

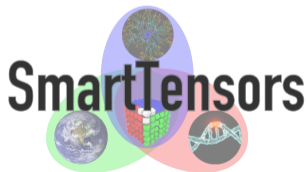
ChemML: Understanding groundwater flow and contaminant transport using machine learning

Velimir V. Vesselinov (monty) (monty@envitrac.com)

Tracy L. Kliphuis (trace) (trace@envitrac.com)

EnviTrace LLC, Santa Fe, NM, USA

<https://EnviTrace.com>

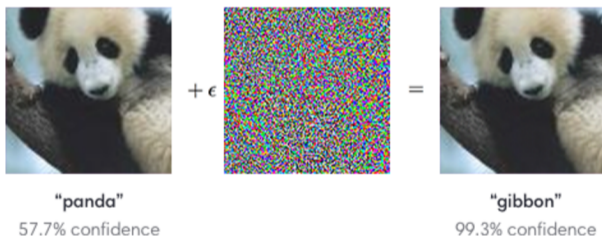




- ▶ Supervised
- ▶ Unsupervised (Self-supervised)
- ▶ Physics (Science) Informed



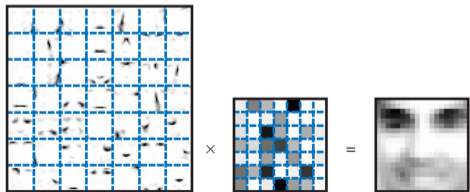
- ▶ requires labeling and huge training (labeled) datasets
- ▶ introduces subjectivity through the labeling process
- ▶ labeling for science application is challenging
- ▶ black box: we do not know why it works
- ▶ cannot discover something that we do not know already
- ▶ can be severely impacted by data noise: **adversarial examples**



⇒ major limitations of the **supervised** ML methods for **science** applications



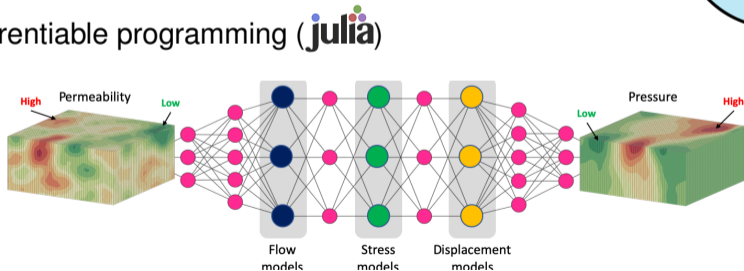
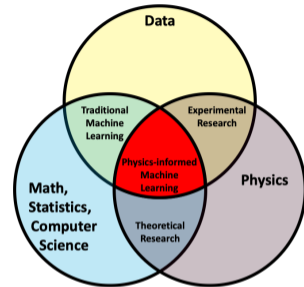
- ▶ extracts hidden features (signals, signatures) in the processed data automatically without any prior information
- ▶ applicable for both categorization and prediction
- ▶ produces unbiased analyses not impacted by data labeling, subject-matter-expert (SME) opinions, and physics assumptions
- ▶ identifies features that distinguish images of animals (e.g., cats, dogs, horses, etc.) or geothermal features
- ▶ categorizes data and SME can identify (“label”) animals (or geologic features)
- ▶ SME needed after ML is performed

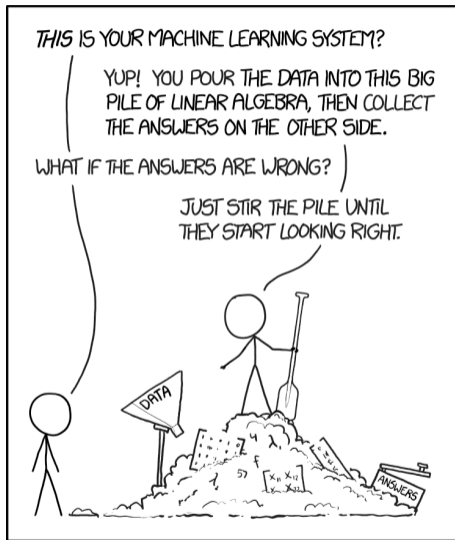


Physics (Science) Informed Machine Learning (PIML)



- ▶ learns from data but includes preconceived science knowledge
- ▶ physics/science information embedded in the ML framework or added as penalties
- ▶ PIML neural networks are problem specific
- ▶ needs SME inputs related to the analyzed problem
- ▶ SME needed before and after ML is performed
- ▶ increases efficiency, accuracy, and robustness
- ▶ requires differentiable programming (**julia**)





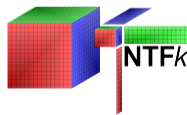
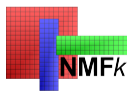
SO MUCH OF "AI" IS JUST FIGURING OUT WAYS TO OFFLOAD WORK ONTO RANDOM STRANGERS.



