# Accelerating and Explaining Earth System Process Models with Machine Learning

*David John Gagne*

*National Center for Atmospheric Research*
*Boulder, CO USA*
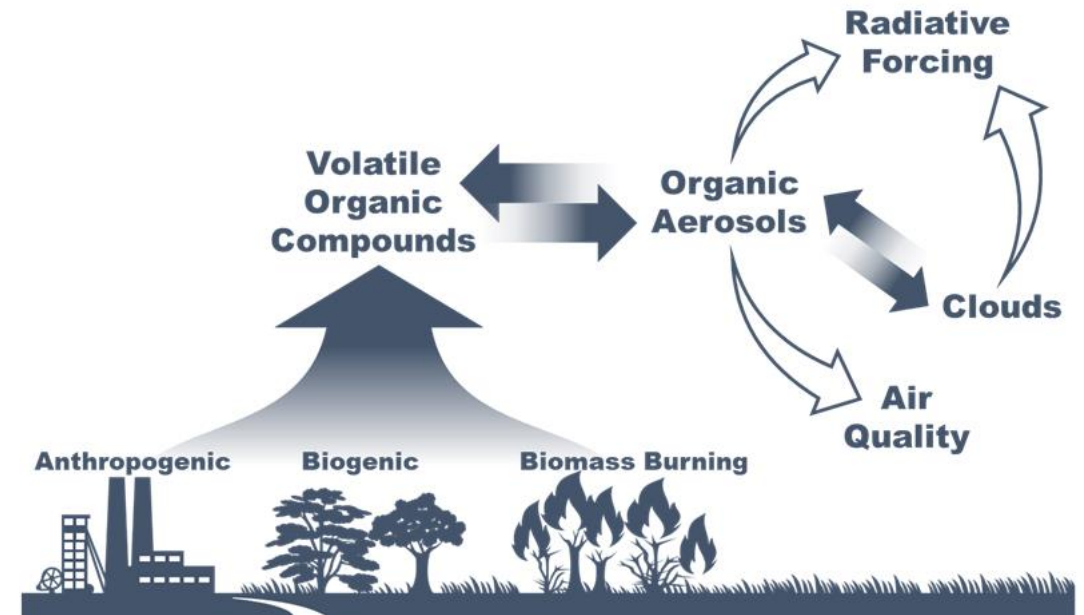
**October 5, 2020**

- Models of small particles in the atmosphere can model bulk properties or small-scale interactions
- Interaction models produce significantly different results from bulk counterparts but are too computationally expensive to run within weather/climate simulations
- Machine learning emulators trained on limited runs from the complex models can approximate them at a far smaller computational cost.
- **Goals**
  – **Develop machine learning emulators for**
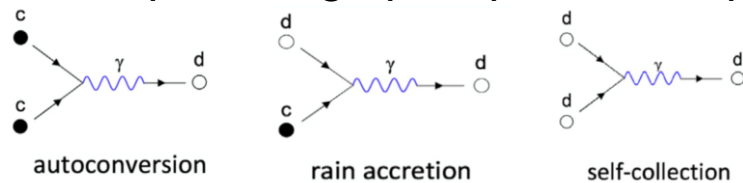
# Machine Learning the Warm Rain Process

Andrew Gettleman, David John Gagne, Chih-Chieh Chen, Matthew Christensen, Zachary Lebo, Hugh Morrison, Gabrielle Gantos

Precipitation formation is a critical uncertainty for weather and climate models.

Different sizes of drops interact to evolve from small cloud drops to large precipitation drops.
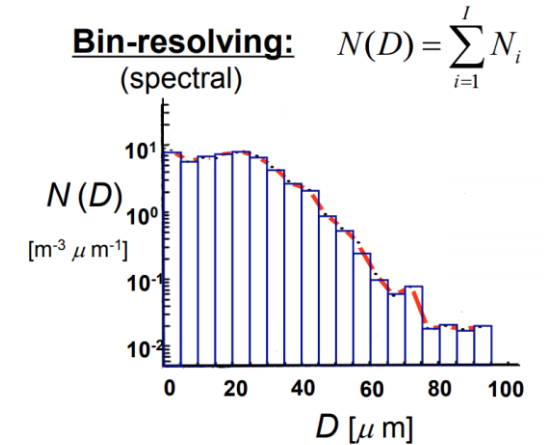


autoconversion    rain accretion    self-collection

Detailed bin codes are too expensive for large scale models, so empirical functions are used instead.
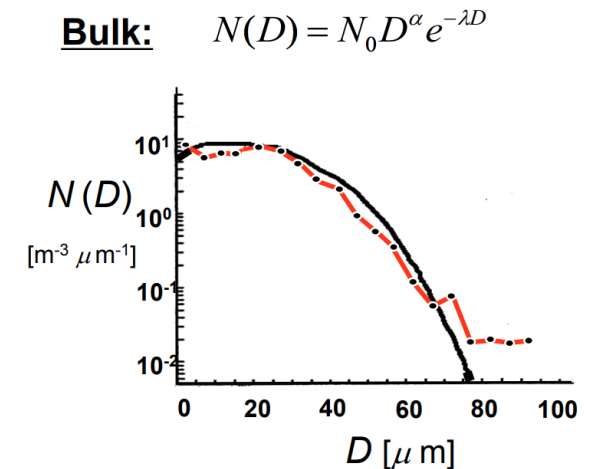
Can a machine learning approach provide a more accurate emulation of precipitation formation processes without a significant increase in computation?

**Goal**: Put a detailed bin process into a global general circulation model and emulate it using ML techniques.
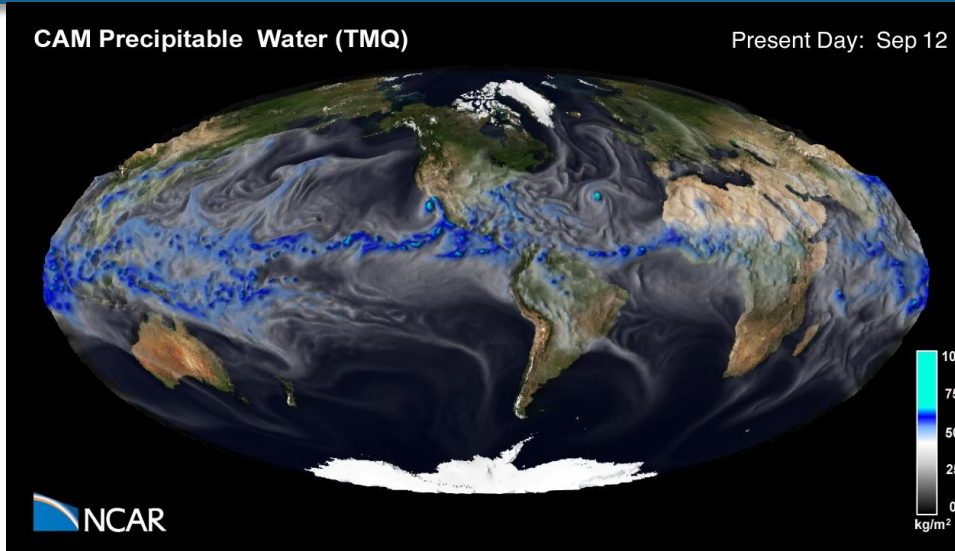
**Bin Scheme (Tel Aviv University (TAU) in CAM6):** Divide particle sizes into bins and calculate evolution of each bin separately.

