

# Cyberinfrastructure Tools for Precision Agriculture in the 21st Century

Michela Taufer

The University of Tennessee Knoxville



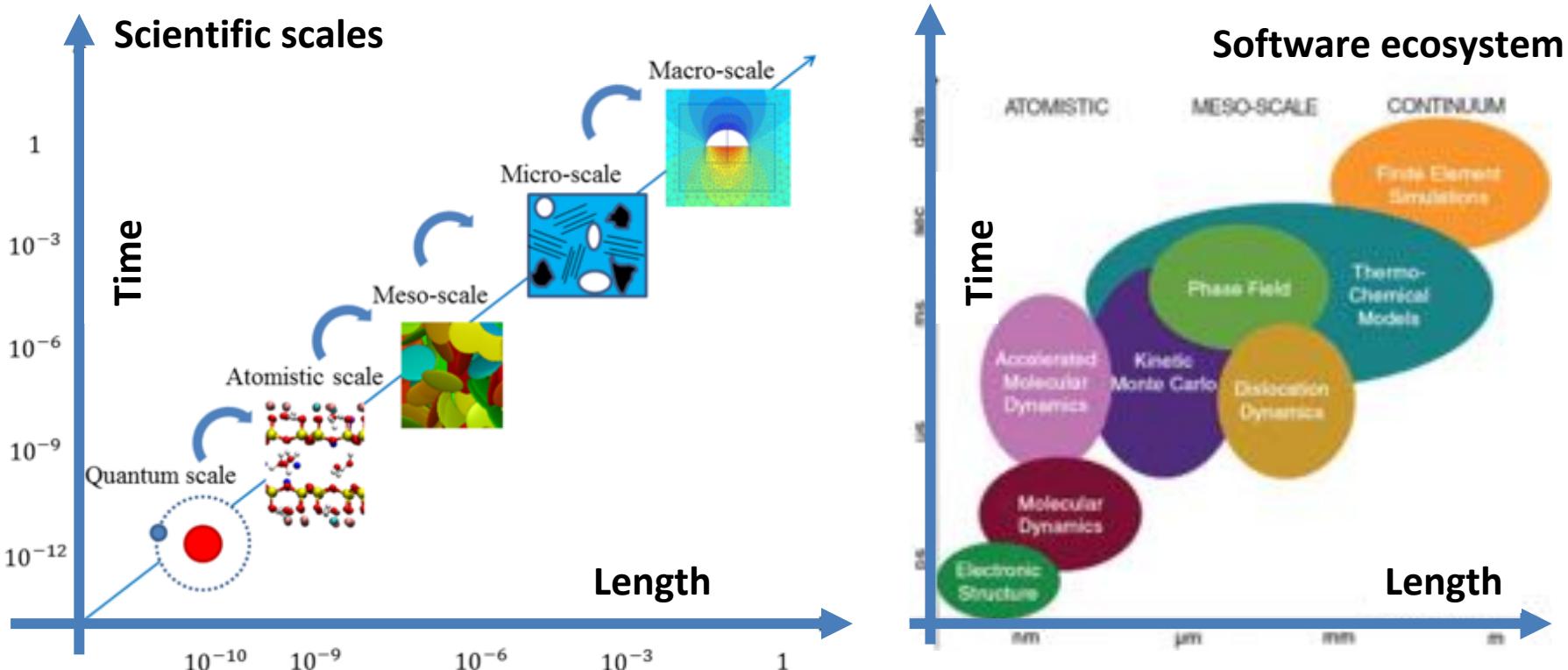
# Contributors and collaborators

- U. Delaware: Ricardo Llamas, Mario Guevara, and Rodrigo Vargas
- UTK: Danny Rorabaugh, Kae Suarez, Leobardo Valera, Ria Patel, and David Icove
- ORNL: Jimmy Landmesser
- UIUC: Craig Willis and Victoria Stodden

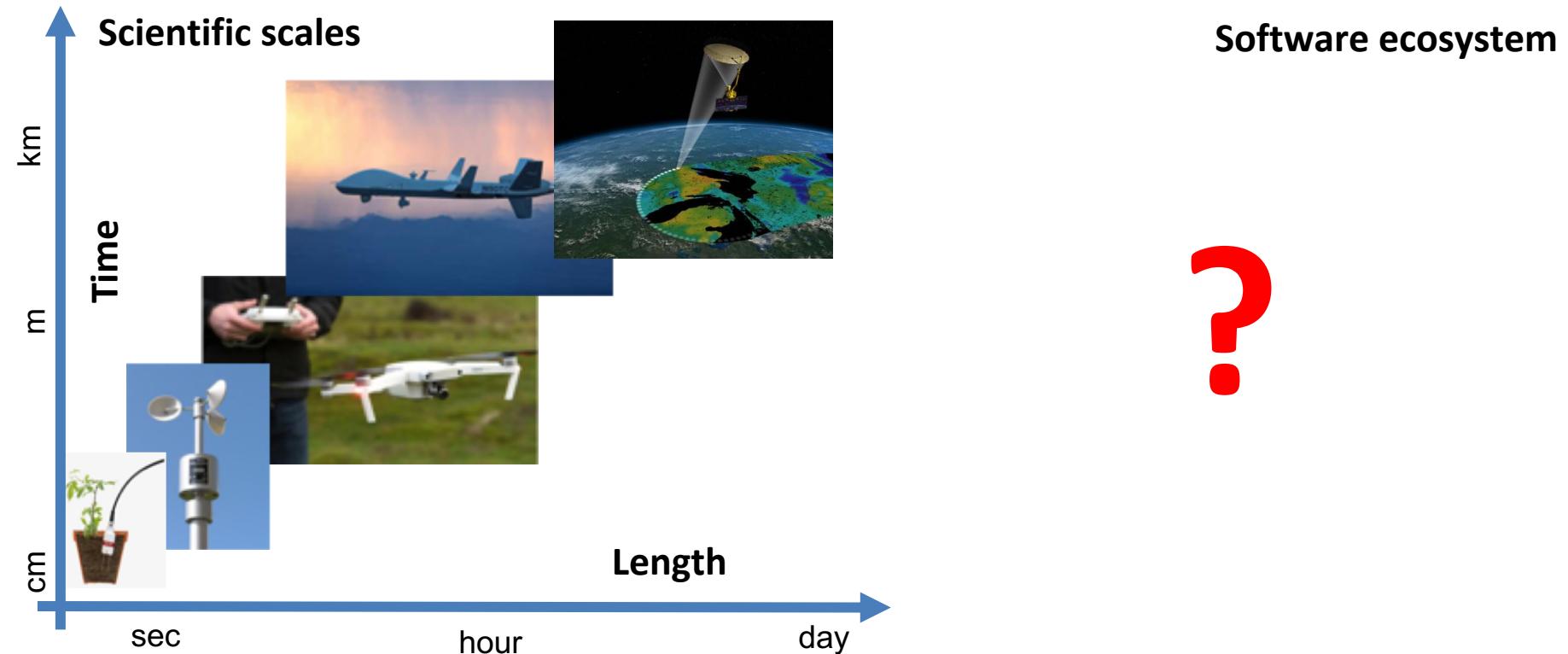
# Sponsors and supporters

- NSF OAC 1854312 CIF21 DIBBs: PD: Cyberinfrastructure Tools for Precision Agriculture in the 21st Century (PIs: Taufer and Vargas)
- NSF OAC 1941443 EAGER: Reproducibility and Cyberinfrastructure for Computational and Data-Enabled Science (PIs: Stodden and Taufer)
- IBM Shared University Research (SUR) Award
- NSF XSEDE JetStream: Allocations EAR180011 and TRA180041
  - Many thanks to Jeremy Fischer, IU

