

# Accelerated simulation of air pollution using NVIDIA RAPIDS

**Christoph A. Keller**<sup>1,2</sup>, Thomas L. Clune<sup>1</sup>,  
Matthew A. Thompson<sup>1,3</sup>, Matthew A. Stroud<sup>1,4</sup>,  
Mat J. Evans<sup>5</sup>, Zahra Ronaghi<sup>6</sup>

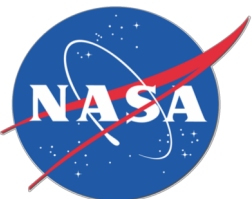
<sup>1</sup>NASA Global Modeling and Assimilation Office (GMAO)

<sup>2</sup>Universities Space Research Association (USRA)

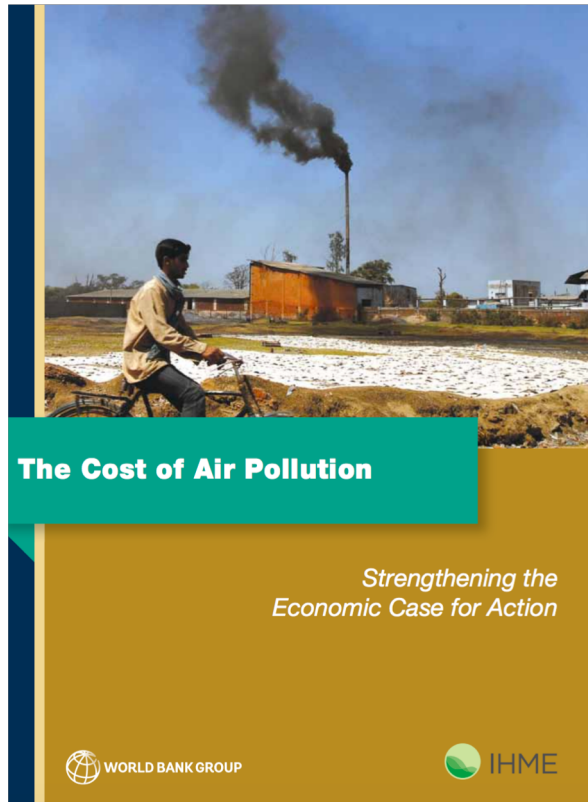
<sup>3</sup>Science Systems and Applications, Inc (SSAI)

<sup>4</sup>ASRC Federal Inuteq, <sup>5</sup>University of York, <sup>6</sup>NVIDIA

GPU Technology Conference :: 4-6 November 2019 :: Washington, DC



# Air pollution is a global problem - mitigating it is a big opportunity



World Bank: ~\$5 trillion in welfare losses in 2013

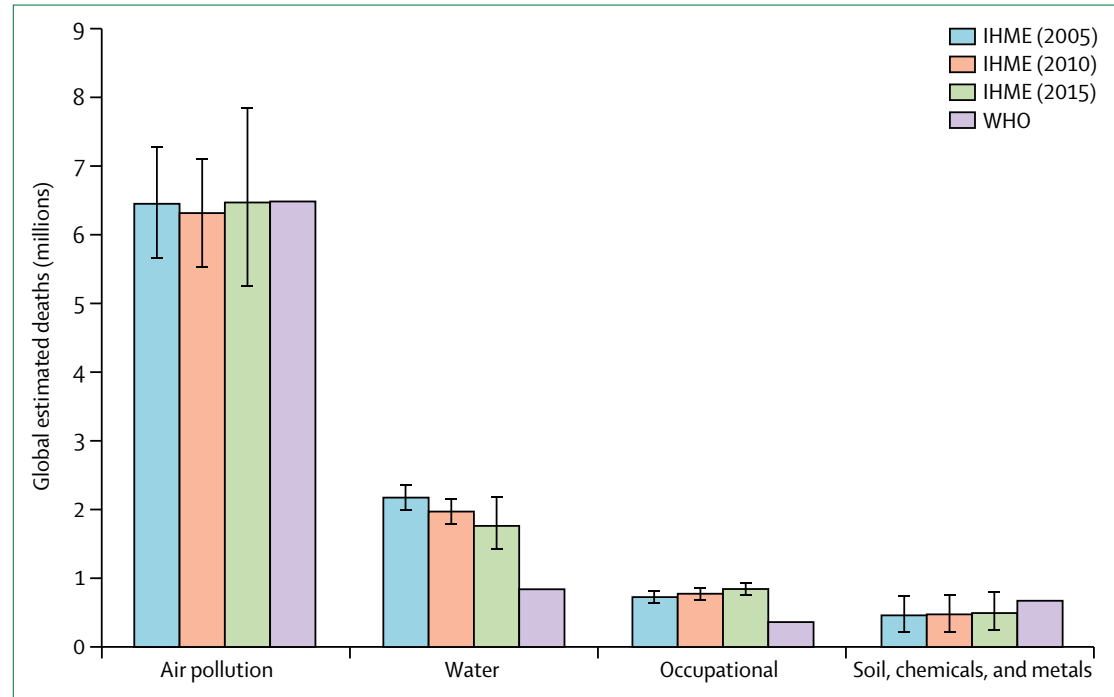
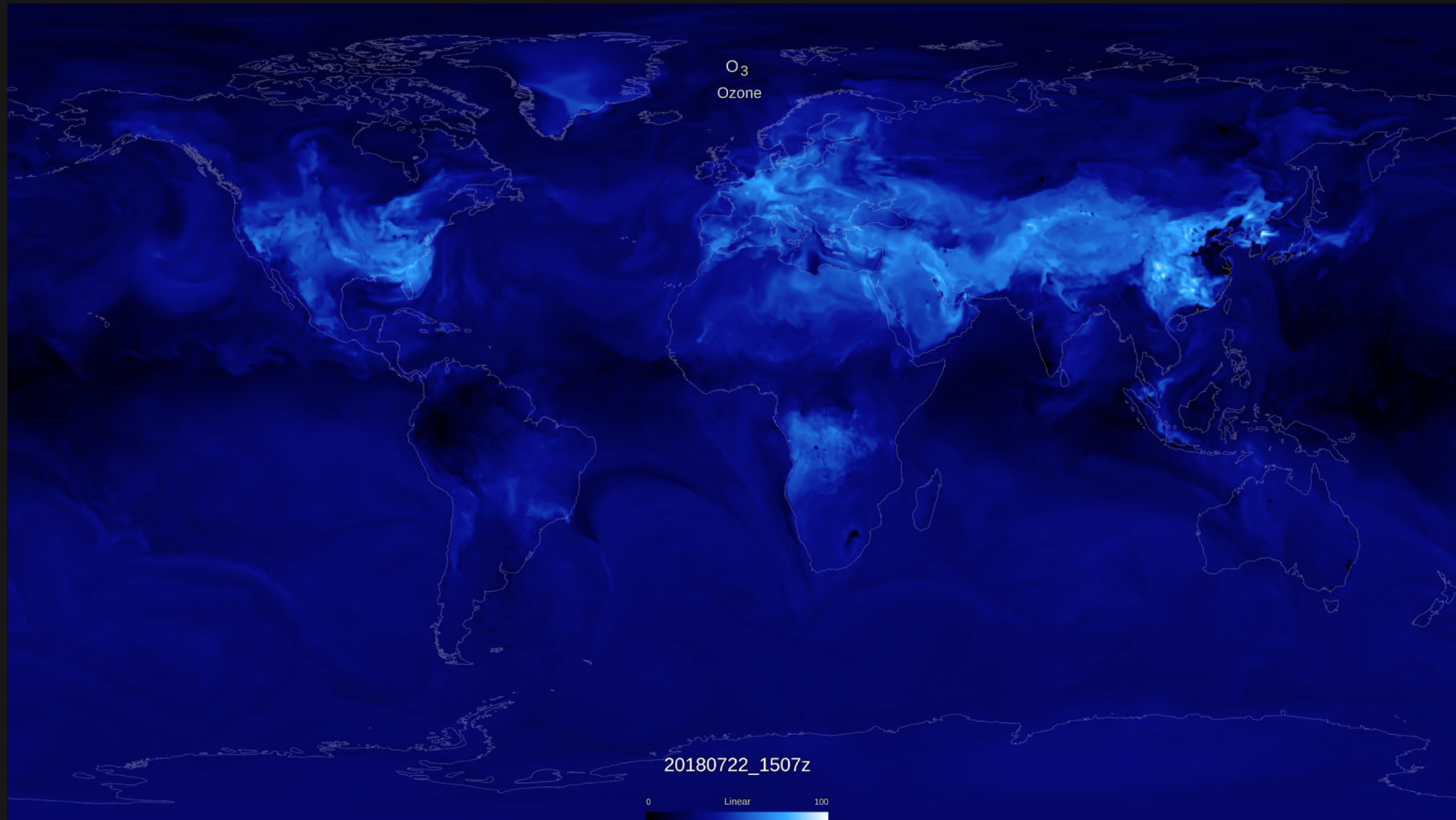


Figure 4: Global estimated deaths (millions) by pollution risk factor, 2005-15  
Using data from the GBD study<sup>42</sup> and WHO.<sup>99</sup> IHME=Institute for Health Metrics and Evaluation.

The Lancet (2017): Air pollution is responsible for 6-7 million death / year

# Numerical simulation of atmospheric chemistry



➤ 56 million grid cells (25x25 km<sup>2</sup>, 72 levels), 250 chemical species



# NASA's GEOS composition forecast (GEOS-CF) model conducts global air quality simulations in near real-time

<https://fluid.nccs.nasa.gov/cf/>

<https://portal.nccs.nasa.gov/datashare/gmao/geos-cf/v1/>

**GrADS Data Server - info for /gmao/geos-cf/assim/chm\_tavg\_1hr\_g1440x721\_v1 : [dds](#) [das](#)**

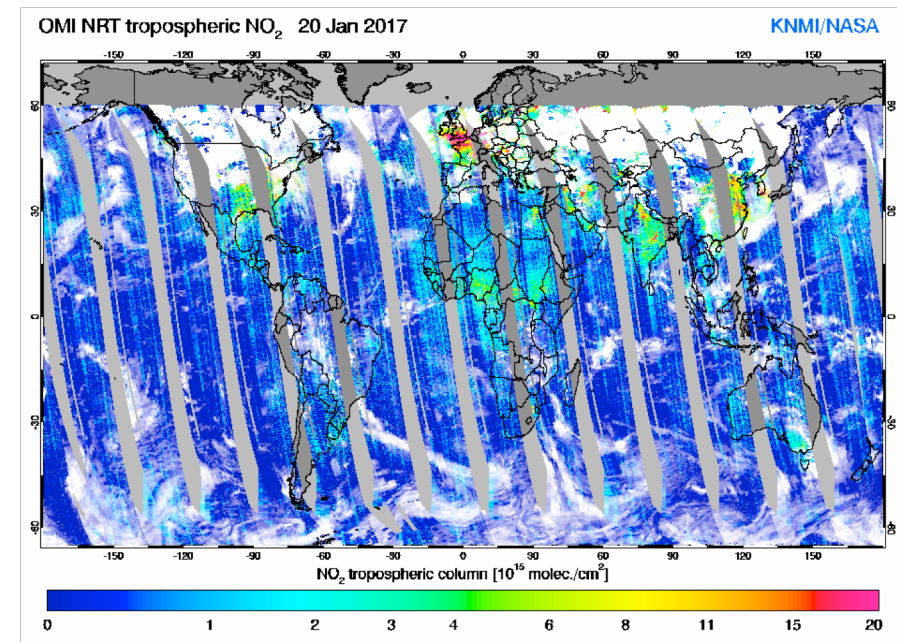
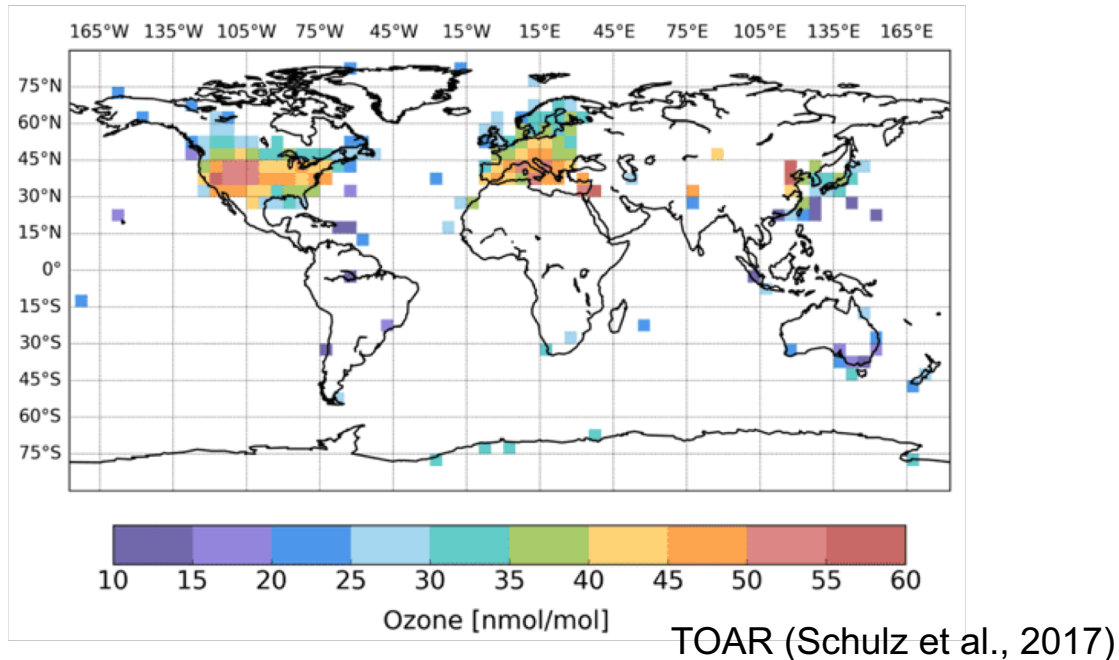
**OPeNDAP/DODS Data URL:** [https://opendap.nccs.nasa.gov/dods/gmao/geos-cf/assim/chm\\_tavg\\_1hr\\_g1440x721\\_v1](https://opendap.nccs.nasa.gov/dods/gmao/geos-cf/assim/chm_tavg_1hr_g1440x721_v1)