- ▶ Feature extraction (**FE**)
- ▶ Blind source separation (**BSS**)
- ▶ Detection of disruptions / anomalies
- ▶ Separate physics processes
- ▶ Discover unknown dependencies and phenomena
- ▶ Develop reduced-order / surrogate models
- ▶ Identify interrelationships between model inputs and outputs
- ▶ Guide development of physics models representing the data
- ▶ Optimize data acquisition
- ▶ Make predictions



- ▶ Novel LANL-patented, open-source, unsupervised and physics-informed Machine Learning (ML) methods and computational techniques
- ▶ Based on matrix/tensor factorization coupled with custom k -means clustering and nonnegativity/sparsity/physics-informed constraints
- ▶ Developed in **julia**
 - <http://tensors.lanl.gov>
 - <https://github.com/SmartTensors>
- ▶ Capable to efficiently process large datasets (GB/TB's)
- ▶ Applied in numerous geoscience projects (SBIR, DOE, ARPA-E, etc.)





▶ ChemML:

- <https://envitrace.com/projects/chemml.html>
- Physics-informed AI/ML for characterization, parameterization, and prediction of contaminant transport and remediation processes
- SBIR DOE funded

▶ GeoTGO:

- <https://envitrace.com/projects/geothermal.html>
- Equitable and inclusive tool for community-based geothermal development
- SBIR DOE funded



▶ **GeoThermalCloud:**

- <https://github.com/SmartTensors/GeoThermalCloud.jl>
- Cloud Fusion of Big Data and Multi-Physics Models using Machine Learning for Discovery, Exploration and Development of Hidden Geothermal Resources
- DOE funded

▶ **ML4Geo:**

- <https://github.com/SmartTensors/ML4Geo.jl>
- Machine Learning based Well Design to Enhance Unconventional Energy Production
- ARPA-E funded



► Field Data:

- Contamination
- Climate
- Seismic
- Geothermal
- Oil/gas production
- CO₂ sequestration
- Wildfires
- COVID-19

► Lab Data:

- X-ray Spectroscopy
- UV Fluorescence Spectroscopy
- Microbial population growth
- Fracture development
- Isotope fractionation

► Operational Data:

- Neutron Accelerator (LANSCE)
- Oil/gas production
- Geothermal energy production
- CO₂ sequestration

► Model Outputs:

- Geothermal
- Watershed
- CO₂ sequestration
- Oil/gas production
- Climate
- Reactive mixing $A + B \rightarrow C$
- Co-polymers Phase separation
- Protein Molecular Dynamics



- ▶ Fleming et al., Machine Learning in Earth and Environmental Science Requires Education and Research Policy Reforms, **Nature Geoscience**, 10.1038/s41561-021-00865-3, 2021.
- ▶ Siler et al., Machine learning to identify geologic factors associated with production in geothermal fields: A case-study using 3D geologic data, Brady geothermal field, Nevada, **Geothermal Energy**, 10.1186/s40517-021-00199-8, 2021.
- ▶ Ahmmed et al., Machine Learning to Discover Mineral Trapping Signatures due to CO2 Injection, **Journal of Greenhouse Gas Control**, 10.1016/j.ijggc.2021.103382, 2021.
- ▶ Ahmmed et al., A comparative study of machine learning models for predicting the state of reactive mixing, **Journal of Computational Physics** 10.1016/j.jcp.2021.110147, 2021.
- ▶ Mehana et al., Machine-Learning Predictions of the Shale Wells? Performance, **Journal of Natural Gas Science and Engineering**, 10.1016/j.jngse.2021.103819, 2021.
- ▶ Vesselinov et al., Unsupervised Machine Learning Based on Non-Negative Tensor Factorization for Analyzing Reactive-Mixing, **Journal of Computational Physics**, Special issue: Machine Learning, 2019.
- ▶ Stanev et al., Unsupervised Phase Mapping of X-ray Diffraction Data by Nonnegative Matrix Factorization Integrated with Custom Clustering, **Nature Computational Materials**, 2018.
- ▶ Vesselinov et al., Nonnegative Tensor Factorization for Contaminant Source Identification, **Journal of Contaminant Hydrology**, 2018.
- ▶ O'Malley et al., Nonnegative/binary matrix factorization with a D-Wave quantum annealer, **PLOS ONE**, 2018.
- ▶ Vesselinov et al., Contaminant source identification using semi-supervised machine learning, **Journal of Contaminant Hydrology**, 2017.
- ▶ Alexandrov, Vesselinov, Blind source separation for groundwater level analysis based on nonnegative matrix factorization, **Water Resources Research**, 2014.



- ▶ Characterization and remediation of contaminants in groundwater is challenging
- ▶ Geologic and geochemical data are massive, complex, and difficult to interpret
- ▶ Data are often not fully used
- ▶ Building geochemical models representing these data and predicting future behavior is also challenging
- ▶ These models also rely on numerous assumptions and are time-consuming to build and execute
- ▶ Assumptions can lead to errors
- ▶ **ChemML** provides an alternative which develops models that are:
 - ▶ fast and simple
 - ▶ data driven, robust and minimize assumptions
 - ▶ defensible, easy to use and understand



- ▶ Example: 4 buckets representing 4 different water types
- ▶ Buckets have different geochemical concentrations and contaminants





- ▶ Water from the buckets is mixed in unknown way in the subsurface
- ▶ Mixing is caused by many poorly understood subsurface processes