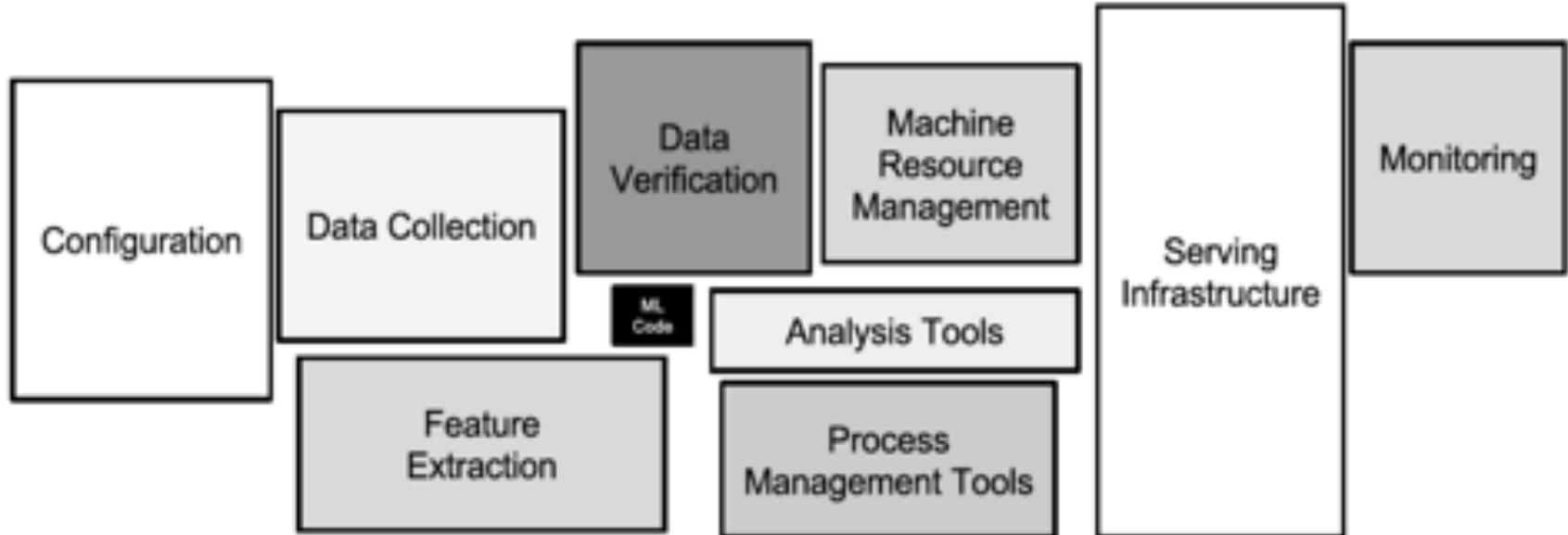
# Multiscale computational modeling



# Multiscale data modeling (MSDM)

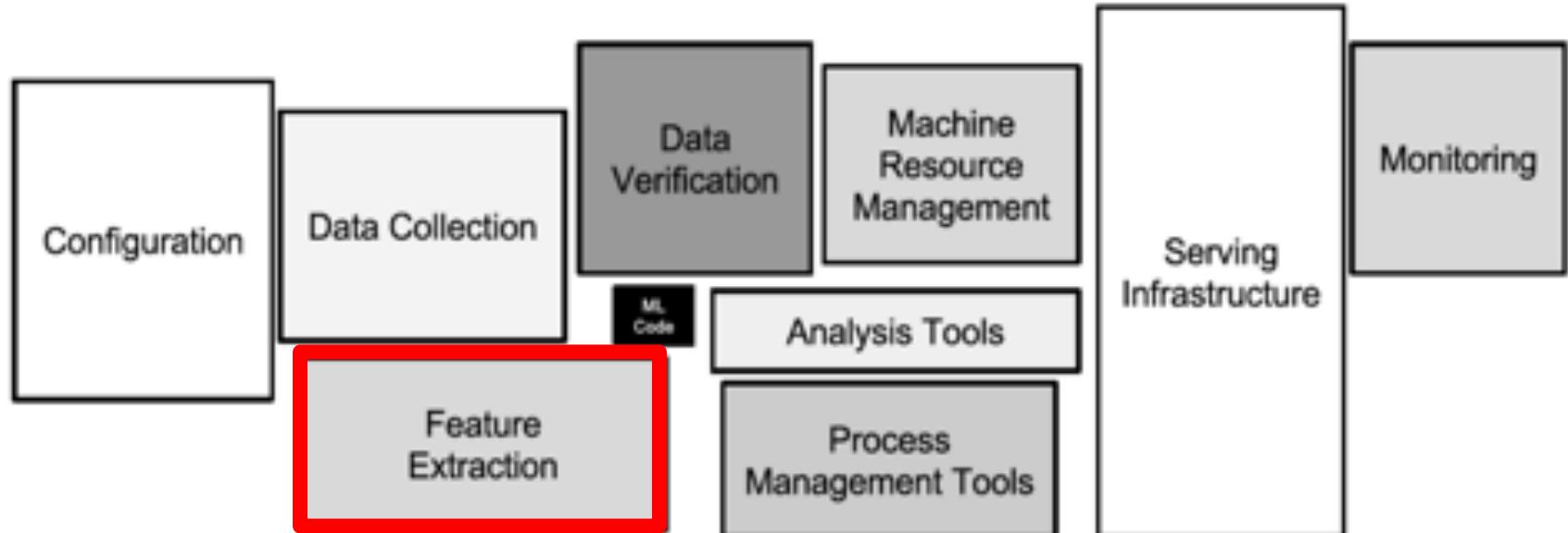


# Hidden (forgotten?) software ecosystem



**“Only a small fraction of real-world ML systems is composed of the ML code”** D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips () Hidden Technical Debt in Machine Learning Systems

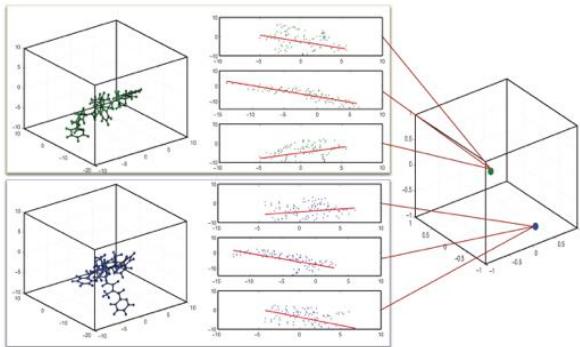
# Hidden (forgotten?) software ecosystem



**"Only a small fraction of real-world ML systems is composed of the ML code"** D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips () Hidden Technical Debt in Machine Learning Systems

# Feature extraction

## Protein-ligand docking

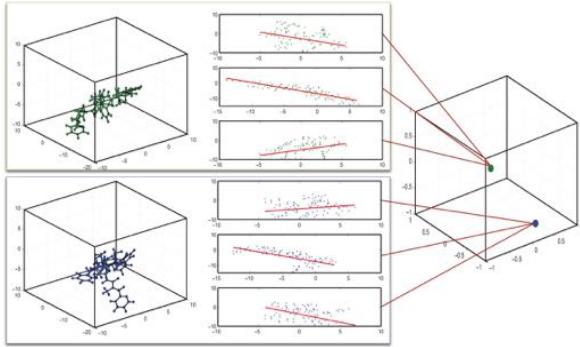


Linear regression: map ligands into **3-D point** representation

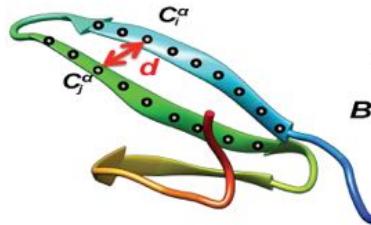
M. Taufer, T. Estrada, and T. Johnston. Algorithms for In Situ Data Analytics of Next Generation Molecular Dynamics Workflows. Numerical algorithms for high-performance computational science. *Issue of Philosophical Transactions A.*, 2019.