**Description:** GEOS CF (Composition Forecast)  
**Documentation:** (none provided)  
**Longitude:** -180.000000000000°E to 179.750000000000°E (1440 points, avg. res. 0.25°)  
**Latitude:** -90.000000000000°N to 90.000000000000°N (721 points, avg. res. 0.25°)  
**Altitude:** 72.000000000000 to 72.000000000000 (1 points)  
**Time:** 00:30Z01JAN2018 to 11:30Z31OCT2019 (16044 points, avg. res. 0.042 days)  
**Variables:** (total of 52)  
**xyle** xylene (c8h10, mw = 106.16 g mol-1) volume mixing ratio dry air  
**dst2** dust aerosol, reff = 1.4 microns (mw = 29.00 g mol-1) volume mixing ratio dry air  
**hno4** peroxyntiric acid (hno4, mw = 79.00 g mol-1) volume mixing ratio dry air  
**pm25su\_rh35\_gcc** sulfate\_particulate\_matter\_with\_diameter\_below\_2.5\_um\_rh\_35

<https://opendap.nccs.nasa.gov/dods/gmao/geos-cf/>



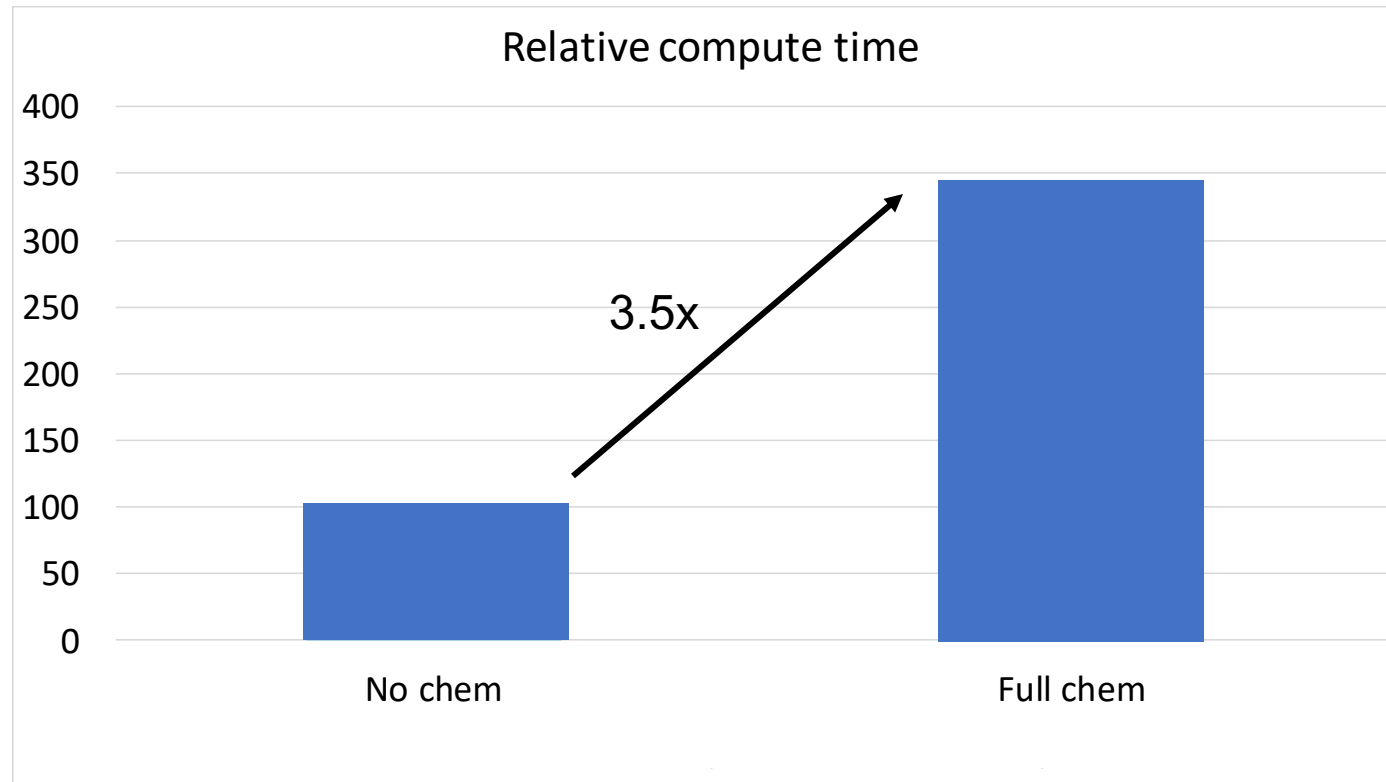
# Need models to fill temporal and spatial gaps in observations



Surface observations are not global

Satellite observations are also discontinuous

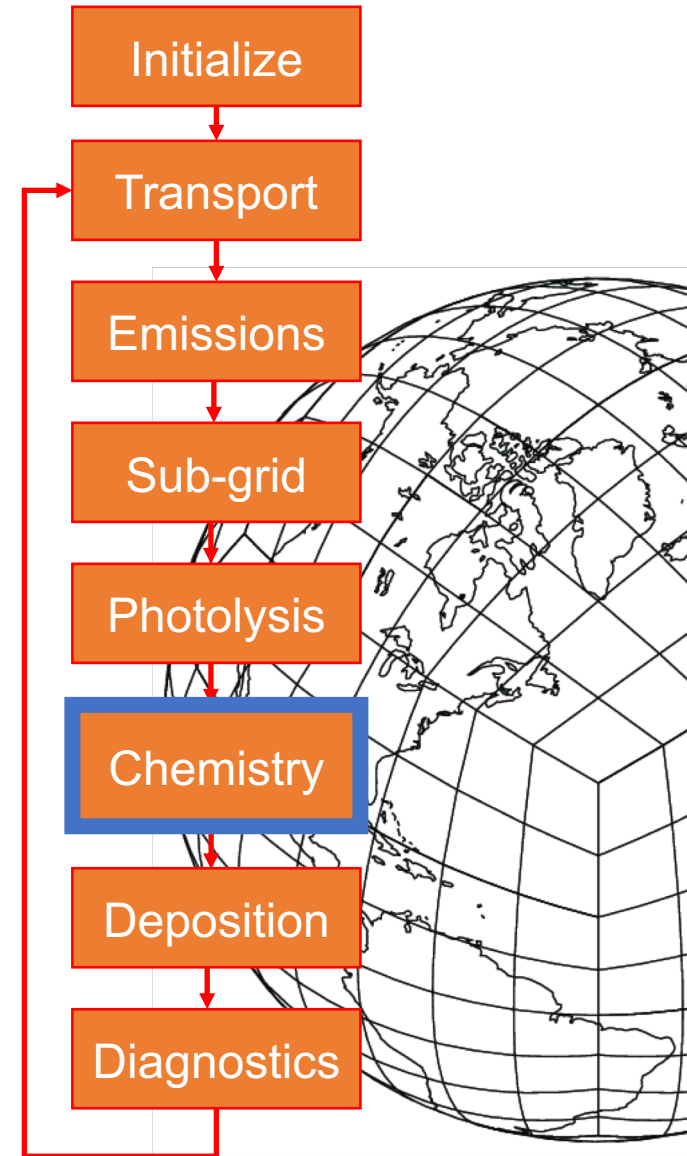
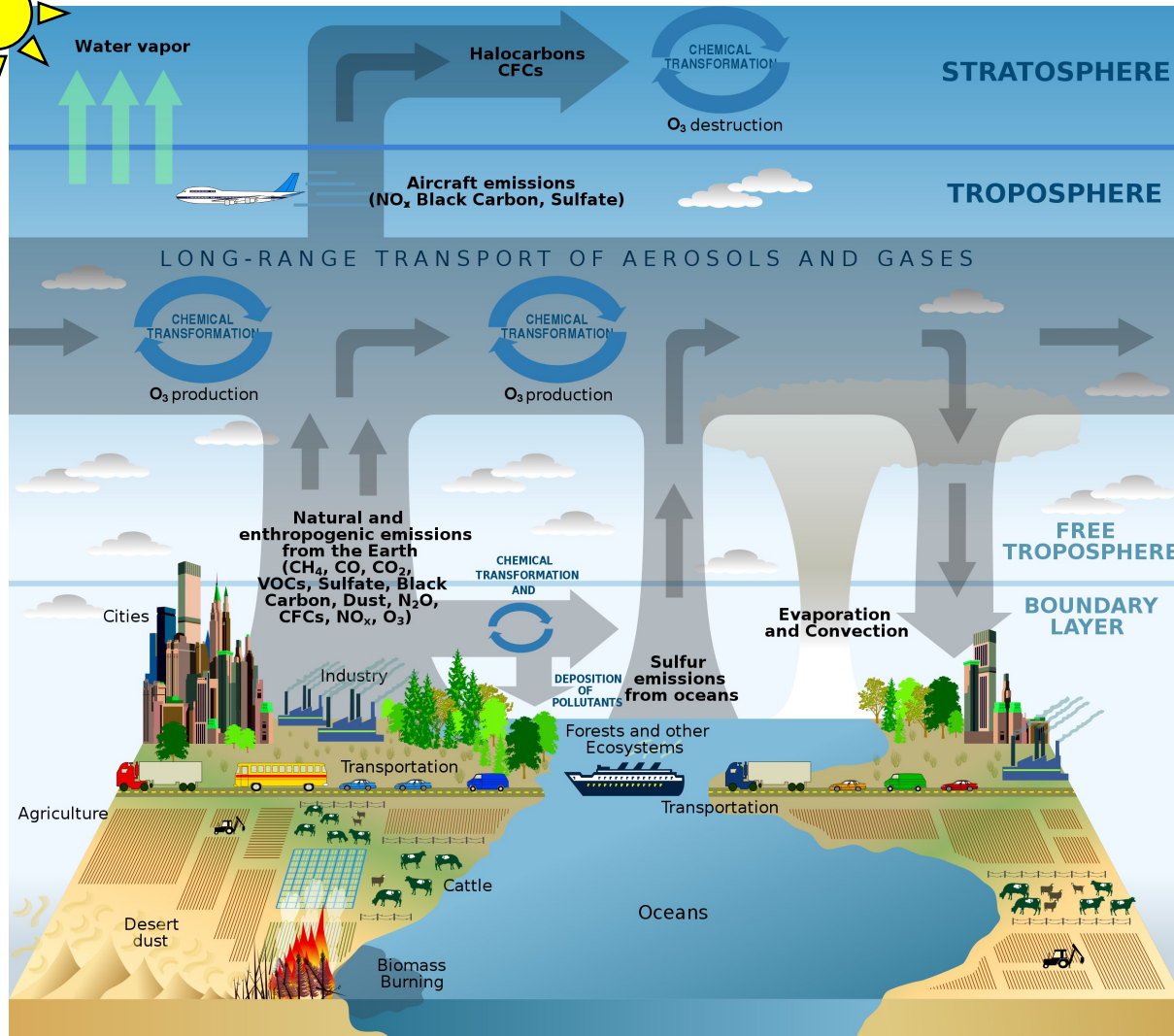
# Atmospheric chemistry models are computationally expensive