- ▶ Water compositions of the original water types (buckets) are typically unknown
- ▶ Only groundwater mixtures observed in monitoring wells are known





- ▶ **ChemML** can estimate bucket composition using only observed mixtures
- ▶ **ChemML** can estimate uncertainties and make predictions





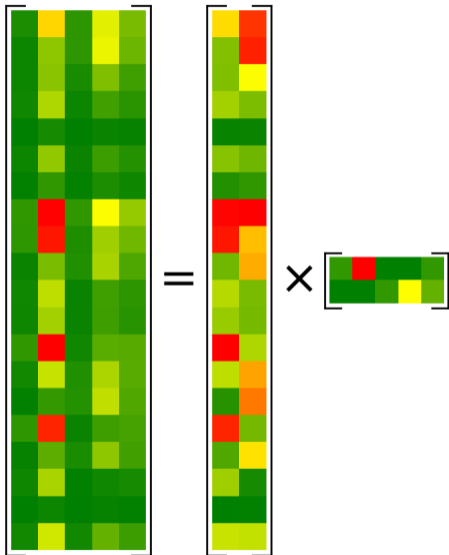
- ▶ **ChemML** uses Machine Learning (ML) to estimate bucket compositions
- ▶ Machine Learning is a form of Artificial Intelligence (AI) that works by exploring data and finding patterns with minimum human intervention



$$\mathbf{X}$$
$$[20 \times 5]$$

\mathbf{X} – data matrix
[geochemical species \times monitoring wells]

\mathbf{X} may have empty cells (data gaps)



$$X = W \times H$$

$$[20 \times 5] = [20 \times 2] \times [2 \times 5]$$

X – **data** matrix

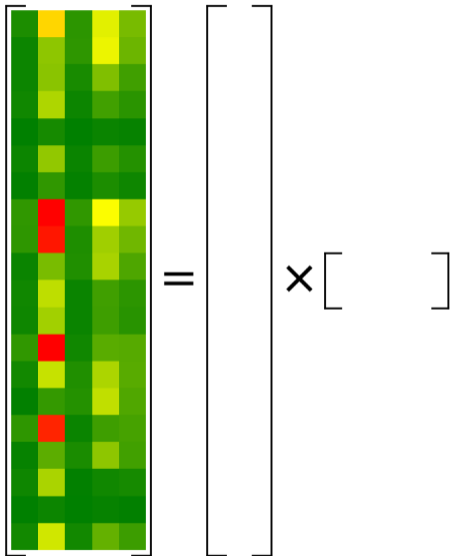
[**geochemical species** \times **monitoring wells**]

W – **bucket** matrix

[**geochemical species** \times **buckets**]

H – **mixing** matrix

[**buckets** \times **monitoring wells**]



$$X = W \times H$$

$$[20 \times 5] = [20 \times ?] \times [? \times 5]$$

⇒ 100 **knowns** (if there are no data gaps)

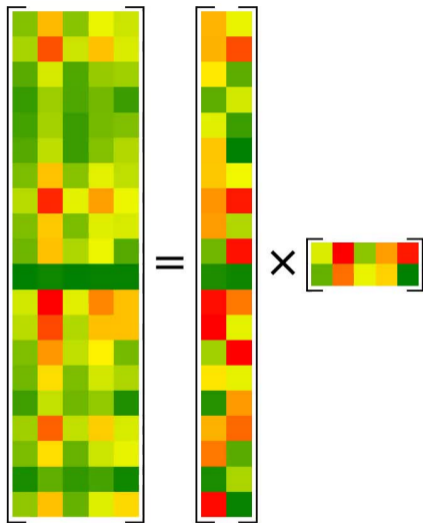
⇒ **unknown** number of buckets (2 or more)

⇒ **unknown** matrix elements of W and H (50 or more)