**Bin-resolving:** (spectral)    $N(D) = \sum_{i=1}^{I} N_i$



**Bulk (MG2 in CAM6):** Calculate warm rain formation processes with a semi-empirical particle size distribution (PSD) based on exponential fit to LES bin microphysics runs.
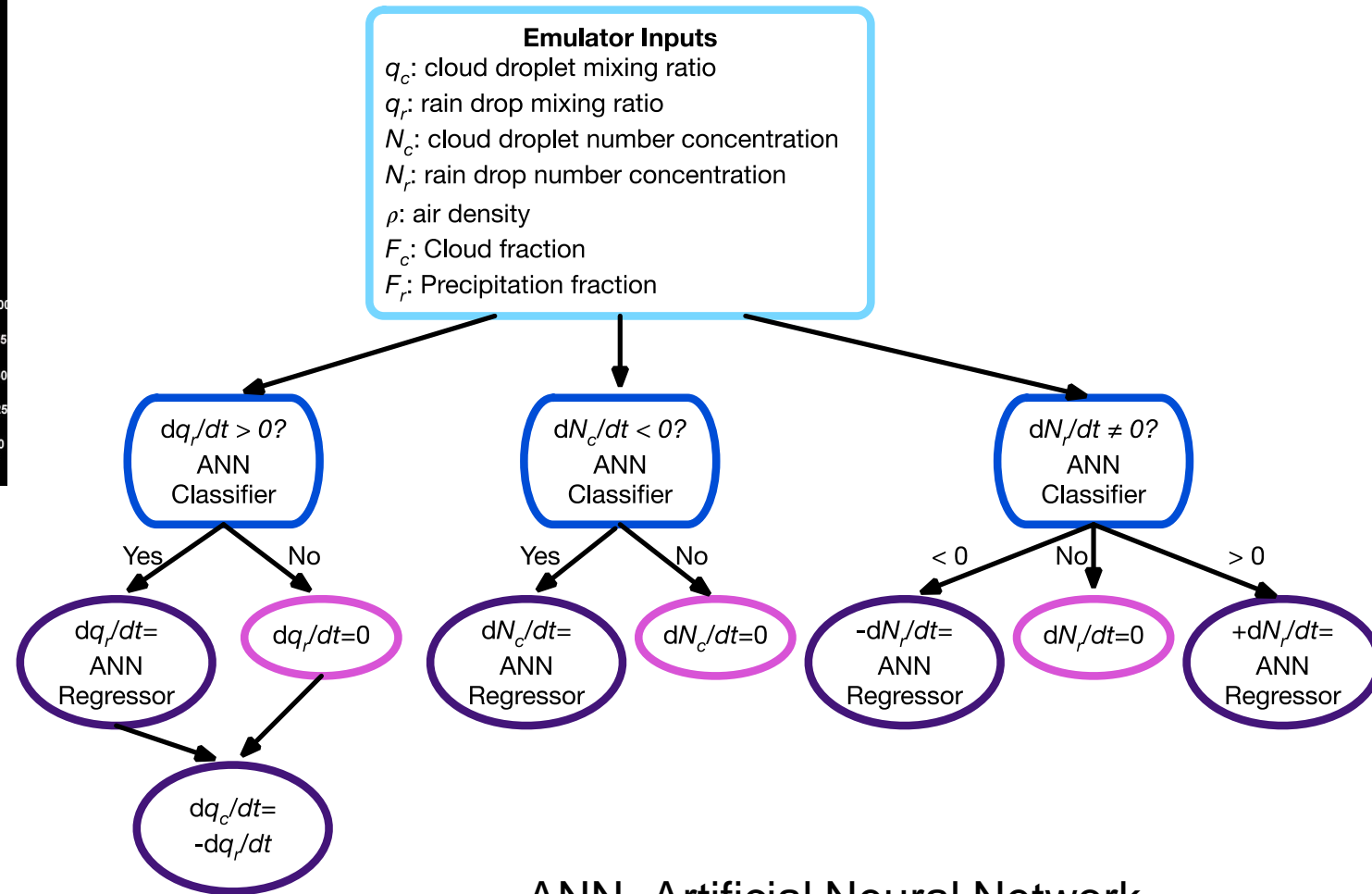
**Bulk:**    $N(D) = N_0 D^\alpha e^{-\lambda D}$

**CAM Precipitable Water (TMQ)**                    Present Day: Sep 12

NCAR

**Emulator Inputs**
$q_c$: cloud droplet mixing ratio
$q_r$: rain drop mixing ratio
$N_c$: cloud droplet number concentration
$N_r$: rain drop number concentration
$\rho$: air density
$F_c$: Cloud fraction
$F_r$: Precipitation fraction

