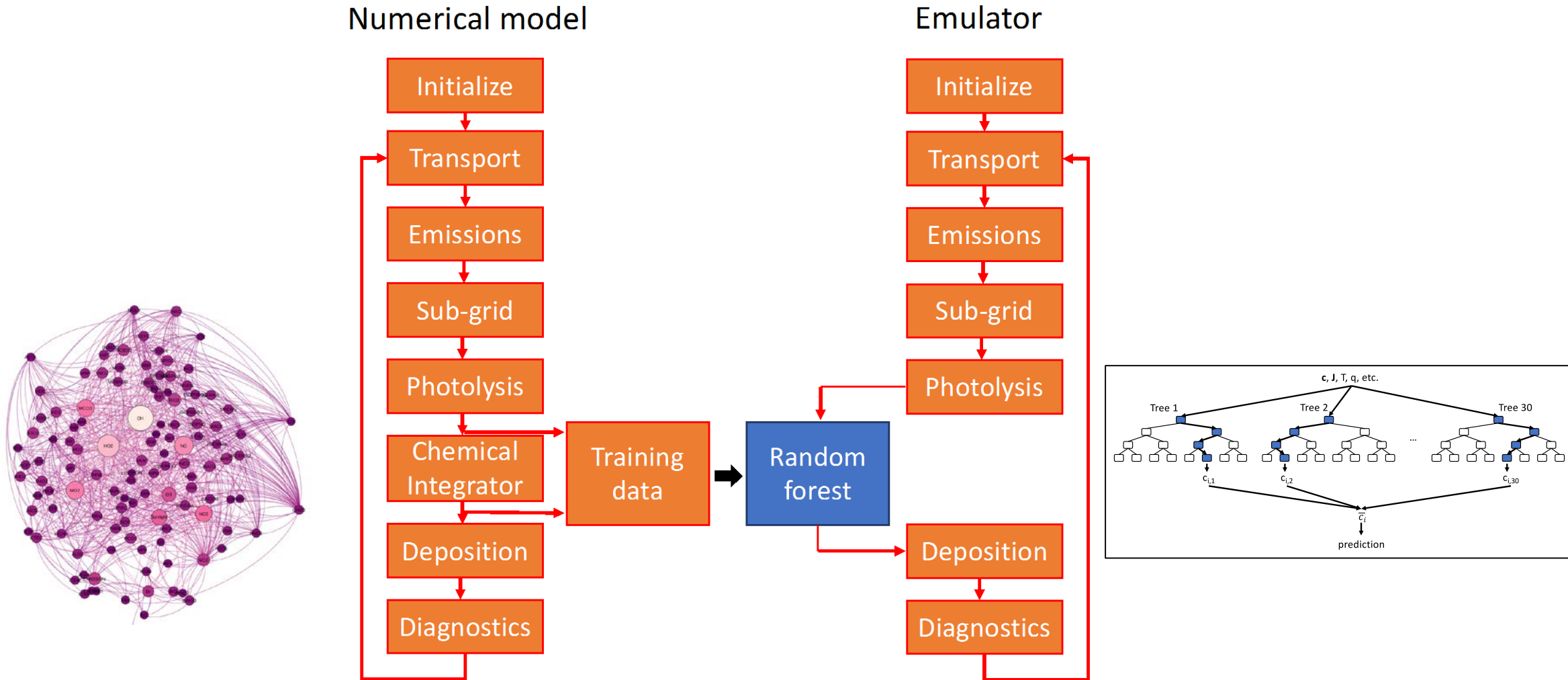
# Random forest / XGBoost can reproduce target concentrations almost perfectly (single-step prediction)