- High-resolution chemistry simulation requires >1000 CPU's
- Throughput: approx. 20 simulation days in 24 hours
- Outputting the full chemical state: ~1.5 TB / simulation day



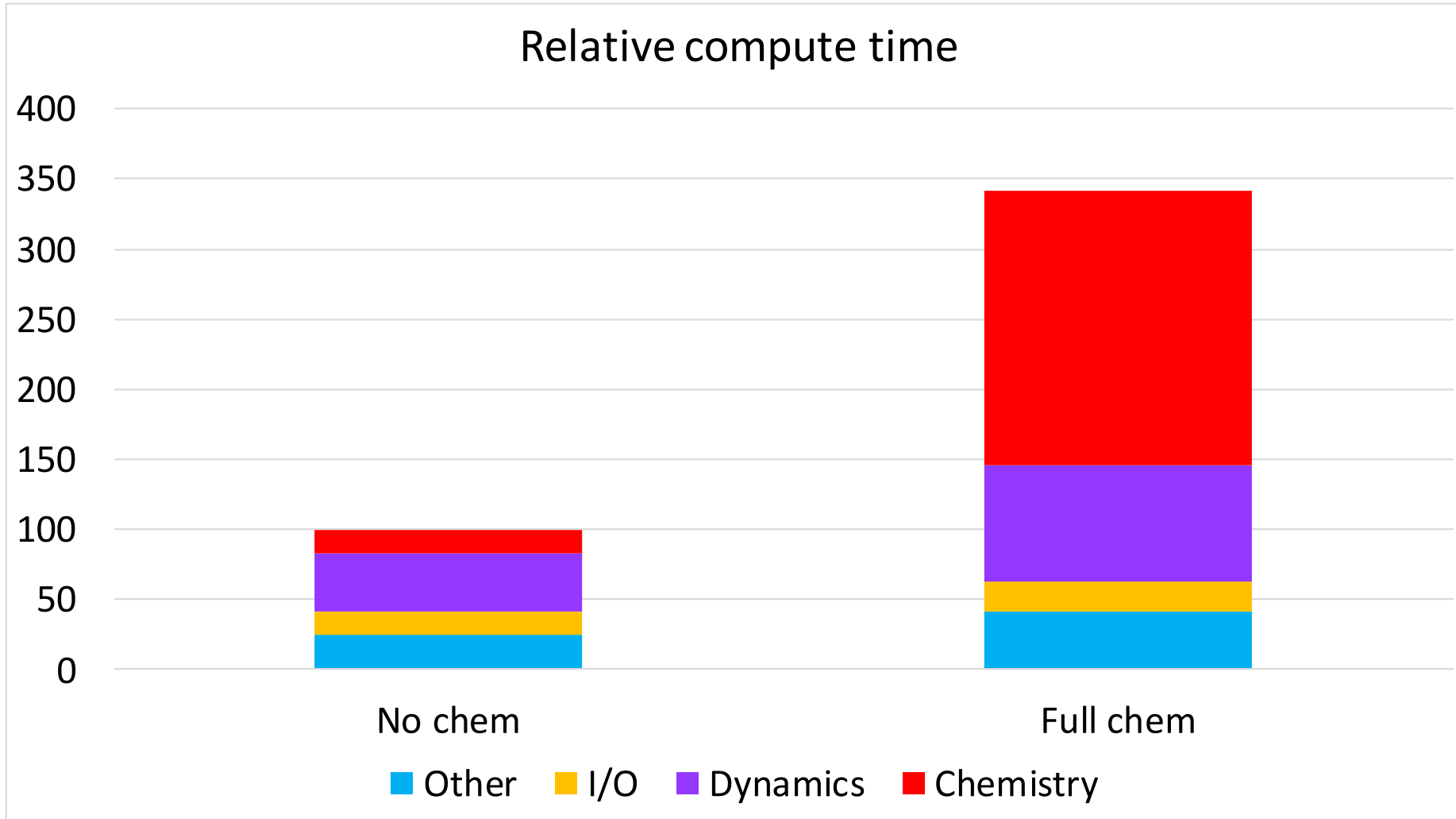
# Numerical simulation of atmospheric composition



<https://digital.library.unt.edu/ark:/67531/metadc11954/m1/37/>

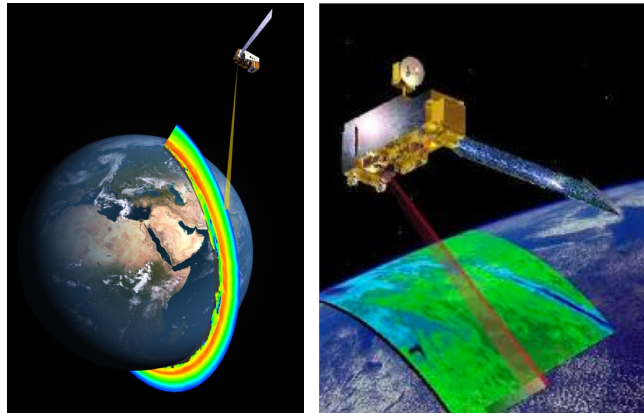


# Chemistry accounts for more than 50% of compute time

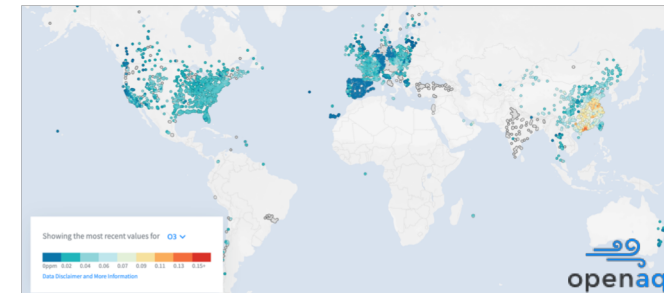
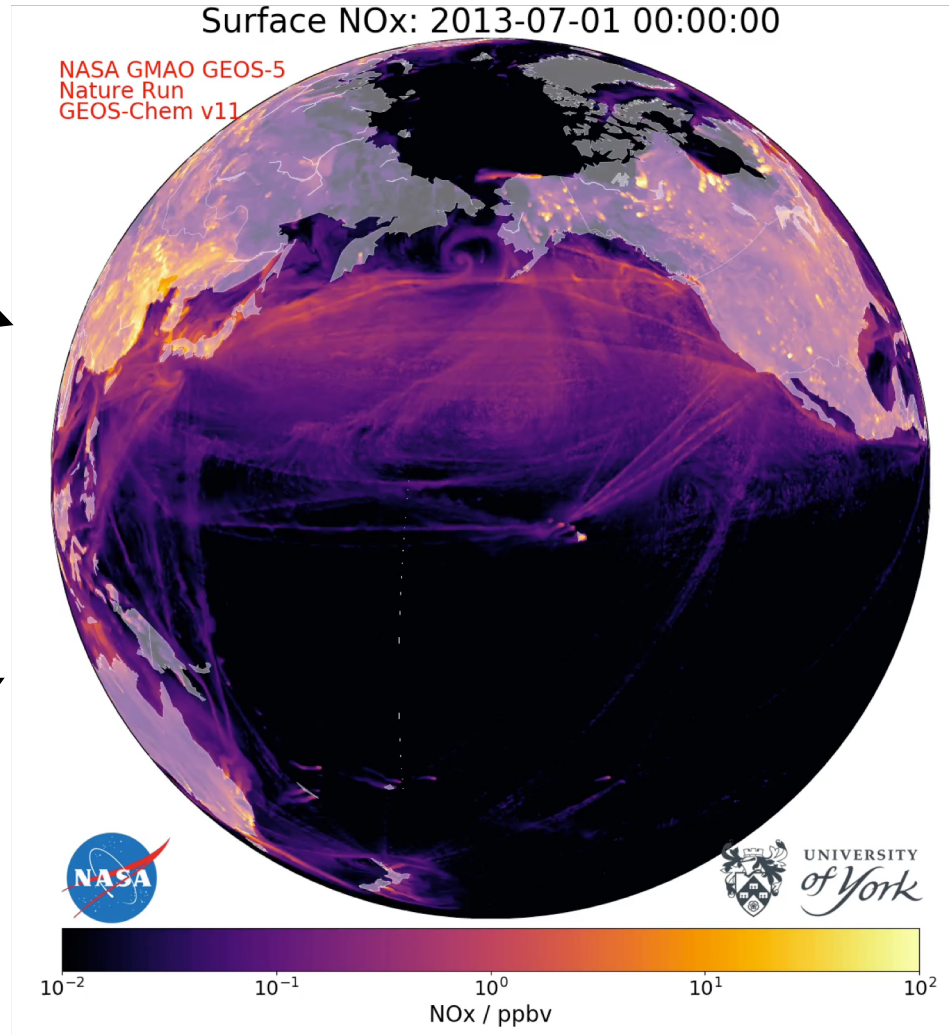




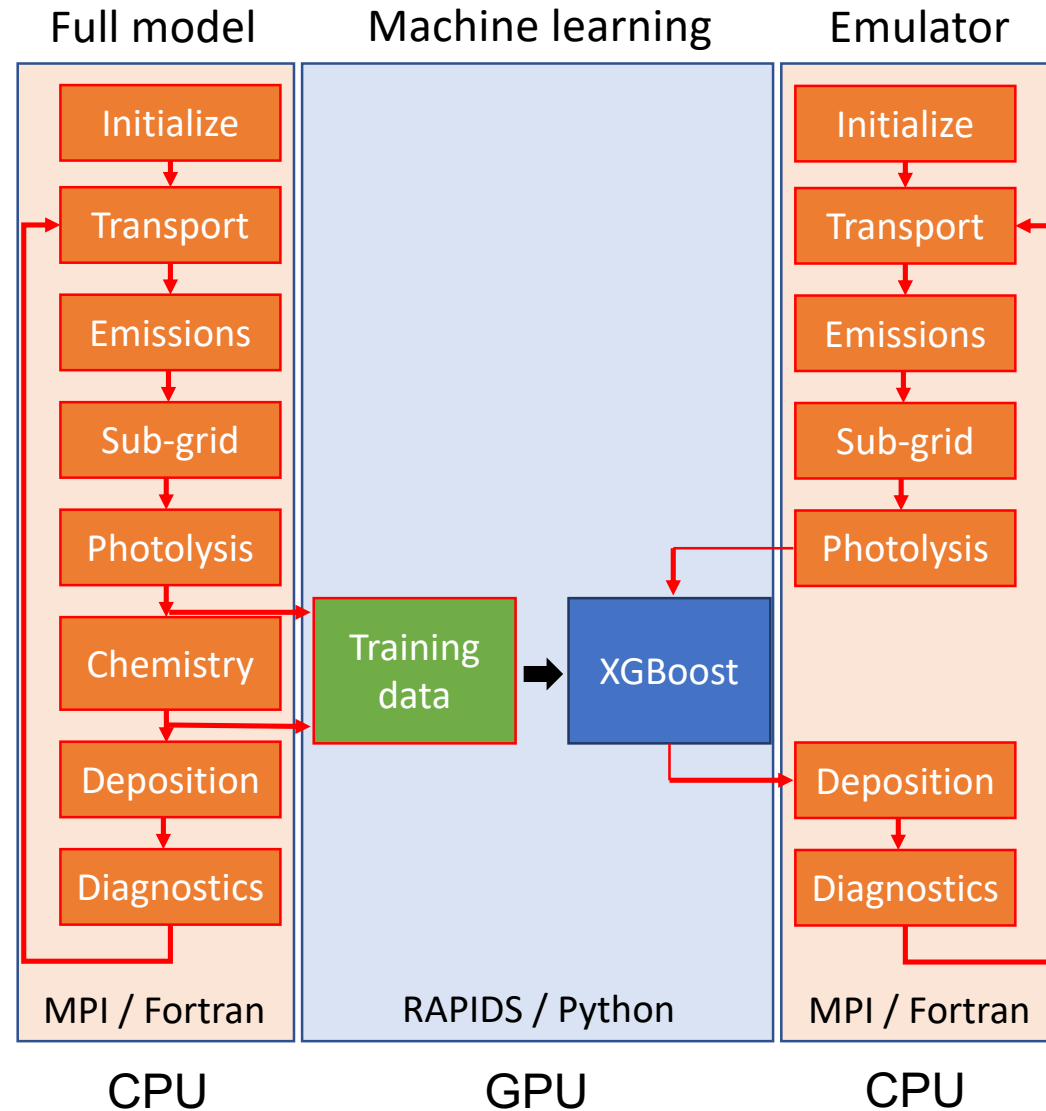
# High computational cost of chemistry currently prevents optimal use of observations



