

Simulation Intelligence in the New Paradigm of Environmental Monitoring

Haruko Wainwright

Nuclear Science and Engineering; Civil and Environmental Engineering

Massachusetts Institute of Technology



Soil and Groundwater Contamination

- **Superfund Sites: >1300 sites (organic/metal/radioactive)**
- **Brownfield Sites: ~450,000**
- **Major environmental justice issues**

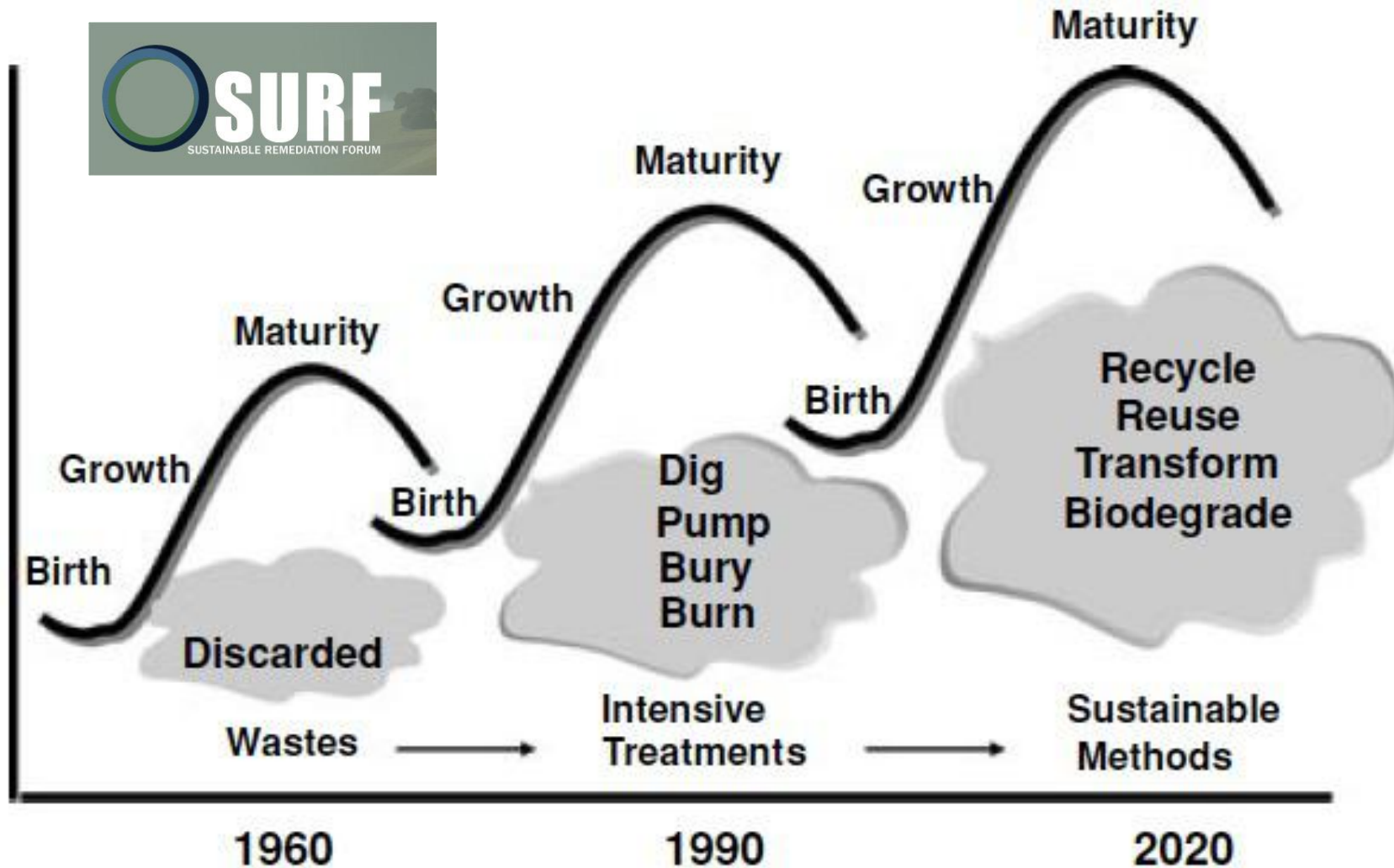


>900 remaining after 30 – 40 years of remediation

→ Challenge of low-concentration large-volume plume



Environmental Remediation: Evolution



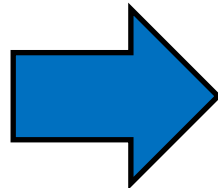
Trade offs:
Contaminant removal
vs

- **Waste**
- **CO2 emission**
- **Energy Use**
- **Ecological Impacts**
- **Noise, Air pollution**

Sustainable Remediation Forum (SURF), "Integrating sustainable principles, practices, and metrics into remediation projects", *Remediation Journal*, 19(3), pp 5 - 114, editors P. Hadley and D. Ellis, Summer 2009

Sustainable Remediation

- **Intensive/invasive clean up → Sustainable methods**
 - **Minimize waste/pollution/energy-use/water-use/ecological damages**
 - **Biodegradation, immobilization**
 - **Monitored natural attenuation**
 - **Longer institutional control with alternative/attractive end-use**
- **Long-term monitoring**



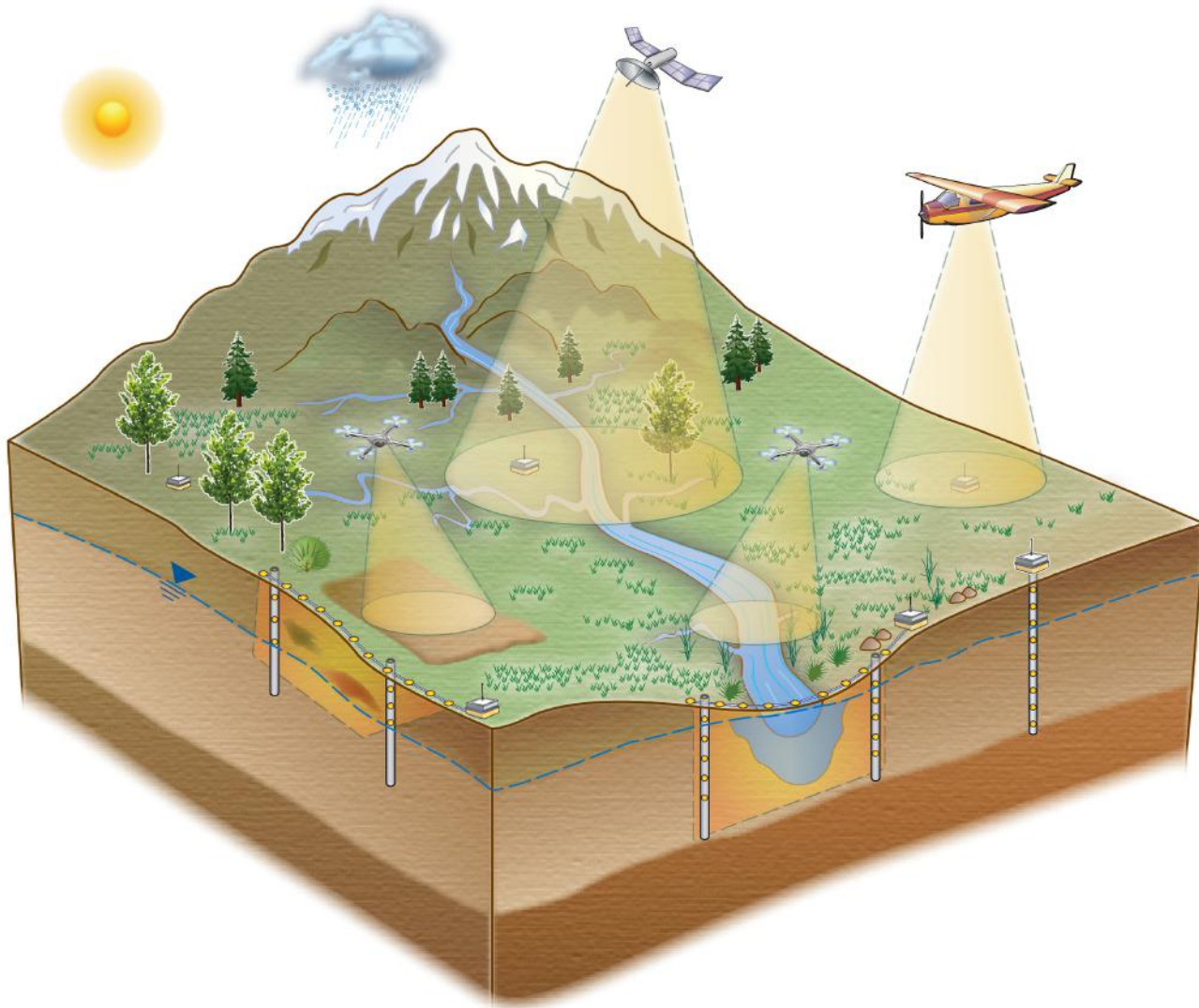
Former Reilly Tar & Chemical Corporation Plant



Rocky Flats National Wildlife Refuge



Earth Systems Monitoring

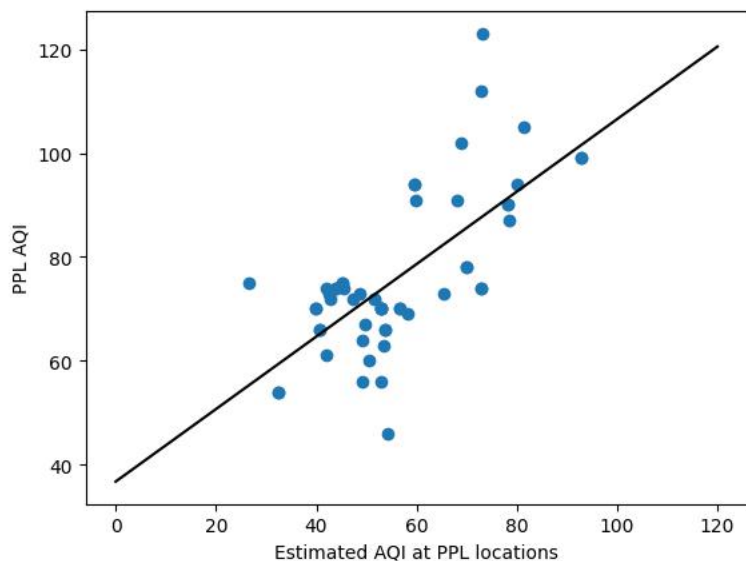


- **Multi-type multi-scale data**
 - Accuracy
 - Coverage
 - Footprints/resolution
- **”Proxy” information**
 - Plants/topography ~ soil
 - Electrical conductivity ~ contaminant concentration
- **Spatial-temporal correlation**
 - Data compression
 - Similar properties in vicinity

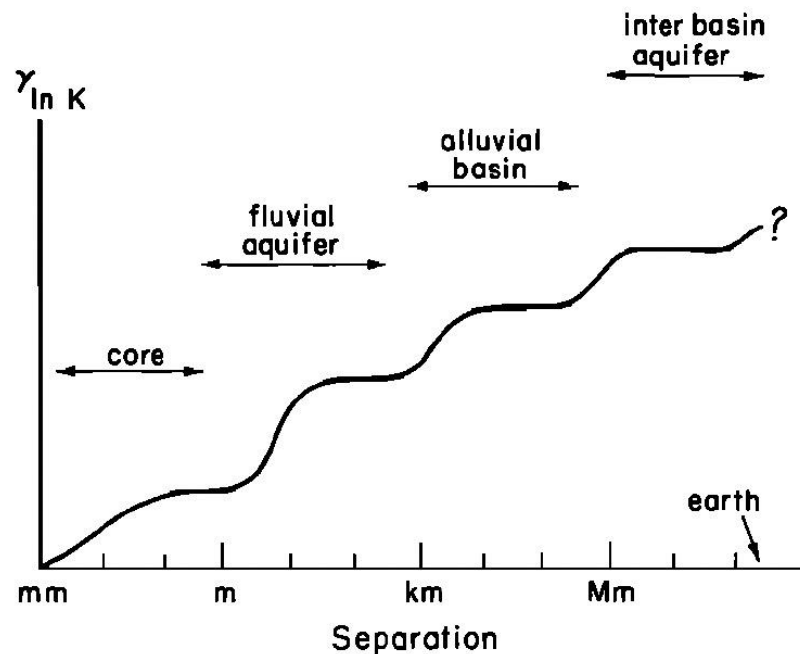


Challenges in ML/AI x Environmental Science

- Lack of training data
- Large uncertainty/variability



- Multiscale heterogeneity



- Large data but little information content

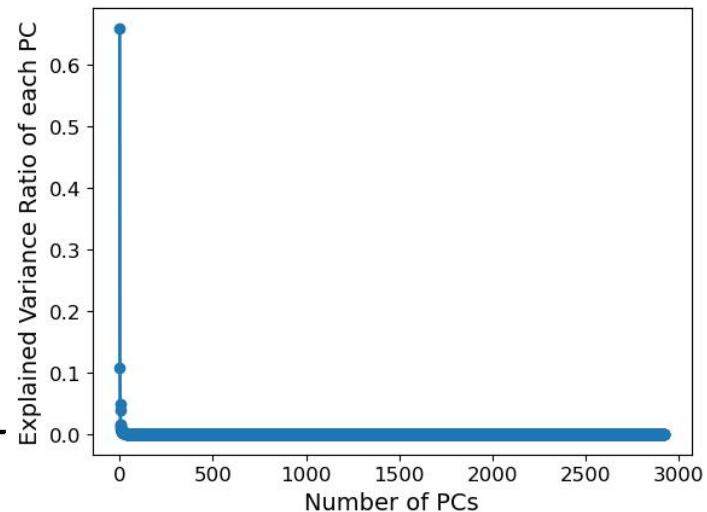
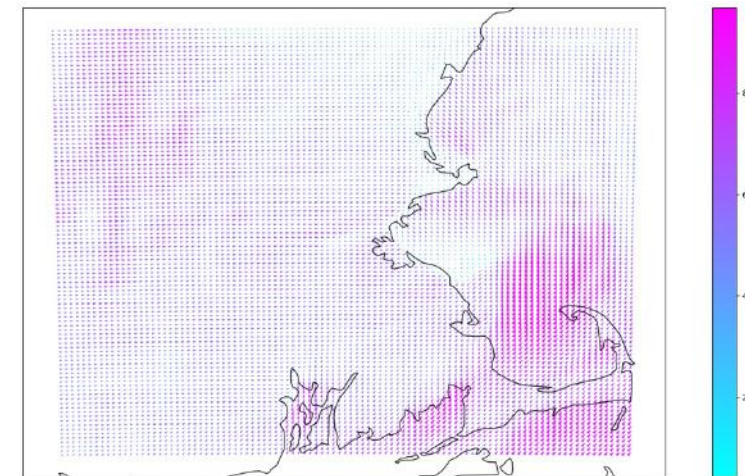
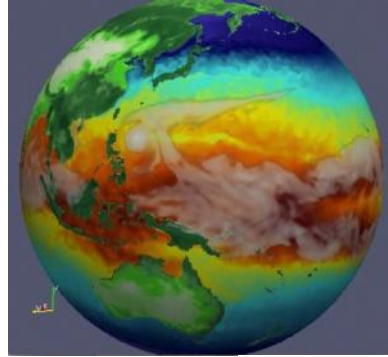
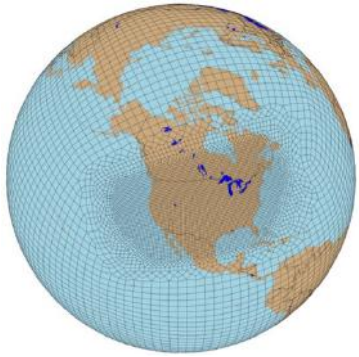


Fig. 8. Hypothetical $\ln K$ variogram illustrating the notion of scale-dependent correlation scales.

Gelhar, 1986

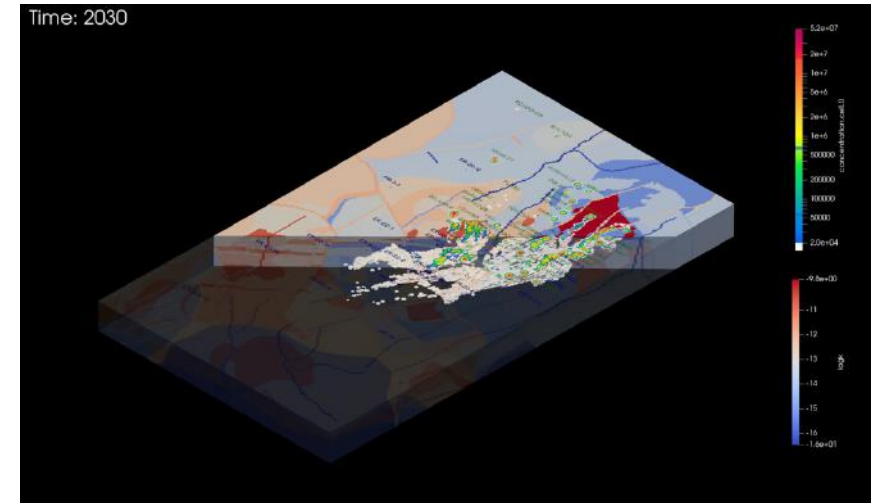
Challenges in Physical Models: Predictability

Global climate models

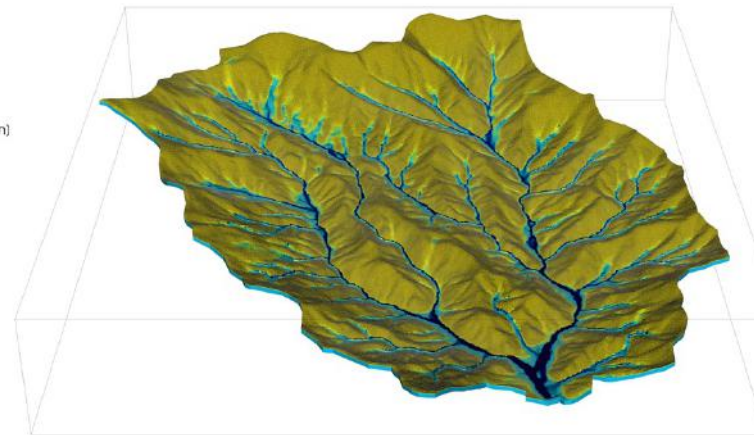


e3sm.org

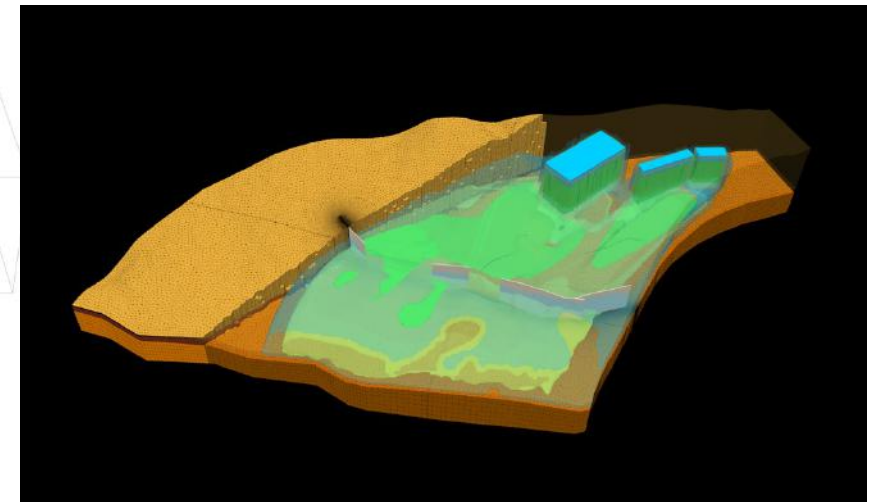
Contaminant transport models



Watershed models



Amanzi-ATS



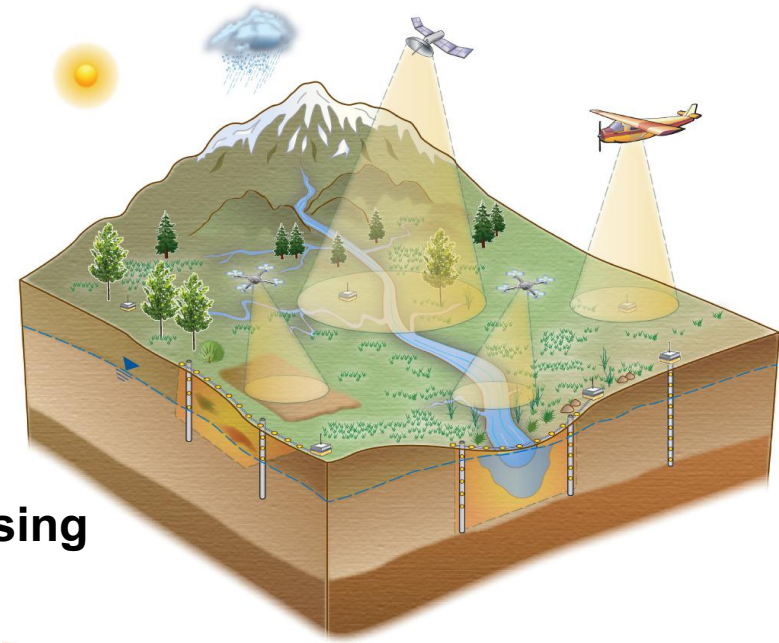
- Parameterization/heterogeneity
- Uncertainty quantification
- Inherit assumptions in models