# Feature extraction

## Protein-ligand docking



Linear regression: map ligands into **3-D point** representation



## Protein folding

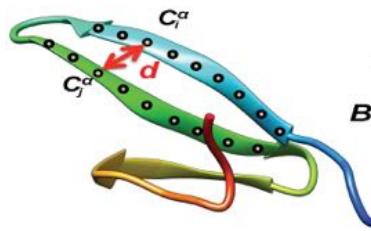
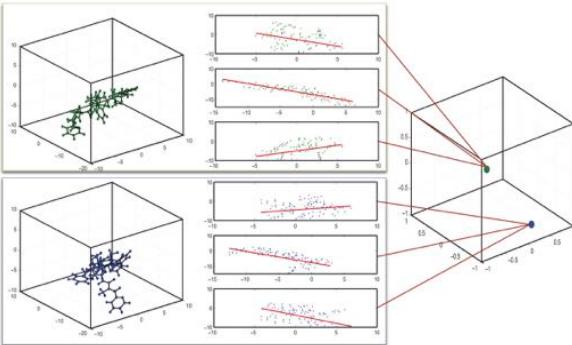
$$B = \begin{bmatrix} 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & d & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ \times & d & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \end{bmatrix}$$

Numerical analyses: map secondary structures into **eigenvalues**

M. Taufer, T. Estrada, and T. Johnston. Algorithms for In Situ Data Analytics of Next Generation Molecular Dynamics Workflows. Numerical algorithms for high-performance computational science. *Issue of Philosophical Transactions A.*, 2019.

# Feature extraction

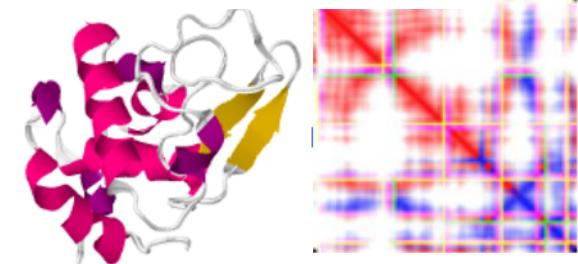
## Protein-ligand docking



## Protein folding

$$B = \begin{bmatrix} 0 & 0 & 0 & C_i^\alpha & x & x & x \\ 0 & 0 & 0 & d & x & x & x \\ 0 & 0 & 0 & x & x & x & \\ x & d & x & 0 & 0 & 0 & \\ x & x & x & 0 & 0 & 0 & \\ x & x & x & 0 & 0 & 0 & \end{bmatrix}$$

## Protein engineering

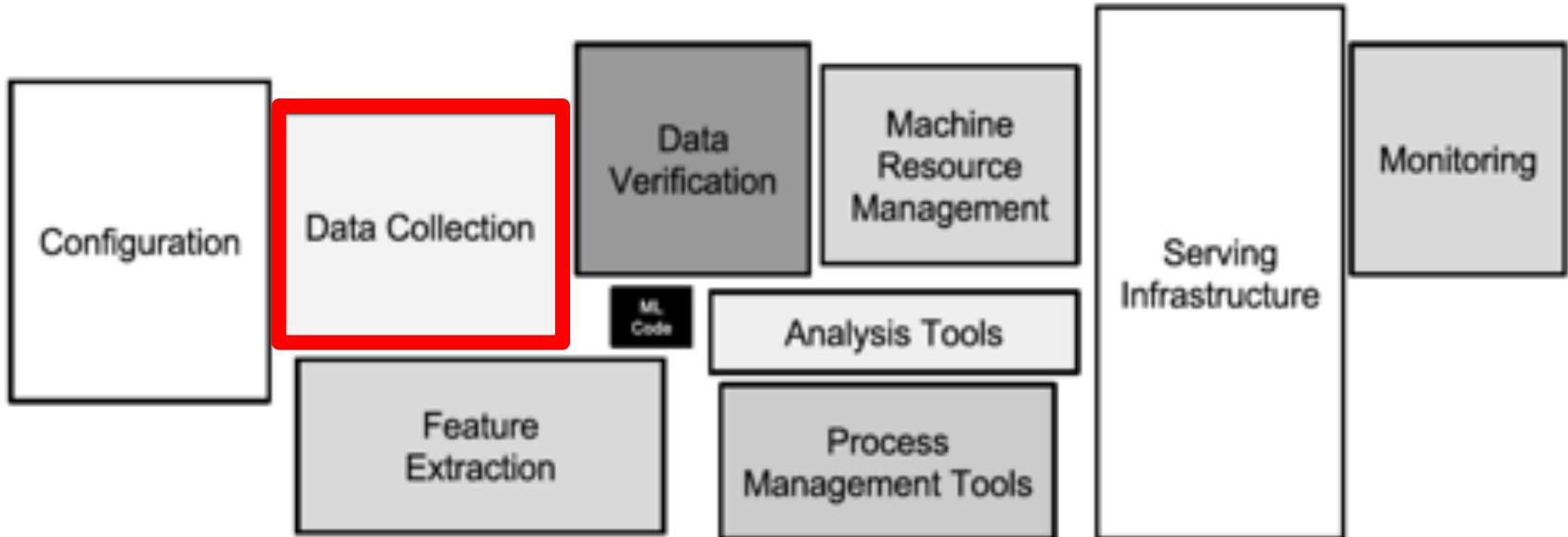


*Deep learning:* map both secondary and tertiary structures into **tensors**

*Numerical analyses:* map secondary structures into **eigenvalues**

M. Taufer, T. Estrada, and T. Johnston. Algorithms for In Situ Data Analytics of Next Generation Molecular Dynamics Workflows. Numerical algorithms for high-performance computational science. *Issue of Philosophical Transactions A.*, 2019.

# Hidden (forgotten?) software ecosystem



**"Only a small fraction of real-world ML systems is composed of the ML code"** D. Sculley, Gary Holt, Daniel Golovin, Eugene Davydov, Todd Phillips () Hidden Technical Debt in Machine Learning Systems