# Random forest / XGBoost solutions reflect known features of chemical kinetics

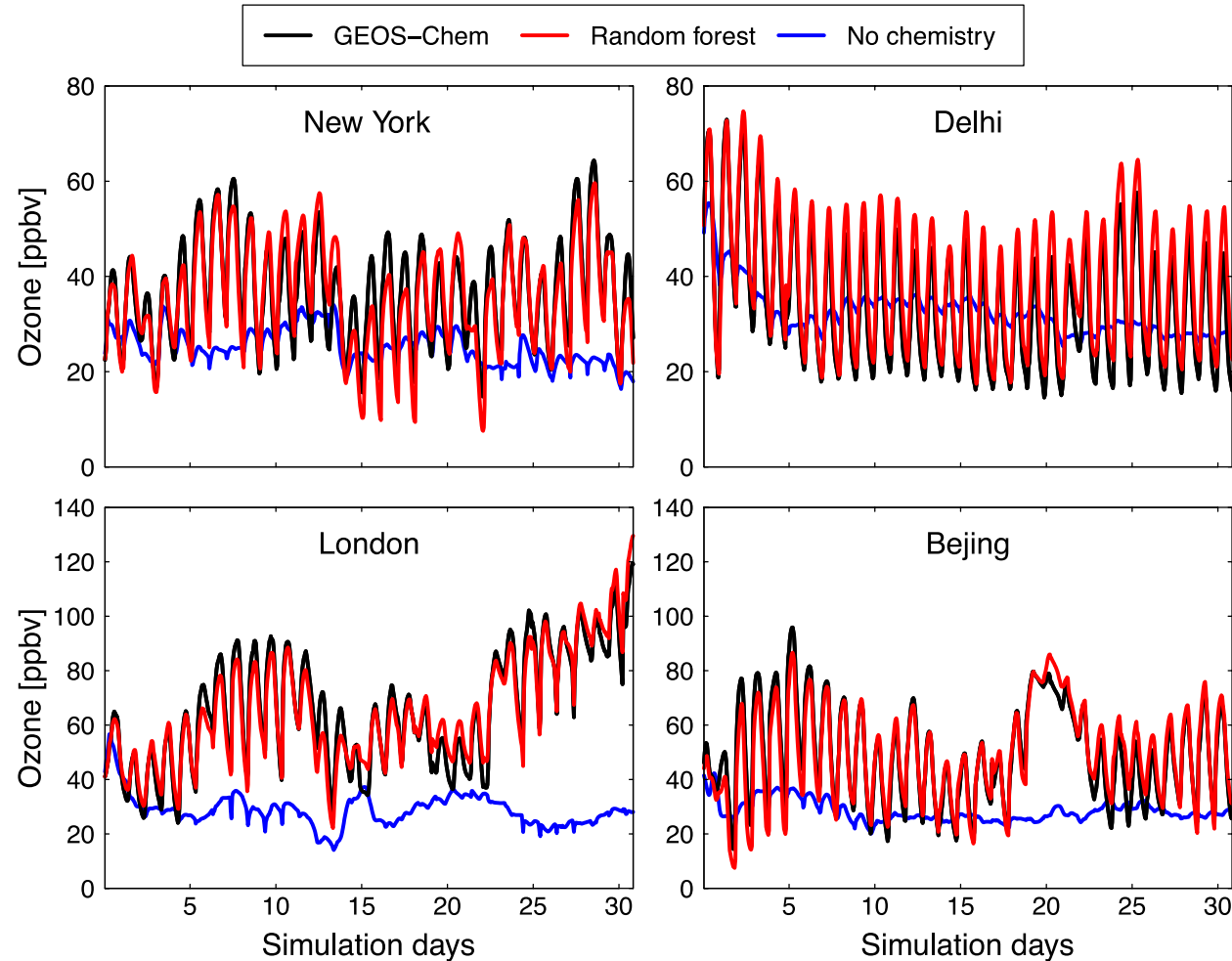


# 1-month simulation with random forest emulator



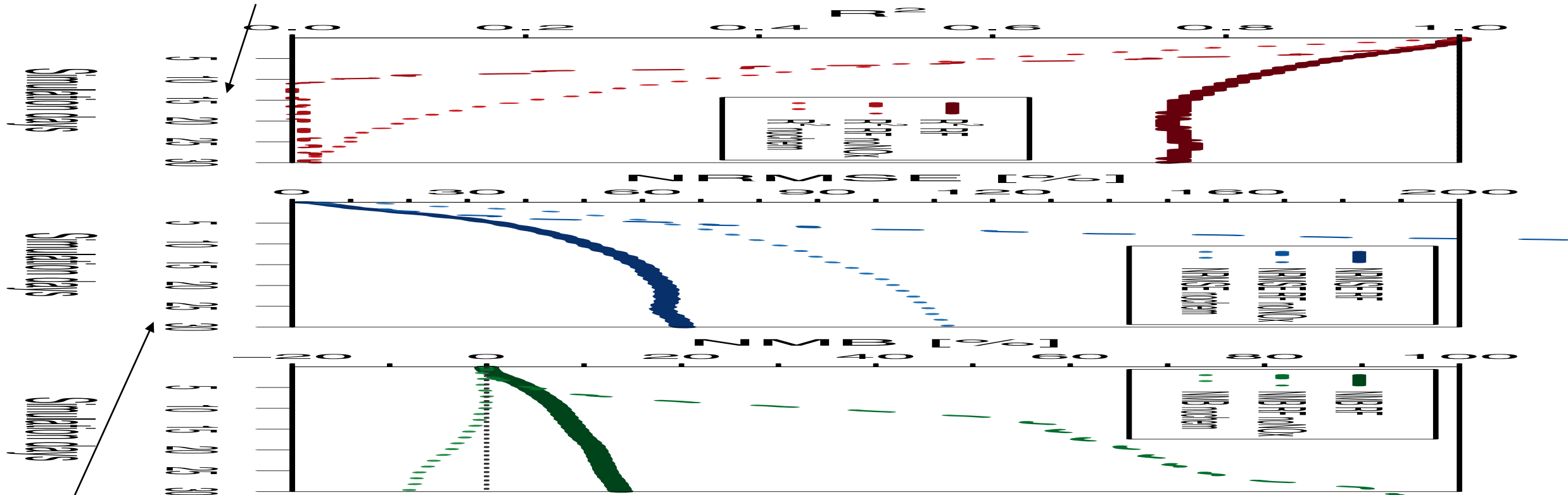