Advanced Long-term Environmental Monitoring Systems



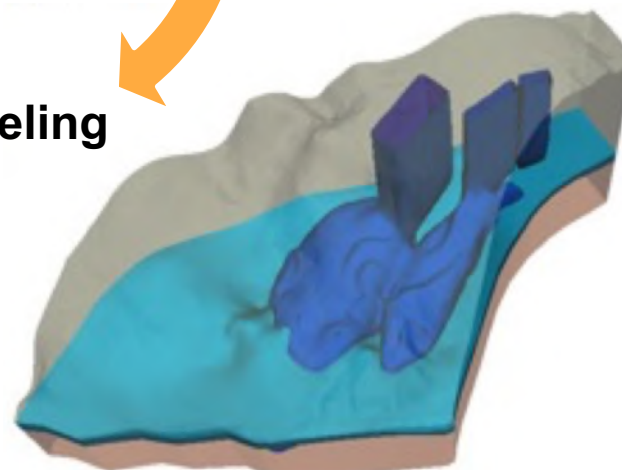
ML/AI

Sensing



EESA20-015

Modeling



ALTEMIS

altemis.lbl.gov

*New paradigm for
long-term monitoring*



Co-Leads



Carol Eddy-Dilek

AI for Soil and Groundwater



Himanshu Upadhyay



Styarth Praveen

Geochemical characterization



Hansell Gonzalez-Raymat



Miles Denham

Reactive Transport Modeling



Haruko Wainwright



Zexuan Xu

Spatial monitoring: Geophysics



Tim Johnson



Baptiste Dafflon



Sebastian Uhlemann

In situ real-time monitoring



Tom Danielson

Spatial monitoring: Radiation



Kai Vetter

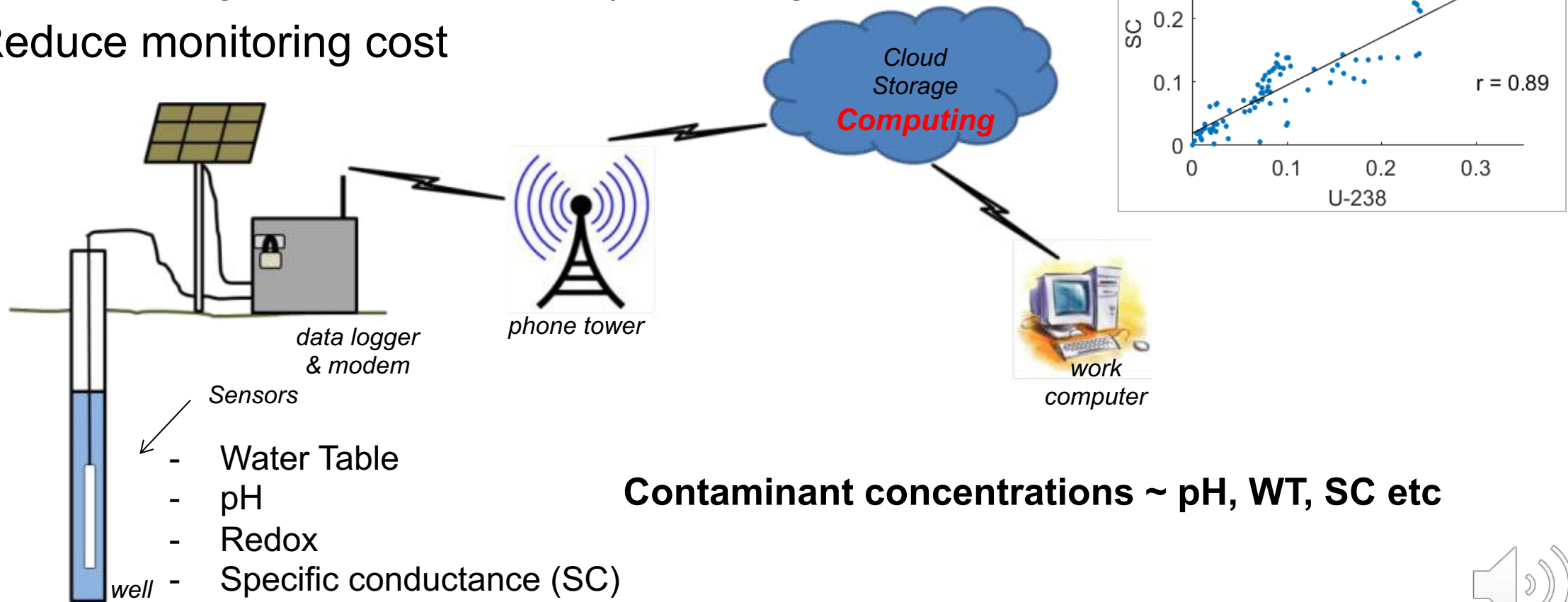


Brian Quiter

In situ Real-time Groundwater Monitoring

- **Low-cost in situ sensors, wireless network, cloud computing**

- Autonomous continuous monitoring
- Detect changes real-time = Early Warning
- Reduce monitoring cost

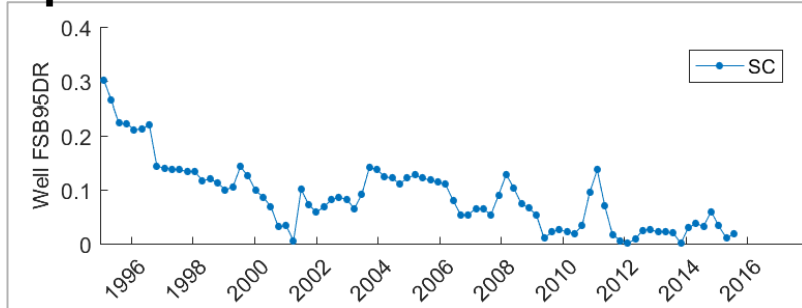


(Schmidt et al., 2018)

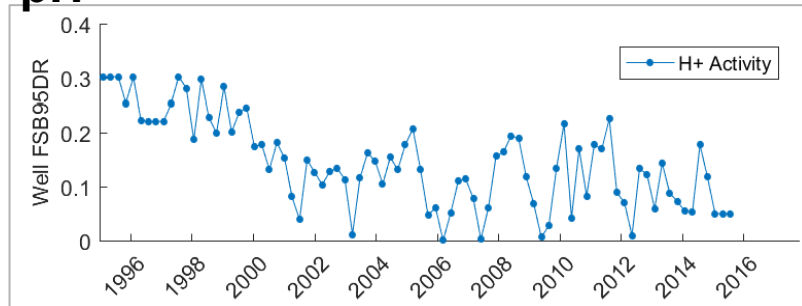


In situ Real-time Groundwater Monitoring

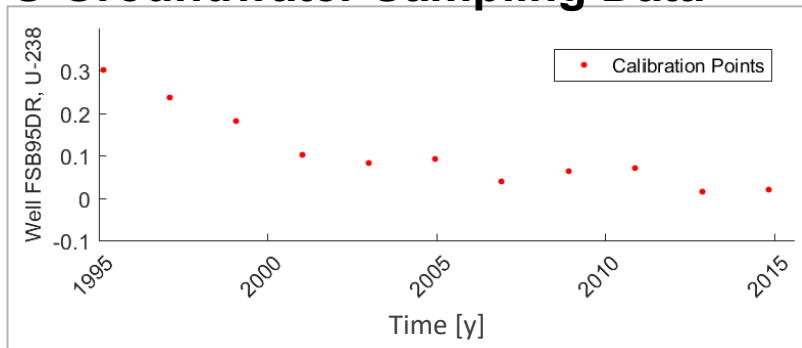
Specific Conductance



pH



U Groundwater Sampling Data



- Estimation of contaminant concentrations: $\{y_1, y_2, \dots, y_T\}$

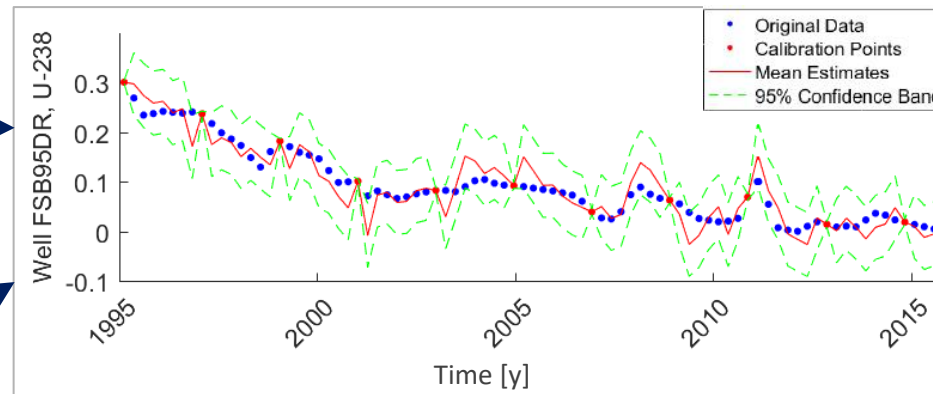
- Time evolution of contaminant concentration

$$y_t = f(y_{t-1}) + \tau \quad \tau \sim N(0, \sigma^2)$$

- Relationship to in situ datasets

$$z_{pH,t} = g_{pH}(y_t) + \varepsilon_{pH}$$

$$z_{SC,t} = g_{SC}(y_t) + \varepsilon_{SC}$$



$$p(y_t | y_{t-1}, z)$$

- Confidence interval captures validation points
- Mean estimate captures the fluctuation
- Reduce #sampling from quarterly to every two years

* Normalized concentrations



Search Citation
Enter search text / DOI
● Environ. Sci. Technol.

Environ. Sci. Technol. Environ. Sci. Technol. Lett.

- Home
- Browse the Journal
- Articles ASAP
- Current Issue
- Submission & Review
- Open Access
- About the Journal

Article

In Situ Monitoring of Groundwater Contamination Using the Kalman Filter

Franziska Schmidt[†], Haruko M. Wainwright^{*‡}, Boris Faybishenko[§], Miles Denham[¶], and Carol Eddy-Dilek[⊥]

[†] Department of Nuclear Engineering, University of California Berkeley, Etcheverry Hall, 2521 Hearst Avenue, Berkeley, California 94709, United States

[‡] Climate and Ecosystem Sciences Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, MS 74R-316C, Berkeley, California 94720-8126, United States

[§] Energy Geosciences Division, Lawrence Berkeley National Laboratory, 1 Cyclotron Road, MS 74R-316C, Berkeley, California 94720-8126, United States

[¶] Panoramic Environmental Consulting, LLC, P.O. Box 9000, Savannah, Georgia 31406, United States

[⊥] Savannah River National Laboratory, Savannah River, South Carolina 29580, United States

Environ. Sci. Technol., 2018, 52 (13), pp 7418–7425

DOI: 10.1021/acs.est.8b00017

Publication Date (Web): June 22, 2018

Copyright © 2018 American Chemical Society

< Previous Article



- HOME
- NEWS
- MULTIMEDIA
- MEETINGS
- PORTALS
- ABOUT

PUBLIC RELEASE: 13-AUG-2018

Algorithm provides early warning system for tracking groundwater contamination

Berkeley Lab researchers devise system to monitor contaminant plumes

DOE/LAWRENCE BERKELEY NATIONAL LABORATORY

- f
- t
- +
- e
- SHARE

PRINT E-MAIL

- News
- Student Resources
- Scholarships
- Student Discounts

New Algorithm Provides Real-Time Monitoring Of Groundwater Pollutants

Sam Bennett 8 Months Ago



AI & Automation Cybersecurity Cloud & Infrastructure Data & Analytics Smart Cities & IoT

Machine learning improves contamination monitoring

BY MATT LEONARD | AUG 14, 2018

Because groundwater is susceptible to pollution from automotive fuel, naturally occurring substances like iron, the Environmental Protection Agency and its state-level counterparts conduct annual or quarterly sampling and analysis.



- NEWS
- EVENTS
- VIDEOS
- TV & PODCASTS
- INDUSTRIES

Efficiency & Environment

Scientists develop new method to track groundwater pollutants in real-time

It is expected to reduce the frequency of manual groundwater sampling and lab analysis and therefore cut the monitoring cost

PyLEnM: A Machine Learning Framework for Long-Term Groundwater Contamination Monitoring Strategies

Aurelien O. Meray, Savannah Sturla, Masudur R. Siddiquee, Rebecca Serata, Sebastian Uhlemann, Hansell Gonzalez-Raymat, Miles Denham, Himanshu Upadhyay, Leonel E. Lagos, Carol Eddy-Dilek, and Haruko M. Wainwright*



Cite This: *Environ. Sci. Technol.* 2022, 56, 5973–5983



Read Online

ACCESS |



Metrics & More

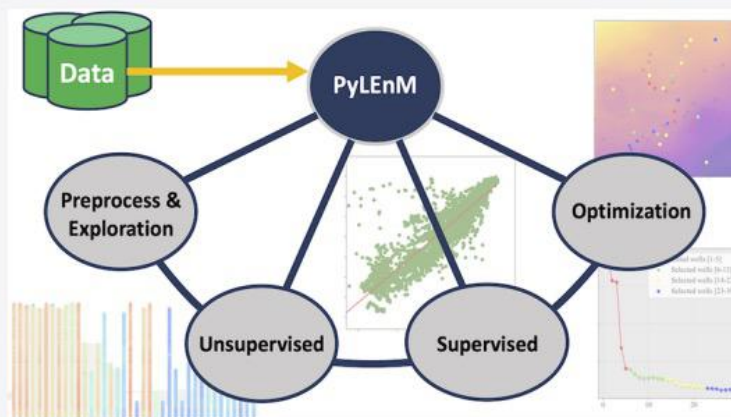


Article Recommendations

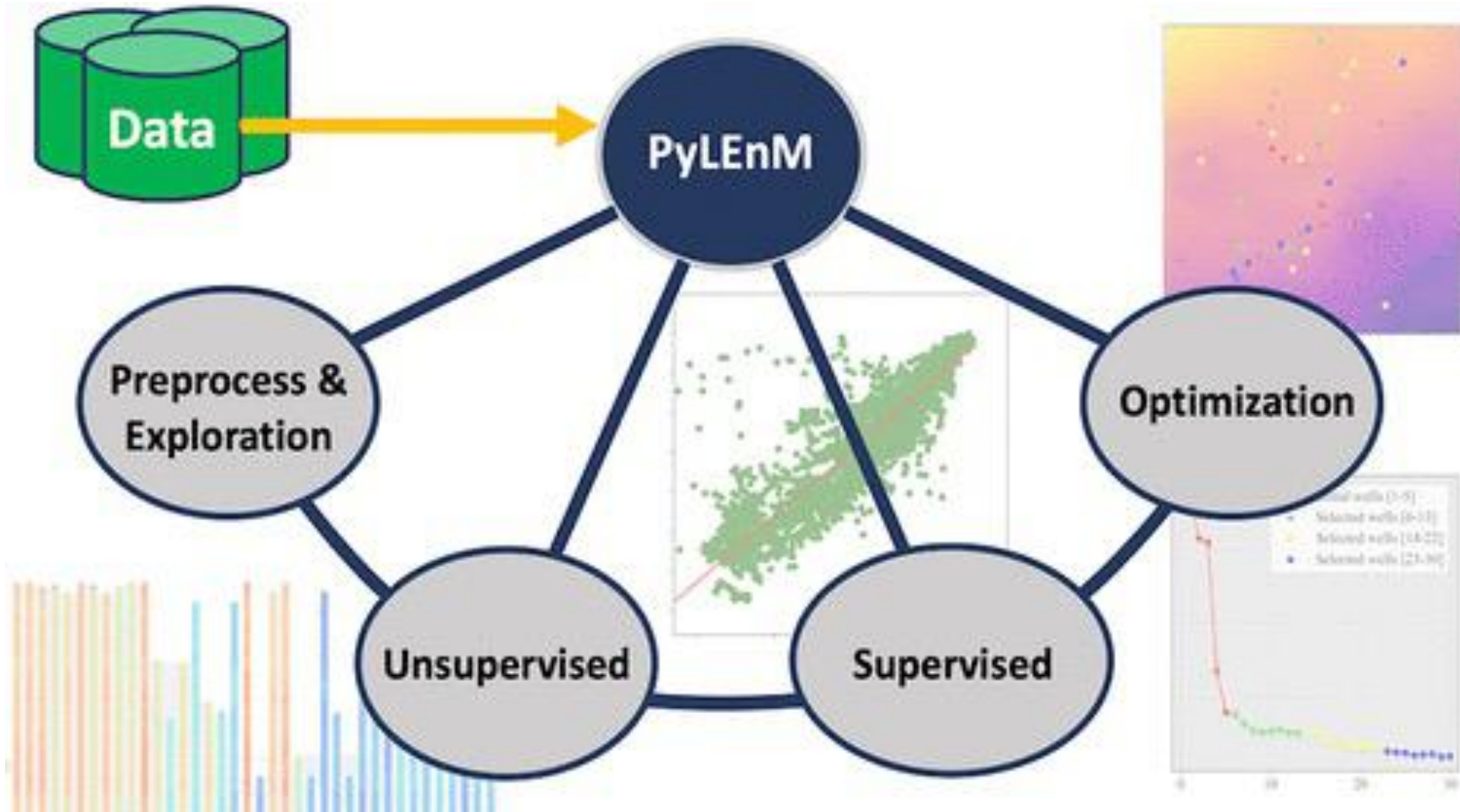


Supporting Information

ABSTRACT: In this study, we have developed a comprehensive machine learning (ML) framework for long-term groundwater contamination monitoring as the Python package PyLEnM (Python for Long-term Environmental Monitoring). PyLEnM aims to establish the seamless data-to-ML pipeline with various utility functions, such as quality assurance and quality control (QA/QC), coincident/colocated data identification, the automated ingestion and processing of publicly available spatial data layers, and novel data summarization/visualization. The key ML innovations include (1) time series/multianalyte clustering to find the well groups that have similar groundwater dynamics and to inform spatial interpolation and well optimization, (2) the automated model selection and parameter tuning, comparing multiple regression models for spatial interpolation, (3) the proxy-based spatial interpolation method by including spatial