www.nasa.gov



# Replace slow chemical integrator with machine learning model



# Use machine learning to emulate chemical transformations in the atmosphere

Inputs

Meteorology:  
- 7 variables

Chemistry:  
- 143 chemical species  
- 91 photolysis rates

ML



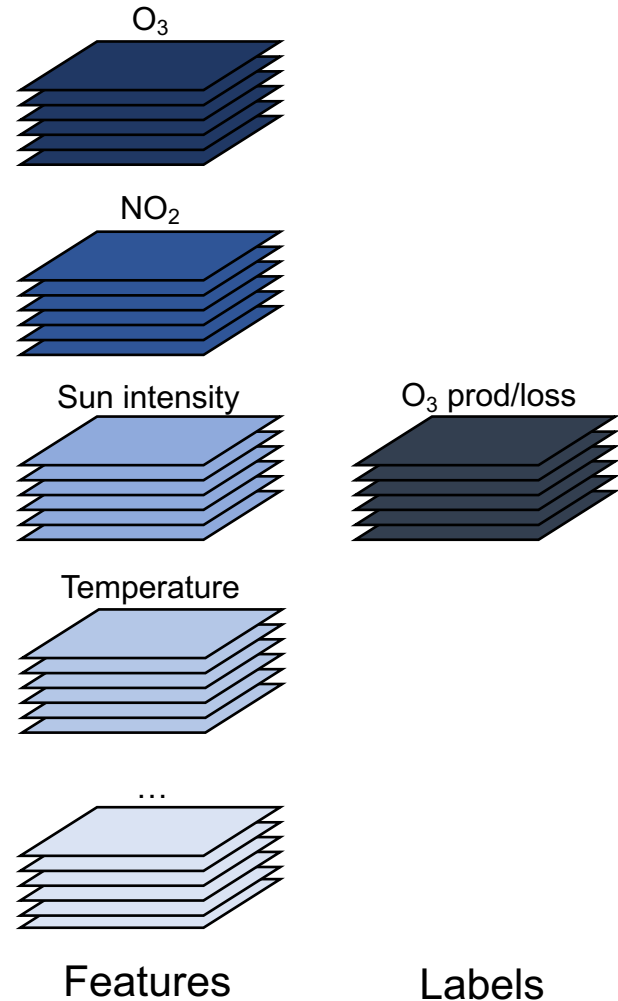
Output

Chemical production / loss

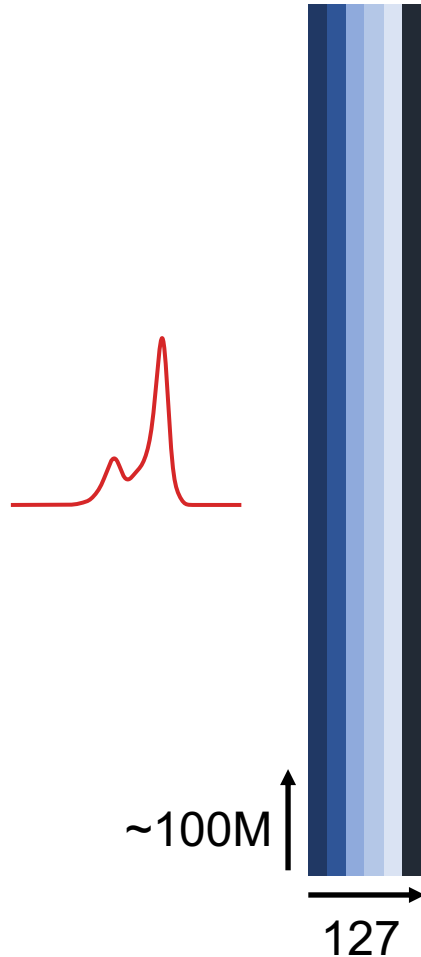
- Algorithm: extreme gradient boosted decision trees (XGBoost)
- Train separate algorithm for each species

# Machine learning workflow

Training data (netCDF format)



Subsample & flatten,  
write to csv (xarray)



Train (XGBoost):

- Read csv, convert to DMatrix
- Train

## Setup 1

Read on CPU (Intel Haswell)  
Train on CPU

## Setup 2:

Read on CPU  
Train on GPU (V100)

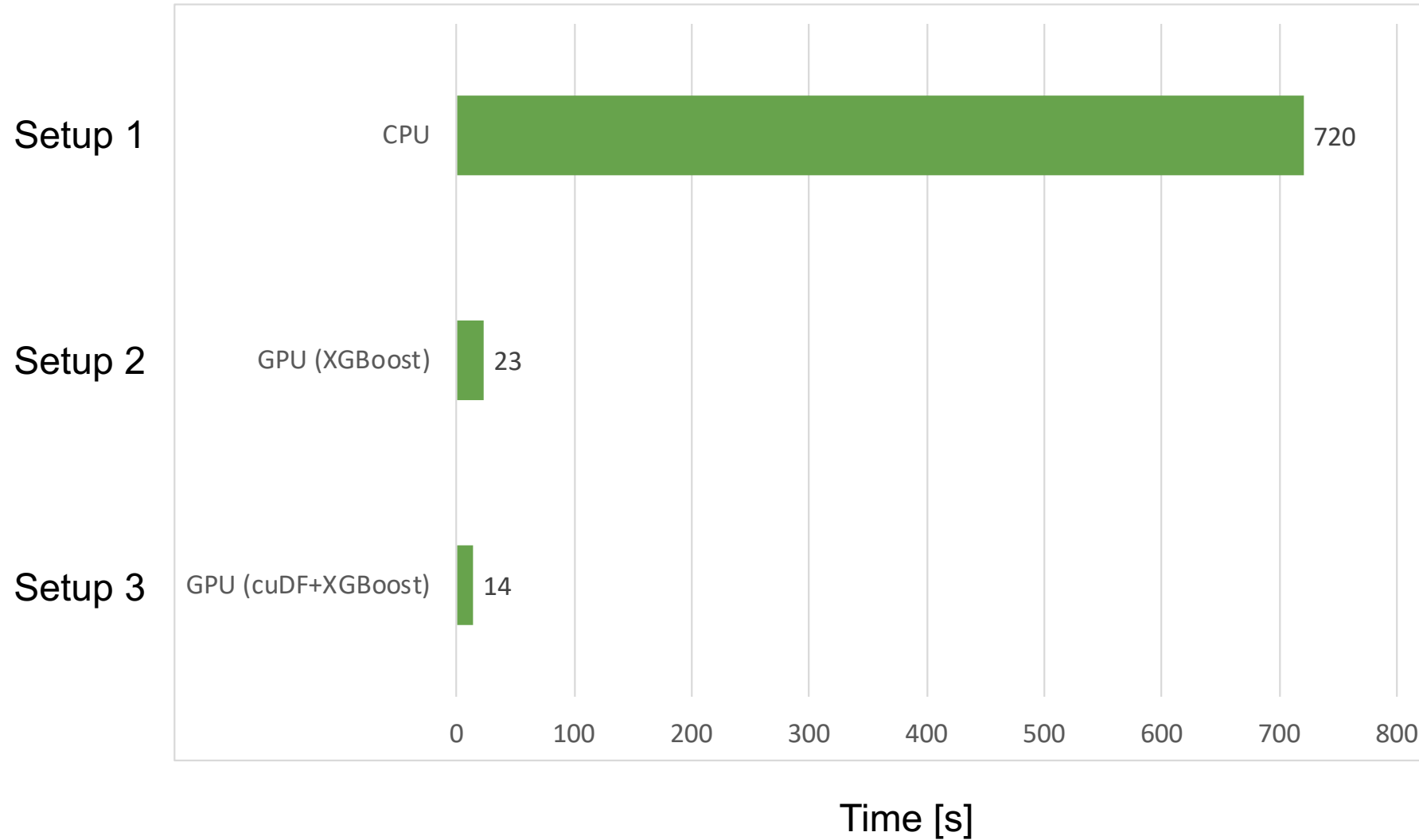
## Setup 3:

Read on GPU (cuDF/cuIO)  
Train on GPU (dask-XGBoost)