⇒ **Physics constraints:**

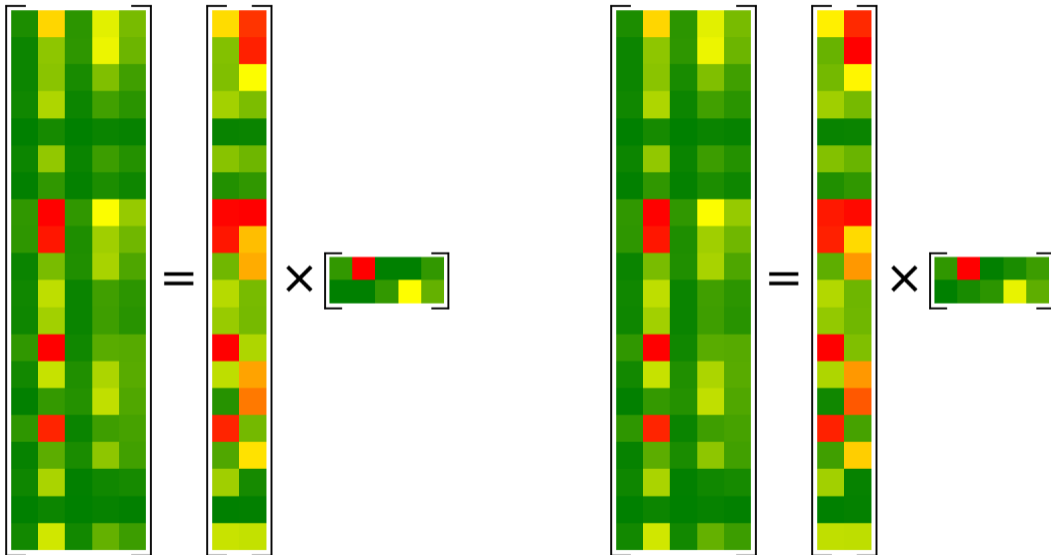
▶ all elements of W and $H \geq 0$

$$\sum_{k=1}^K H_{k,j} = 1 \quad \forall j$$



0001

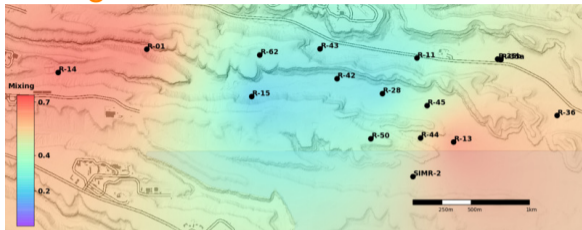
ChemML: true vs. estimated matrix factorization