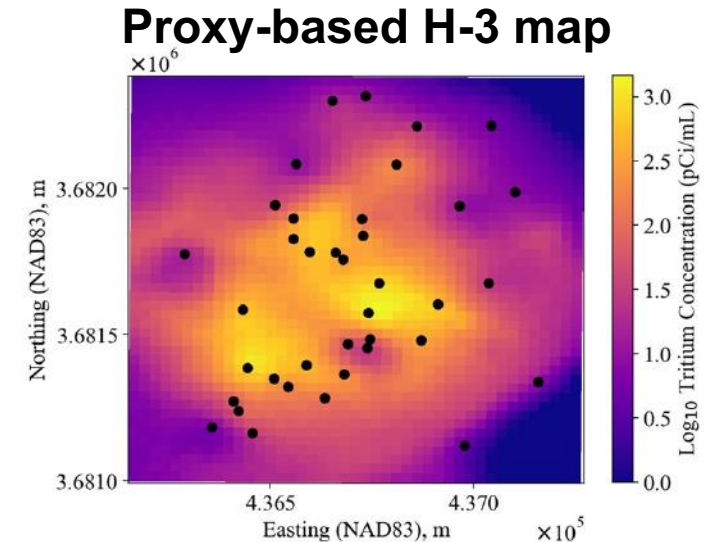
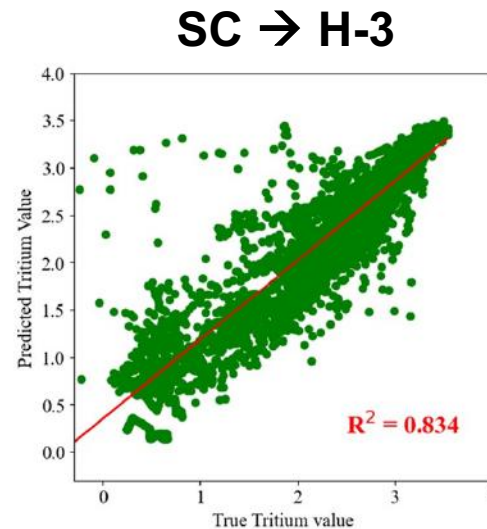
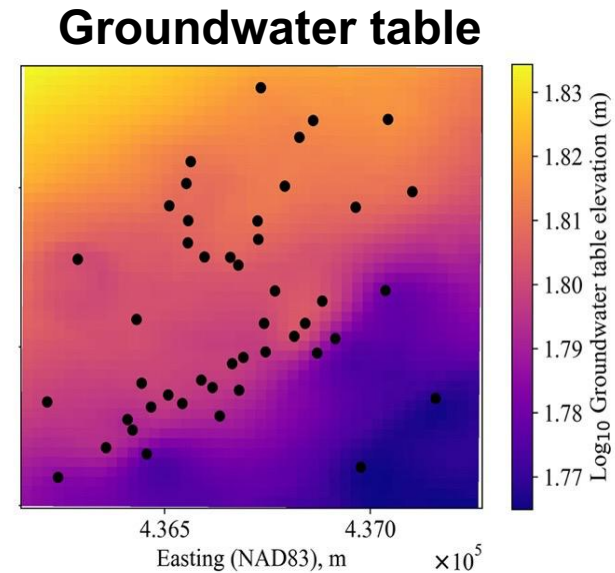
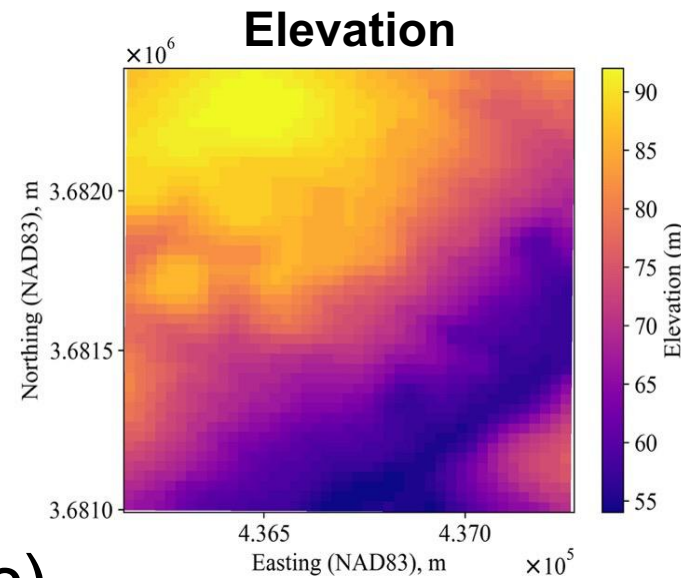


PyLenM: Python for Long-term Env. Monitoring



PyLenM: Supervised Learning

- **Spatiotemporal Interpolation**
 - Groundwater table
 - Contaminant concentration
- **Proxy variables**
 - LiDAR elevation data
 - Topographic metrics (slope etc)
 - Distance from the source
 - In situ measurable SC
→ tritium concentration
- **Comparison of multiple regression methods**



PyLenM: Well Placement Optimization

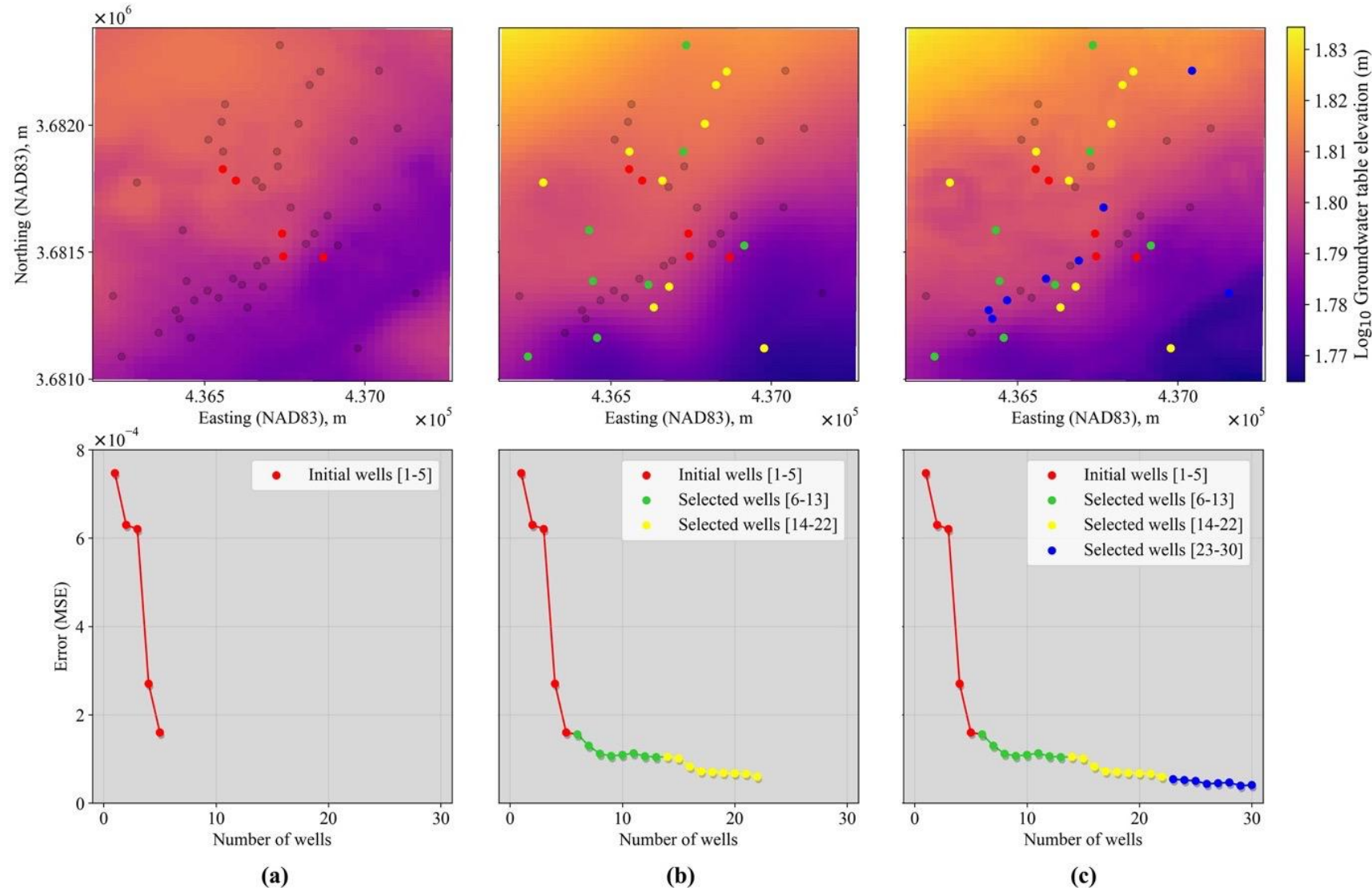
Sub-selection of wells for long-term monitoring

Greedy algorithm

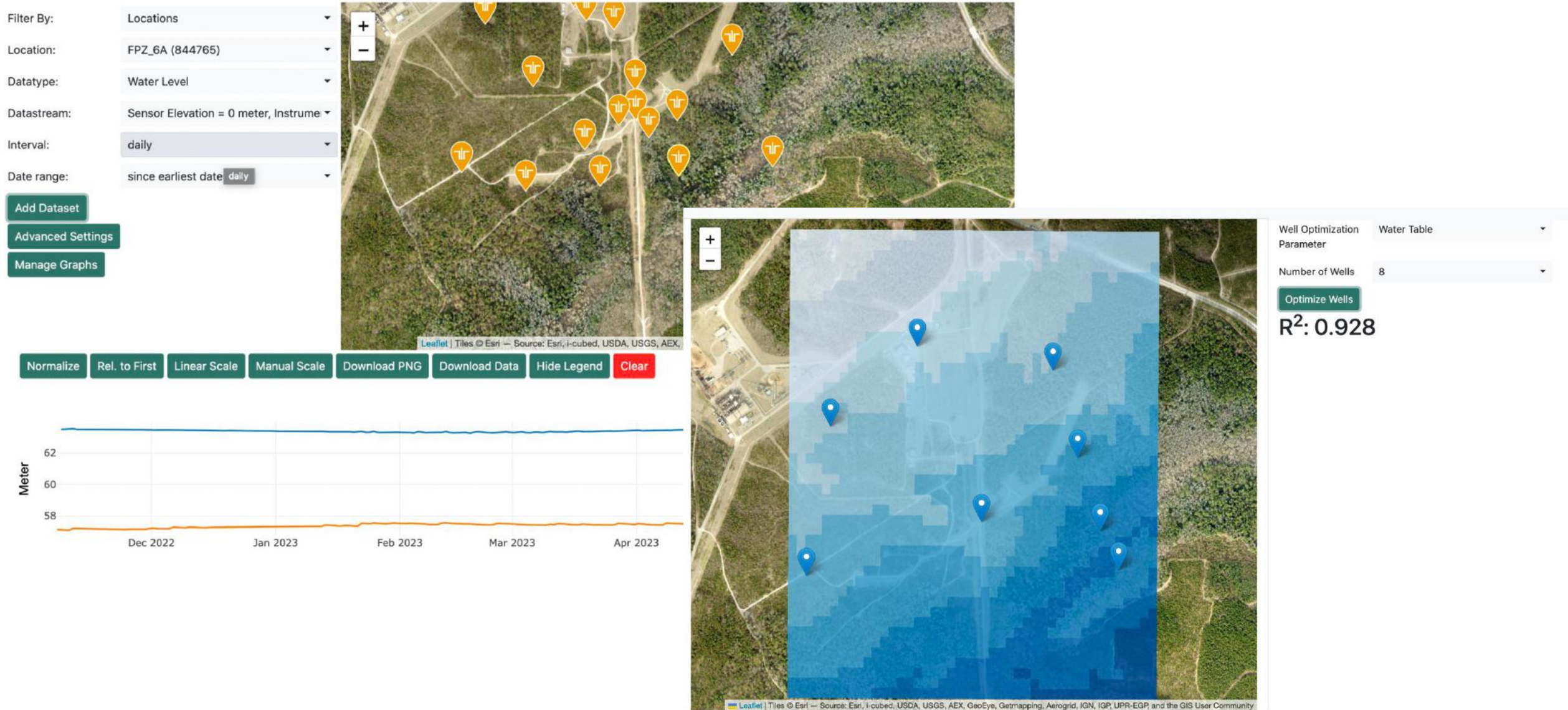
- Reference map created using all the wells
- Interpolation with one additional well at a time
- Find the well that minimize the overall error

Minimum-but-sufficient # wells

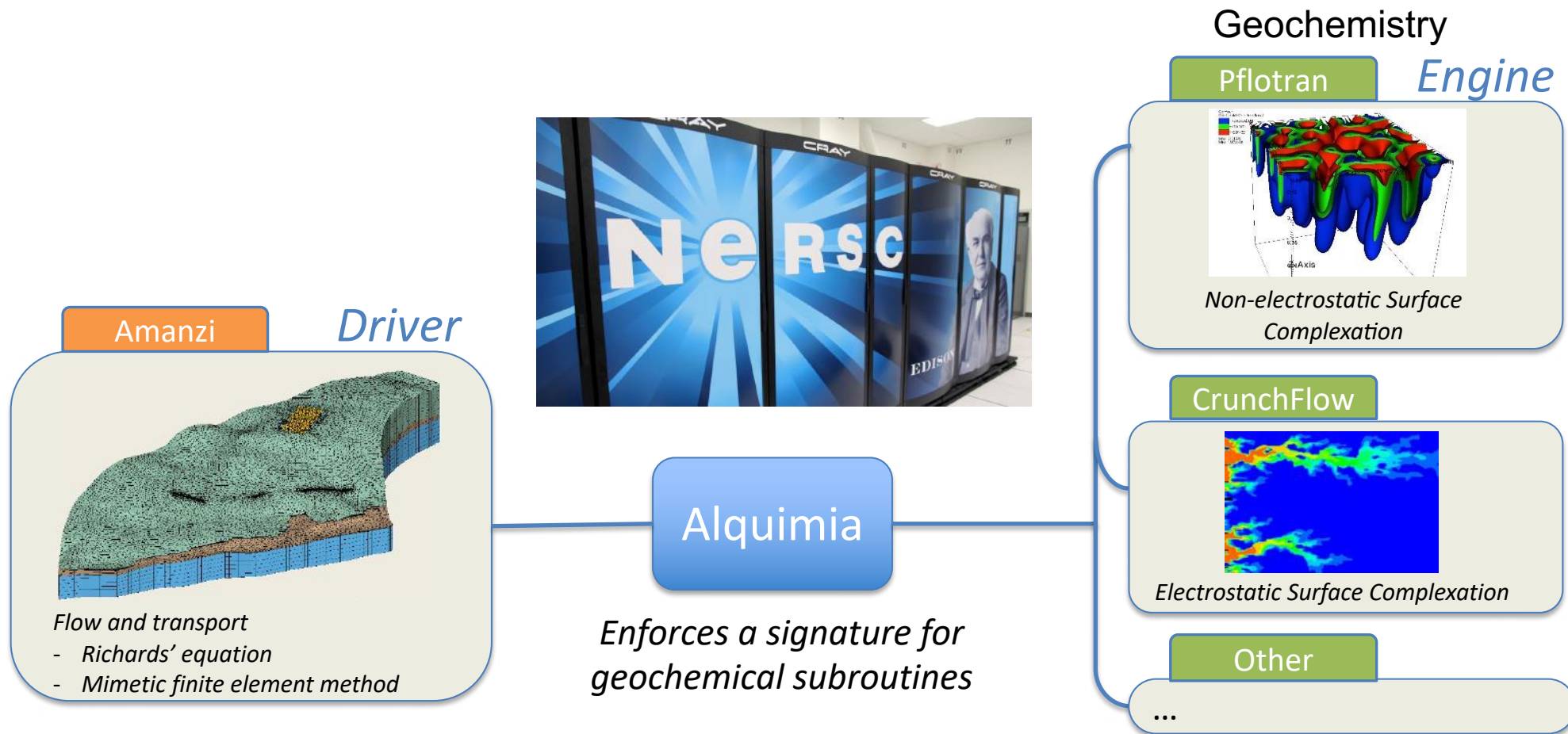
- Error convergence



Web-interface: Implementation



Contaminant Transport Modeling: Amanzi



- **Complex flow**
- **Complex geochemical reactions**



SRS F-Area: Geochemical Model

Reaction	log K (25 °C)
Aqueous species¹	
$(\text{UO}_2)_2 (\text{OH})_2^{+2} \leftrightarrow 2\text{UO}_2^{+2} + 2\text{H}_2\text{O} - 2\text{H}^+$	5.62
$(\text{UO}_2)_2 \text{CO}_3 (\text{OH})_3^- \leftrightarrow 2\text{UO}_2^{+2} + 3\text{H}_2\text{O} + \text{HCO}_3^- - 4\text{H}^+$	11.18
$(\text{UO}_2)_2 \text{OH}^{+3} \leftrightarrow 2\text{UO}_2^{+2} + \text{H}_2\text{O} - \text{H}^+$	2.7
$(\text{UO}_2)_3 (\text{CO}_3)_6^{-6} \leftrightarrow 3\text{UO}_2^{+2} + 6\text{HCO}_3^- - 6\text{H}^+$	7.97
$(\text{UO}_2)_3 (\text{OH})_4^{+2} \leftrightarrow 3\text{UO}_2^{+2} + 4\text{H}_2\text{O} - 4\text{H}^+$	11.9
$\text{UO}_2 (\text{OH})_4^{-2} \leftrightarrow \text{UO}_2^{+2} + 4\text{H}_2\text{O} - 4\text{H}^+$	32.4
$(\text{UO}_2)_3 (\text{OH})_5^+ \leftrightarrow 3\text{UO}_2^{+2} + 5\text{H}_2\text{O} - 5\text{H}^+$	15.55
$(\text{UO}_2)_3 (\text{OH})_7^- \leftrightarrow 3\text{UO}_2^{+2} + 7\text{H}_2\text{O} - 7\text{H}^+$	32.2
$(\text{UO}_2)_3 \text{O} (\text{OH})_2 (\text{HCO}_3)^+ \leftrightarrow 3\text{UO}_2^{+2} + 3\text{H}_2\text{O} + \text{HCO}_3^- -$	9.68
$(\text{UO}_2)_4 (\text{OH})_7^+ \leftrightarrow 4\text{UO}_2^{+2} + 7\text{H}_2\text{O} - 7\text{H}^+$	21.9
$\text{UO}_2 \text{NO}_3^+ \leftrightarrow \text{UO}_2^{+2} + \text{NO}_3^-$	-0.3
$\text{UO}_2 (\text{OH})^+ \leftrightarrow \text{UO}_2^{+2} + \text{H}_2\text{O}$	5.25
$\text{UO}_2 (\text{OH})_2 (\text{aq}) \leftrightarrow \text{UO}_2^{+2} + 2\text{H}_2\text{O} - 2\text{H}^+$	12.15
$\text{UO}_2 (\text{OH})_3^- \leftrightarrow \text{UO}_2^{+2} + 3\text{H}_2\text{O} - 3\text{H}^+$	20.25
$\text{UO}_2 \text{CO}_3 (\text{aq}) \leftrightarrow \text{UO}_2^{+2} + \text{HCO}_3^- - \text{H}^+$	0.39
$\text{UO}_2 (\text{CO}_3)_2^{-2} \leftrightarrow \text{UO}_2^{+2} + 2\text{HCO}_3^- - 2\text{H}^+$	4.05
$\text{UO}_2 (\text{CO}_3)_3^{-4} \leftrightarrow \text{UO}_2^{+2} + 3\text{HCO}_3^- - 3\text{H}^+$	9.14
$\text{CaUO}_2 (\text{CO}_3)_3^{-2} \leftrightarrow \text{Ca}^{+2} + \text{UO}_2^{+2} + 3\text{HCO}_3^- - 3\text{H}^+$	3.8
$\text{Ca}_2\text{UO}_2 (\text{CO}_3)_3 (\text{aq}) \leftrightarrow 2\text{Ca}^{+2} + \text{UO}_2^{+2} + 3\text{HCO}_3^- - 3\text{H}^+$	0.29
$\text{MgUO}_2 (\text{CO}_3)_3^{-2} \leftrightarrow \text{Mg}^{+2} + \text{UO}_2^{+2} + 3\text{HCO}_3^- - 3\text{H}^+$	5.19
$\text{UO}_2 \text{SiO} (\text{OH})_3^+ \leftrightarrow \text{SiO}_2 (\text{aq}) + \text{UO}_2^{+2} + 2\text{H}_2\text{O} - \text{H}^+$	2.48