# Surface concentrations over polluted regions are well reproduced by ML model



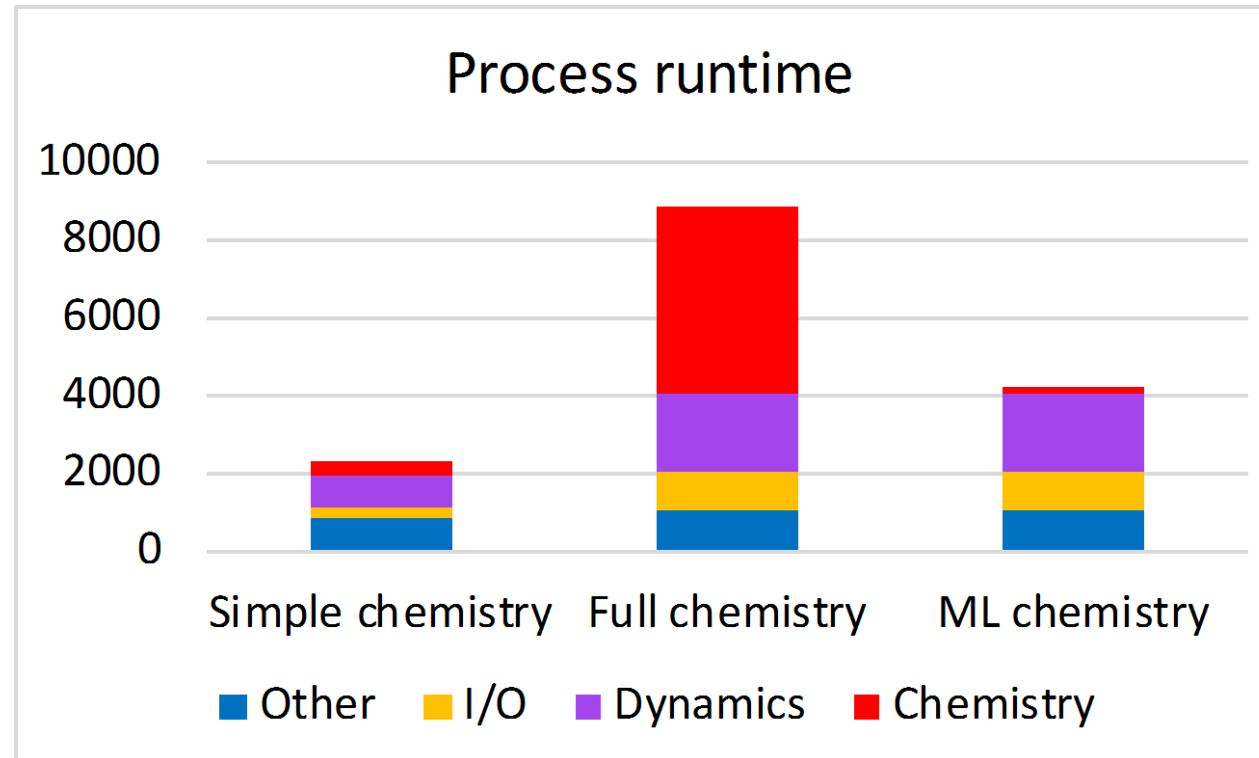
# Machine learning model remains stable over the long-term (but only if NOx is predicted as a family)