d$q_r$/dt > 0?
ANN
Classifier

d$N_c$/dt < 0?
ANN
Classifier

d$N_r$/dt ≠ 0?
ANN
Classifier

Yes                    No

Yes                    No

< 0        No        > 0

d$q_r$/dt=
ANN
Regressor

d$q_r$/dt=0

d$N_c$/dt=
ANN
Regressor

d$N_c$/dt=0

-d$N_r$/dt=
ANN
Regressor

d$N_r$/dt=0

+d$N_r$/dt=
ANN
Regressor

d$q_c$/dt=
-d$q_r$/dt

ANN=Artificial Neural Network

## Data Generation

1. Run CAM6 for 2 years with fixed forcing from other CESM components
2. Output global microphysics input and output fields every 25 hours
3. Identify grid cells with non-negligible cloud and rain water mixing ratios
4. Save filtered data to csv files
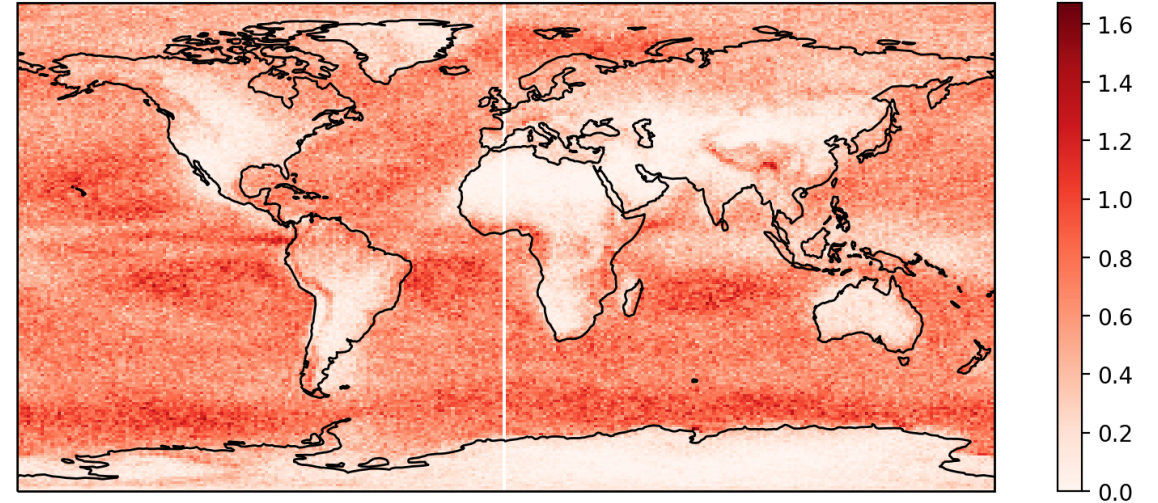5. Logarithmic transform and normalize input and tendency fields

NCAR
UCAR          Contact: dgagne@ucar.edu, @DJGagneDos

# Classifier Results

## Classifier Results

|  | TAU QR 1 | TAU QR 0 | Total |
|---|---|---|---|
| NN QR 1 | **41.7%** | 0.7% | 42.4% |
| NN QR 0 | 0.8% | **56.8%** | 57.6% |
| Total | 42.5% | 57.5% | **98.4%** |

|  | TAU NC 1 | TAU NC 0 | Total |
|---|---|---|---|
| NN NC 1 | **52.9%** | 0.5% | 53.4% |
| NN NC 0 | 0.2% | **46.3%** | 46.5% |
| Total | 53.1% | 46.8% | **99.3%** |

|  | TAU NR -1 | TAU NC 0 | TAU NR 1 | Total |
|---|---|---|---|---|
| NN NR -1 | **35%** | 0.0% | 0.4% | 35.4% |
| NN NR 0 | 0.1% | **43.1%** | 0.3% | 43.5% |
| NN NR 1 | 0.2% | 0.5% | **20.4%** | 21.1% |
| Total | 35.3% | 43.6% | 21.1% | **98.5%** |



$dq_r/dt$ Classifier False Positive Relative Frequency (%)



$dq_r/dt$ Classifier False Negative Relative Frequency (%)

# Regressor Results

| Output | $R^2$ | MAE | Hellinger |
|--------|-------|-----|-----------|
| $dq_r/dt$ | 0.991 | 0.095 | 4.53e-4 |
| $dn_c/dt$ | 0.995 | 0.112 | 1.49e-3 |
| $dn_r/dt < 0$ | 0.995 | 0.081 | 6.04e-4 |
| $dn_r/dt > 0$ | 0.978 | 0.178 | 1.18e-3 |

- CAM6: Control
- TAU or TAU-bin: Stochastic Collection Kernel
- TAU-ML: Machine learning Emulator for TAU code

- For each, global 0.9°x1.25° simulation, 9 years, 1st year high frequency instantaneous output
  - Base (2000 Climatology)
  - Pre-Industrial (1850) aerosols. (For aerosol cloud interactions)
  - SST+4K (For Cloud Feedbacks)

NCAR
UCAR

A) Frequency of Mass Fixer Activation

B) Annual Mean Mass Fixer Frequency 936hPa

How often does mass fixer kick in and where?
- Low altitudes and tropical high altitudes (cirrus)
- Low altitude (below is 936hPa), mostly in sub-tropical strato-cumulus regions, edge of stratus regions. Mostly SH.
- Also a tropical peak at 800hPa

# Precipitation Feedbacks

# Cloud Feedbacks



- ACI are similar between control and TAU code.
- Slightly lower LWP change, but forcing is similar, a bit higher in S. Hemisphere.
- Emulator reproduces TAU results.

NCAR
UCAR

Goal: understand average sensitivities of input fields while accounting for nonlinear interactions within model

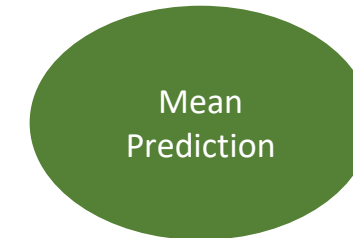1. Set all instances for one variable in a dataset to a single value

| Temperature | Dewpoint | Pressure |
|---|---|---|
| **280** | 10 | 986 |
| **280** | 14 | 1014 |
| **280** | 2 | 992 |
| **280** | 25 | 1025 |
| **280** | 6 | 950 |

2. Feed fixed data through model

Machine Learning or Physical Model

3. Calculate mean prediction for fixed value

Mean Prediction

4. Repeat process for range of input values

potential temperature 10 m

Outside the range of the training data (red) the neural network extrapolates mostly linearly

- Neural network emulator set largely replicates the behavior of the TAU bin microphysics warm rain processes
- Successfully ran in CAM6 in training climate
- Both tendencies and feedbacks from emulator closely match original scheme

**Challenges**

- Running in future and pre-industrial climates results in more calls to mass fixer
- Linear extrapolation behavior may not be appropriate for certain variables. How to constrain?
- Training superdroplet scheme emulator

NCAR
UCAR

# Machine Learning Emulation of the GECKO-A Chemistry Model

David John Gagne, Charlie Becker, John Schreck, Keely Lawrence, Siyuan Wang, Alma Hodzic
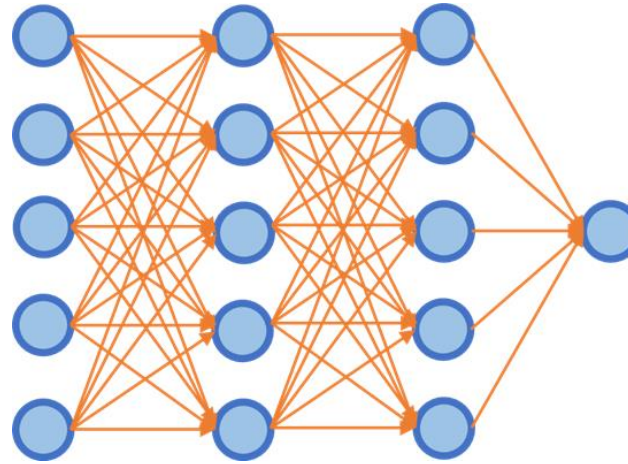
NCAR
UCAR