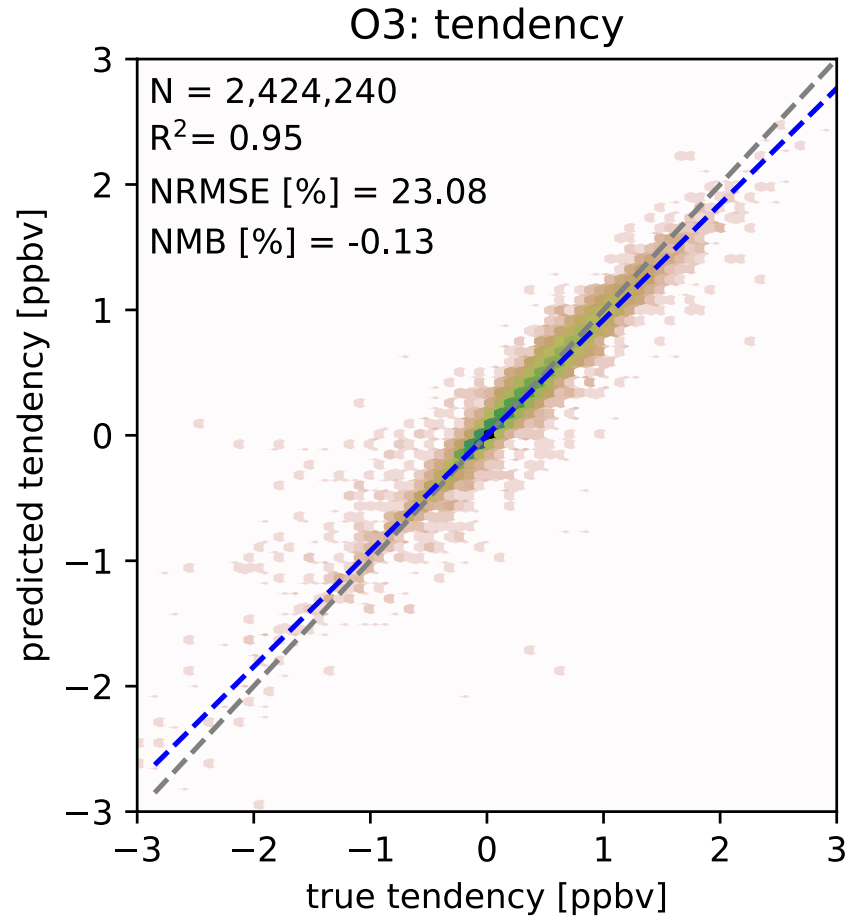




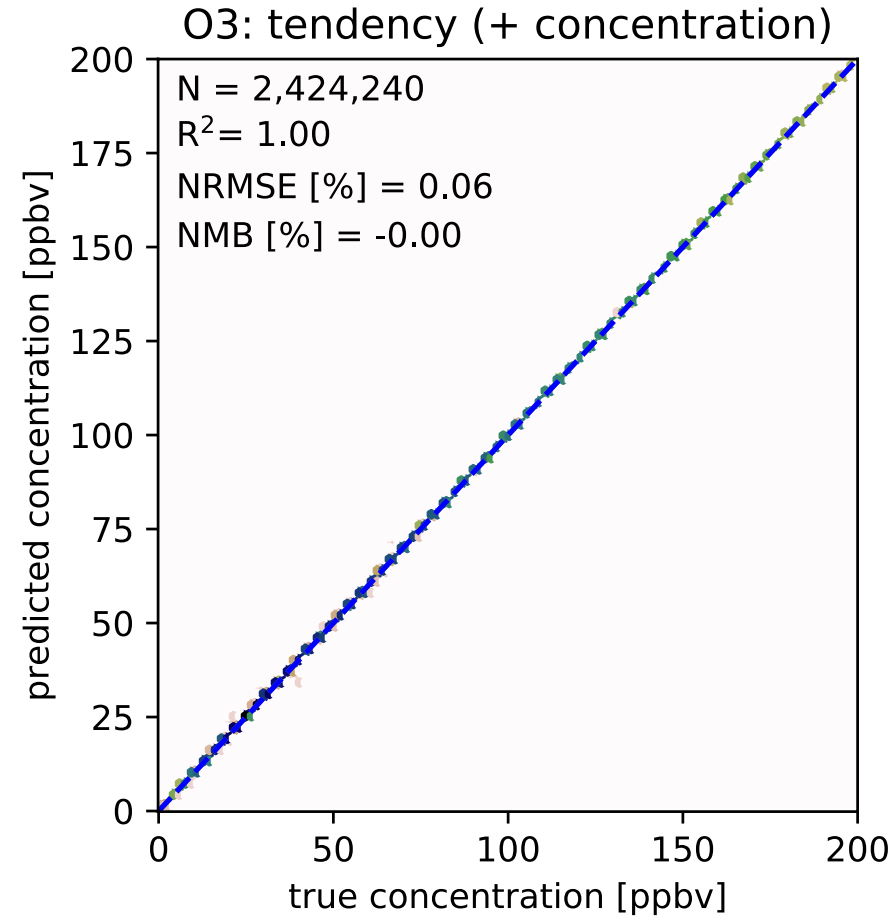
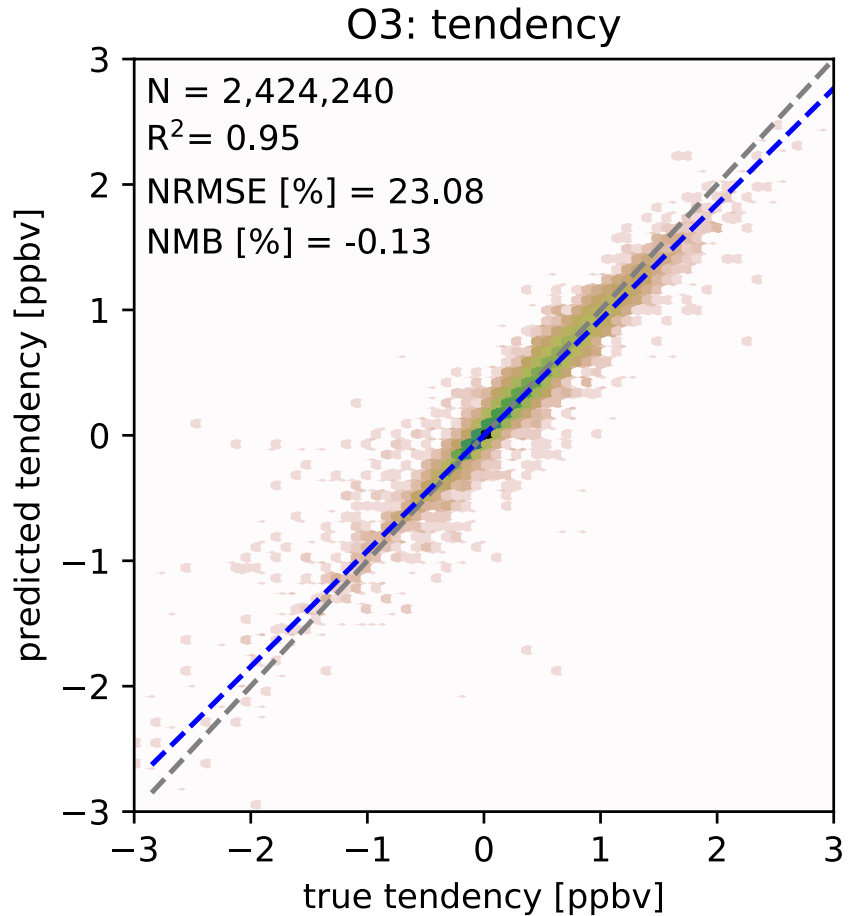
# XGBoost training benchmarks



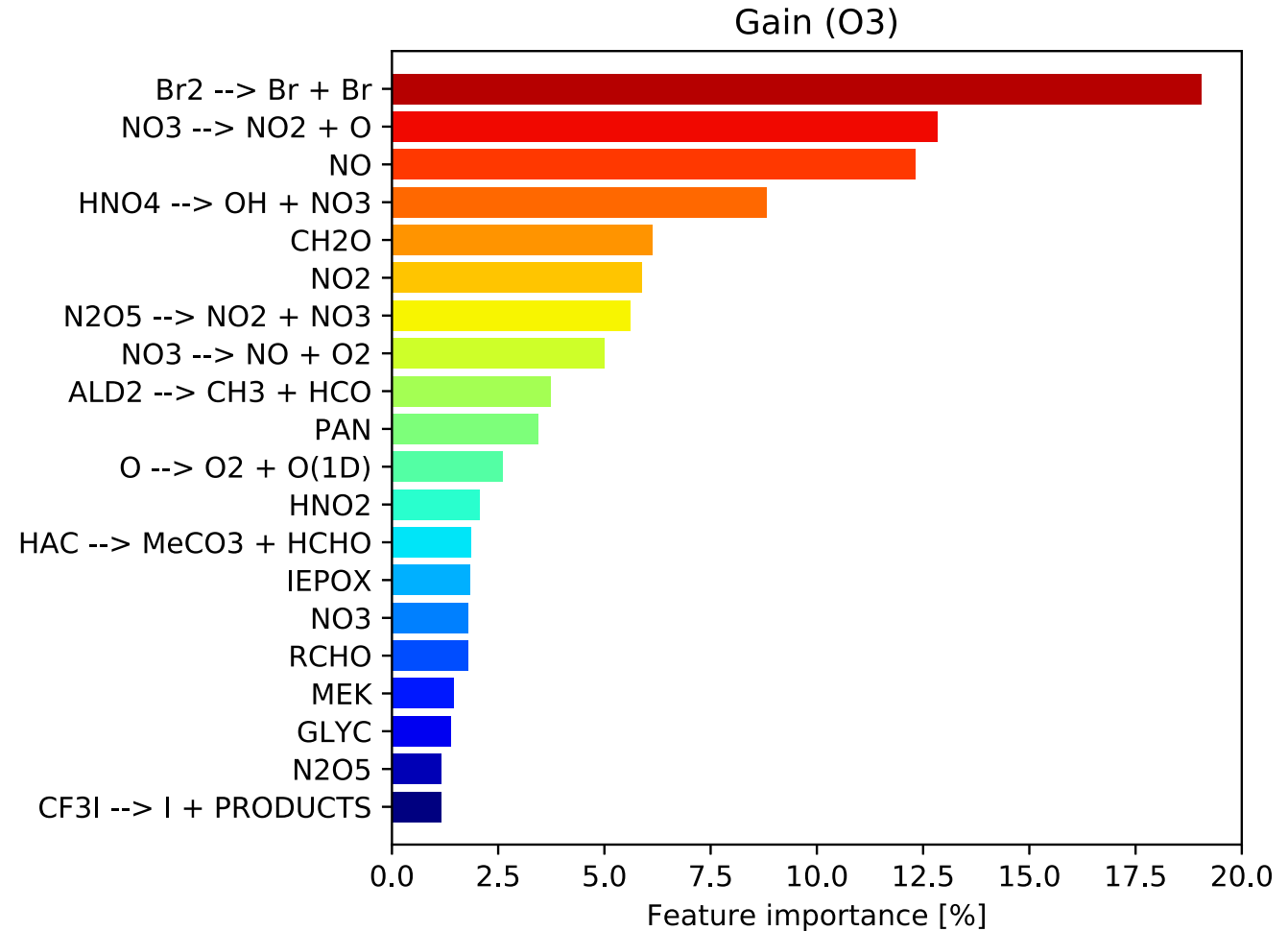
# XGBoost reproduces target concentrations well (single-step prediction)