Model with NOx family prediction



Model without NOx family prediction

# Speedup potential



- Offline evaluation of one forest is 1000x faster than numerical integration
- Current implementation is very inefficient (2x slower than full chemistry)
- Currently working on seamless integration of XGBoost



# Summary

- Tree models do a good job at simulating atmospheric chemistry
- Adding constraints (e.g., chemical families) to the machine learning model is critical
- Potential applications:
  - Chemical data assimilation
  - Air quality forecasting
- Issues:
  - Train on very large data sets (>1 TB)
  - Dynamics for >200 chemical species is still slow

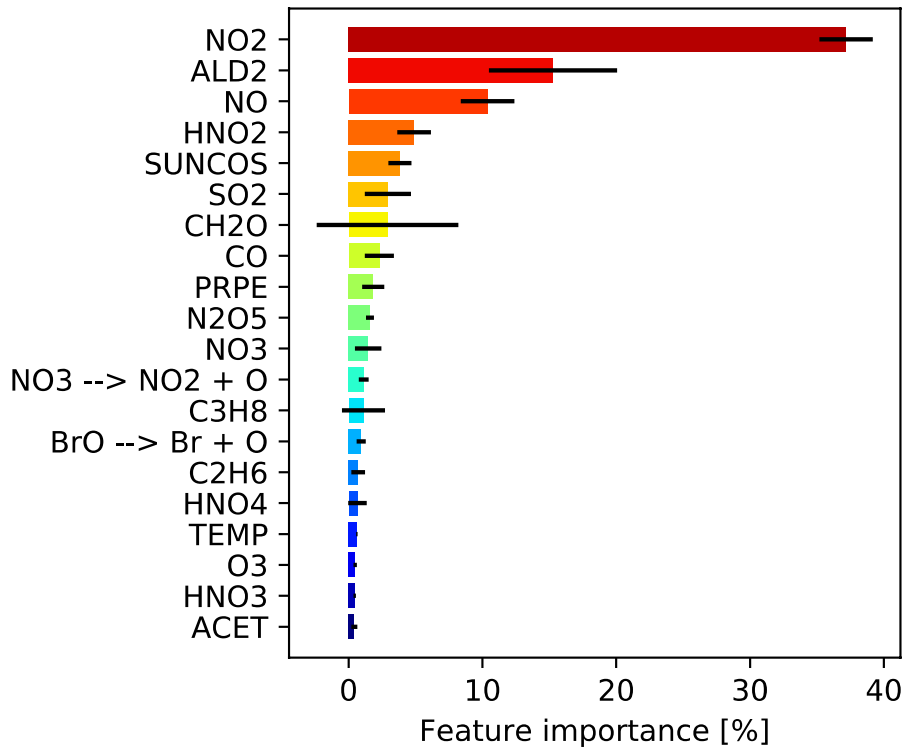
Keller and Evans: Application of random forest regression to the calculation of gas-phase chemistry within the GEOS-Chem chemistry model v10, GMD, 2019.