# GECKO-A Challenge: Build An Emulator For 3-D Models?
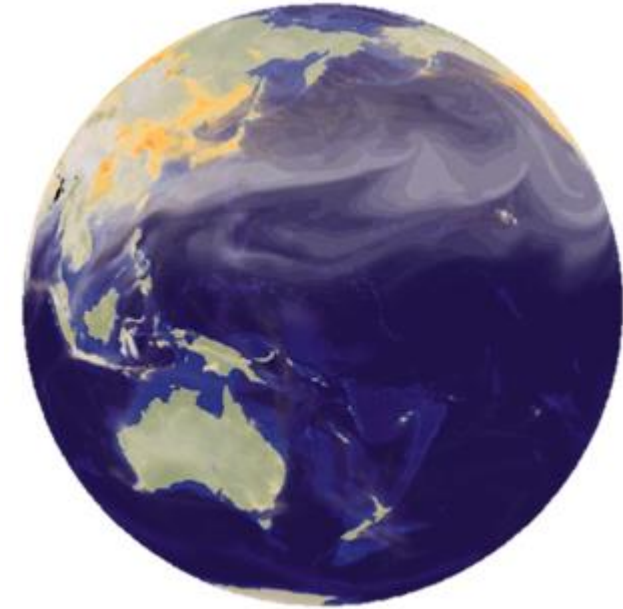
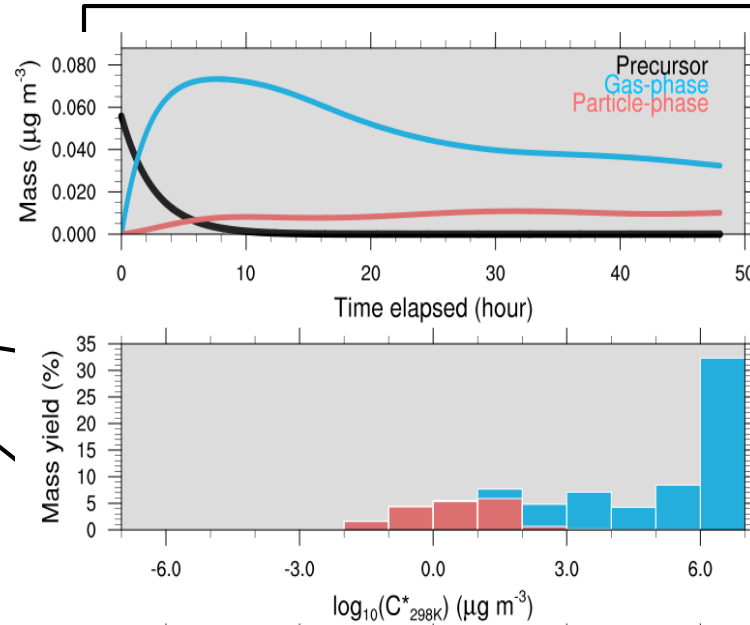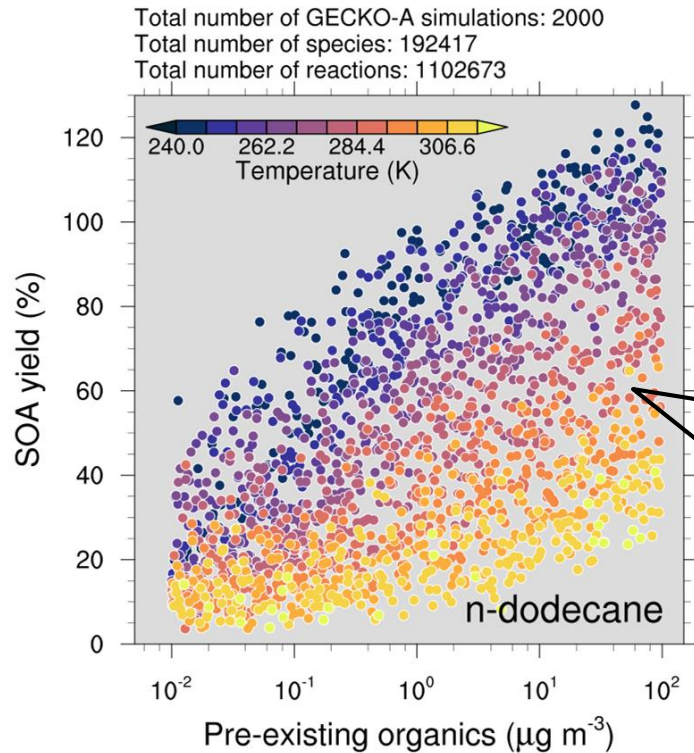**GECKO-A Training Library**



**Machine-Learning Emulator**



**3-D Models**



- Many inspiring applications out there: machine-learning emulators using explicit/process-level models, and implementing the trained emulators into large-scale models. Such explicit/process-level models are otherwise too expensive for large-scale models.

- The goal of this project is to train the machine-learning emulator using the "library" generated by the hyper-explicit chemical mechanism, GECKO-A.

# Goal: Build Emulator to Predict the Total Organic Aerosol



Total number of GECKO-A simulations: 2000
Total number of species: 192417
Total number of reactions: 1102673

**Demo: what the data looks like**

**GECKO-A Library:**
- 2000 GECKO-A simulations: in each run, we run GECKO-A under certain condition for 5 days
- 2000 input files (csv).
- Each file contains: (i) mass of precursors; (ii) mass of products in the gas-phase; and (iii) mass of products in the particle-phase. All (i)-(iii) as a function of time.

NCAR
UCAR

# GECKO Data

**Metadata**

| Metadata | Units | Label |
|---|---|---|
| Number Experiments | 2000 | id |
| Total Timesteps | 1440 | Time |
| Timestep Delta | 300 seconds | - |

**Potential Input Variables**

| Variable Name | Units | Type |
|---|---|---|
| Precursor | ug/m3 | Varies |
| Gas | ug/m3 | Varies |
| Aerosol | ug/m3 | Varies |
| Temperature | K | Static |
| Solar Zenith Angle | degree | Static |
| Pre-existing Aersols | ug/m3 | Static |
| o3 | ppb | Static |
| nox | ppb | Static |
| oh | 10^6 molec/cm3 | Static |

**Potential Output Variables**

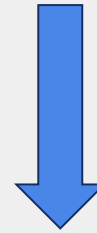| Variable Name | Units | Type |
|---|---|---|
| Precursor (at t+1) | ug/m3 | Varies |
| Gas (at t+1) | ug/m3 | Varies |
| Aerosol (at t+1) | ug/m3 | Varies |