# XGBoost reproduces target concentrations well (single-step prediction)

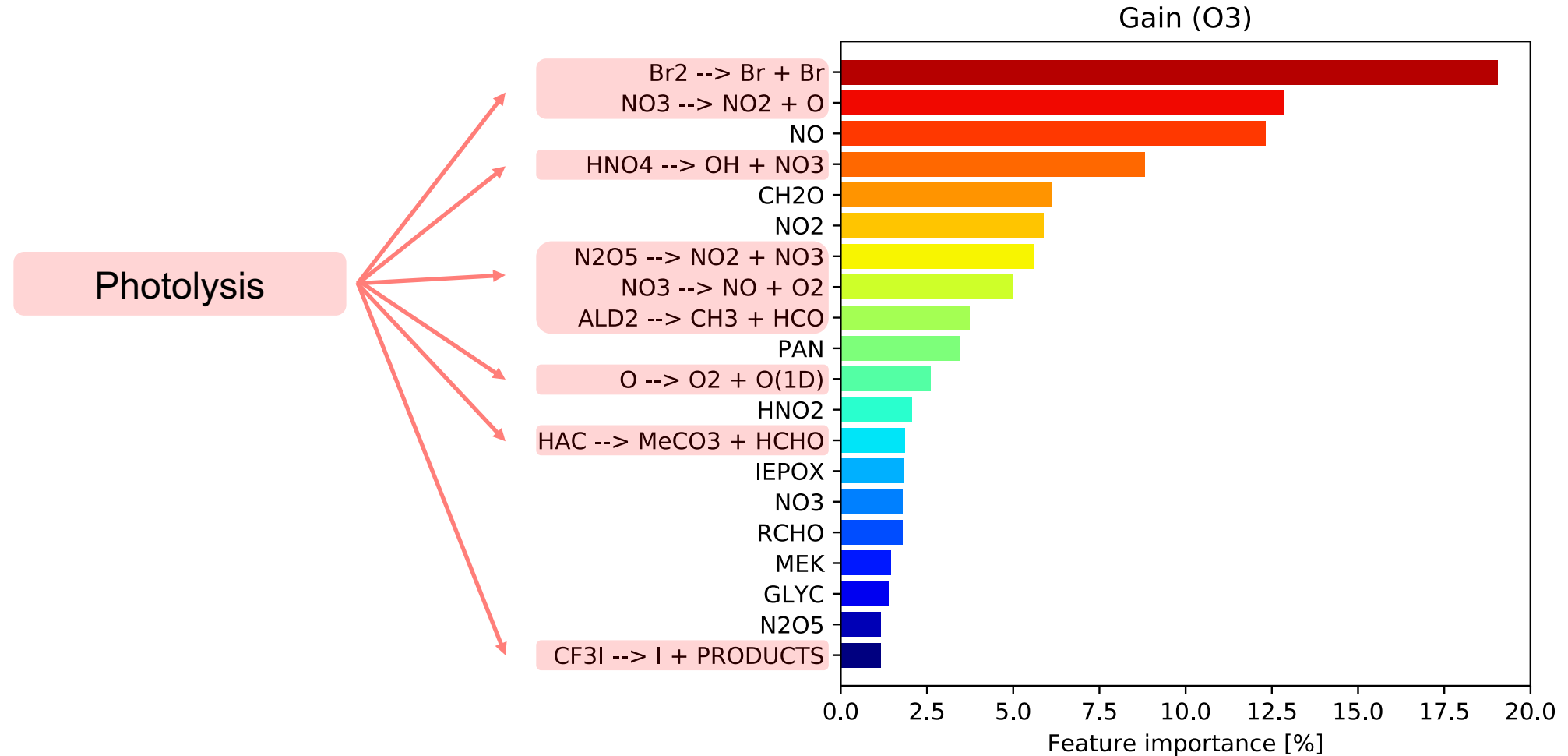


# XGBoost solution reflects known features of chemical kinetics

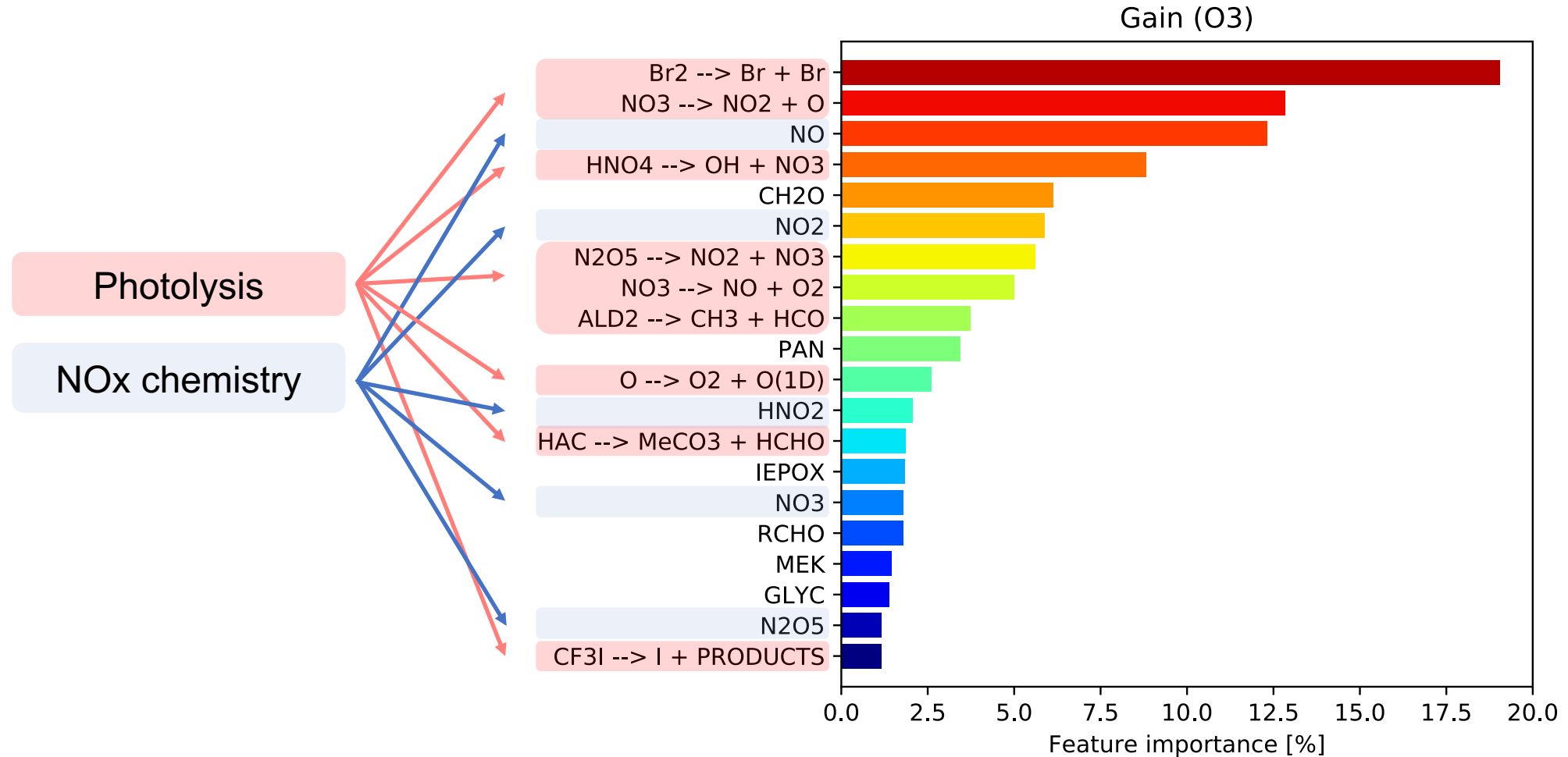




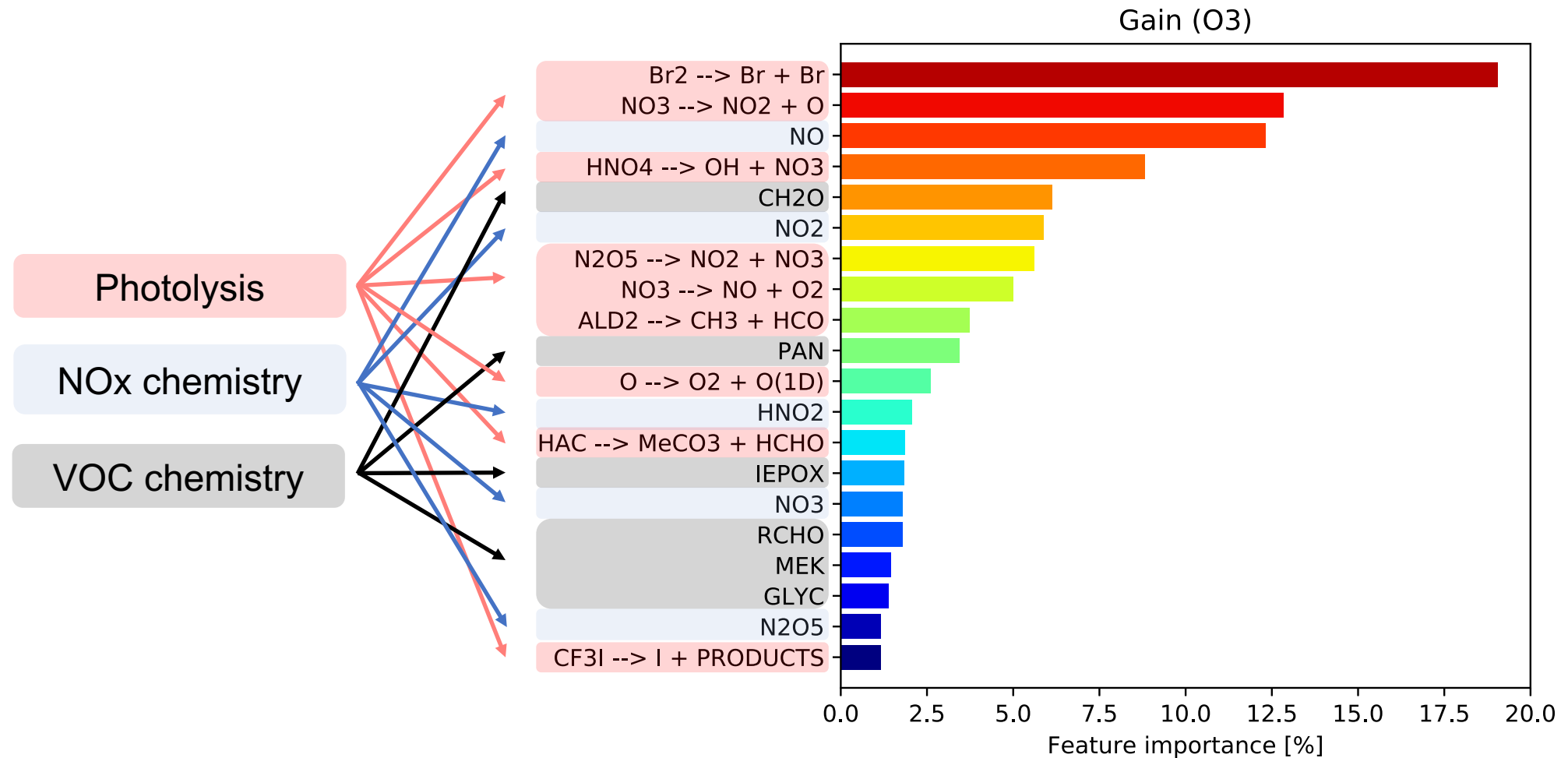
# XGBoost solution reflects known features of chemical kinetics



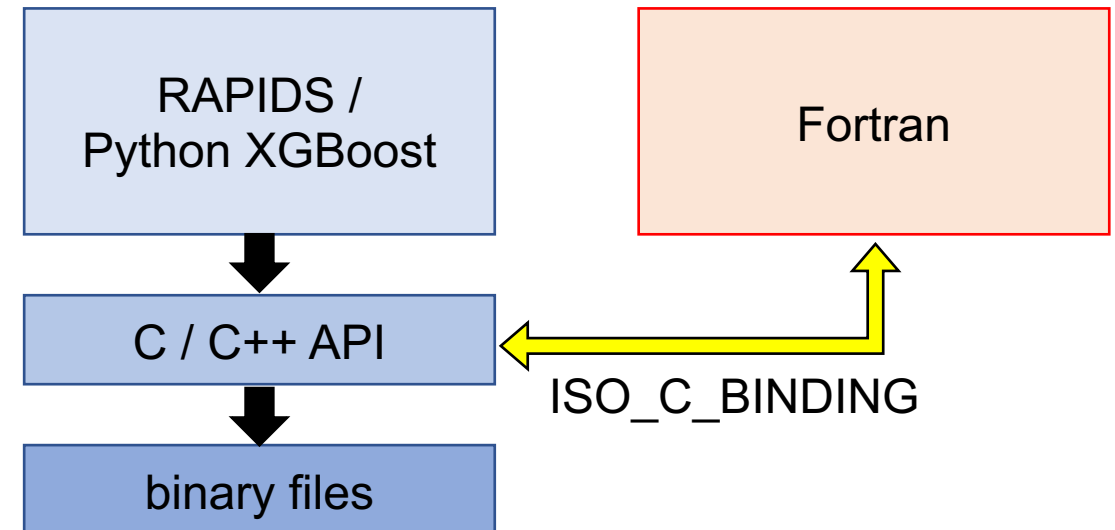
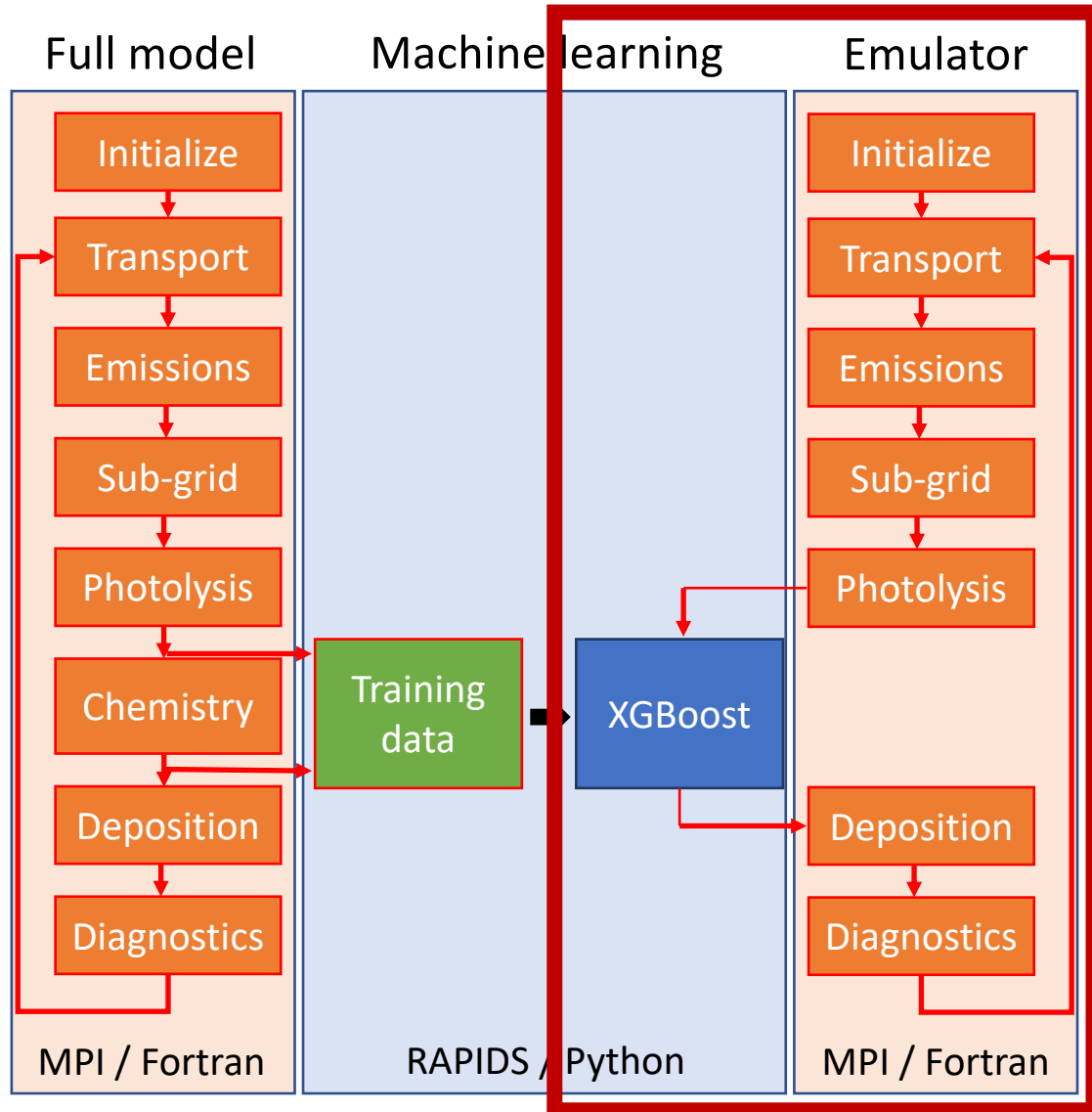
# XGBoost solution reflects known features of chemical kinetics