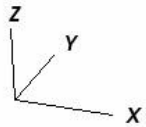
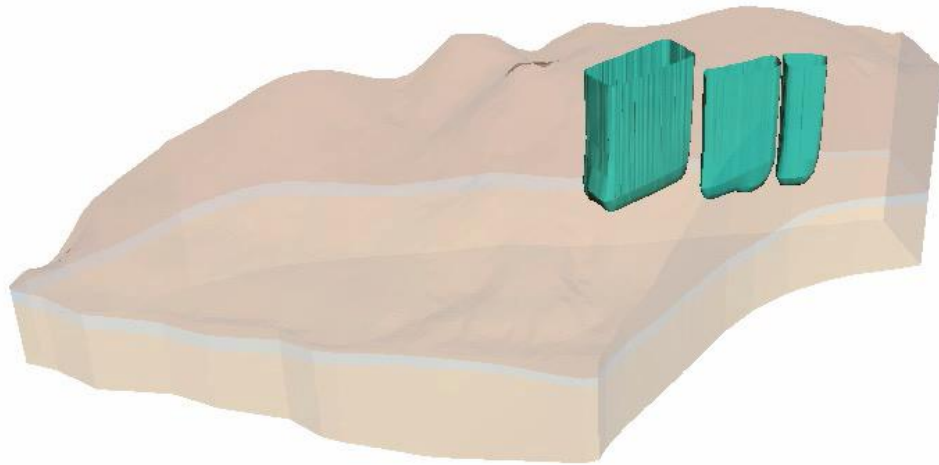
Surface and exchange species²

$(> \text{k-OH})_2 \text{UO}_2^+ \leftrightarrow 2 > \text{k-OH}^{-0.5} + \text{UO}_2^{+2}$	-5.3
$(> \text{k-OH})_2 \text{UO}_2 \text{CO}_3^- \leftrightarrow 2 > \text{k-OH}^{-0.5} + \text{UO}_2^{+2} + \text{HCO}_3^- - \text{H}^+$	-6.2
$> \text{k-OH}_2^{+0.5} \leftrightarrow > \text{k-OH}^{-0.5} + \text{H}^+$	-4.9
$> \text{k-OHNa}^{+0.5} \leftrightarrow > \text{k-OH}^{-0.5} + \text{Na}^+$	2.1
$> \text{k-OH}_2 \text{NO}_3^{-0.5} \leftrightarrow > \text{k-OH}^{-0.5} + \text{H}^+ + \text{NO}_3^-$	-4.9
$> \text{k}_2\text{UO}_2 \leftrightarrow 2 > \text{k}^- + \text{UO}_2^{+2}$	-7.1
$> \text{kNa} \leftrightarrow > \text{k}^- + \text{Na}^+$	-2.9
$> \text{kH} \leftrightarrow > \text{k}^- + \text{H}^+$	-4.5
$> \text{k}_2\text{Ca} \leftrightarrow 2 > \text{k}^- + \text{Ca}^{+2}$	-6.8
$> \text{k}_3\text{Al} \leftrightarrow 3 > \text{k}^- + \text{Al}^{+3}$	-8
$(> \text{Fe-OH})_2 \text{UO}_2^+ \leftrightarrow 2 > \text{Fe-OH}^{-0.5} + \text{UO}_2^{+2}$	-14.11
$(> \text{Fe-OH})_2 \text{UO}_2 \text{CO}_3^- \leftrightarrow 2 > \text{Fe-OH}^{-0.5} + \text{UO}_2^{+2} + \text{HCO}_3^- - \text{H}^+$	-4.35
$> \text{Fe-OH}_2^{+0.5} \leftrightarrow > \text{Fe-OH}^{-0.5} + \text{H}^+$	-9.18
$(> \text{Fe-OH})_2 \text{CO}_3^- \leftrightarrow 2 > \text{Fe-OH}^{-0.5} + \text{H}^+ + \text{HCO}_3^- - 2\text{H}_2\text{O}$	-12.23
$> \text{Fe-OCO}_2 \text{Na}^{-0.5} \leftrightarrow > \text{Fe-OH}^{-0.5} + \text{Na}^+ + \text{HCO}_3^- - \text{H}_2\text{O}$	-3.28
$> \text{qz-OH}_2^+ \leftrightarrow > \text{qz-OH} + \text{H}^+$	1.1 ³
$> \text{qz-O}^- \leftrightarrow > \text{qz-OH} - \text{H}^+$	8.1 ³
$> \text{qz-ONa} \leftrightarrow > \text{qz-OH} - \text{H}^+ + \text{Na}^+$	6.8 ⁴



3D Plume Modeling and Simulations

DB: plot_data.VisIt.xmf
Time: 1956



user: user
Sun Apr 14 10:34:15 2019

Helpful for understanding

- Residual contaminants under the basins and within the clay layer
- Climate change impact (Libera et al., 2019; Xu et al., 2022)

High computational burden

- >1M grid blocks
- Complex geochemical reactions
- Up to 100s of simulations for UQ

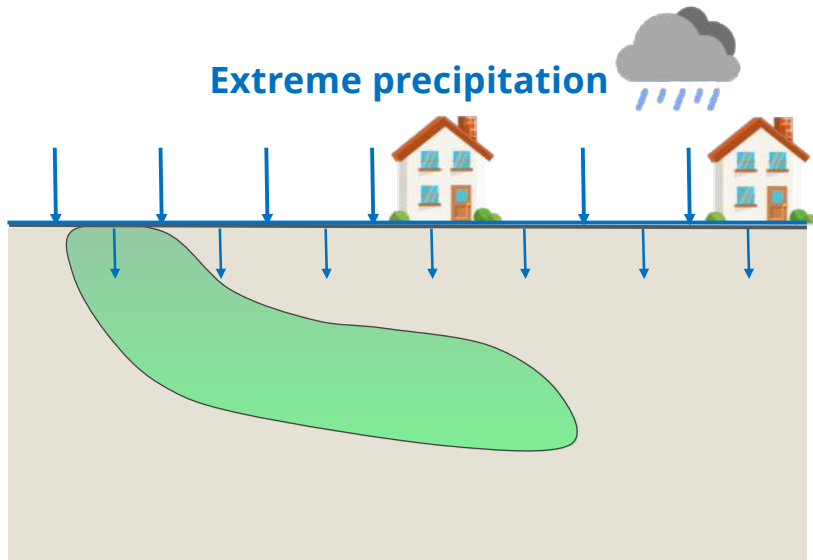
Challenging to fit at all the points

- Geological heterogeneity
- Conceptual model error

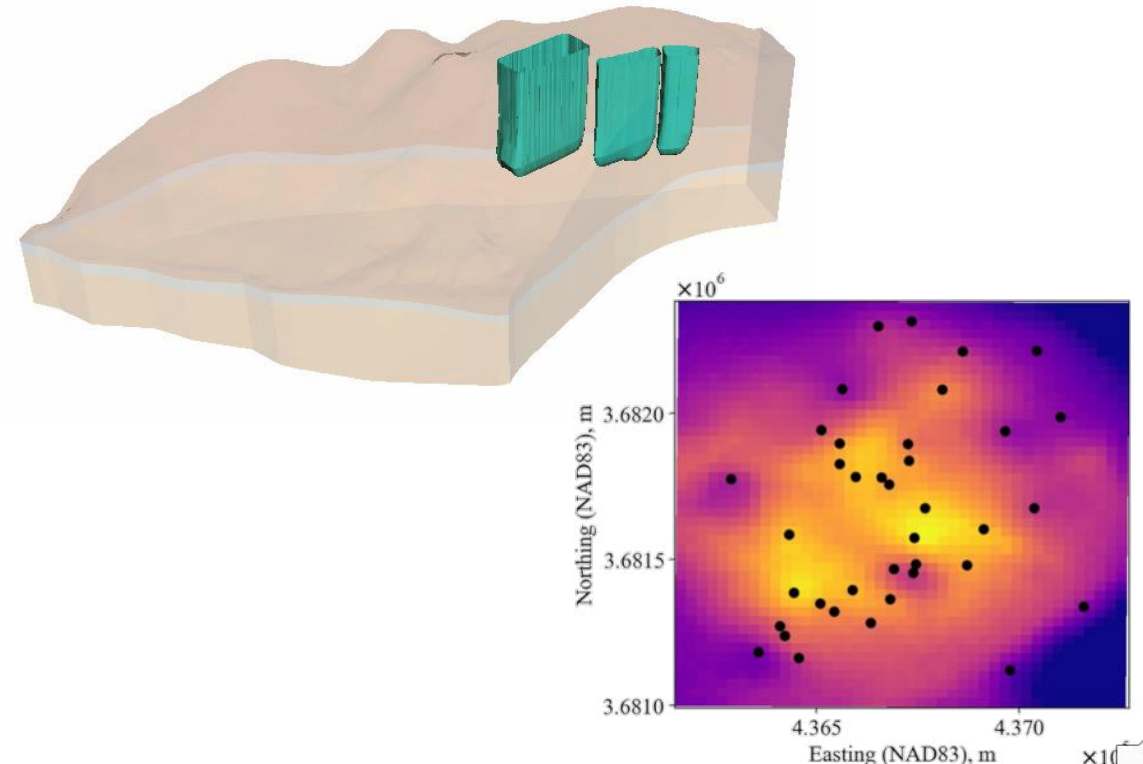


Simulation Intelligence: Simulations x ML/AI

Climate Change Impact on Groundwater contamination
→ Emulator with Fourier Neural Operator



Physics-informed interpolation
→ Model-data integration with Bayesian hierarchical model

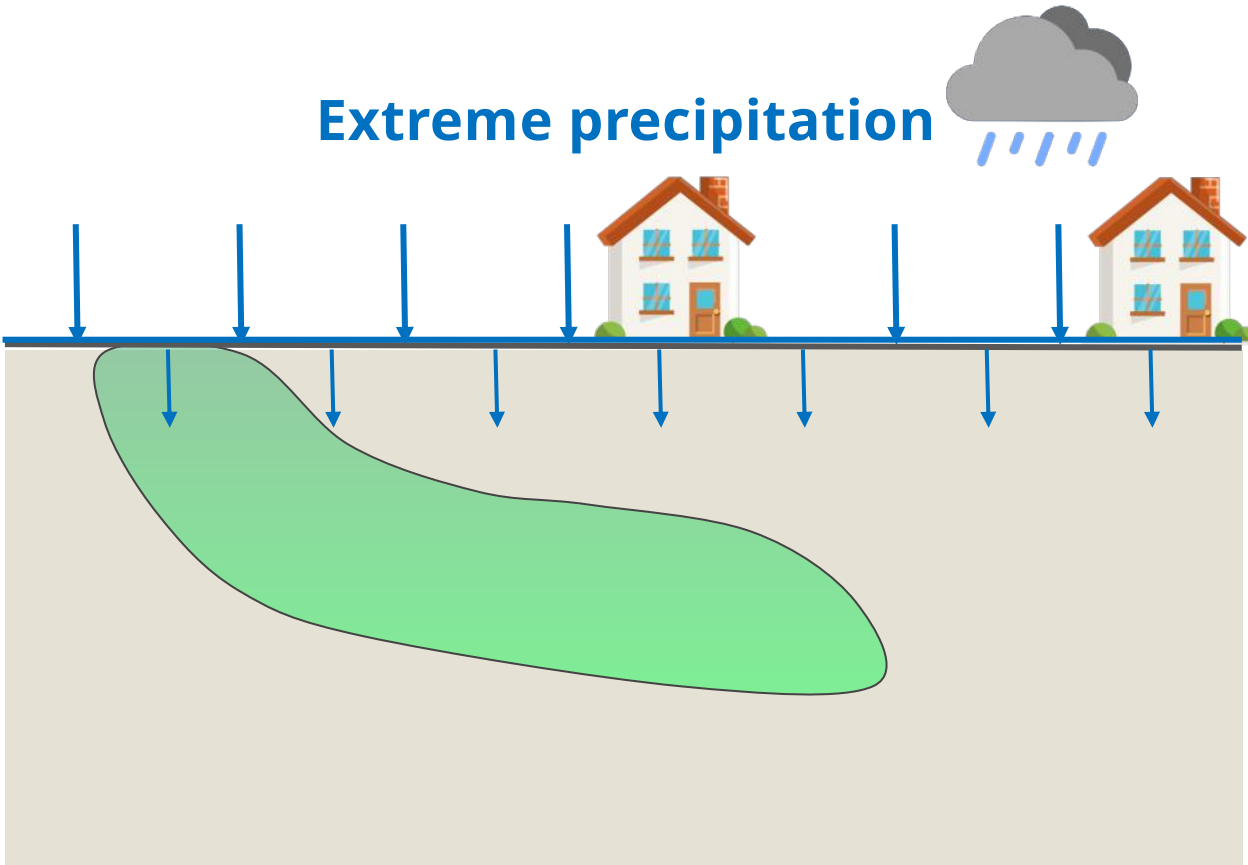


In collaboration with NASA Frontier Development Lab

Lavin, A., Zenil, H., Paige, B., Krakauer, D., Gottschlich, J., Mattson, T., ... & Pfeffer, A. (2021). Simulation intelligence. Towards a new generation of scientific methods. *arXiv preprint arXiv:2112.03235*.

Climate Change Impacts on Contamination

Extreme precipitation



Higher precipitation

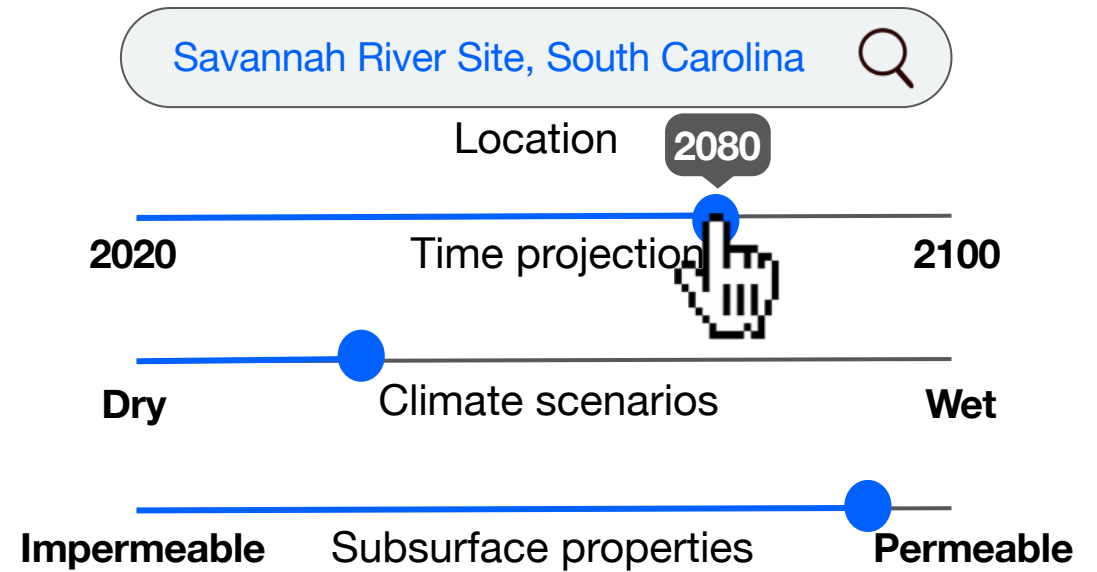
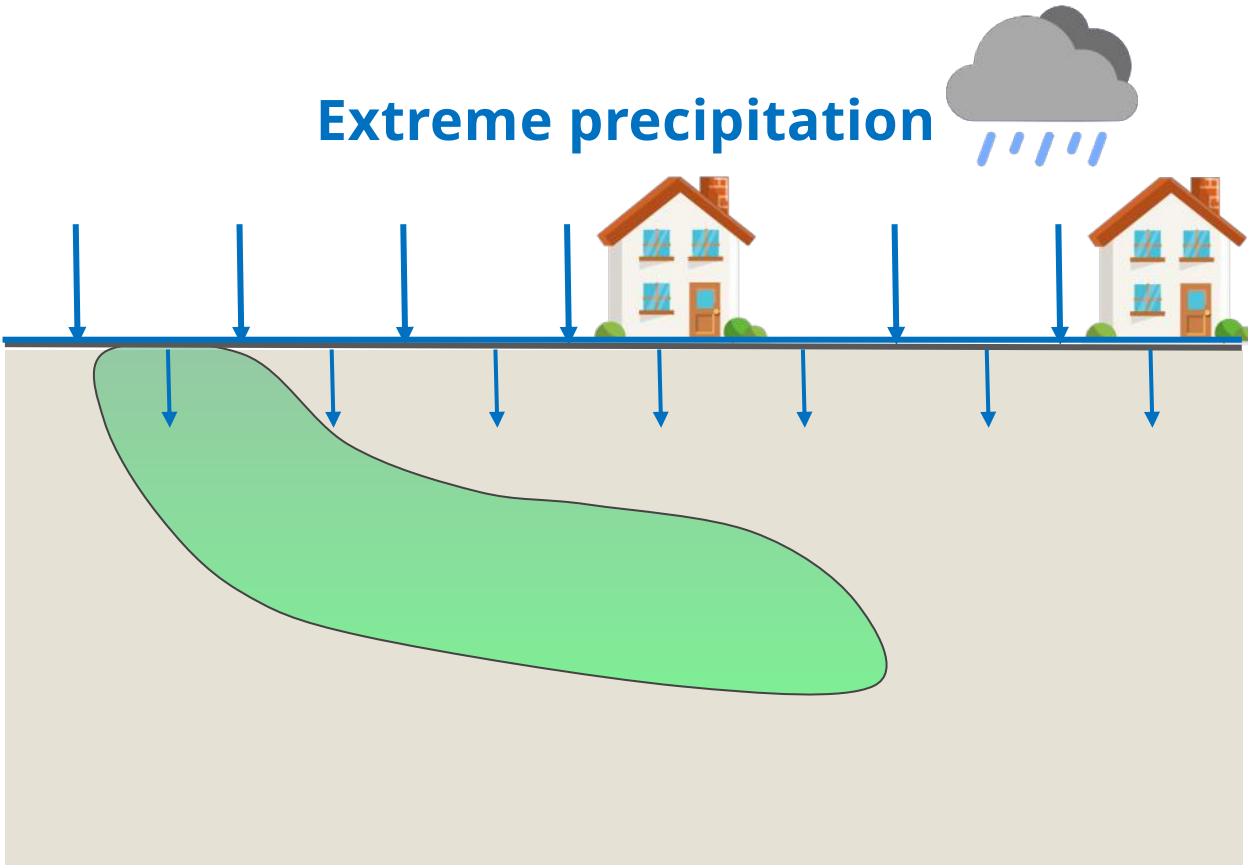
- Re-mobilize residual contaminants?
- Dilute concentrations?
- Change management strategies?
- Change monitoring well configuration?



Climate Change Impacts on Contamination

When and where to make modification?

Extreme precipitation



- But computation is pretty heavy
- We can't run simulation on laptops

Wang, L., Kurihana, T., Meray, A., Mastilovic, I., Praveen, S., Xu, Z., ... & Wainwright, H. (2022). Multi-scale Digital Twin: Developing a fast and physics-informed surrogate model for groundwater contamination with uncertain climate models. *arXiv preprint arXiv:2211.10884*.