- Use fixed environmental conditions and concentration of precursor, gas and aerosol for a given precursor type
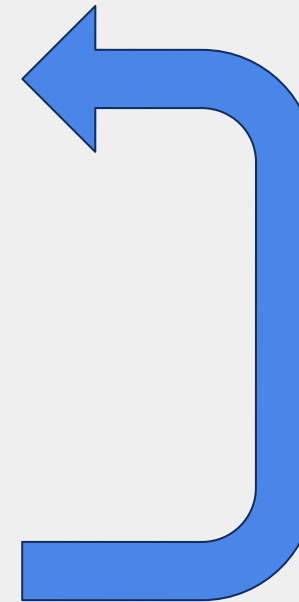- Generated data for toluene, dodecane, and alpha-pinene

# Base Model

INPUT$_{(t)}$

$\downarrow$

OUTPUT$_{(t+1)}$

# Box Emulator Model

STARTING CONDITIONS

$\downarrow$

BASE MODEL
INPUT$_{(t)}$

$\downarrow$

BASE MODEL
PREDICTION$_{(t+1)}$

Loop for length
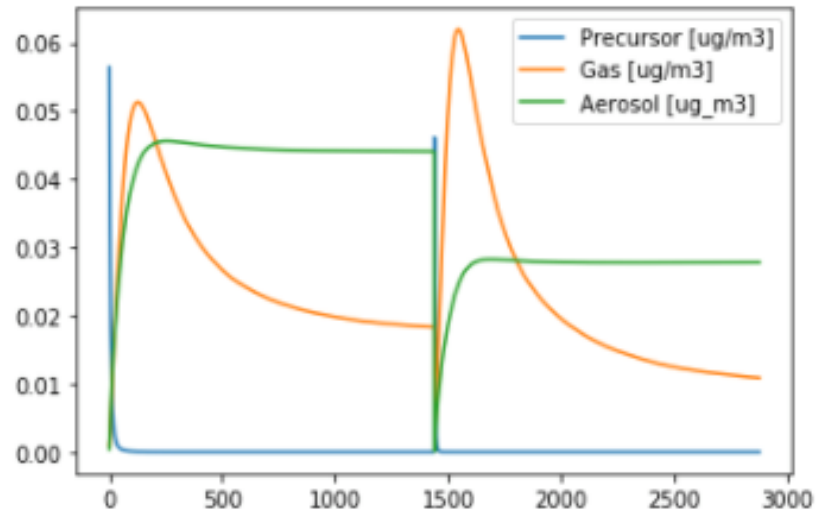of experiment

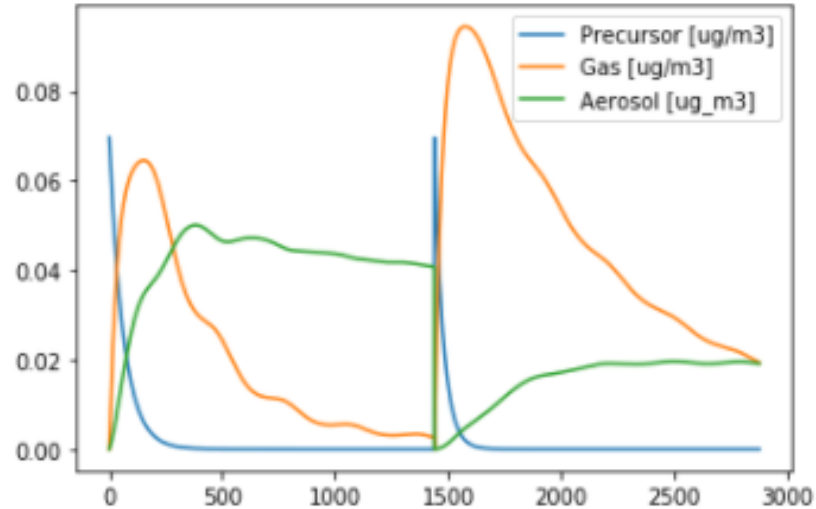Ensemble Runs - dodecane



LSTM FRAMEWORK:
- Recurrent neural network combined with 1D convolutions through time (depends on previous 20 time steps to predict single future time step)
- Trained on 1600, 5-day Experiments (300s time steps) - validated on 200 experiments

CHALLENGES:
- Recurrent networks tend to prevent runaway error propagation but have major challenges incorporating them into a 3D transport model
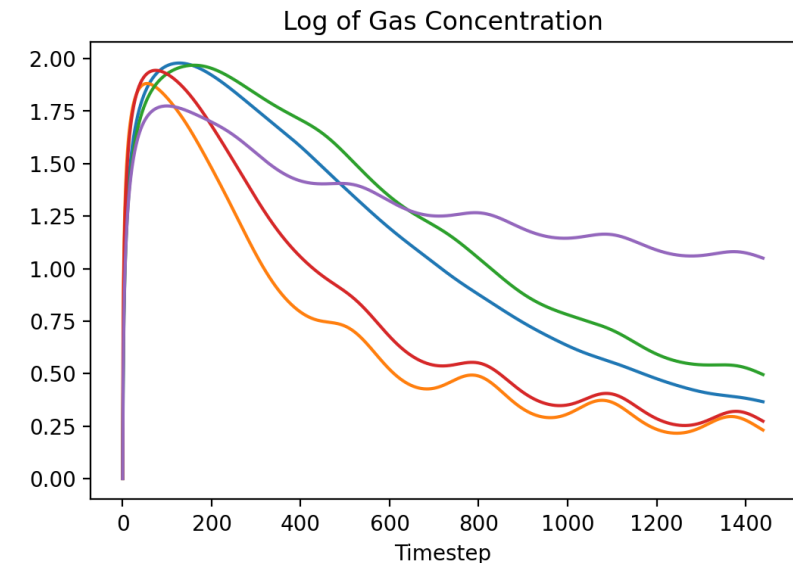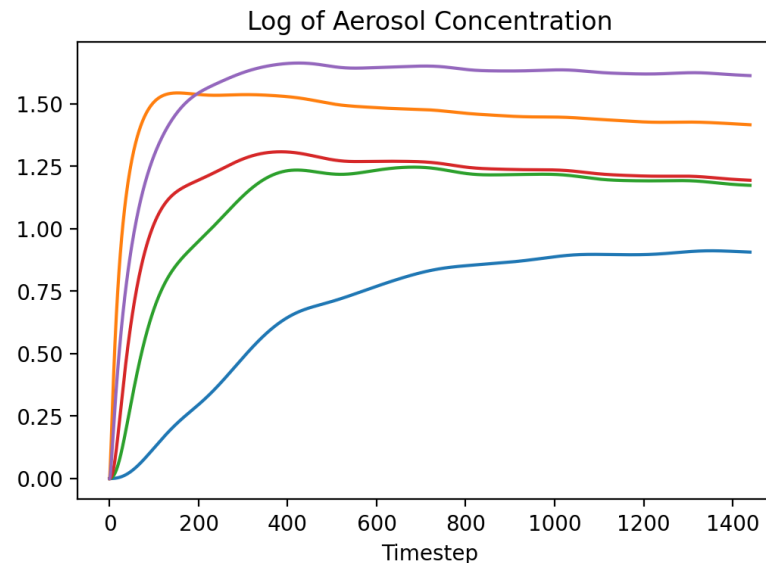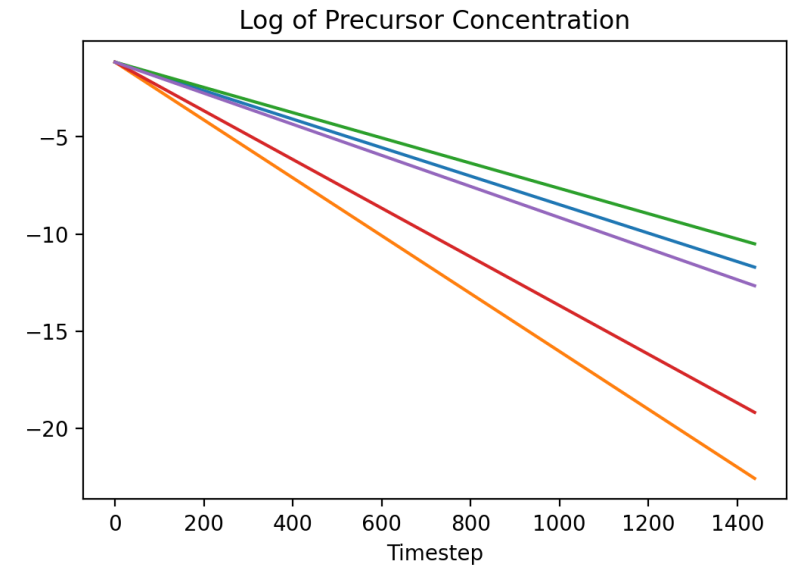
NCAR
UCAR