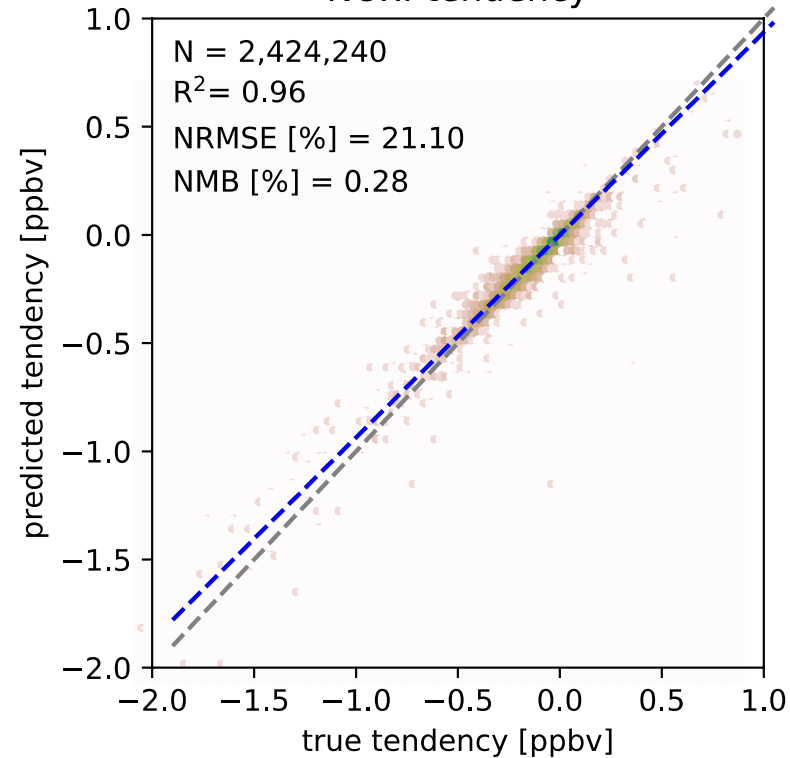


# Prediction of NOx

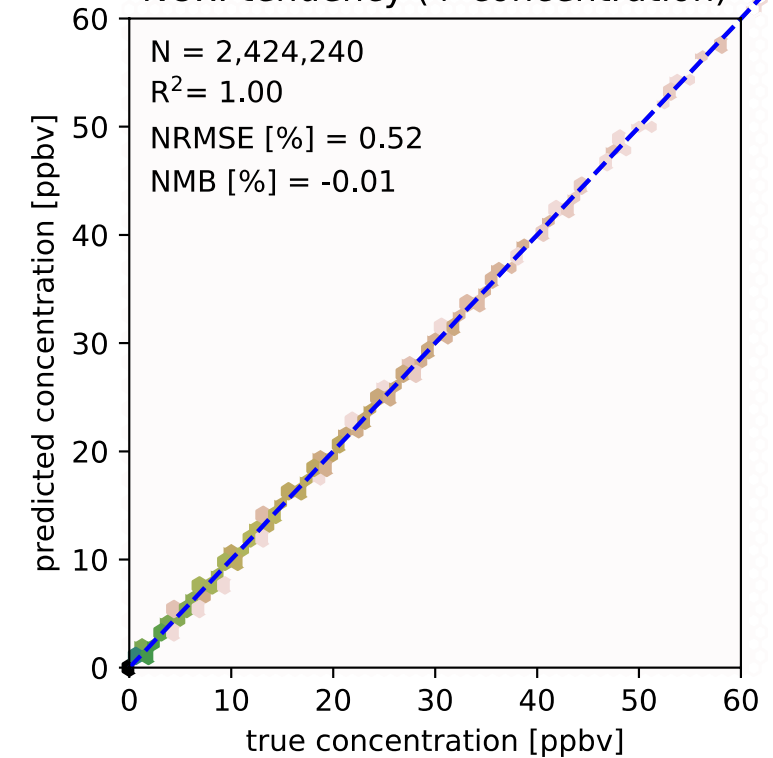
### NOx - feature importances



### NOx: tendency



### NOx: tendency (+ concentration)



# Surface concentrations over polluted regions are well reproduced by ML model