Emulator/Surrogate Modeling



HPC
Clusters



Parameters
from PDF

$$\{p_1, p_2, \dots, p_N\}$$

Simulations:

$$\phi = f(p)$$
$$\{\phi_1, \phi_2, \dots, \phi_N\}$$

Regressions:

$$\phi \sim f(p)$$

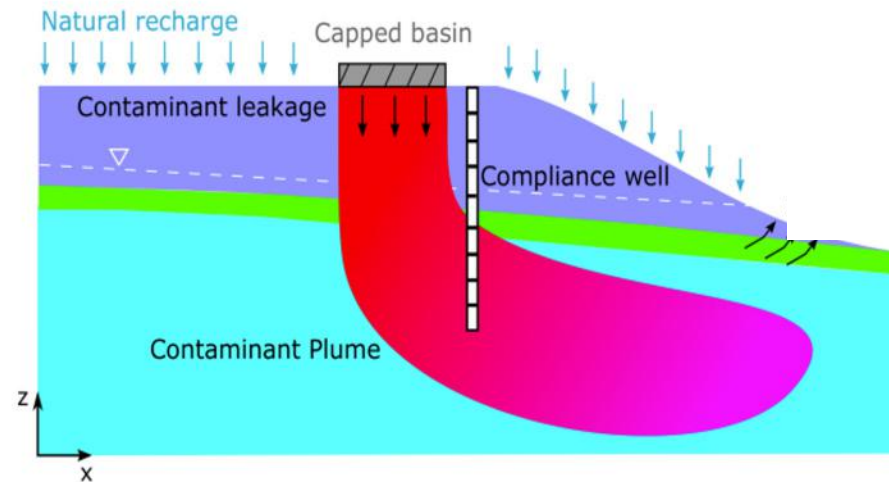
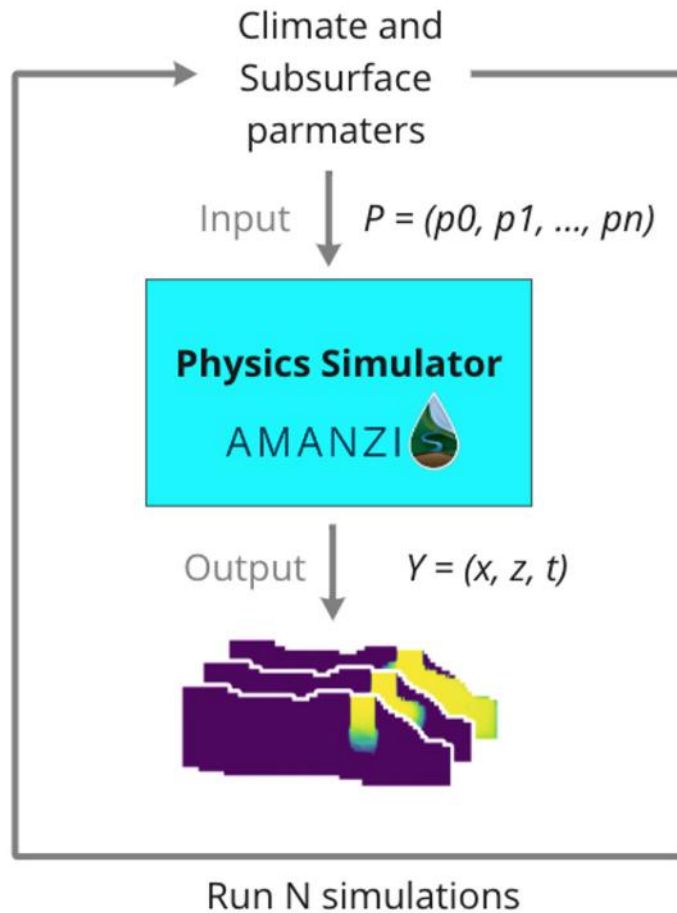
Emulator
predictions

$$\phi = f_{emulate}(p)$$

Statistical representation of physical models



2D Flow and Transport simulator



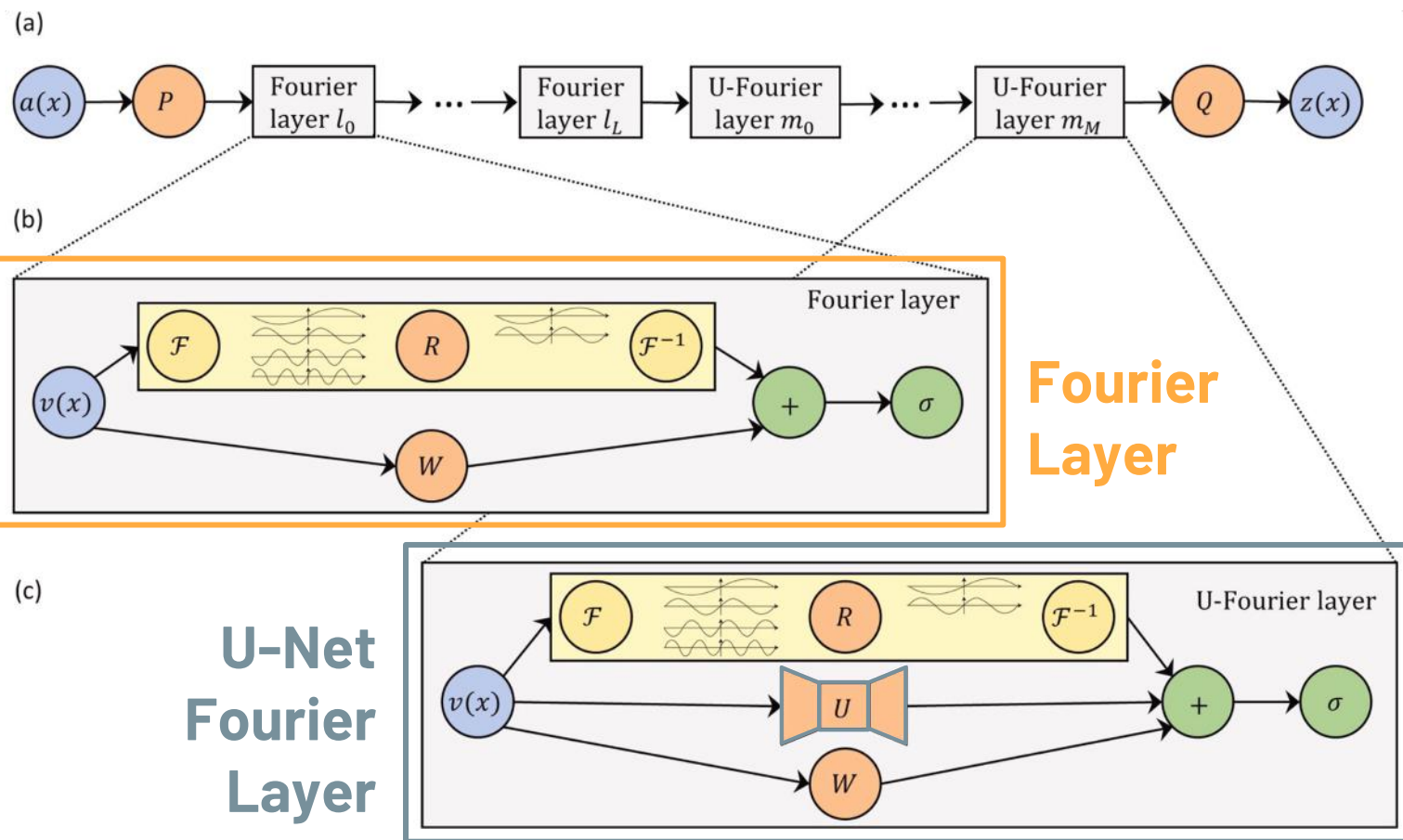
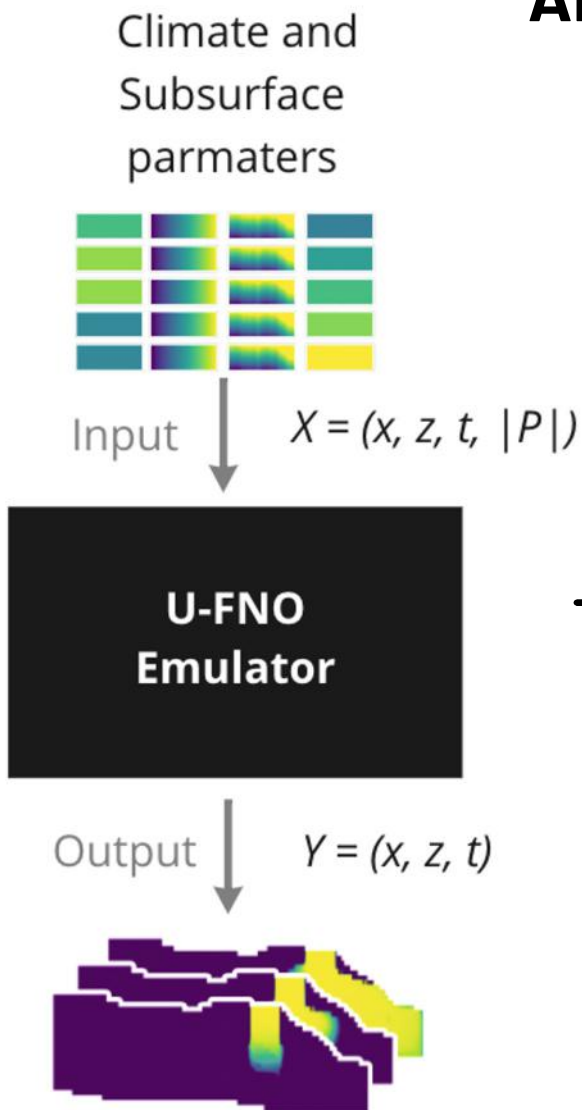
Parameter	Details
Permeability	upper layer
Porosity	upper layer
Alpha	inverse air entry suction
Sr	residual water content
m	$m = 1-1/n$, a measure of the pore-size distribution
Source concentration	Initial contaminant concentration
Discharge rate (Cap at 1988)	Source/cap discharge rate in volumetric water
Time-varying recharge	Climate data (precip. & ET) (history, mid-century, late-century)

Flow and reactive transport model in 2D



Deep learning architecture: UFNO

An enhanced Fourier Neural Operator (Wen et al. 2022)



Data-driven & Physics-informed Loss

A combined loss function to minimize multiple errors → UFNOB

$$\mathcal{L}(y, \hat{y}) = \mathcal{L}_{MRE}(y, \hat{y}) + \beta_1 \mathcal{L}_{der}(y, \hat{y}) + \beta_2 \mathcal{L}_{plume}(c', \hat{c}') + \beta_3 \mathcal{L}_{BC}(\hat{y})$$

Mean Relative Error

$$\mathcal{L}_{MRE}(y, \hat{y}) = \frac{\|y - \hat{y}\|_2}{\|y\|_2}$$

Derivatives

$$\mathcal{L}_{der}(y, \hat{y}) = \frac{\|\partial y / \partial x - \partial \hat{y} / \partial x\|_2}{\|\partial y / \partial x\|_2} + \frac{\|\partial y / \partial z - \partial \hat{y} / \partial z\|_2}{\|\partial y / \partial z\|_2}$$

Contaminant boundary

$$\mathcal{L}_{plume}(c', \hat{c}') = \frac{\|\partial c' / \partial x - \partial \hat{c}' / \partial x\|_2}{\|\partial c' / \partial x\|_2} + \frac{\|\partial c' / \partial z - \partial \hat{c}' / \partial z\|_2}{\|\partial c' / \partial z\|_2}, \text{ where } c' = \begin{cases} 0, & c < MCL \\ 1, & c \geq MCL \end{cases} \quad \hat{c}' = \begin{cases} 0, & \hat{c} < MCL \\ 1, & \hat{c} \geq MCL \end{cases}$$

MCL: maximum contaminant level

$$\mathcal{L}_{BC}(\hat{y}) = \|\hat{q}_x|_{\partial D}\|_2 + \|\hat{q}_z|_{\partial D}\|_2 + \|\partial \hat{h}|_{\partial D}\|_2$$

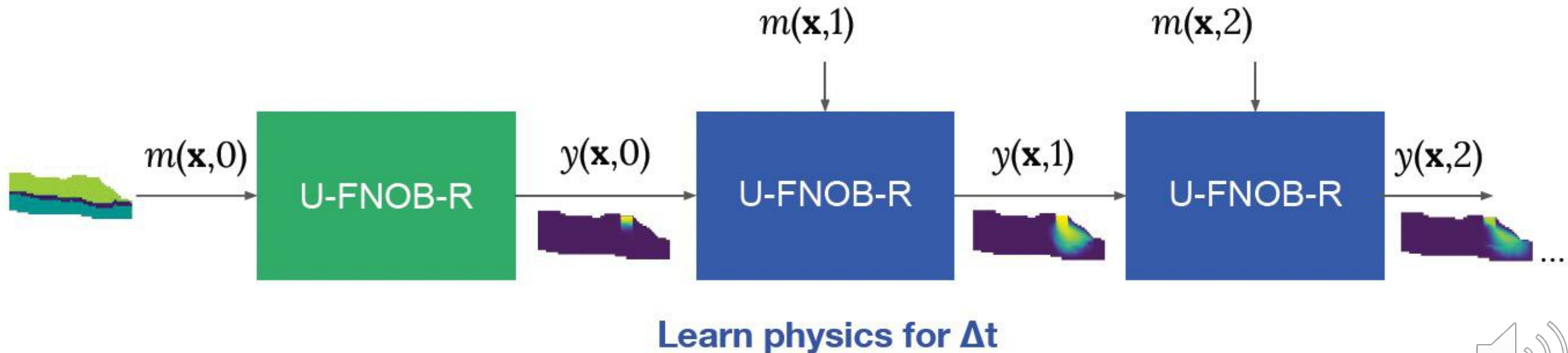


U-FNOB and U-FNOB-Recurrent

a) Architecture 1: U-FNOB



b) Architecture 2: U-FNOB-R (Recurrent)



Emulator-based Plume Prediction

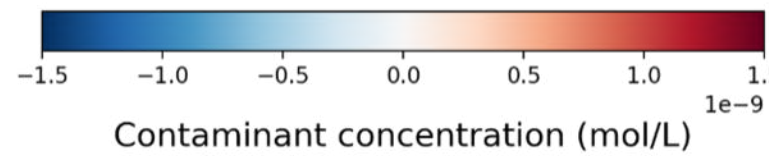
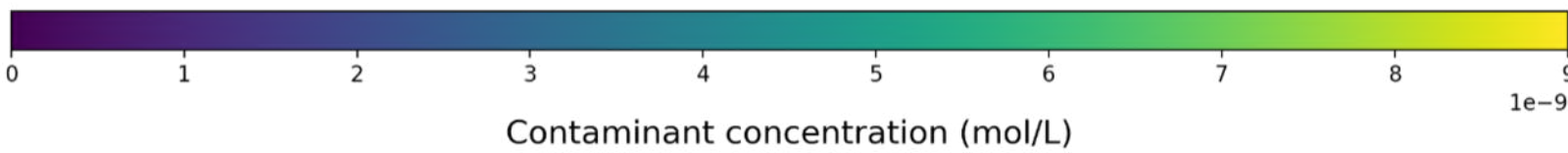
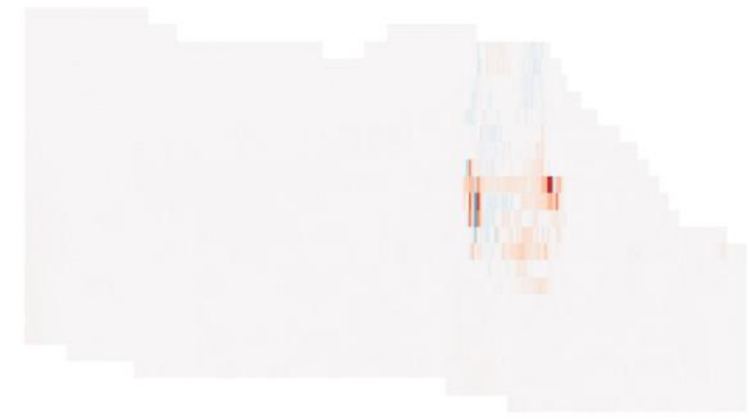
1955