"Background" bucket

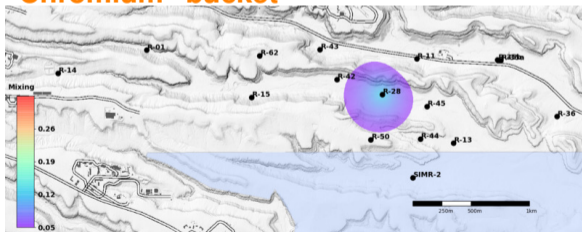
2005



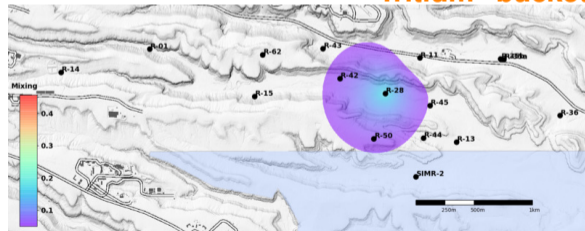
"Nitrate" bucket



"Chromium" bucket



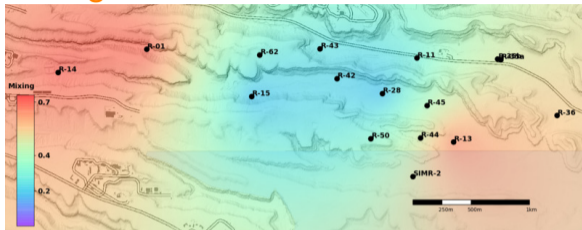
"Tritium" bucket



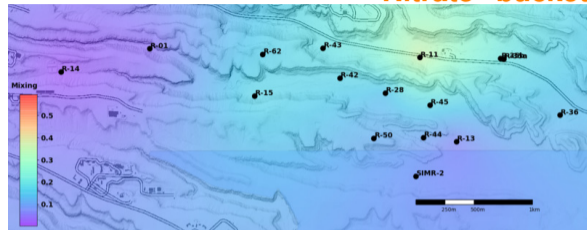


"Background" bucket

2006



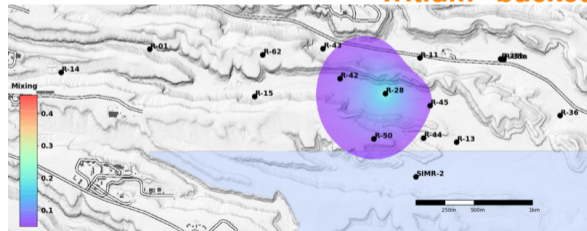
"Nitrate" bucket



"Chromium" bucket



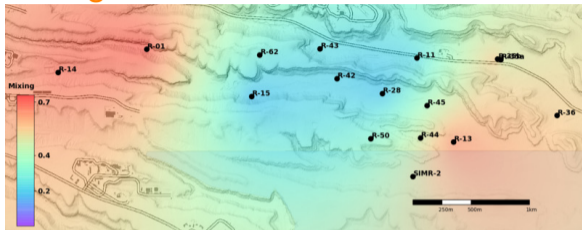
"Tritium" bucket



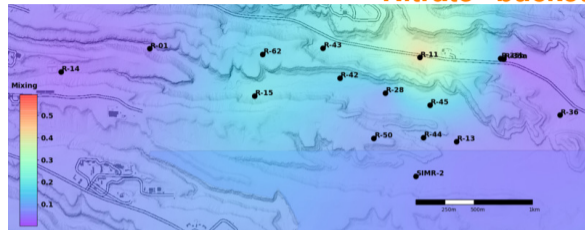


"Background" bucket

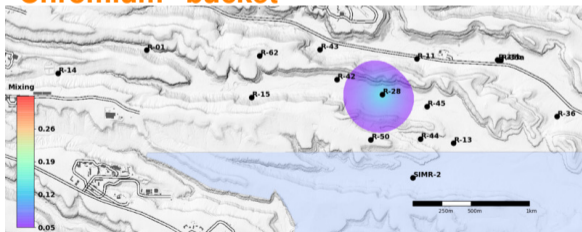
2007



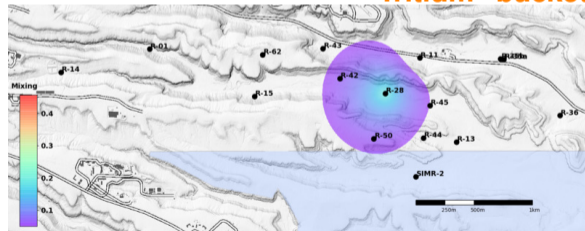
"Nitrate" bucket



"Chromium" bucket



"Tritium" bucket

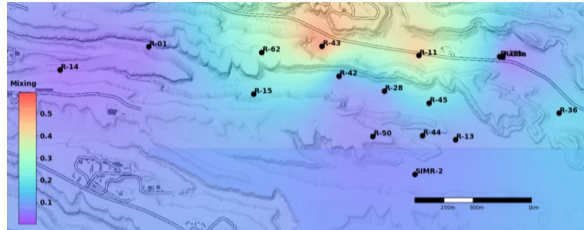
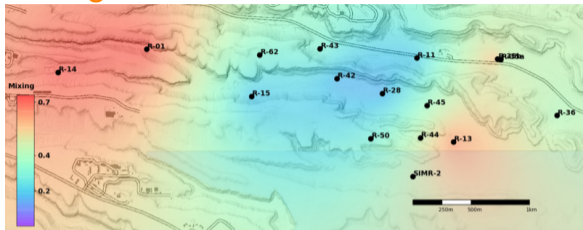




"Background" bucket

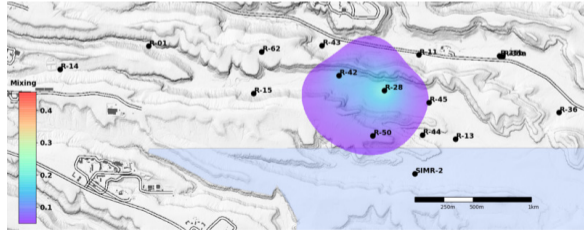
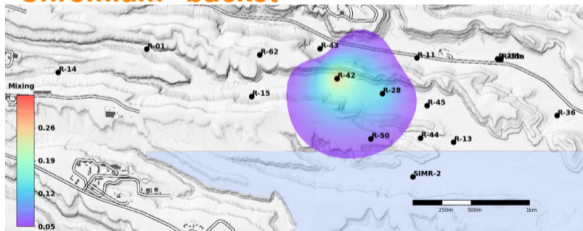
2008