# Data collection at the edge

Point Field  
Measurements

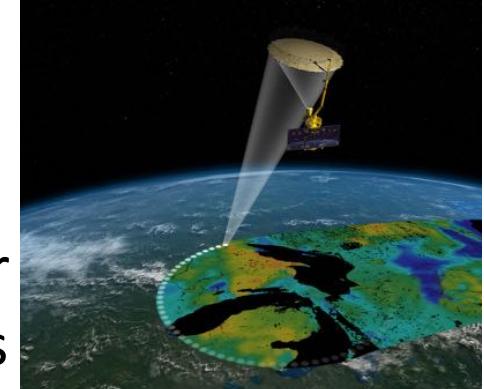


# Data collection at the edge

## Point Field Measurements



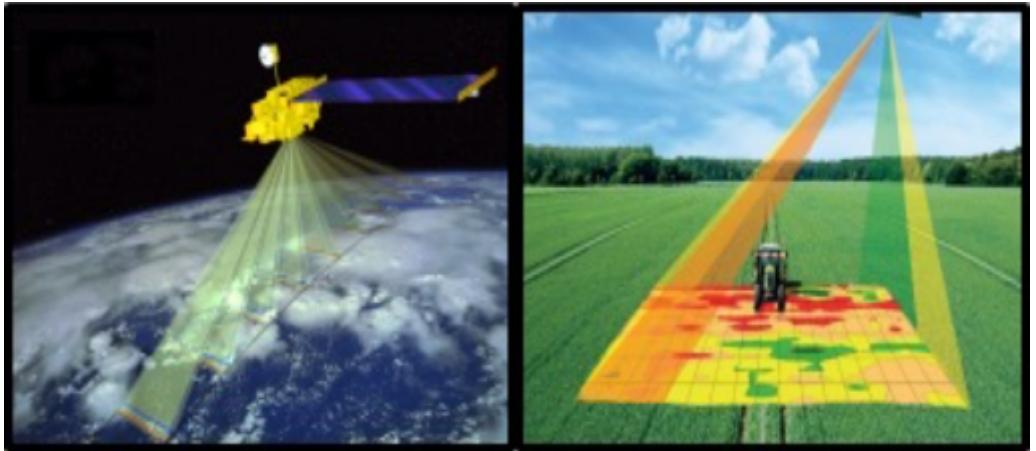
## Remote Sensor Measurements



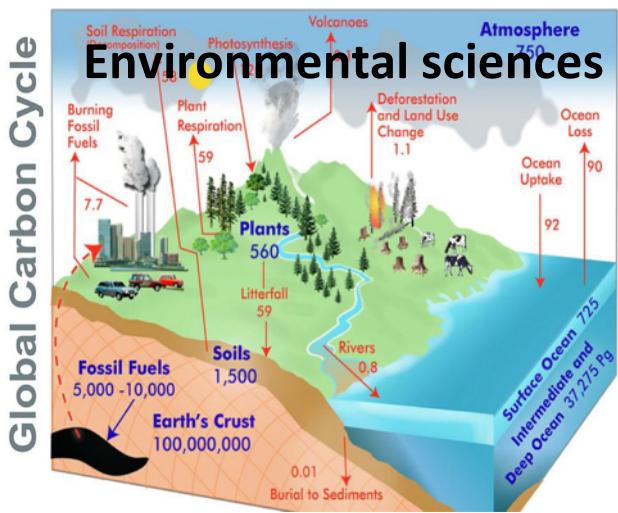
# Challenges in MMDM

- Design and implement robust and sustainable software ecosystems
- Combine analytics and computing across heterogenous platforms (i.e., HPC, Cloud, and edge computing)
- Build trust in results through reproducibility, replicability, and transparency (RRT)

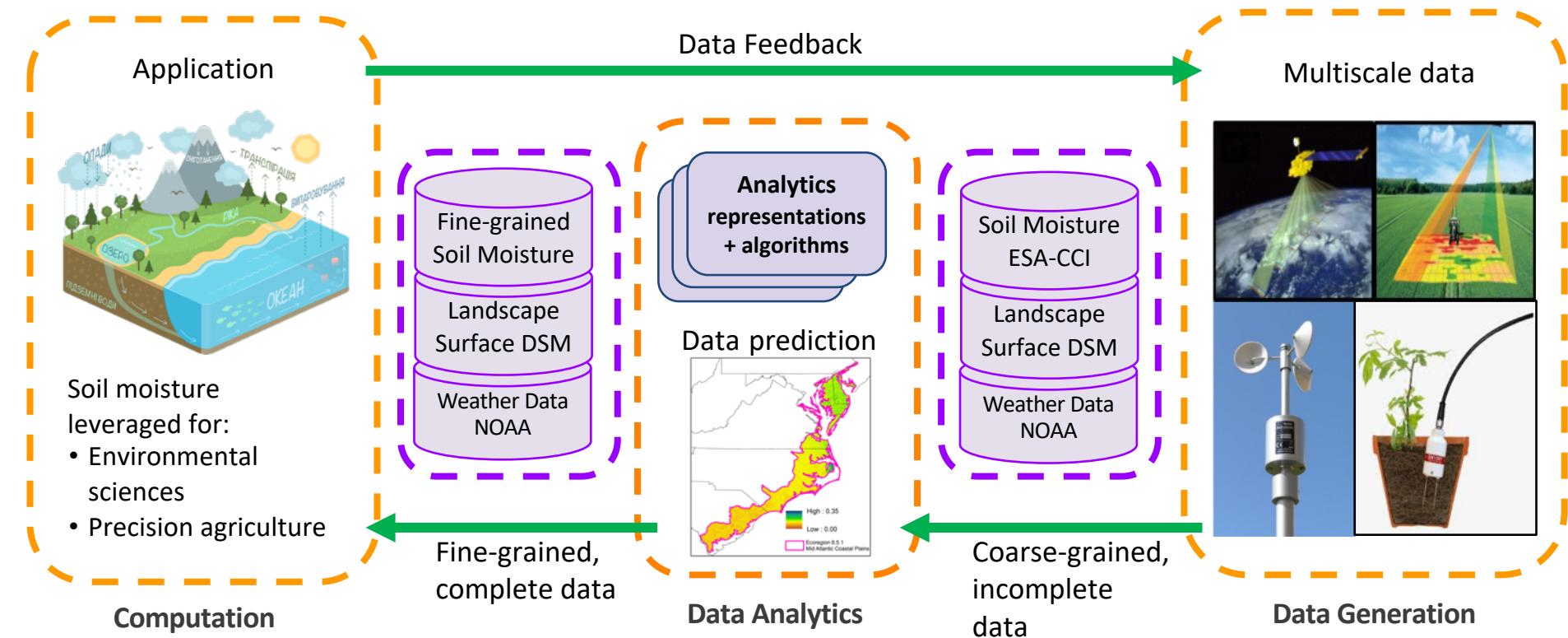
# Relevance of soil moisture data



- Satellite-borne remote sensing technology
  - Infrared to radio
  - Active and passive



# Workflows for precision agriculture



# **Design and implement a software ecosystem for precision agriculture**

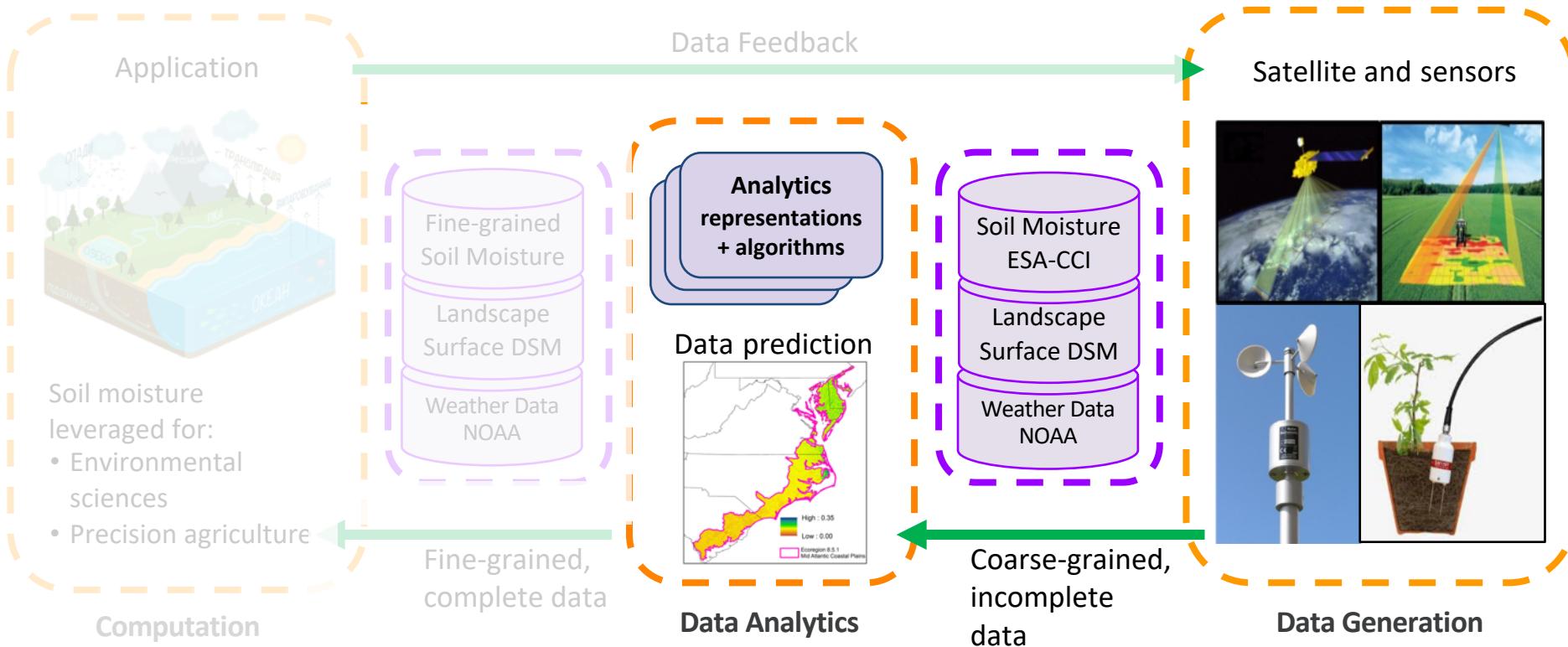
**Collaborators: Rodrigo Varga's Group (UD)**