Contaminant Concentration

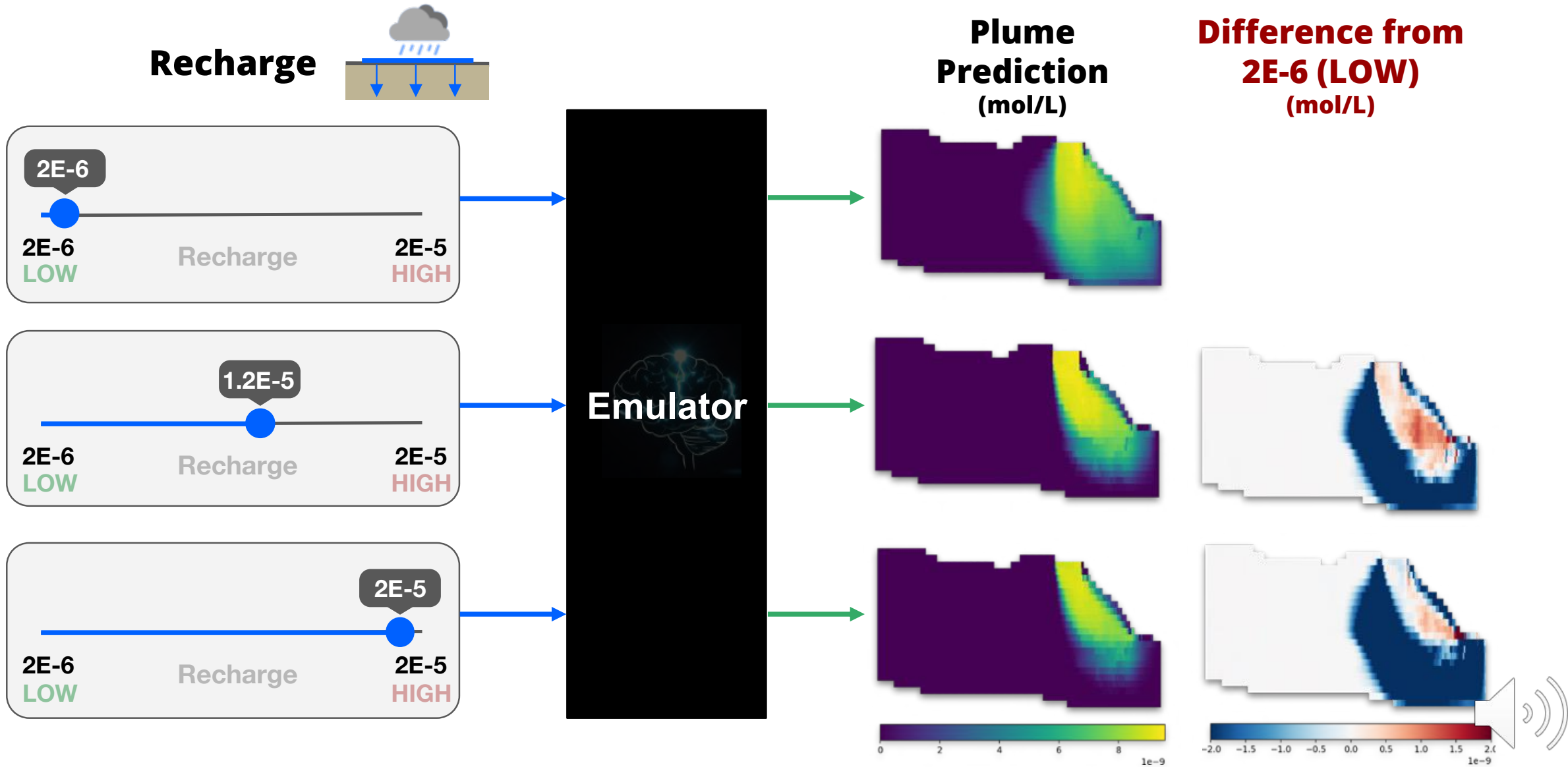
TRUTH

PREDICTION

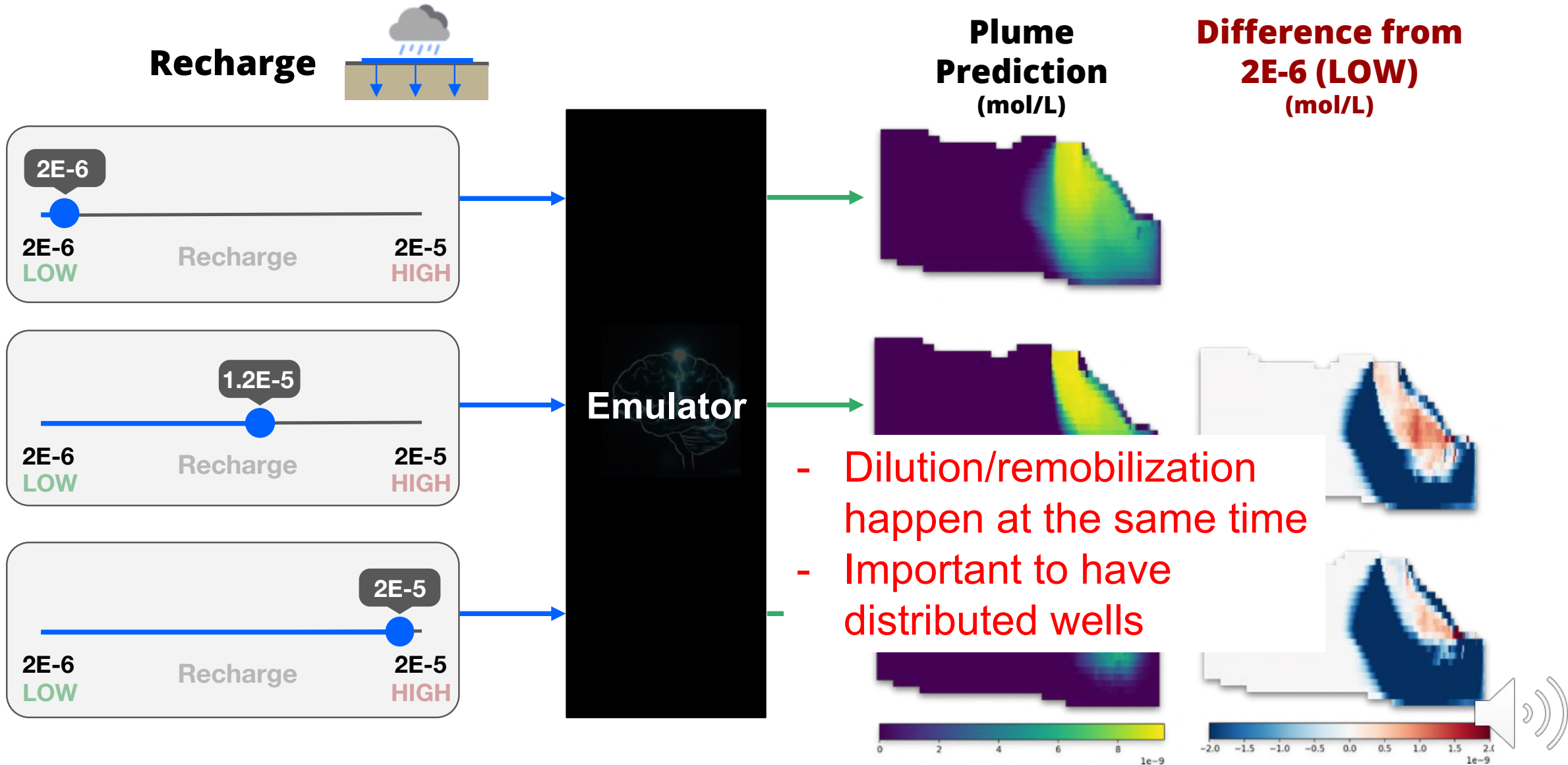
DIFFERENCE



Off-Line Climate Change Assessment



Off-Line Climate Change Assessment



Comparison of Different Strategies

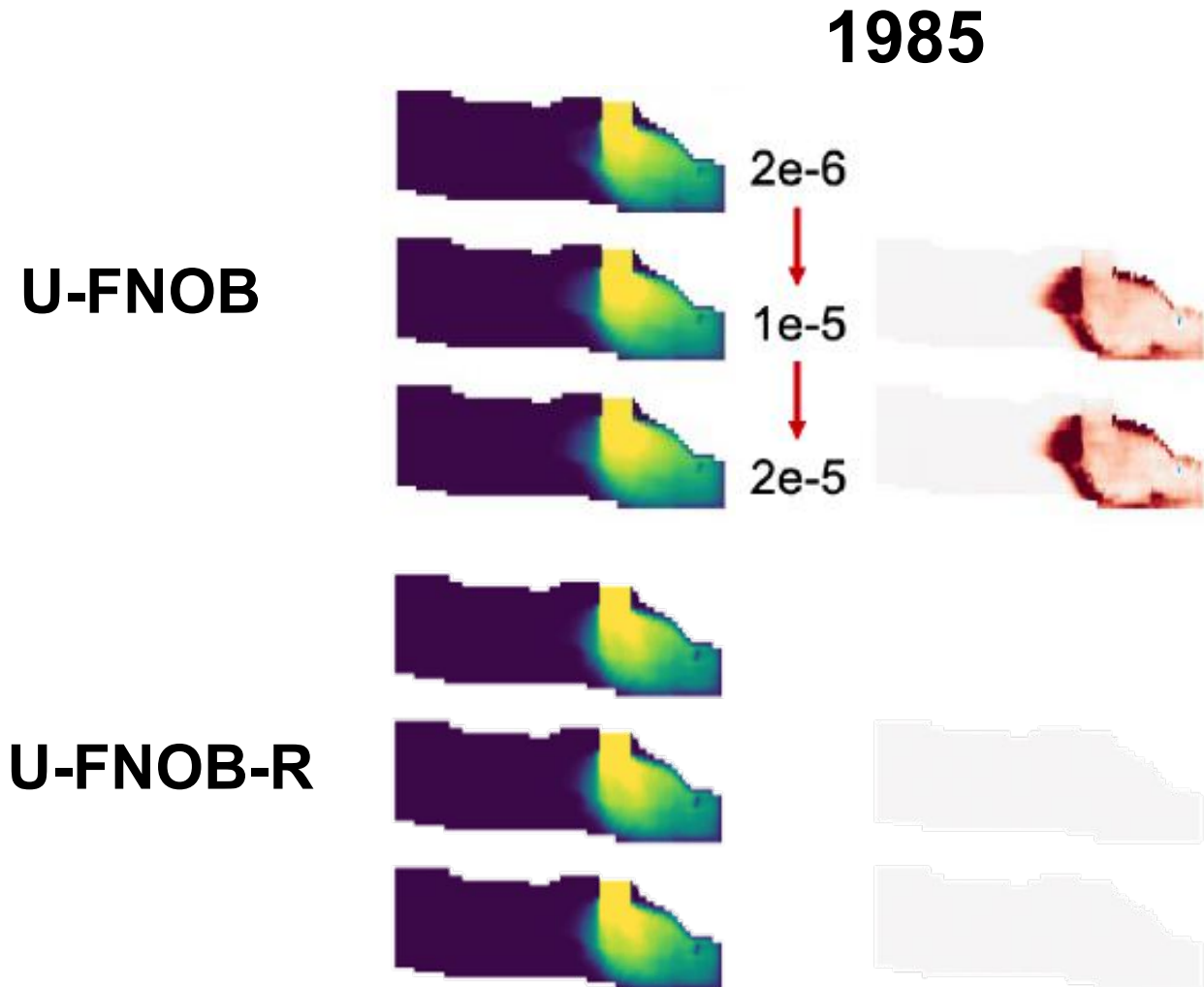
Index	Architectures	Epochs	Loss ($\beta_1, \beta_2, \beta_3$)	MRE	MSE
1	FNO- R	30	(0,0,0)	0.051	2.71e-4
2	FNO	30	(0,0,0)	0.055	2.98e-4
3	U-FNOB-R	30	(0,0,0)	0.035	1.51e-4
4	U-FNOB	30	(0,0,0)	0.037	1.29e-4
5	U-FNOB	30	(0.1,0,0)	0.029	8.83e-5
6	U-FNOB	30	(0,0.1,0)	0.033	1.10e-4
7	U-FNOB	30	(0,0,0.1)	0.034	1.27e-4
8	U-FNOB	30	(0.1,0.1,0.1)	0.028	8.14e-5
9	U-FNOB-R	150	(0.1,0.1,0.1)	0.020	4.49e-5
10	U-FNOB	150	(0.1,0.1,0.1)	0.014	2.44e-5

- FNOB is better than FNOB-Recurrent



Full Time Series vs Recurrent

Changing the mid-century precipitation (kg-water m⁻²s⁻¹) between 2020 and 2060

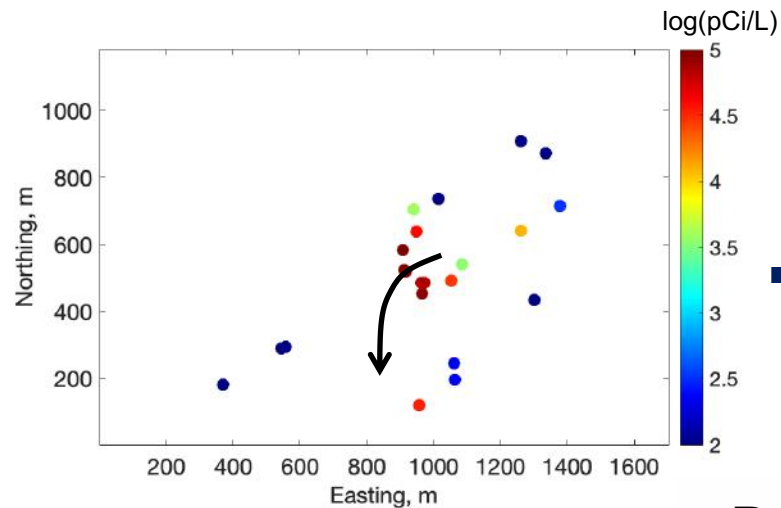


- FNOB changes plumes in the past
- Regression with full time-series does not know past/future
- Recurrent one is more realistic

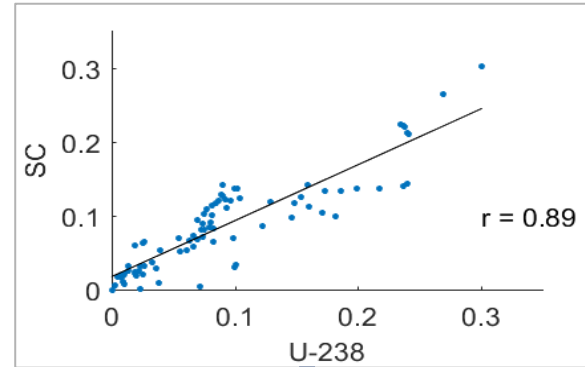
→ Different strategies for time-dependent parameters?

Physics-informed Spatiotemporal Interpolation

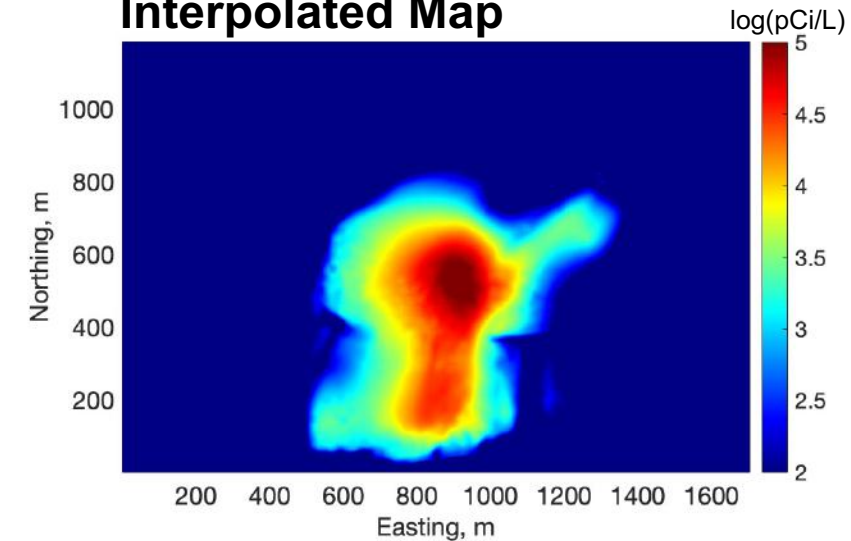
Concentrations at Wells (2015)



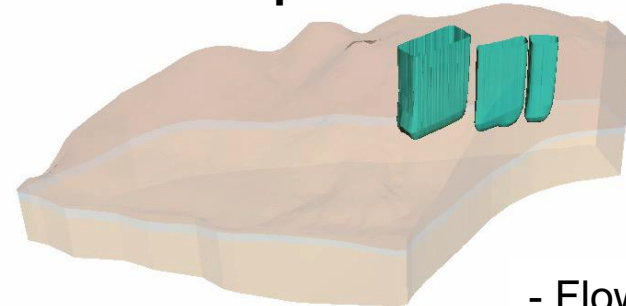
In situ Sensor Data



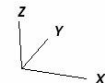
Interpolated Map



Reactive Transport Model



- Flow direction
- Plume source



- Ensemble simulations

$$\{y_1, y_2, \dots, y_T\}$$



Advances in GP for Large Datasets

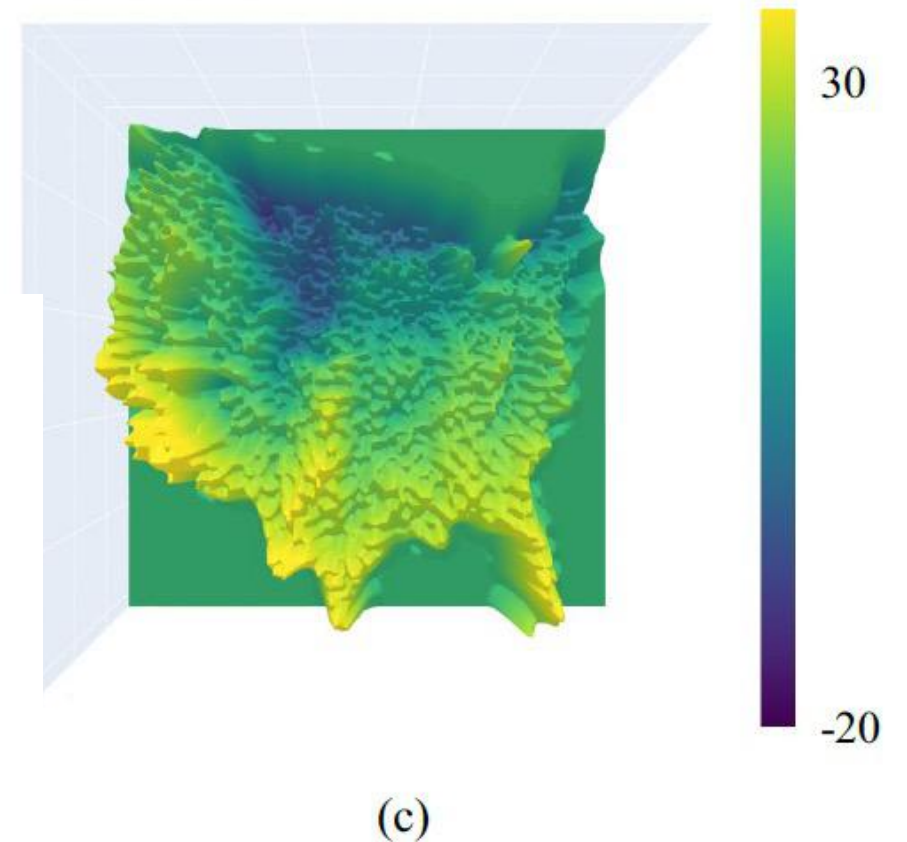
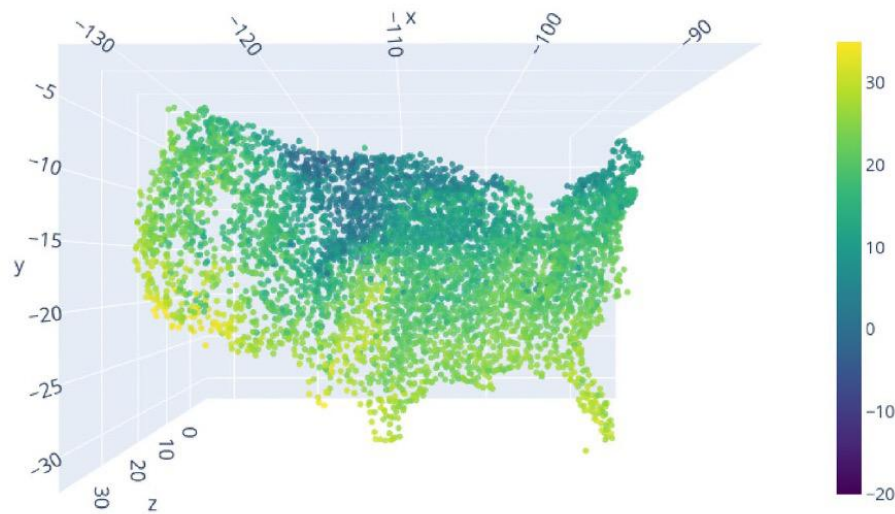
scientific reports

OPEN

Exact Gaussian processes for massive datasets via non-stationary sparsity-discovering kernels

Marcus M. Noack¹, Harinarayan Krishnan¹, Mark D. Risser² & Kristofer G. Reyes³

Check for updates



Bayesian Hierarchical Model for Physics-informed Interpolation

- Estimate the spatiotemporal distribution of contaminant concentrations y (ppm) conditioned on groundwater sampling data (z_G) and in situ sensor (z_S)

$$p(\mathbf{y}|\mathbf{z}_G, \mathbf{z}_S)$$

- Ensemble simulations of plume and concentrations: $\phi = f(\mathbf{p})$
- **Address the bias and errors of simulations: $y = g(\phi) + \varepsilon$**
 - Probably not good for long-term prediction/extrapolation
 - Good for improving the current interpolation
- **Gaussian Process Model: $y \sim N(g(\phi), C)$, C = spatially correlated covariance**

Bayesian Hierarchical Model for Physics-informed Interpolation

- Posterior distribution

$$p(\mathbf{y}|\mathbf{z}_G, \mathbf{z}_S) \propto \int p(\mathbf{z}_S|\mathbf{y})p(\mathbf{z}_G|\mathbf{y})p(\mathbf{y}|\phi(\mathbf{p}), \boldsymbol{\theta})d\boldsymbol{\theta} d\mathbf{p}$$

- Data models:

- $p(\mathbf{z}_S|\mathbf{y})$: the correlation between sensor data and concentrations

$$p(\mathbf{z}_S|\mathbf{y}) = N(h(\mathbf{y}), \tau^2), \tau^2 \text{ is the measurement error}$$

- \mathbf{z}_G : Concentrations plus i.i.d errors

- Prior model:

- $p(\boldsymbol{\theta}) \rightarrow$ GP parameters
- $p(\mathbf{p}) \rightarrow$ model parameters



Bayesian Hierarchical Model for Physics-informed Interpolation

- Ensemble simulations: $\{\phi_1, \phi_2, \dots, \phi_N\}$
- Jeffrey's prior for θ