"Nitrate" bucket



"Chromium" bucket

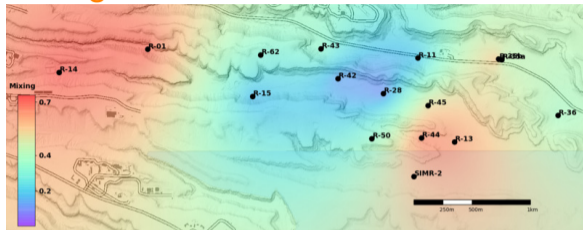
"Tritium" bucket



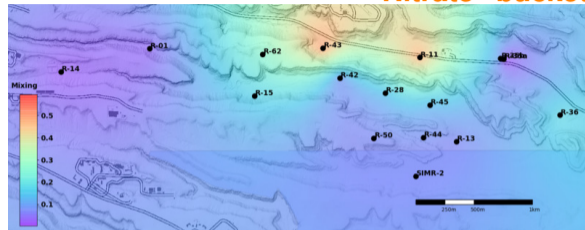


"Background" bucket

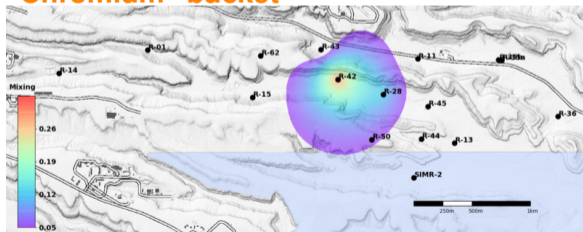
2009



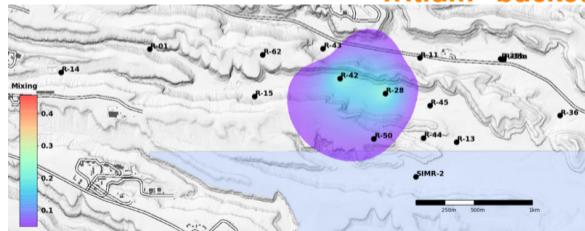
"Nitrate" bucket



"Chromium" bucket



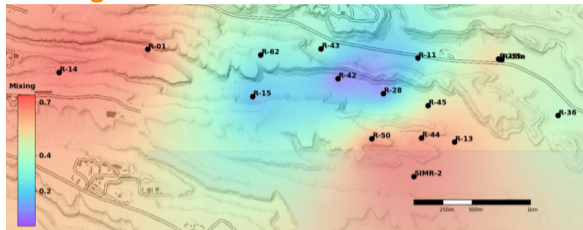
"Tritium" bucket



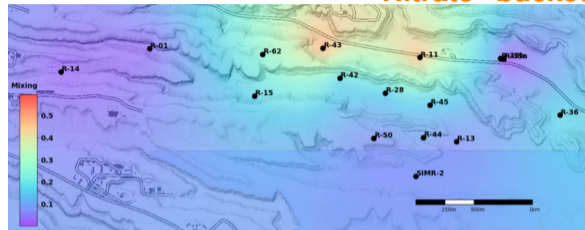


"Background" bucket

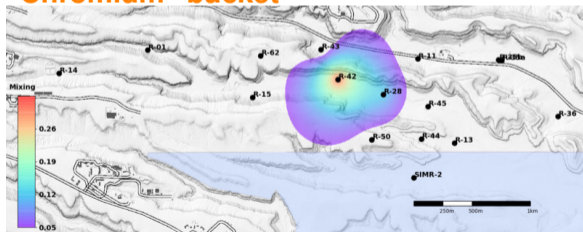
2010



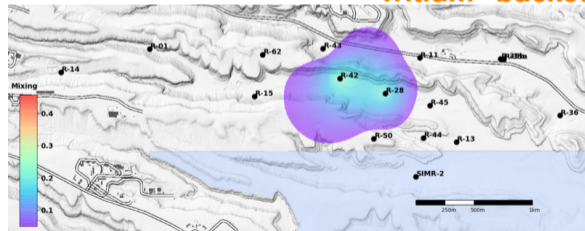
"Nitrate" bucket



"Chromium" bucket



"Tritium" bucket

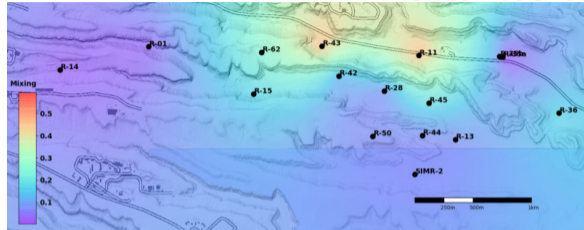
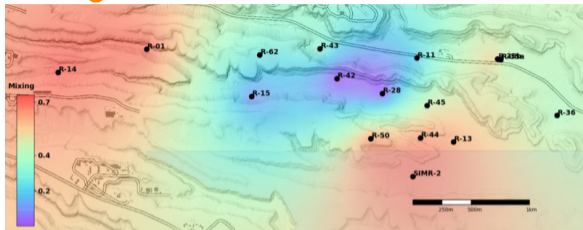




"Background" bucket

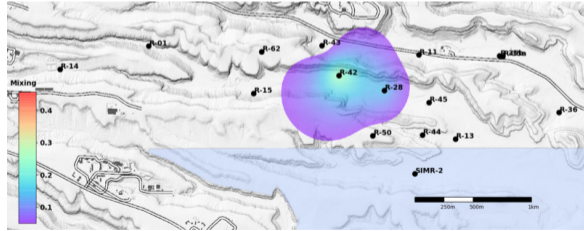
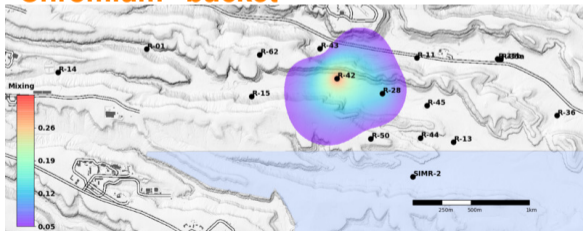
2011

