

# High-Accuracy Module Emulators from Physically-Constrained AI Algorithms

## Authors

Anthony S. Wexler and Paul Ullrich, University of California, Davis; Qi Tang, Lawrence Livermore National Laboratory; Manishkumar Shrivastava, Pacific Northwest National Laboratory

## Focal Area(s)

How do we use AI tools to integrate observations, simulated data and physical and chemical fundamentals (Focal Area 3) into model components (Focal Area 2) that have high accuracy and stability and low computational burden to improve Earth System Predictability?

## Science Challenge

Earth system modeling of the hydrological cycle involves compute-intensive modules representing complex chemical and physical process. Recently, AI tools that are far less compute intensive have been developed that emulate these modules, but many of these efforts are not yet sufficiently accurate or even stable. We know a lot about the physics and chemistry of earth system processes. The Science Challenge is developing AI tools that not only incorporate observations and simulated data, but also incorporate the physics and chemistry of the process, while still maintaining the compute efficiency.

## Rationale

Predicting extreme events in the hydrological cycle will involve high spatial resolution earth system models. The grand challenge is running global earth system models at a horizontal grid-spacing of a few hundred meters comparable to vertical resolution (rather than several hundred kms, used currently), so that fine-scale processes of aerosol and cloud microphysics, and convection and entrainment can be explicitly resolved. In many earth system models, the compute limitation resides with three modules, namely atmospheric chemistry, aerosol dynamics and radiative transfer, while the wall-clock limitation resides with dynamics and transport. Increasing horizontal spatial resolution to be comparable to the vertical resolution will dramatically increase the compute resources consumed by all of these modules. One path forward is to replace the three modules with machine-learned emulators. Current modules are based on the physics and chemistry of processes, even if they are highly parameterized. Machine-learned module replacements typically memorize data generated by the original module, achieving dramatic computational efficiency improvements, but frequently at the expense of accuracy, stability or both. New machine learning frameworks are needed that integrate all of our knowledge, including fundamental physics and chemistry along with observations or generated data within earth system models. The first such frameworks have recently been developed<sup>1</sup>.

## Narrative

## High-Accuracy Module Emulators from Physically-Constrained AI Algorithms

Atmospheric chemistry, aerosols, clouds and their interactions play a crucial role in regulating the energy and water cycles of the Earth system. The DOE Energy Exascale Earth System Model version 1 (E3SMv1) shows a strong aerosol effective radiative forcing (ERF,  $-1.65 \text{ W m}^{-2}$ )<sup>2</sup>. This ERF is substantially larger than the IPCC AR5's best estimate ( $-0.9 \text{ W m}^{-2}$ ) and falls outside the likely range ( $-1.5$  to  $-0.4 \text{ W m}^{-2}$ ). Such strong aerosol ERF requires a high equilibrium climate sensitivity (ECS, 5.3 K) (attributable to an unusually large positive shortwave cloud feedback) to simulate the global surface temperature observations. Similar to E3SMv1, several other Coupled Model Intercomparison Project phase 6 (CMIP6) models also show high ECS values exceeding the upper limit of the expected range<sup>3</sup>. These recent findings suggest the pressing need for better representations of chemistry-aerosol-cloud related processes in the future Earth system models to achieve improved predictability.

There is a chain of processes that lead to clouds. Atmospheric chemistry forms condensable organic and inorganic compounds that either nucleate to form new particles or condense on existing nuclei growing them to cloud condensation nuclei (CCN). Accurate atmospheric chemistry representations require tens to hundreds of chemical tracers. The aerosol size distribution and chemistry need to be represented to capture the particle dynamics that leads to CCN, which depend on the size of the particles and their chemical composition. Representing the aerosol size distribution (including interstitial and cloud-borne aerosols) requires another tens to hundreds of tracers. Once the aerosol size and chemical distribution are known, front and cell atmospheric dynamics govern their role in cloud formation and precipitation. Uncertainty in representation of aerosol size and composition distribution with earth system models leads to uncertainties in cloud microphysics, CCN activation and precipitation.

A number of recent efforts have used AI/ML tools to develop emulators that memorize the input-output relationship of these compute-intensive modules<sup>4-5</sup> dramatically increasing their speed in some cases, but always at the expense of accuracy and sometimes stability. Although memorizing these input-output relationships is a good first step, scientists in the community know a lot about the physics and chemistry in these modules yet this knowledge is not yet incorporated into the AI/ML emulators.