Algorithm 1: Sampling-Resampling Scheme

1. For k from 1 to N :

1.1. Read the k -th ensemble simulation ϕ_k .

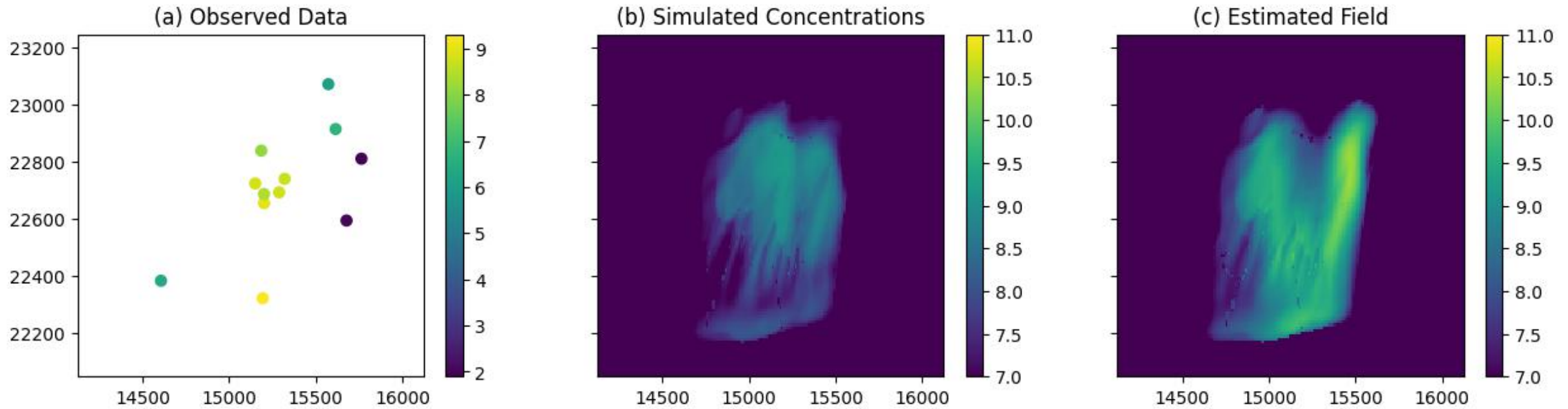
1.2. Sample the hyperparameters θ

1.3. Apply fvGP and estimate \mathbf{y}_k and likelihood

$$L_k = p(\mathbf{z}_G|\mathbf{y})p(\mathbf{z}_S|\mathbf{y})p(\mathbf{y}|\phi, \theta)$$

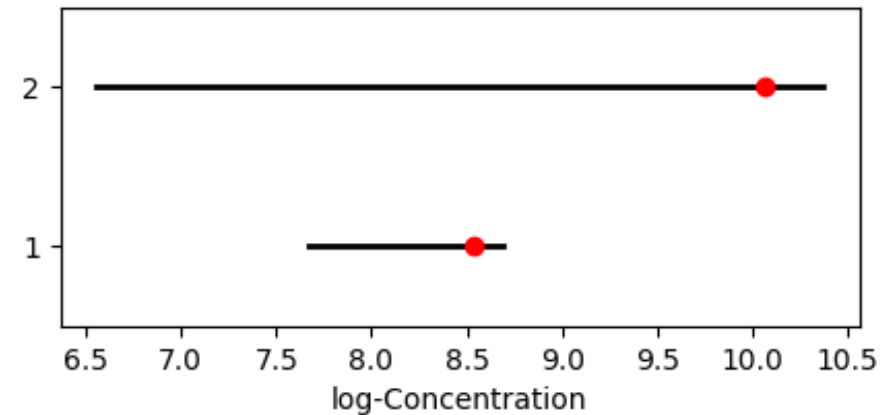
2. Resample \mathbf{y}_k based on the likelihood L_k for the posterior distribution

Physics-informed Spatiotemporal Interpolation



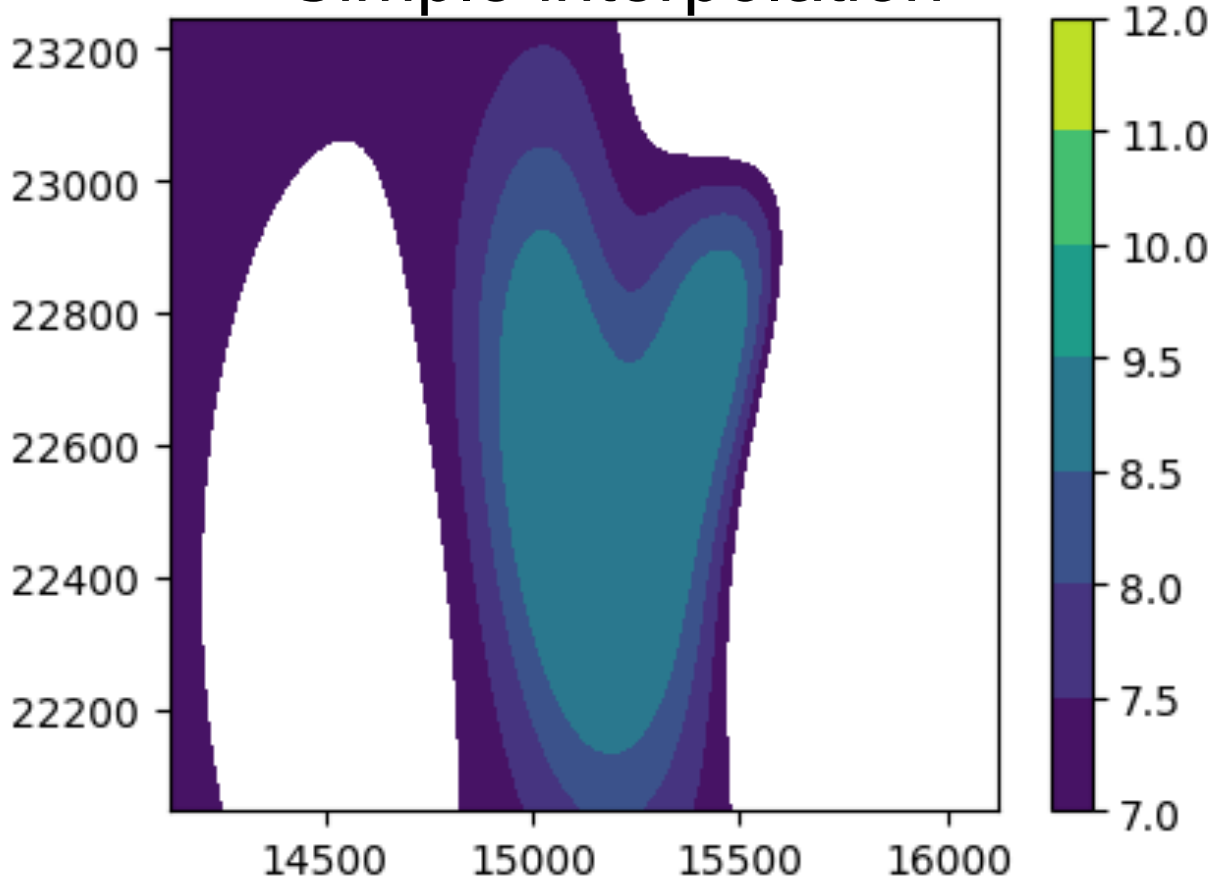
Performance Confirmation

- Confidence intervals
- Points not included in the estimation

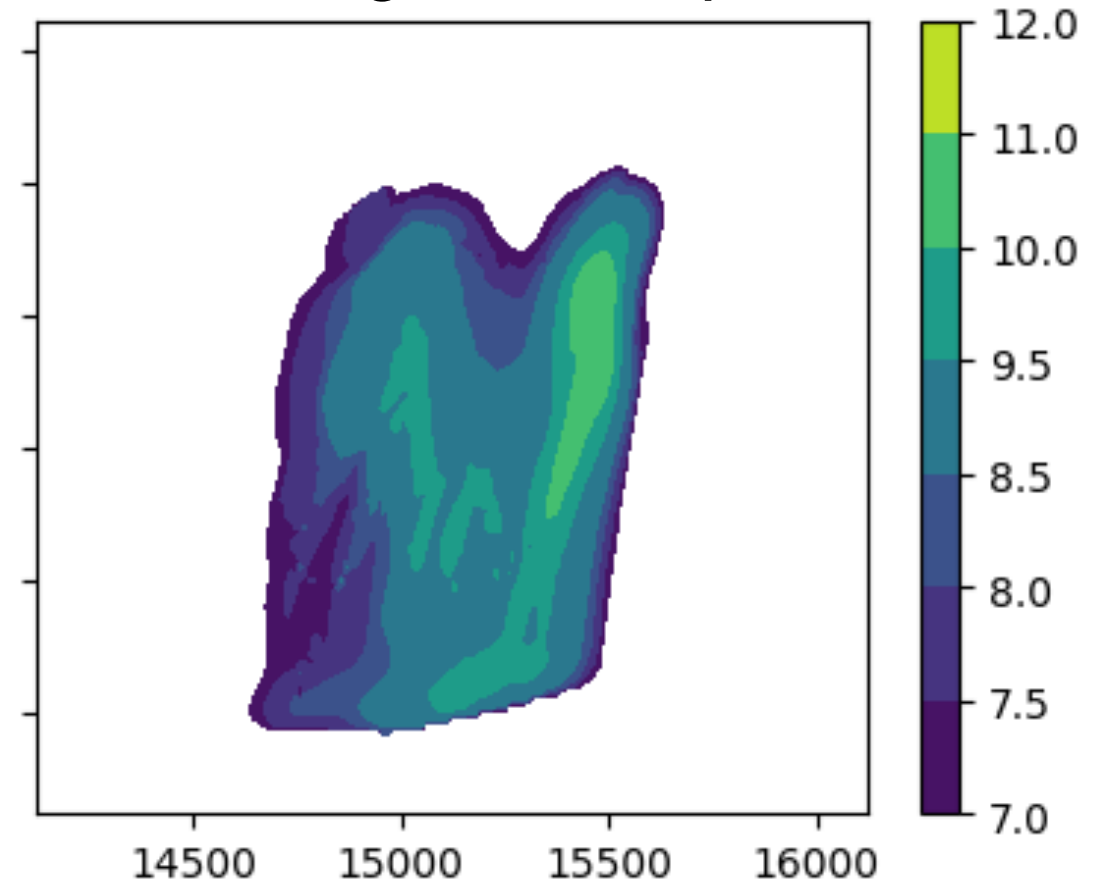


Physics-informed Spatiotemporal Interpolation

Simple interpolation




Integrated map



Impact on Plume Extent: The integrate plume map capture the source locations, plume direction and dispersion

ML Pathway to Adaptation: Challenge

- **Conceptual model development is difficult**
 - Geological heterogeneity, unknowns
- **It is a sensitive topic**
 - Failed weather prediction → 
 - Failed contaminant transport prediction → legal actions...
 - QA/QC of codes
- **Regulations**
 - Processes/paperwork for sensor installment, monitoring modification

ML Pathway to Adaptation: Opportunities

- **Understand regulations**
 - Stepwise implementation: in situ sensor deployment
 - Reducing sampling frequencies is easier
 - Then reducing # variables and reducing # wells
- **Emphasize additional safety assurance**
 - Continuous monitoring → early warning, explaining anomalies
 - Guide monitoring strategies (e.g., climate change)
- **Autonomous/autonomous monitoring → AI-assisted monitoring**
 - Anomaly detection → instrument failure, system changes
 - Realistic plume visualization
 - Digital twin → simulate what can happen in the future

POWER & OPERATIONS

Importance of environmental monitoring for consent-based siting of nuclear facilities

Sat, Nov 19, 2022, 6:04AM | Nuclear News | Haruko Wainwright and Carol Eddy-Dilek



2

Distributed Sources: Agriculture Runoff

- **Runoff water may contain:**
 - Soil
 - Nutrients: nitrogen, phosphorous, trace metals
 - Pesticides: herbicides, insecticides, fungicides
- **Impacts on water quality:**
 - Decreased water quality
 - Harmful algal blooms
 - Fish kills

At least one pesticide was found in about:

- 94 percent of water samples and
- 90 percent of fish samples taken from streams across the Nation
- In nearly 60 percent of shallow wells sampled.

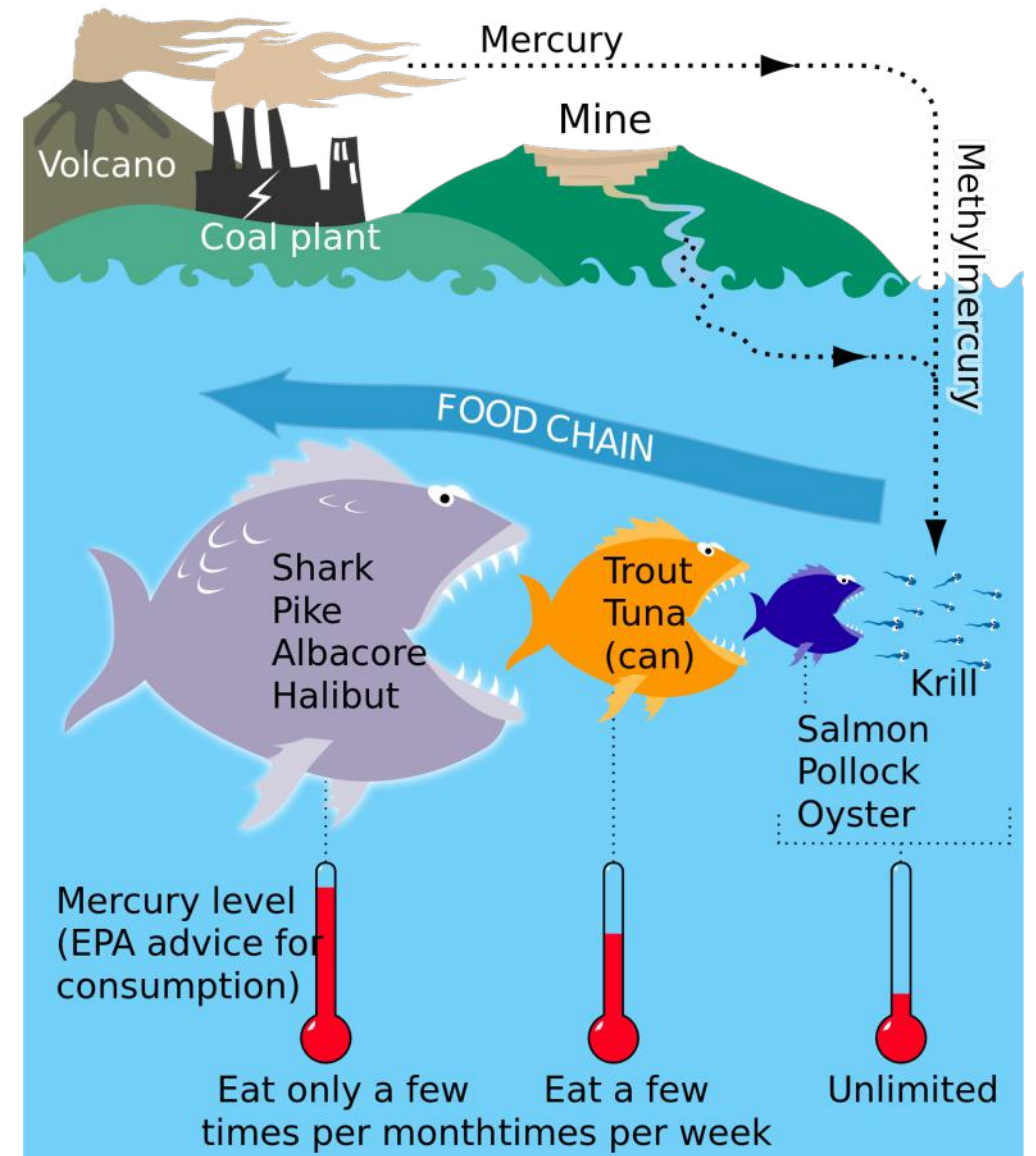
(<https://www.usgs.gov/mission-areas/water-resources/science/agricultural-contaminants>)



<https://www.usgs.gov/media/images/farmers-application-fertilizer>

Distributed Sources: Mercury

The biggest single source of mercury is the **burning of fossil fuels, especially coal**, which releases 160 tons of mercury a year into the air in the US alone. (Woods Hole Oceanographic Institution)



Other Distributed Sources/Contamination

Science

Current Issue First release papers Archive About

PERSPECTIVE TOXICOLOGY

Microplastics and human health

Knowledge gaps should be addressed to ascertain the health risks of microplastics

A. DICK VETHAAK AND JULIETTE LEGLER [Authors Info & Affiliations](#)

SCIENCE • 12 Feb 2021 • Vol 371, Issue 6530 • pp. 672-674 • DOI:10.1126/science.abe5041

20,681 265



Humans are exposed to different types of fibers and particles, including microplastics. The health effects of microplastics are largely unknown. PHOTO: DICK VETHAAK

NEWS CAREERS JOURNALS

Science

NEWS HEALTH

What we don't know about wildfire smoke is likely hurting us

Wildfires may affect our lungs and immune systems long after the blaze dies down

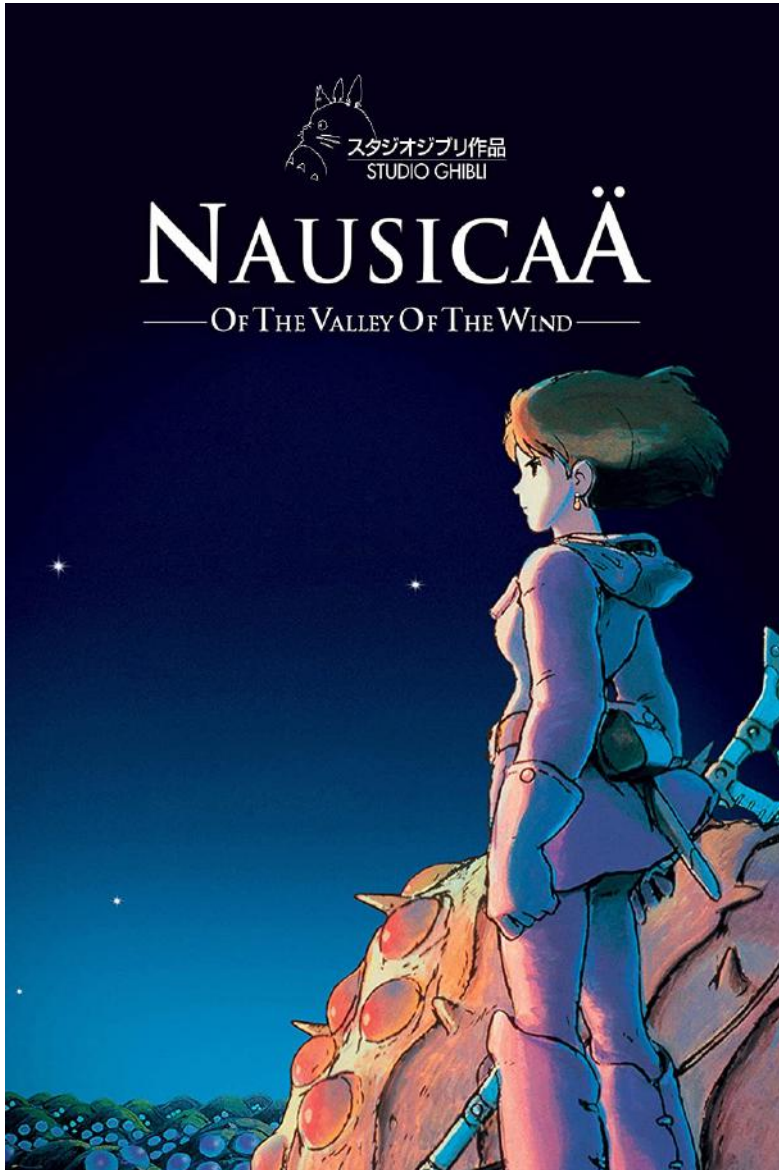
thejapantimes

JAPAN / SCIENCE & HEALTH

PFAS in western Tokyoites' blood more than double national average

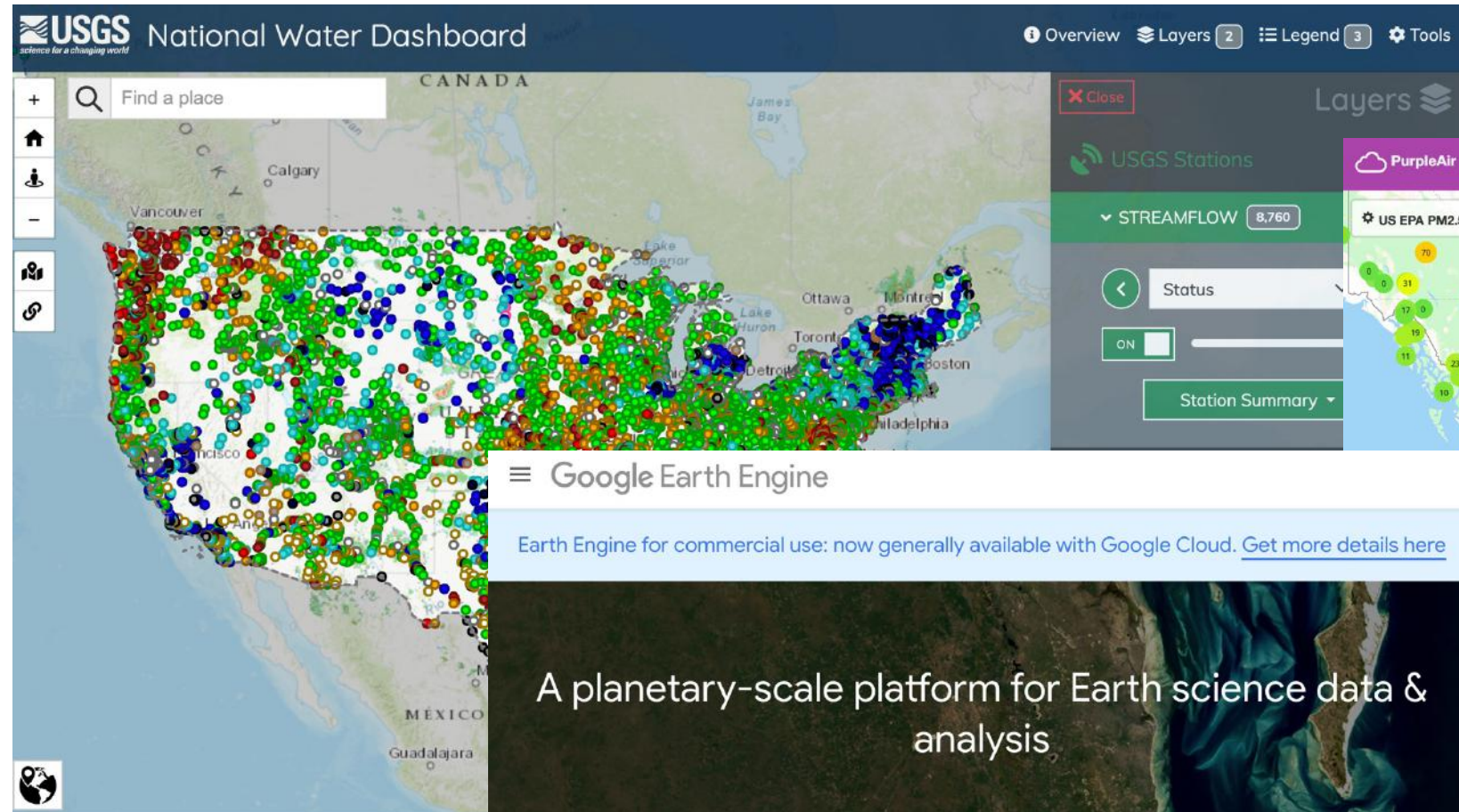


Environmental Science is Critical



- Our environment is more polluted/contaminated than people think
- Often substances that people don't worry too much end up spreading out widely and impacting our life
- Everyone needs to be more aware of pollution issues, and more vigilant to protect our health