"Nitrate" bucket



"Chromium" bucket

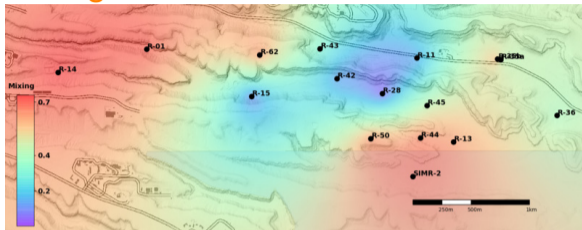
"Tritium" bucket



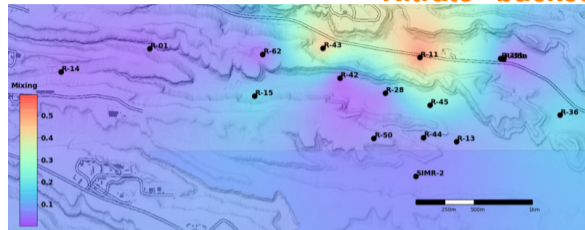


"Background" bucket

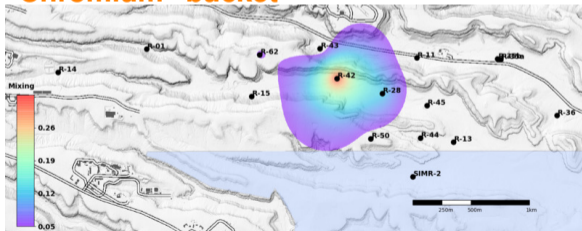
2012



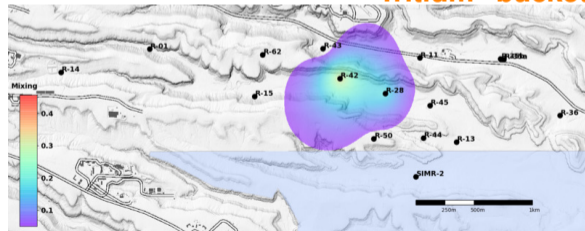
"Nitrate" bucket



"Chromium" bucket



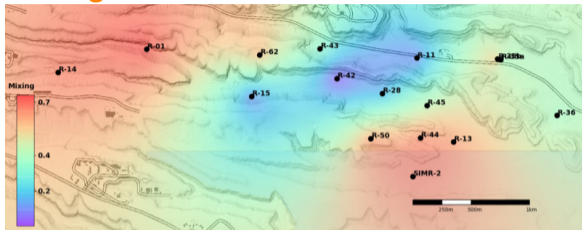
"Tritium" bucket



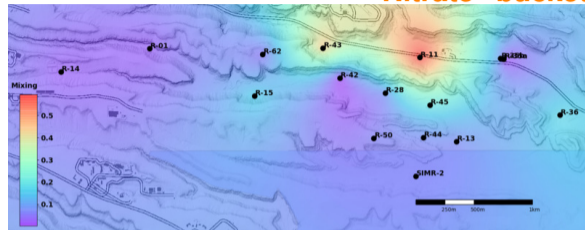


"Background" bucket

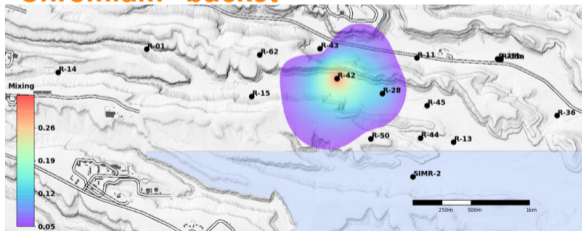
2013



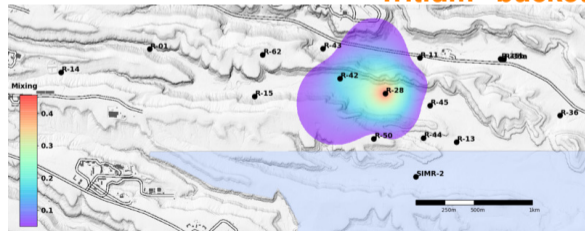
"Nitrate" bucket



"Chromium" bucket



"Tritium" bucket

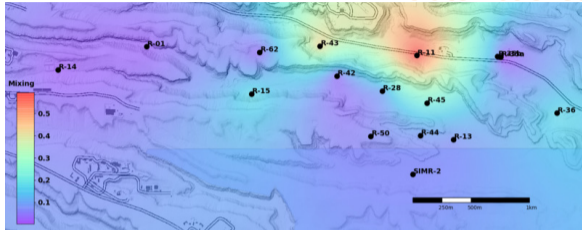
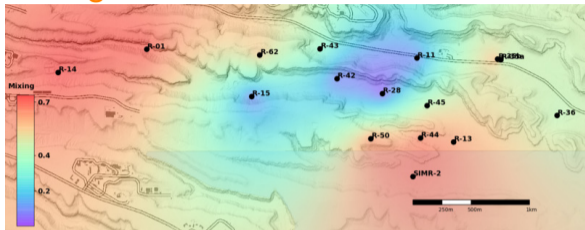




"Background" bucket

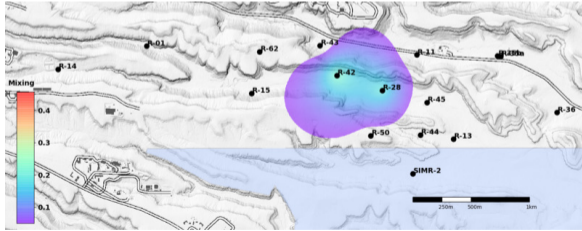
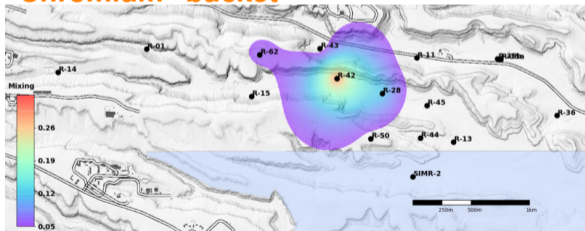
2014

"Nitrate" bucket



"Chromium" bucket

"Tritium" bucket

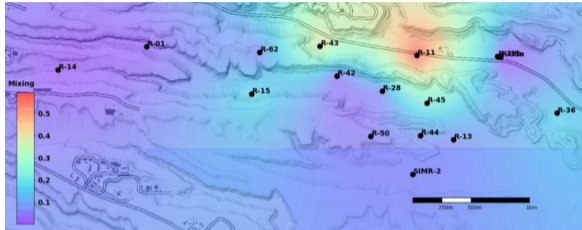
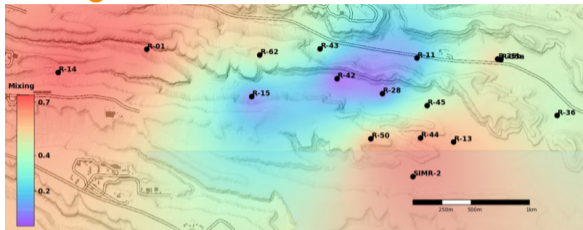