SINGLE TIMESTEP FRAMEWORK:
- Fully-connected single-layer neural network with SELU activation function
- Trained on single timestep input

CHALLENGES:
- Machine learning captures early changes fine but struggles with later parts of run
- Performs well in offline tests but not in emulator box model
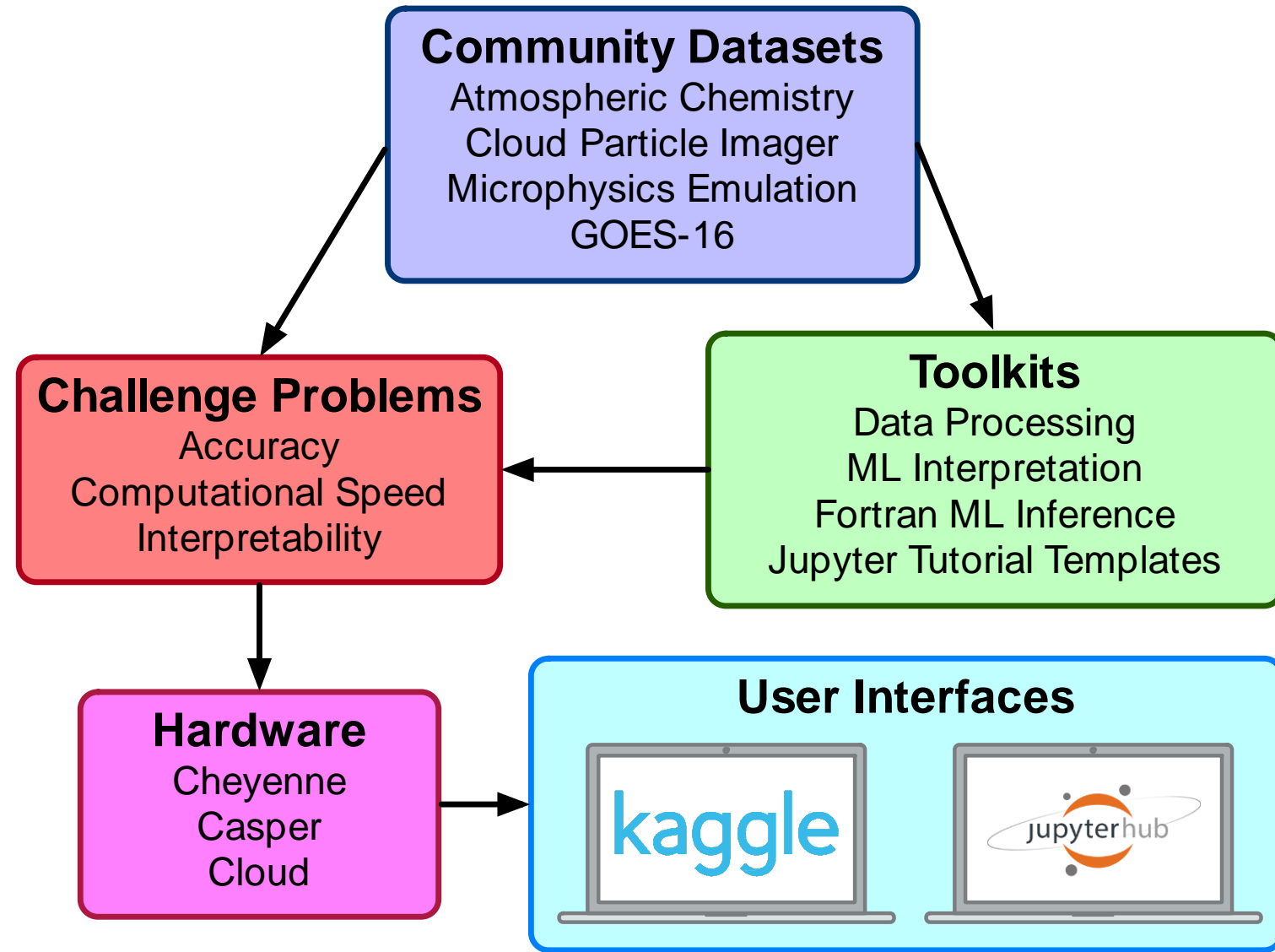
NCAR
UCAR

- Large magnitude difference in both absolute values and tendencies between early and later parts of each simulation
- Machine learning optimization biased toward adjusting initial values more due to larger errors
- Precursor decreases exponentially but gas and aerosol have more complex variability pattern



Log of Precursor Concentration



Log of Aerosol Concentration



Log of Gas Concentration

# AI for Earth System Science Hackathon Motivation

- Interest in AI/ML for weather and climate problems is growing rapidly
- Earth System Science practitioners need help getting started with ML
- Not enough trained ML-ESS experts to work with everyone
- **Solution:** Host a summer school and hackathon!
- Invite AI-ESS experts to lecture about different aspects of AI and ESS
- Create training materials and domain-focused challenge problems