**Platform: NSF XSEDE Jetstream**

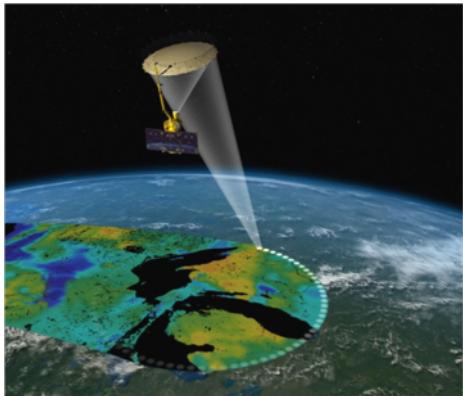
**NSF OAC 1854312 CIF21 DIBBs: PD: Cyberinfrastructure**

**Tools for Precision Agriculture in the 21st Century**

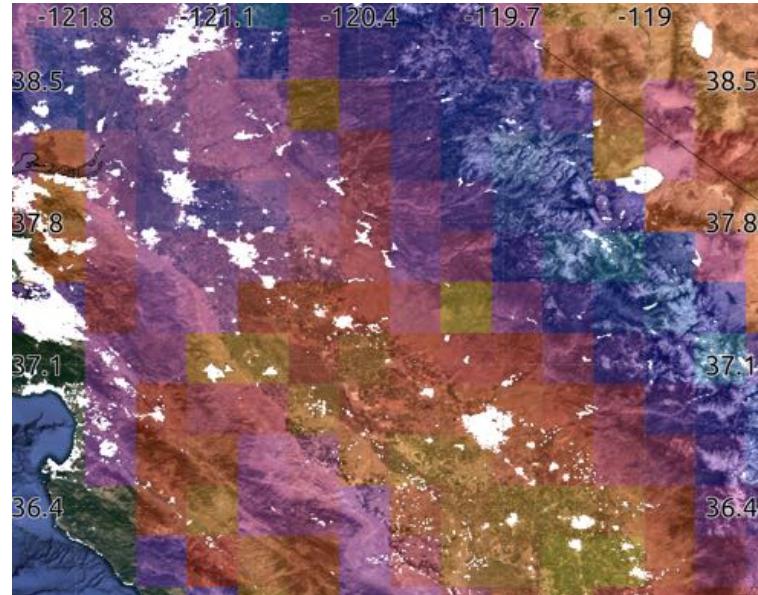
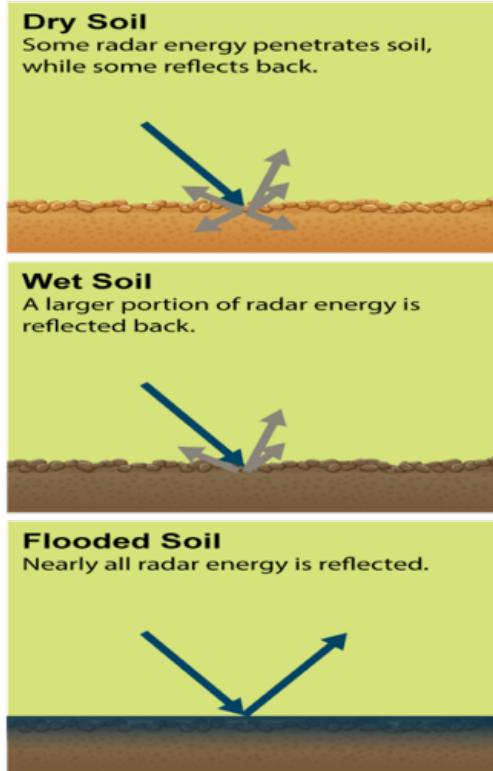
# Data analytics for soil moisture



# Challenge 1: incomplete soil moisture data (I)

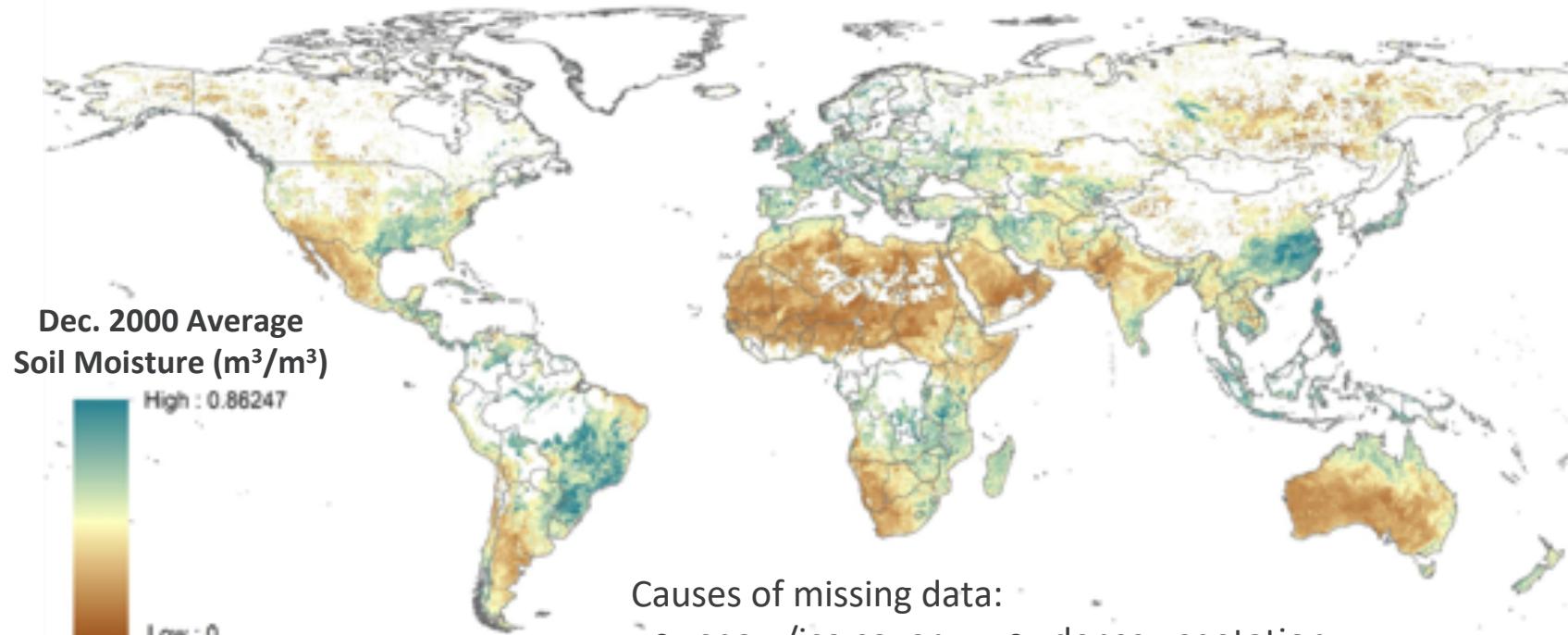


Satellites collect raster data across the surface of the Earth



Visualization example of the ESA-Climate Change Initiative Soil Moisture database with a coarse pixel size of 27x27km