"Background" bucket

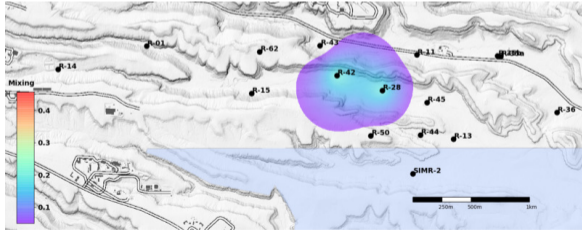
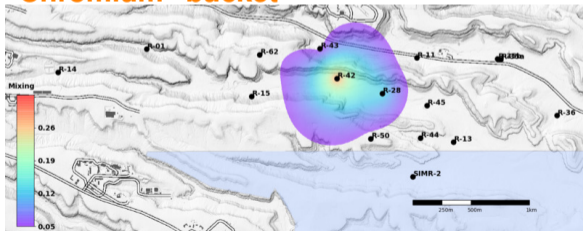
2015

"Nitrate" bucket



"Chromium" bucket

"Tritium" bucket

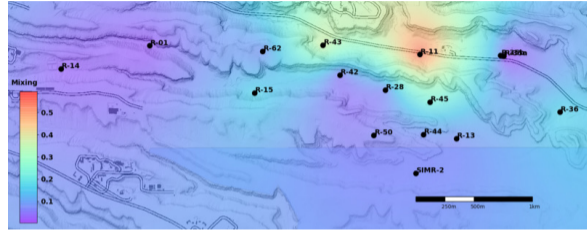
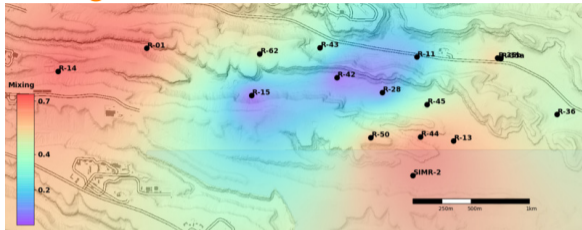




"Background" bucket

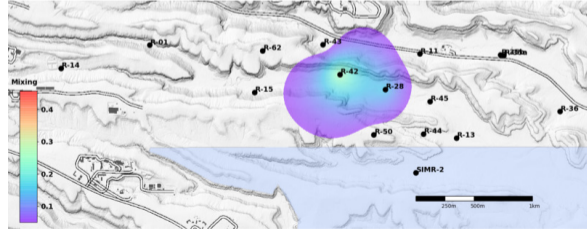
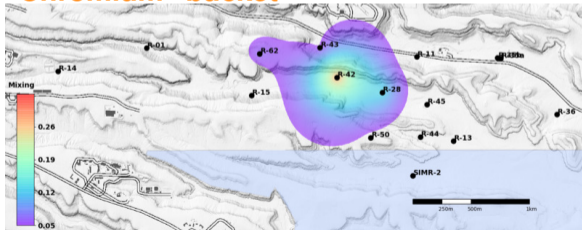
2016

"Nitrate" bucket



"Chromium" bucket

"Tritium" bucket

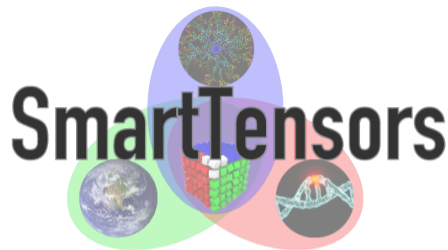




- ▶ All the models are wrong, but some are useful
- ▶ **ChemML** provides effective **data mining solutions**
- ▶ **ChemML** applies novel ML methods developed by our team



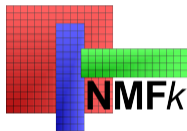
- ▶ Developed **novel** unsupervised and physics-informed ML methods
- ▶ Implemented in **cutting-edge** open-source computational framework called **SmartTensors** based on Nonnegative Factorization (Matrices/Tensors)
- ▶ **SmartTensors** has been used to solve various real-world problems
- ▶ **SmartTensors** deployment as a service on <https://JuliaHub.com> is coming soon
- ▶ **SmartTensors** just received 2 2021 R&D100 awards





► Codes:

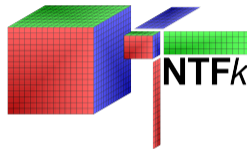
NMF_k



MADS



NTF_k



► Projects, Notebooks, Examples, Tutorials:

<http://tensors.lanl.gov>

<http://SmartTensors.github.io>

<https://github.com/SmartTensors>

<https://github.com/SmartTensors/NMFk.jl/tree/master/notebooks>

<https://github.com/SmartTensors/GeoThermalCloud.jl>

<https://github.com/SmartTensors/SmartTensorsTutorials.jl>