# XGBoost solution reflects known features of chemical kinetics



# 1-month simulation with XGBoost emulator





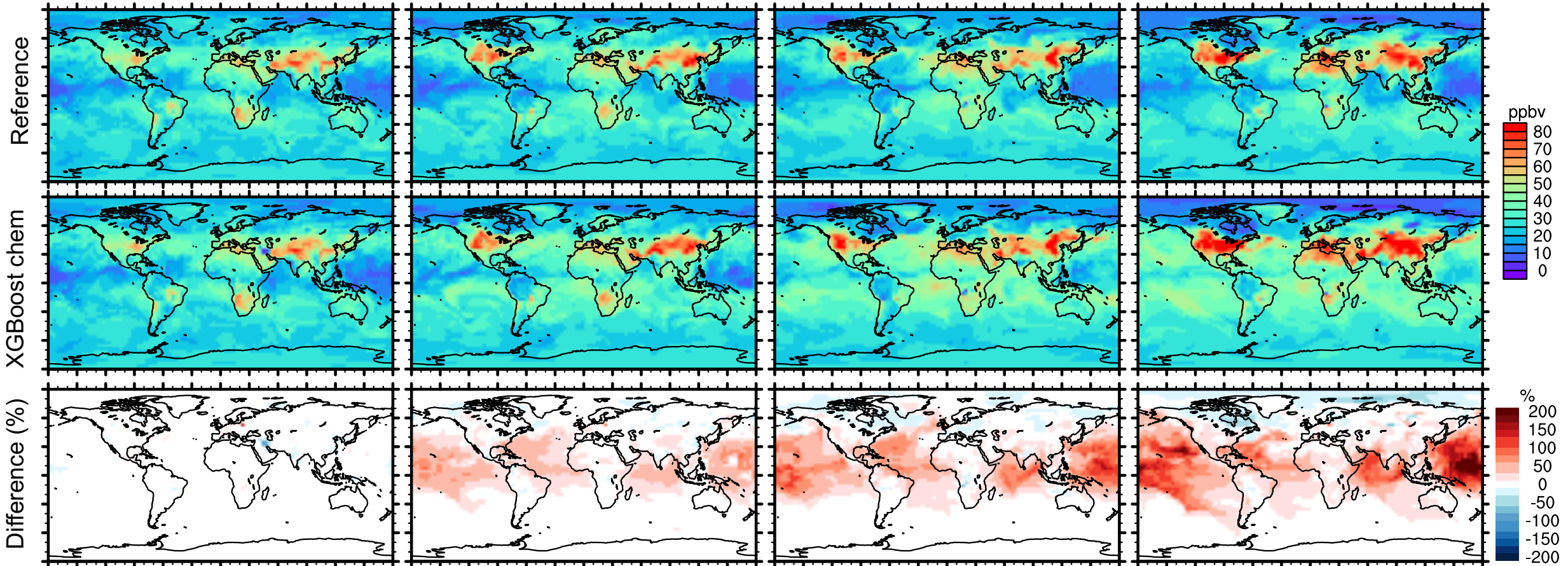
# Emulator model is generally accurate, but overestimates ozone concentrations over remote regions

After 1 day

After 5 days

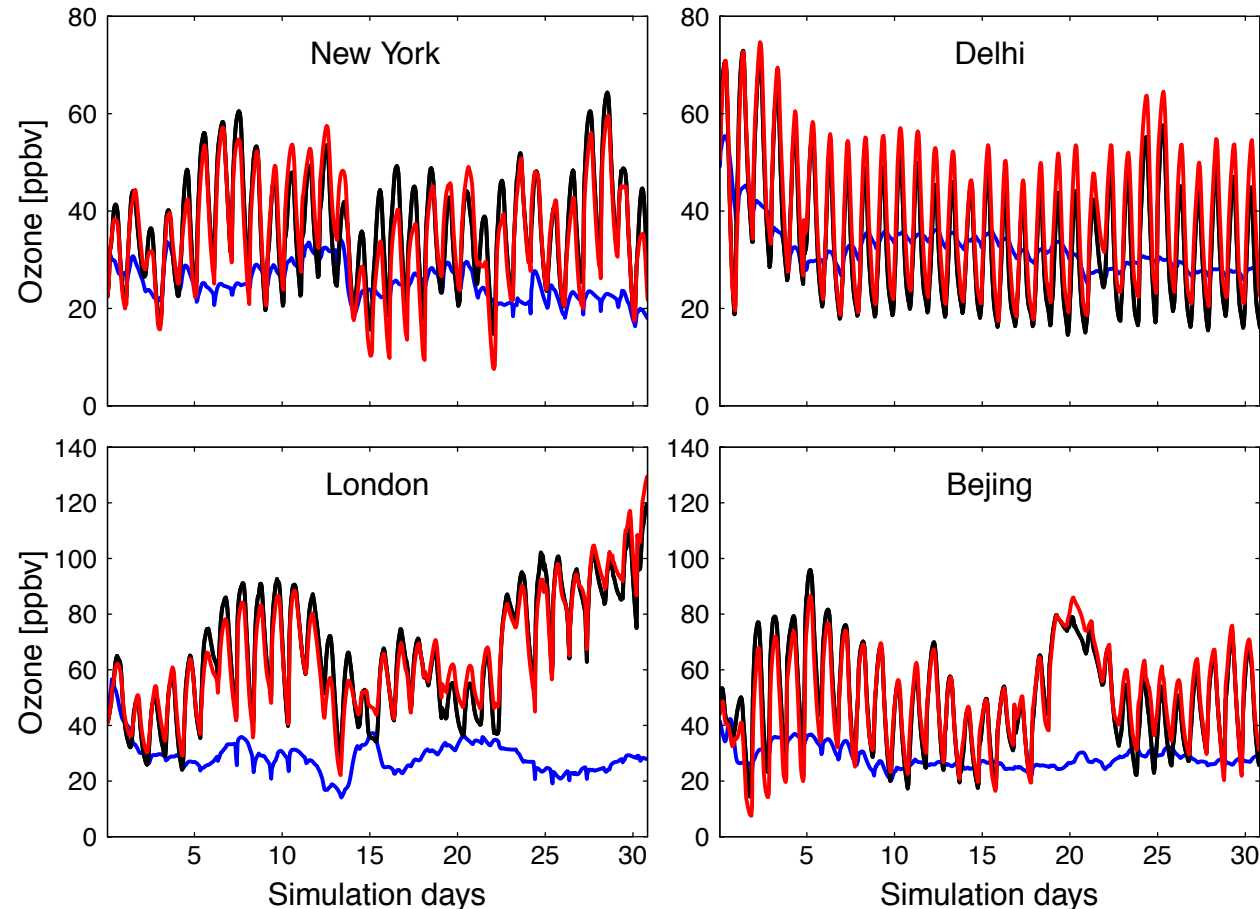
After 10 days

After 30 days

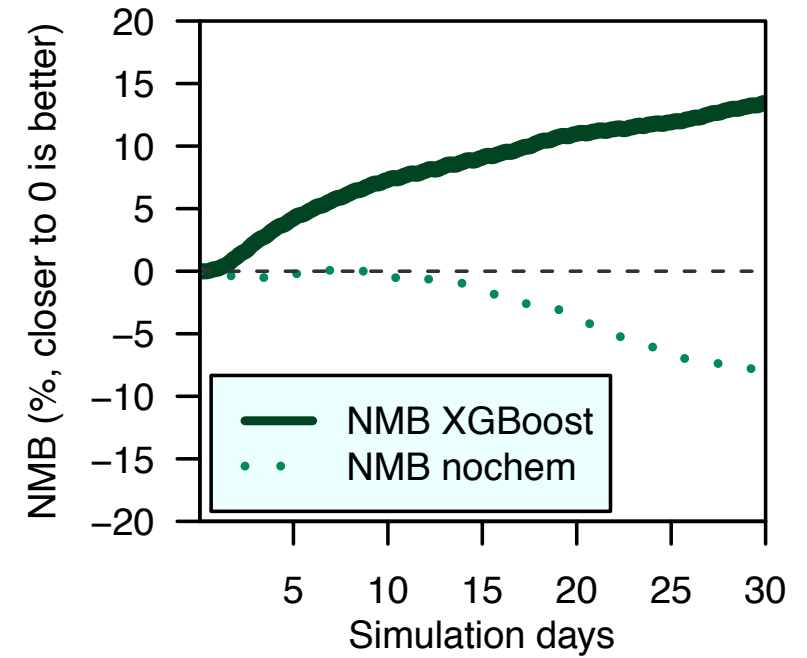
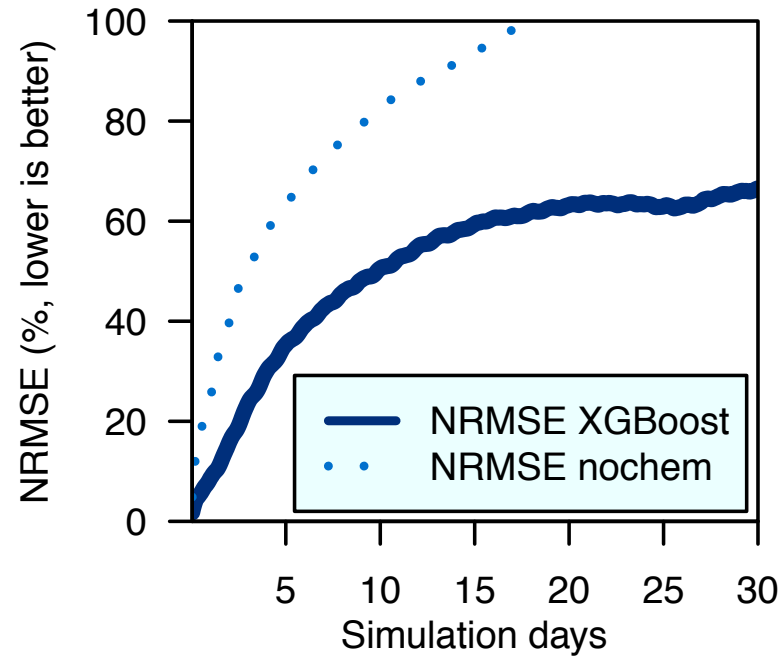
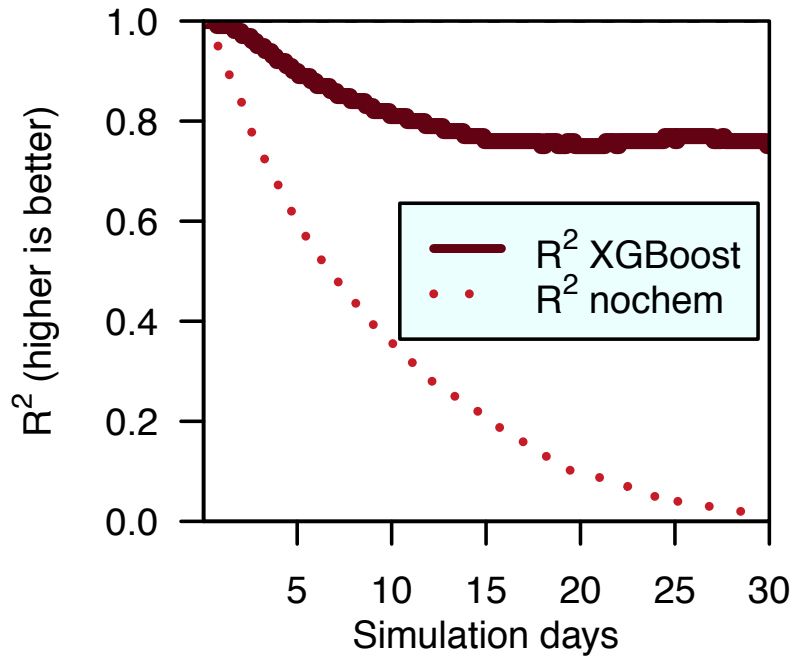