**Community Datasets**
Atmospheric Chemistry
Cloud Particle Imager
Microphysics Emulation
GOES-16

**Toolkits**
Data Processing
ML Interpretation
Fortran ML Inference
Jupyter Tutorial Templates

**Challenge Problems**
Accuracy
Computational Speed
Interpretability

**Hardware**
Cheyenne
Casper
Cloud

**User Interfaces**

kaggle

jupyterhub

- Give participants machine learning experience with realistic ESS data and problems
- Provide them with sufficient computational resources to train more complex models
- Work collaboratively with a new team with diverse backgrounds
- Originally planned to be in person at Mesa Lab, but COVID happened
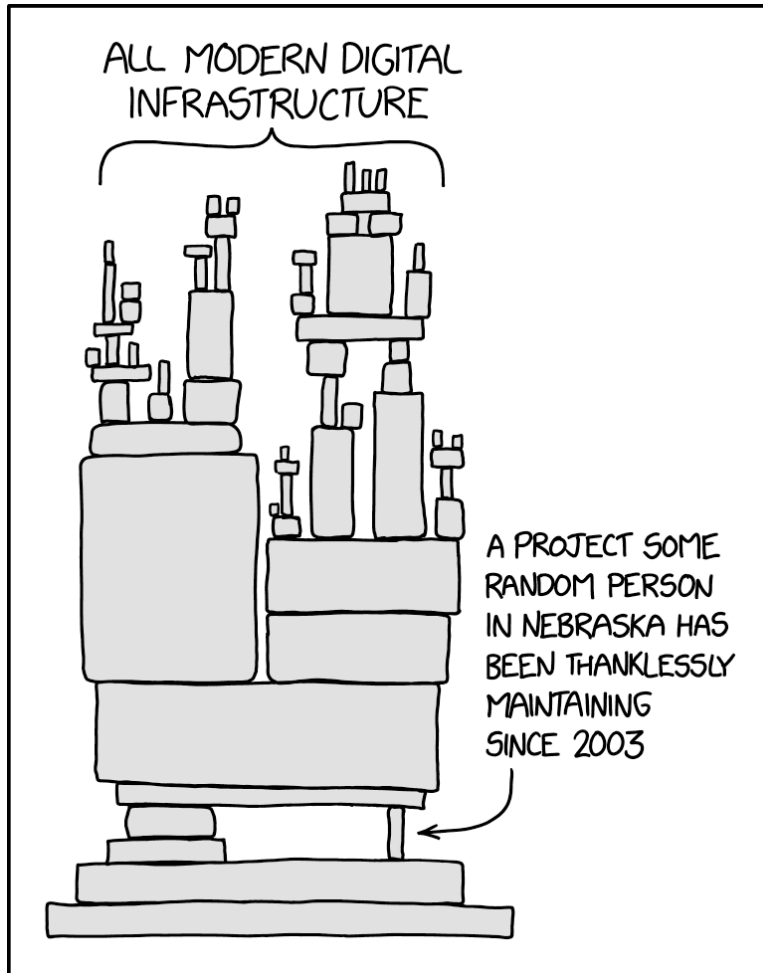- Virtual hackathon allowed greatly expanded participation (80->200)

## User Interface

- Jupyterlab in web browser after logging in with Google credentials.
- Jupyterlab is preinstalled with full python data science environment and access to a GPU.
- User can save data in virtual machine that persists over lifetime of hackathon.
- Users can also run challenge problems through Google Colab notebooks
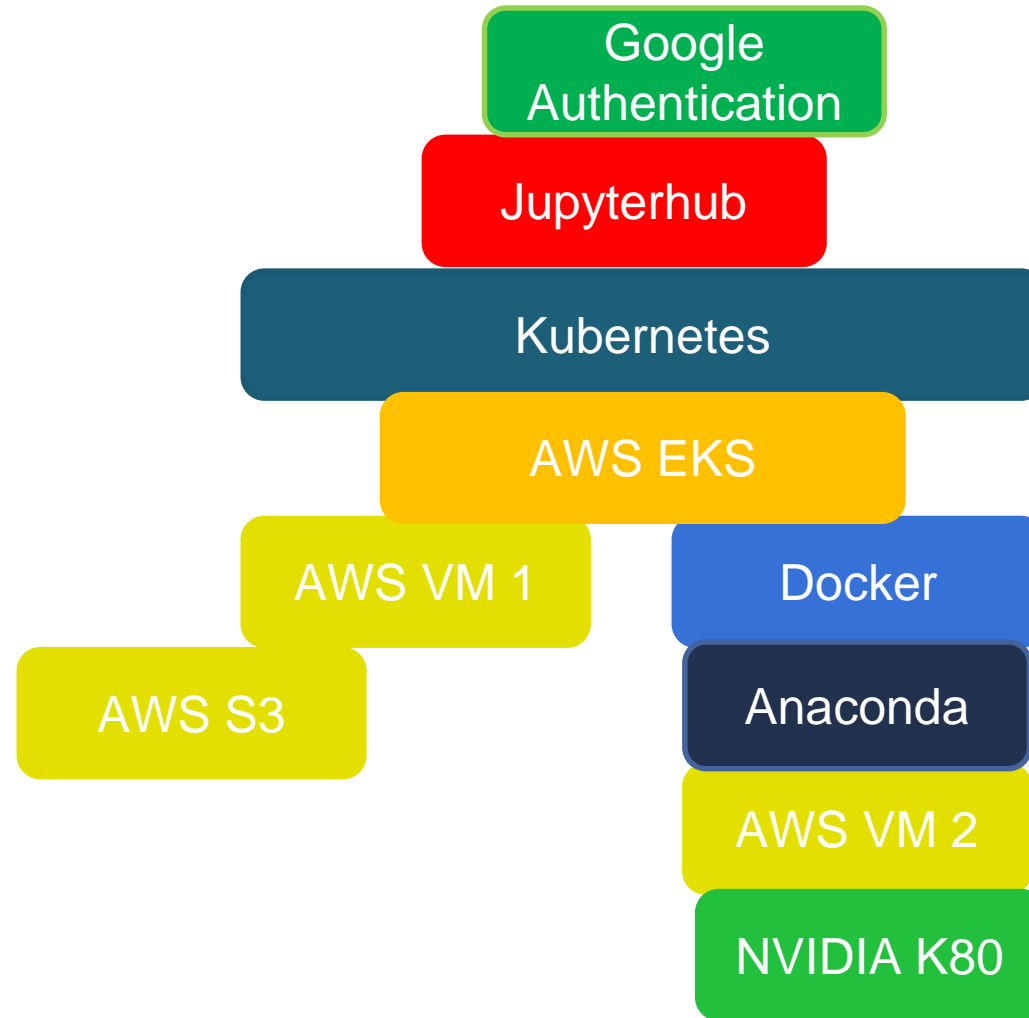
## Team Setup