# Challenge 1: incomplete soil moisture data (II)



Causes of missing data:

- snow/ice cover
- frozen surface
- dense vegetation
- extremely dry surface

## Challenge 2: coarse-grained soil moisture data (I)

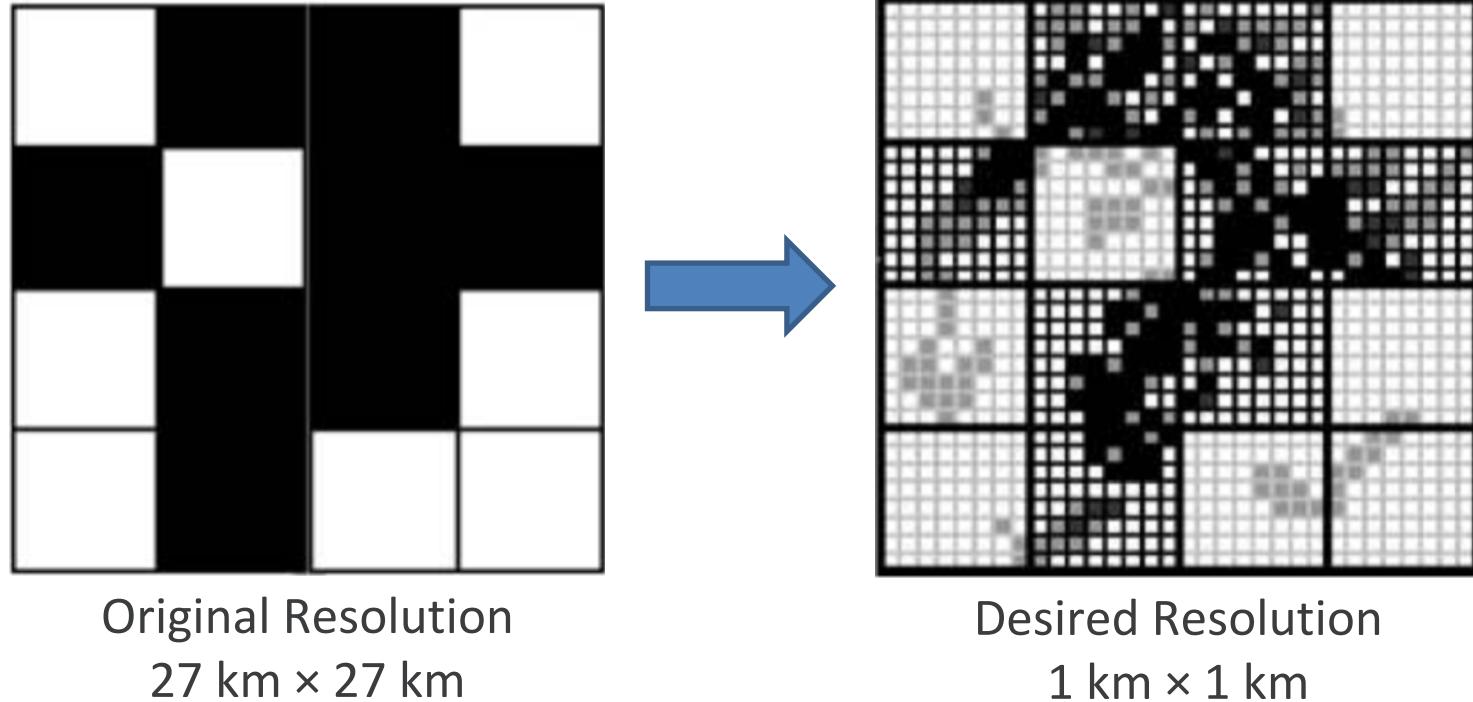
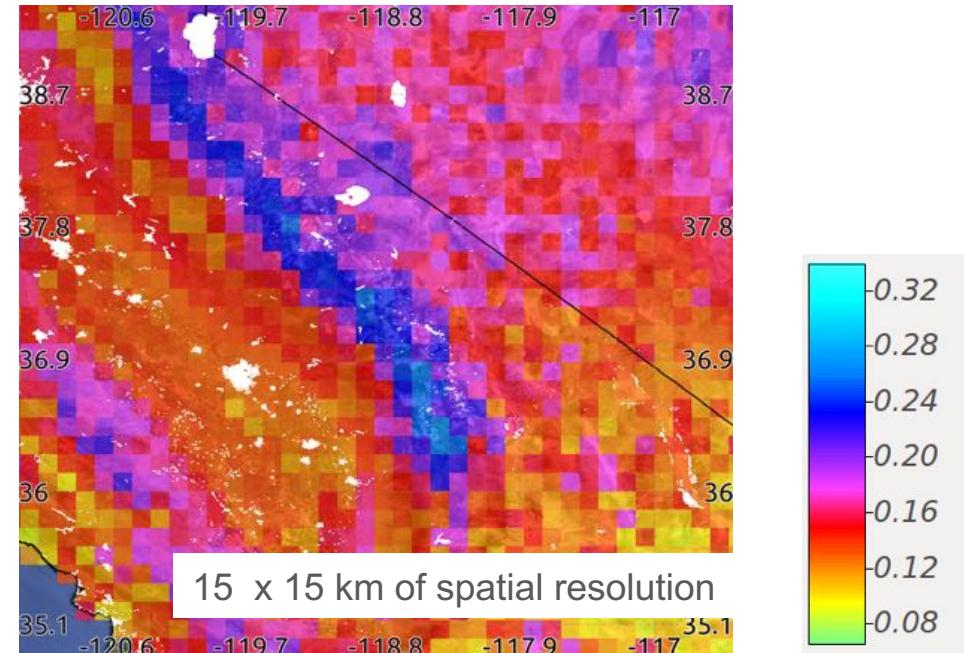
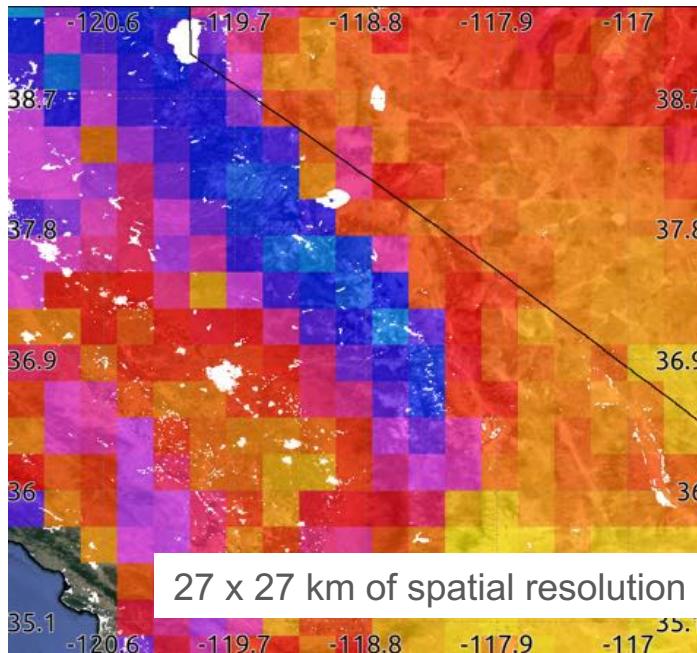


Image source: McPherson et al., *Using coarse-grained occurrence data to predict species distributions at finer spatial resolutions—possibilities and limitations*, Ecological Modeling 192:499–522, 2006.