Recent work took a first step to incorporate fundamental principles and input-output relationships into AI/ML-based emulators<sup>1, 6-7</sup>. In one of these works<sup>1</sup> the AI/ML emulators memorize the fluxes instead of the quantities so that conservation principles are automatically obeyed and applied this to an atmospheric chemistry example. AI/ML emulators do not always conserve mass or energy resulting in one source of error. In general, algorithms and methods are needed whereby emulators not only memorize input-output relationships but also incorporate fundamental physical and chemical principles.

GPUs can provide an additional order of magnitude speed increase, especially for compute-intensive simulations involving large number of tracers. Modules that code atmospheric physics and chemistry are frequently not suitable for computation by GPUs, which is also the case for a number of machine learning tools. But some of these tools, like certain neural networks, are suitable for computation with GPUs. So the goal is that AI/ML emulators must employ GPU-compatible algorithms in addition to memorizing input-output relationships and take advantage of the computational ability of GPUs.

As discussed above, emulators can be validated by comparing their predictions to the computational results of the original modules and by incorporating fundamental physics and chemistry into their development. But each module and the entire model also needs validation against measurements of aerosol size and composition distribution in order to build trust in the model and to identify the model

## **High-Accuracy Module Emulators from Physically-Constrained AI Algorithms**

shortcomings where resources need to be allocated. So, the next step is to incorporate field data into the validation and AI/ML training process. Aircraft, remote sensed and ground-based measurements can all contribute. An integrated analysis of AI/ML emulators of processes and field observations can also be used for uncertainty quantification of processes. For instance, UC Davis operates the IMPROVE network (<https://aqrc.ucdavis.edu/improve>) for the National Parks Service (NPS) and performs chemical analysis for the EPA's Chemical Speciation Network. The IMPROVE network has about 165 sites around the US, with a few in other countries, that have been measuring PM<sub>2.5</sub> mass and chemical composition for decades. One difficulty using data to validate models is that the data is collected at a point whereas the models are making predictions over a domain. But most of the sites in the IMPROVE network are in rural and remote areas, usually state and national parks and forests, where there are few or no local sources, other than biogenic ones, and these biogenic sources are often continuous for a large spatial extent. As a result, the point source measurements in IMPROVE better represent spatial domain extents that are more comparable to those in regional and global models than data collected in urban areas where the gradients are steep so the extents small. In October 2019, DOE partnered with NPS to install a sampler at the ARM Southern Great Plains site, so data from this site can also be used in synergy with other ARM data streams.

### **Suggested Partners/Experts (Optional)**

Intentionally left blank.

# High-Accuracy Module Emulators from Physically-Constrained AI Algorithms

## References (Optional)

1. Sturm, P. O.; Wexler, A. S., A mass-and energy-conserving framework for using machine learning to speed computations: a photochemistry example. *Geoscientific Model Development* **2020**, *13* (9), 4435-4442.
2. Golaz, J. C.; Caldwell, P. M.; Van Roekel, L. P.; Petersen, M. R.; Tang, Q.; Wolfe, J. D.; Abeshu, G.; Anantharaj, V.; Asay-Davis, X. S.; Bader, D. C., The DOE E3SM coupled model version 1: Overview and evaluation at standard resolution. *Journal of Advances in Modeling Earth Systems* **2019**, *11* (7), 2089-2129.
3. Zelinka, M. D.; Myers, T. A.; McCoy, D. T.; Po-Chedley, S.; Caldwell, P. M.; Ceppi, P.; Klein, S. A.; Taylor, K. E., Causes of higher climate sensitivity in CMIP6 models. *Geophysical Research Letters* **2020**, *47* (1), e2019GL085782.
4. Kelp, M. M.; Tessum, C. W.; Marshall, J. D., Orders-of-magnitude speedup in atmospheric chemistry modeling through neural network-based emulation. *arXiv preprint arXiv:1808.03874* **2018**.
5. Keller, C. A.; Evans, M. J., Application of random forest regression to the calculation of gas-phase chemistry within the GEOS-Chem chemistry model v10. *Geoscientific Model Development* **2019**, *12* (3), 1209-1225.
6. Beucler, T.; Pritchard, M.; Rasp, S.; Ott, J.; Baldi, P.; Gentine, P., Enforcing analytic constraints in neural-networks emulating physical systems. *arXiv preprint arXiv:1909.00912* **2019**.
7. Beucler, T.; Rasp, S.; Pritchard, M.; Gentine, P., Achieving conservation of energy in neural network emulators for climate modeling. *arXiv preprint arXiv:1906.06622* **2019**.