- Users communicate with other team members through Slack
- Teams of 5 were assigned randomly from the registration info.
- Scheduled hackathon period from 2 to 6 PM Mountain Time each day

# Hackathon Cloud Infrastructure

From xkcd.com

ALL MODERN DIGITAL INFRASTRUCTURE

A PROJECT SOME RANDOM PERSON IN NEBRASKA HAS BEEN THANKLESSLY MAINTAINING SINCE 2003

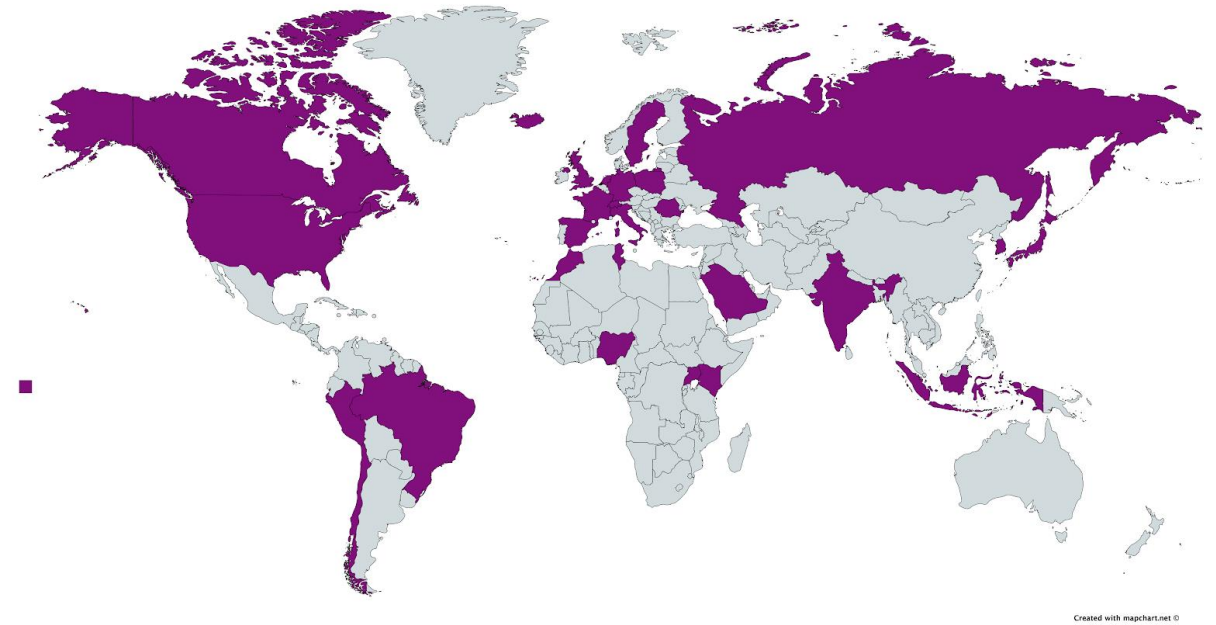Google Authentication
Jupyterhub
Kubernetes
AWS EKS
AWS VM 1
Docker
Anaconda
AWS S3
AWS VM 2
NVIDIA K80

- Over 250 initial participants from 5 continents
- ~150 participate throughout the week
- 72 teams
- Over 33K Slack Messages
- ~$36K in AWS for Earth compute credits used



Created with mapchart.net ©

# Administration Challenges

- Setting up Jupyterhub and Kubernetes on AWS
  - Lots of trial and error in the last week
  - Could not get autoscaling to work
- Code and VM bugs/failures
  - Users discovered coding bugs and needed more libraries/extensions in Python environment
  - Using more RAM than available crashes the VM instead of going to swap
- Slack communication
  - Too many notifications turned on by default for admins
  - Receiving questions through mix of official channels and PMs
  - Could have used more people helping answer team questions
- Team management
  - People dropped out throughout the week
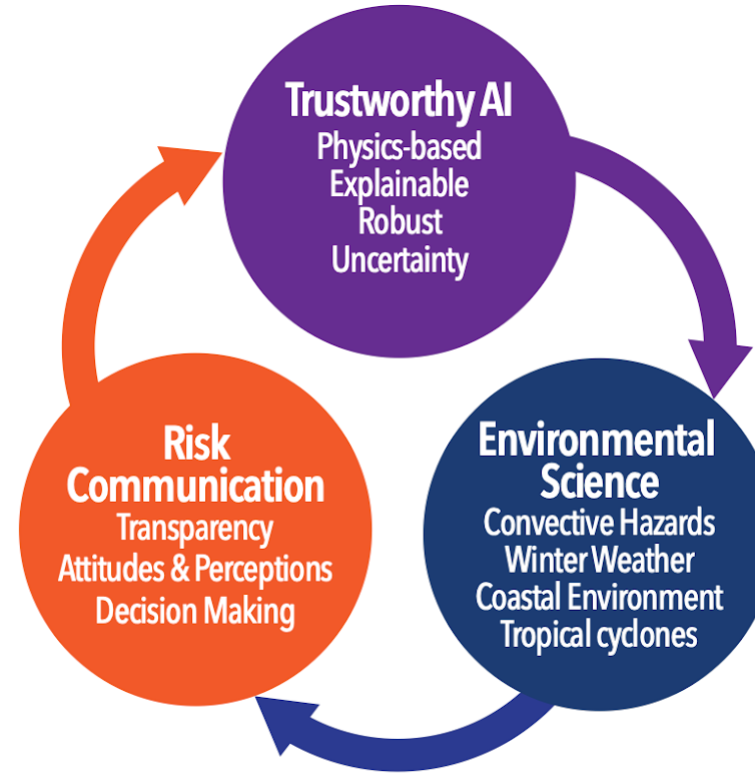  - Some people requested transfers to different teams/problems

# Lessons Learned

- Have fewer, more robustly tested and documented challenge problems
- Create challenge problems targeted at different experience levels
- Have synchronous work periods targeted at different time zones
- Satellite admin sites to support different time zones
- Provide clearer guidance up front about tasks, goals, best practices, and expectations
- Provide regular feedback on team submissions
- Charge a registration fee to incentivize participation

Hackathon notebooks:
https://github.com/NCAR/ai4ess-hackathon-2020

NSF AI Institute for Research on Trustworthy AI in Weather, Climate, and Coastal Oceanography

**Microphysics Emulation**

**GECKO Emulation**

**AI4ESS Hackathon**

AI4ESS Presentations, Notebooks and Data Links at ai4ess.org

**Contact Me**
Email: dgagne@ucar.edu
Twitter: @DJGagneDos
Github: djgagne