Open Data, Open Science



USGS National Water Dashboard

Overview Layers 2 Legend 3 Tools

Find a place

USGS Stations

STREAMFLOW 8,760

Status

ON

Station Summary

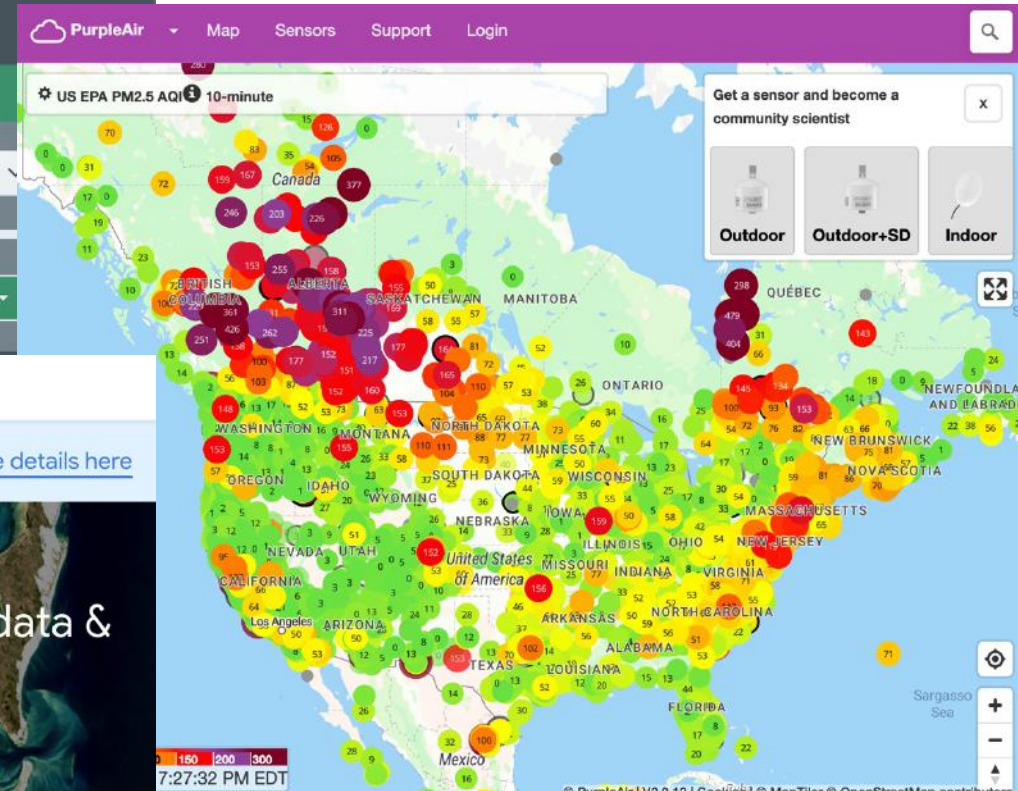
Google Earth Engine

Earth Engine for commercial use: now generally available with Google Cloud. [Get more details here](#)

A planetary-scale platform for Earth science data & analysis

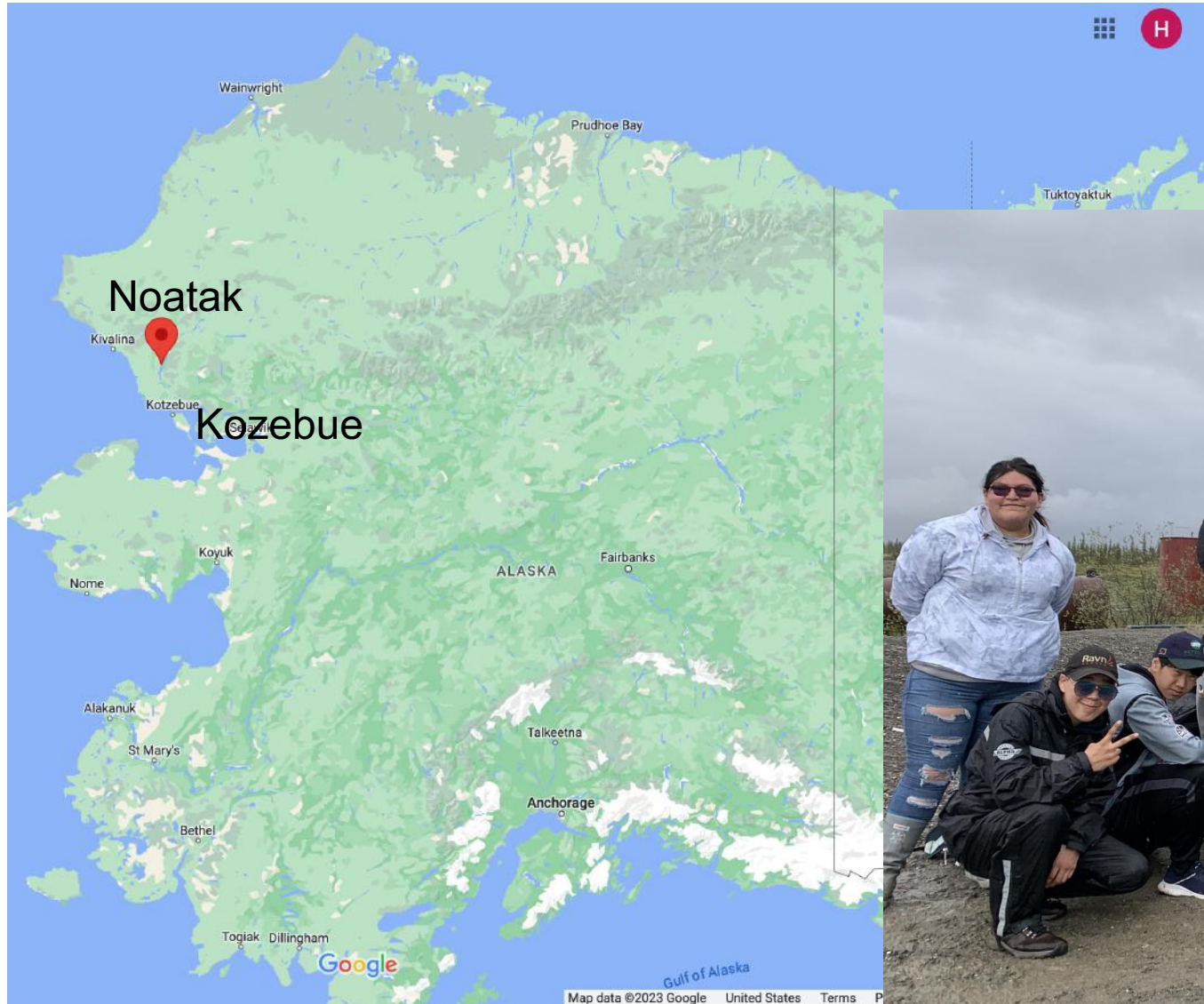
Powered by Google's cloud infrastructure

▶ Watch Video



- Public database
- API

Toward Citizen Science: K-12 Education



T³ Teaching Through Technologies



Hands on Activities with Sensors



Names: Mathys, Jeremy, T Date/Location:

#	pH	TDS	Temp.	Answer	Note	Strip
1	4.0	2000	70.5	SEE	strip color	
2					Obs. color = ac strange	
3	4.0	862	71	SEE	clear	
4	6.0	352	70.5	SEE	clear	
5	9.0	0	100	SO	water	
6	6.0	200	70.5	SEE	AT L	
7	10	780	70.5	VIN	smells vinegar	
8	4	240	69.5	WASH	10/20/20	

Sensor Technologies for Teaching



GP Data Integration for Air Quality

GP4AQ

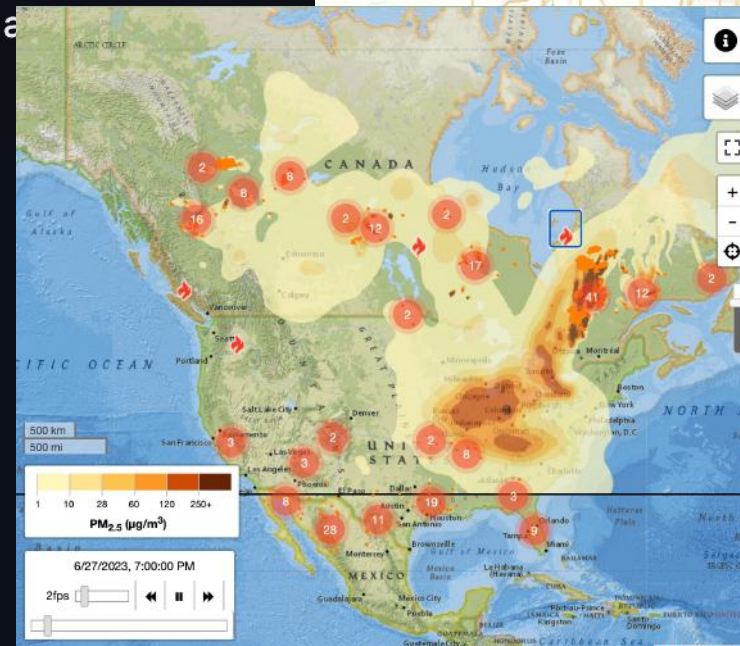
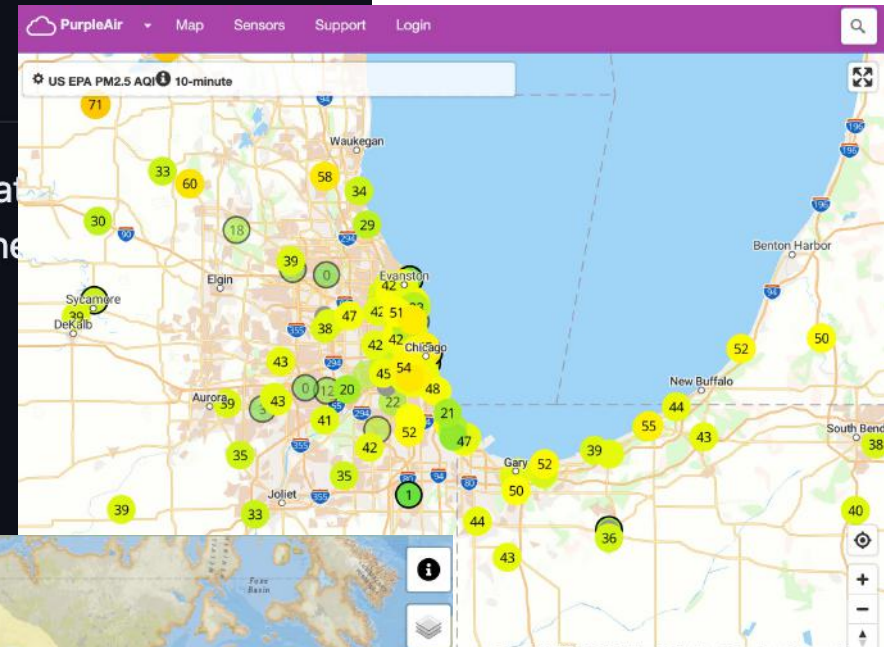
github.com/hmwainw/GP4AQ

These Jupyter notebooks demonstrate the data integration of EPA and Purple air sensor data, simulation results, using the Gaussian Process regression, for interpolating and mapping the data over space.

Note:

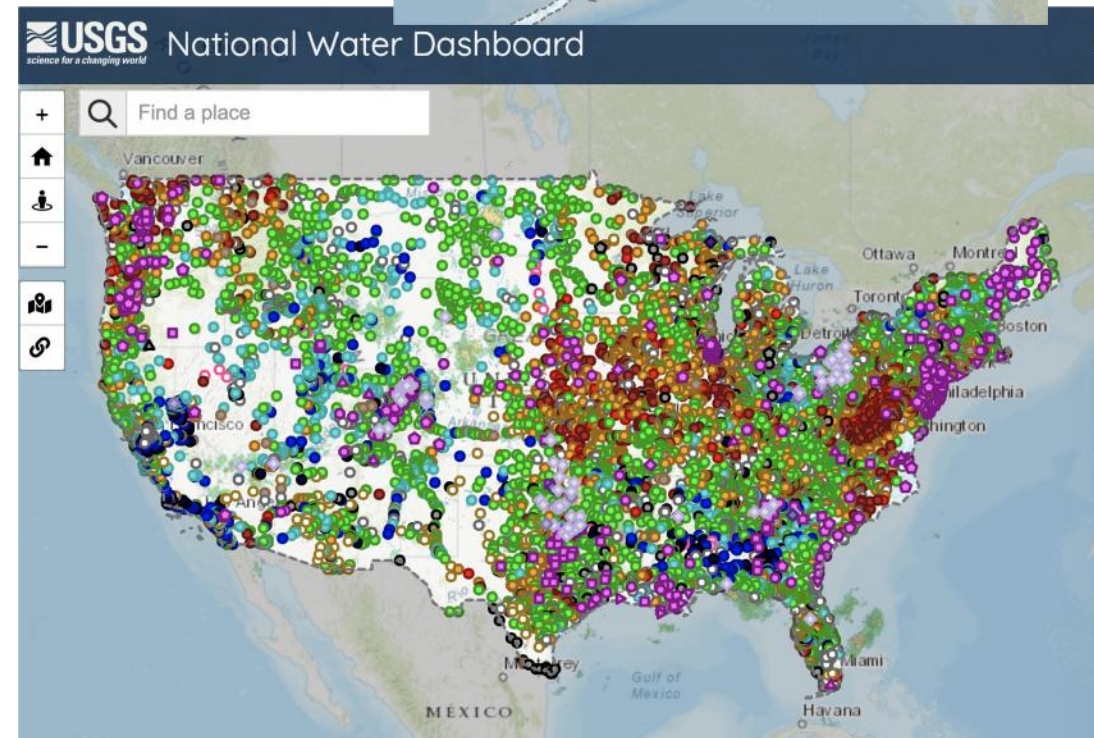
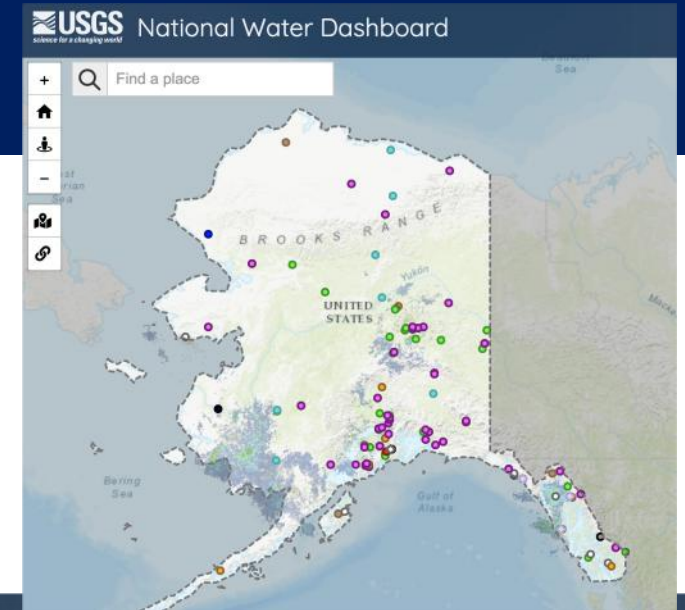
- You have to have purpleair_chicago.txt in the Google Drive root directory.
- The same functions appear in multiple notebooks. These notebooks are for a specific purpose.

1. Define the grid and domain for spatial estimation
2. Download EPA air quality data
3. Download Purple air quality data
4. Download plume simulation maps
5. Interpolate EPA air quality data
6. Interpolate Purple air quality data
7. Integrate EPA and Purple air quality data
8. Integrate EPA, Purple air, and plume simulation data



Next Step

- **Environmental Monitoring Network for rural America**
- **USGS Database does not cover rural regions**
 - Too remote
 - Data quality concern
- **Can high schools be the base for environmental network in rural regions?**
- **Improve STEM education**
 - **More college/PhD from rural regions!**



Challenges.... Tech/AI in Environment/Climate

- **Pollution monitoring is not exciting when nothing happens**
 - **Attributes more relevant to daily life?**
 - River temperature for fishing?
 - Soil moisture sensors for gardening?



- **Students good at math/science are not interested in the environment and climate**
 - **Environmental data in math/statistics education?**
 - Open data and problem sets?

Summary

- **Long-term monitoring of soil and groundwater contamination**
 - Sustainable remediation: long-term institutional control
 - Ensure the stability/safety of contaminated sites and detect anomalies
- **ALTEMIS: Multiscale multi-type data integration**
 - Integration of proxy information (e.g., spatial data, in situ sensors)
 - PyLenM: Framework from various data to ML and decision making
 - Model simulations to inform monitoring and management
- **Simulation Intelligence: Simulations x ML/AI**
 - U-FNO for emulating simulation results to understand climate change impact on residual contamination
 - Bayesian hierarchical models with GP for physics-informed spatial interpolation (physics-informed monitoring)
- **General contaminations: Democratizing environmental science**
 - Citizen science for water/air quality, tackling environmental justice issues
 - AI/ML for environmental science

Thank You!

Contact

Haruko Wainwright

HMWainw@MIT.EDU

Acknowledgment

DOE Office of Environmental Management

DOE Office of Science