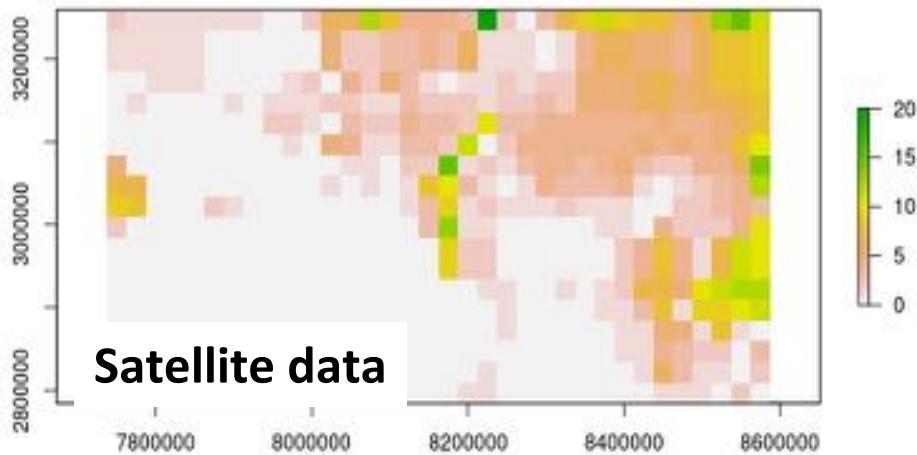
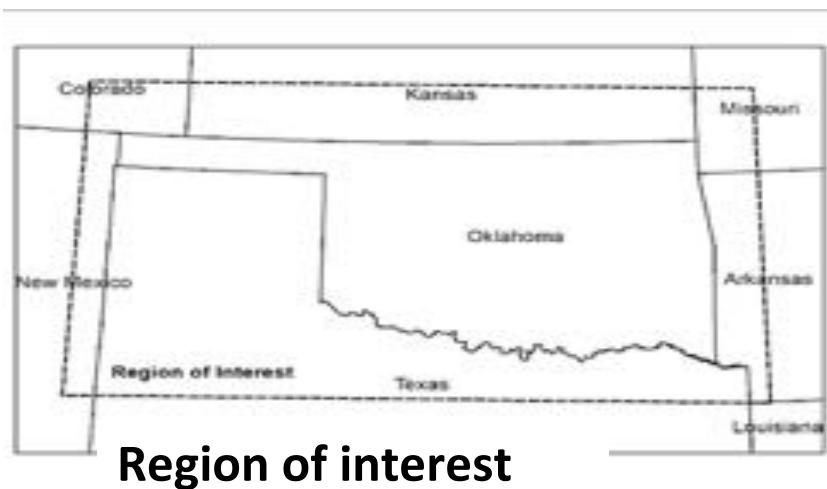
# Challenge 2: coarse-grained soil moisture data (II)

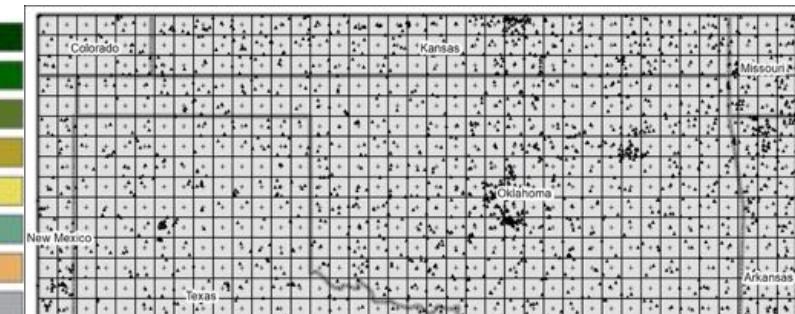
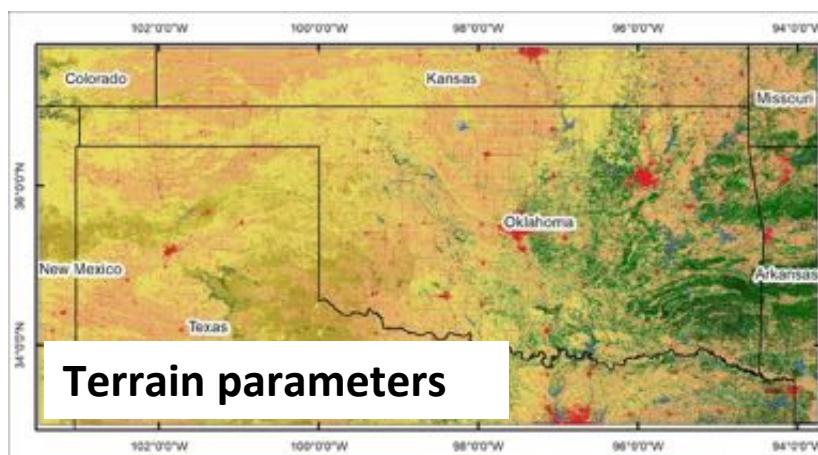
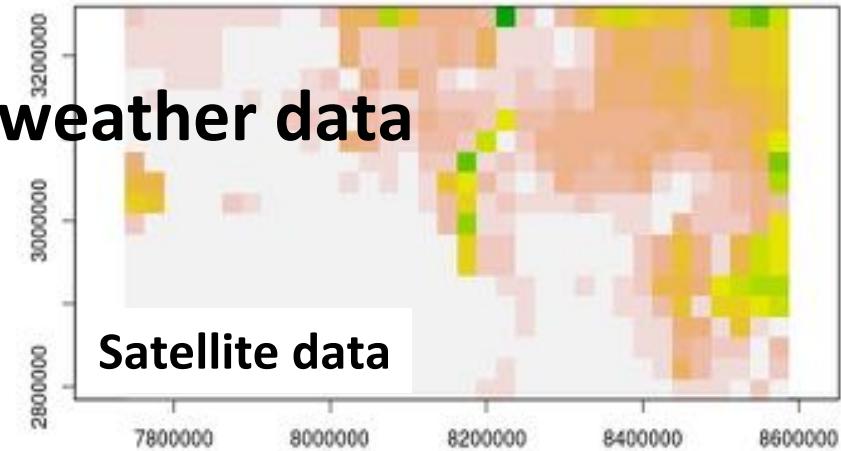
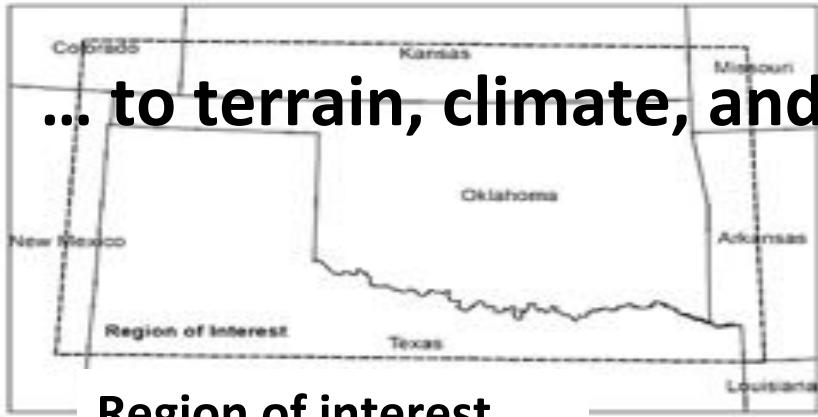
Original product ESA CCI ( $\text{m}^3 \text{ m}^{-3}$ , mean 2013)



M. Guevara , M. Taufer, and R. Vargas. Gap-Free Annual Soil Moisture Global across 15km Grids: 1991-2016. Earth System Science Data, 2019.