# Surface concentrations over polluted regions are well reproduced by ML model

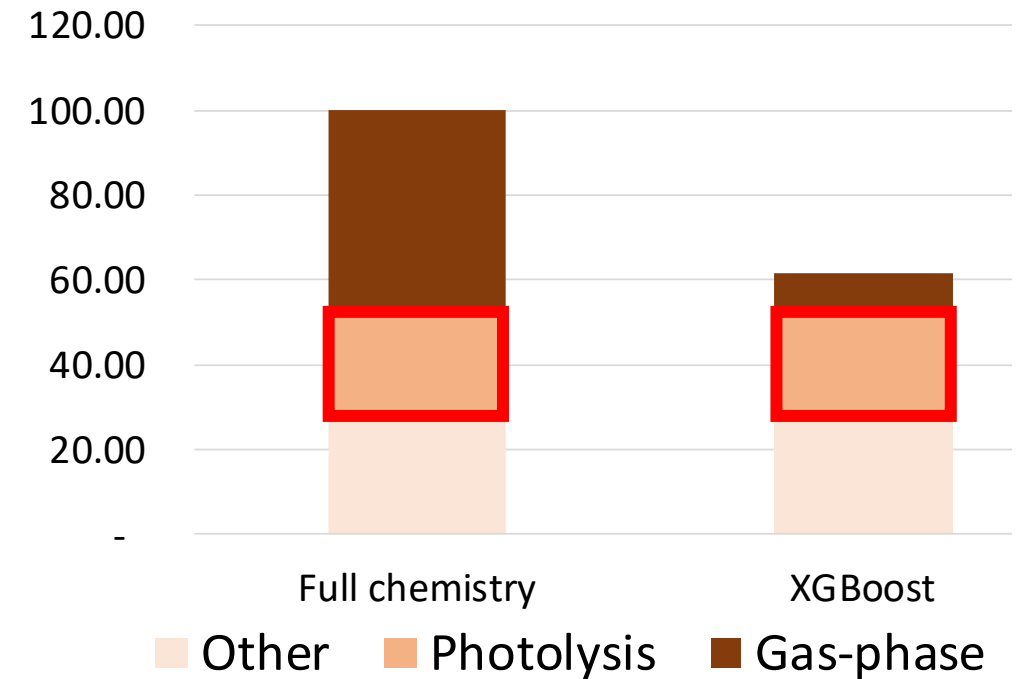
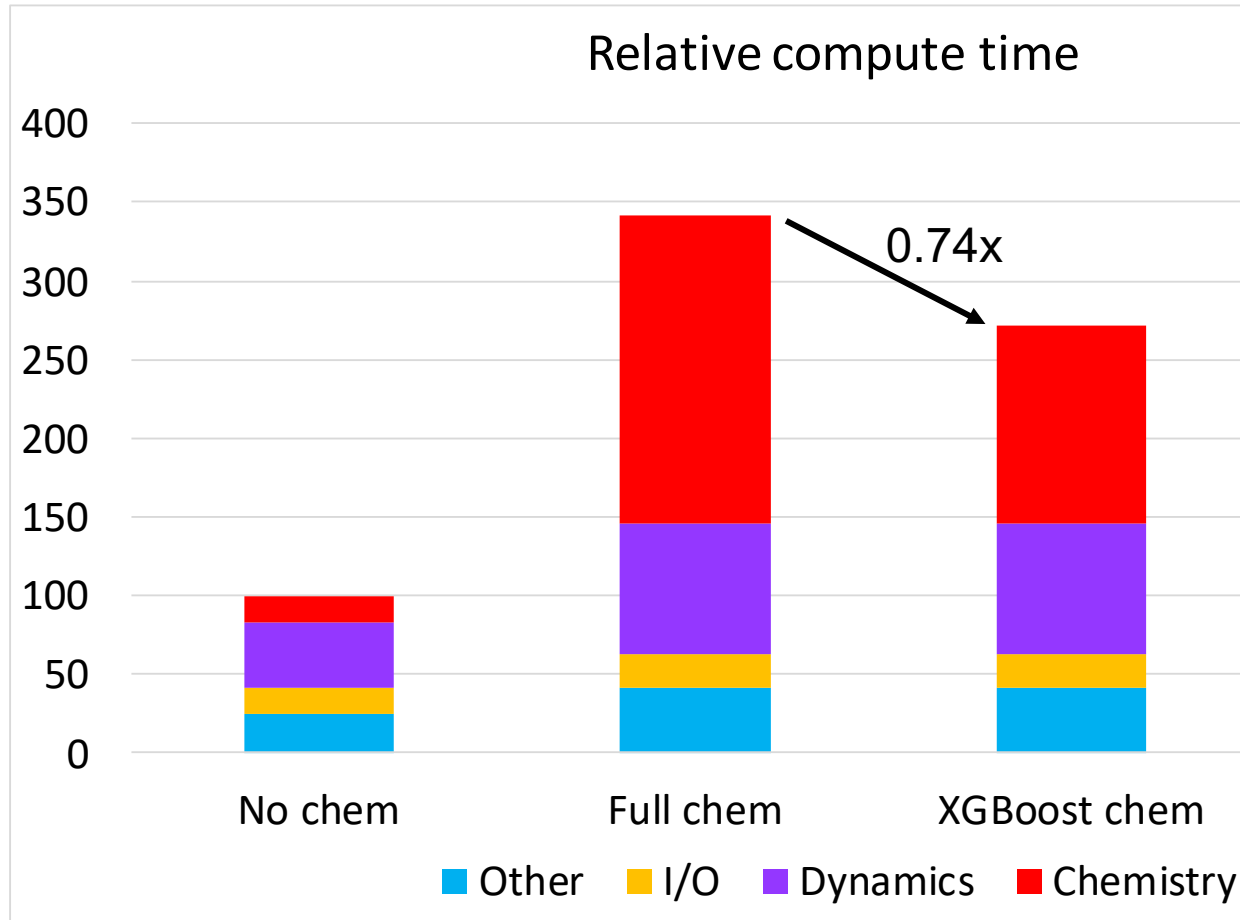
— Reference — XGBoost — No chemistry



# Machine learning model remains stable over the long-term



# Model speedup



- XGBoost model is ~25% faster than reference model
- Chemistry is still slowest part of the model

# Incorporating photolysis calculation into the ML algorithm

Inputs

Meteorology:  
- 7 variables

Chemistry:  
- 143 chemical species  
- 91 photolysis rates

ML

Output

Chemical production / loss

- Original ML algorithm uses as input 91 photolysis rates

# Incorporating photolysis calculation into the ML algorithm

## Inputs

Meteorology:  
- 7 10 variables

Chemistry:  
- 143 chemical species  
- 91 photolysis rates  
- 14 aerosol columns

ML

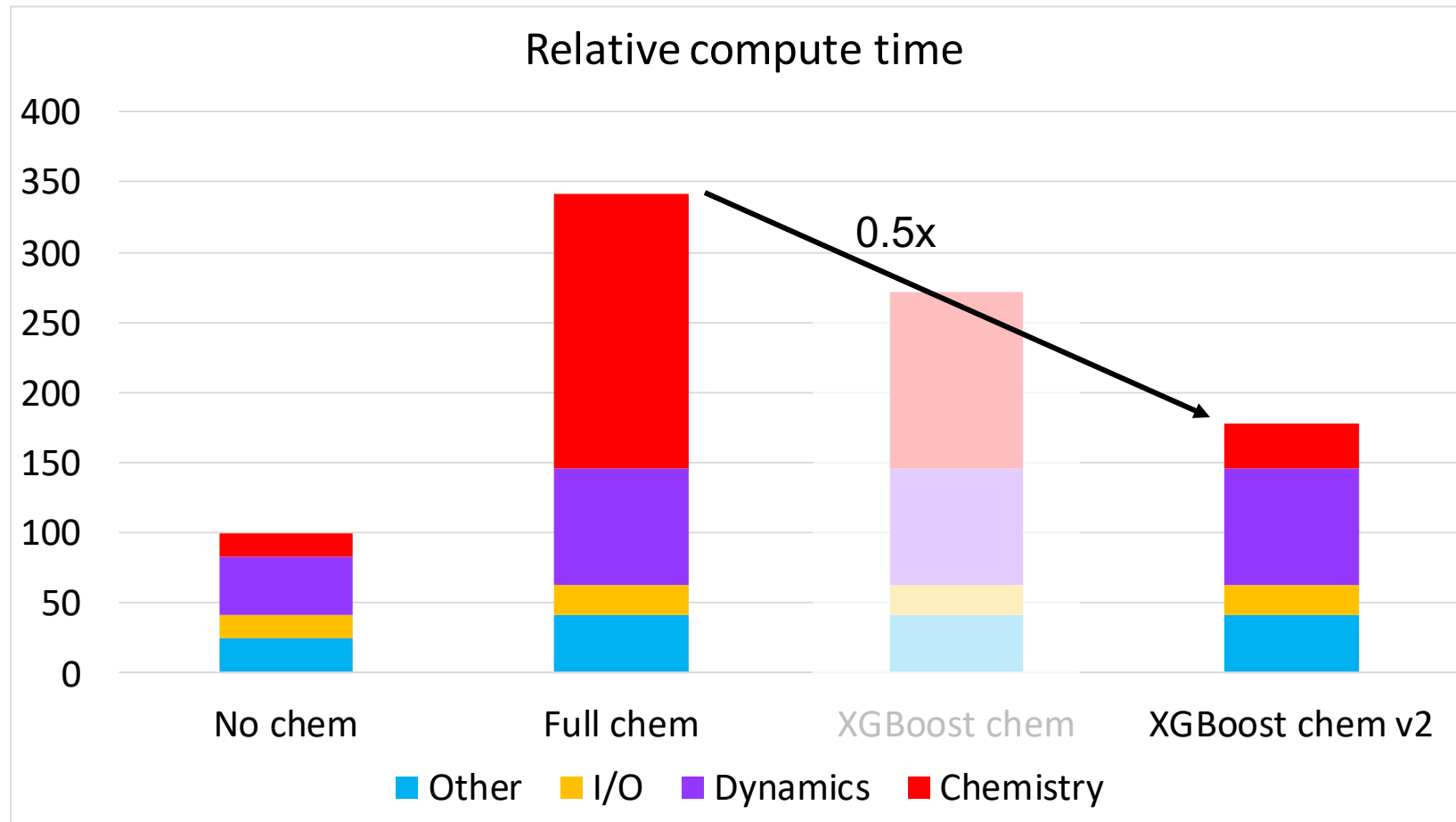


## Output

Chemical production / loss

- Original ML algorithm uses as input 91 photolysis rates
- Replace photolysis rates with quantities needed to compute photolysis

# Model speedup with optimized XGBoost model



- XGBoost chemistry model is now ~2 times faster than reference model
- Chemistry >6x faster than before, dynamics becomes bottleneck



# Summary

- Machine learning can help speed up air quality models by at least 2-5x
- Benefits:
  - Better use of satellite observations
  - Improve (short to medium-term) air quality forecasts
- Ongoing work:
  - Train on very large data sets (>1 TB)
  - Better coupling between CPU and GPUs (model side)
  - 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.