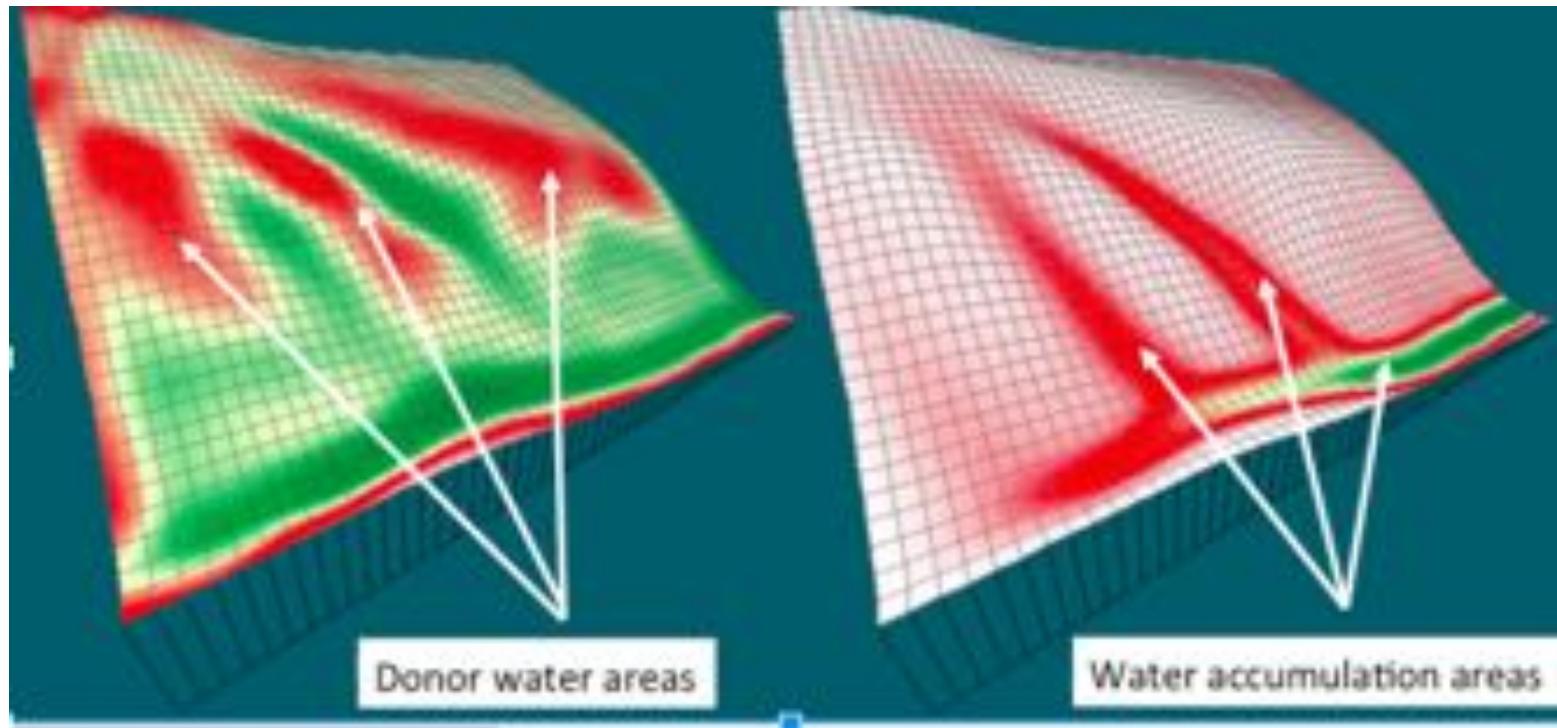
# Integration of multiscale data: from satellites ...



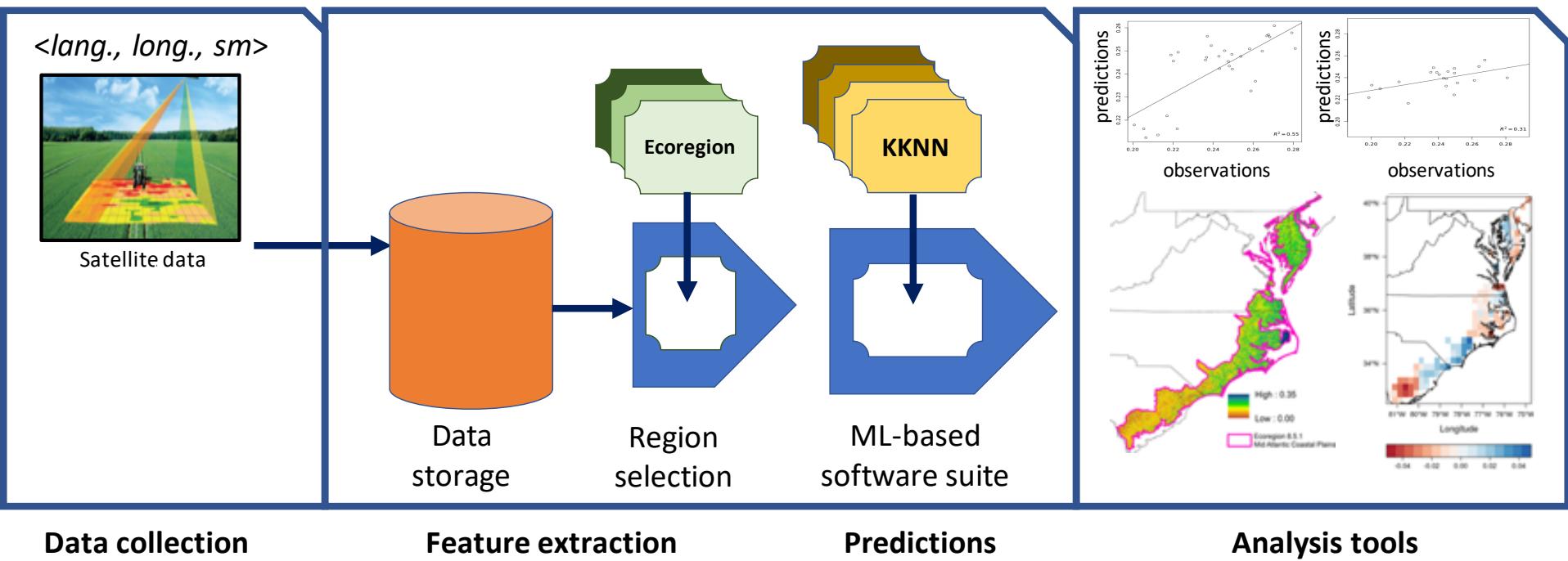


Global Historical Climatology Network (GHCN)  
and other local data (field measurements)

# Example of terrain parameters: water wetness index



# SOMOSPIE: SOil MOisture SPatial Inference Engine



Data collection

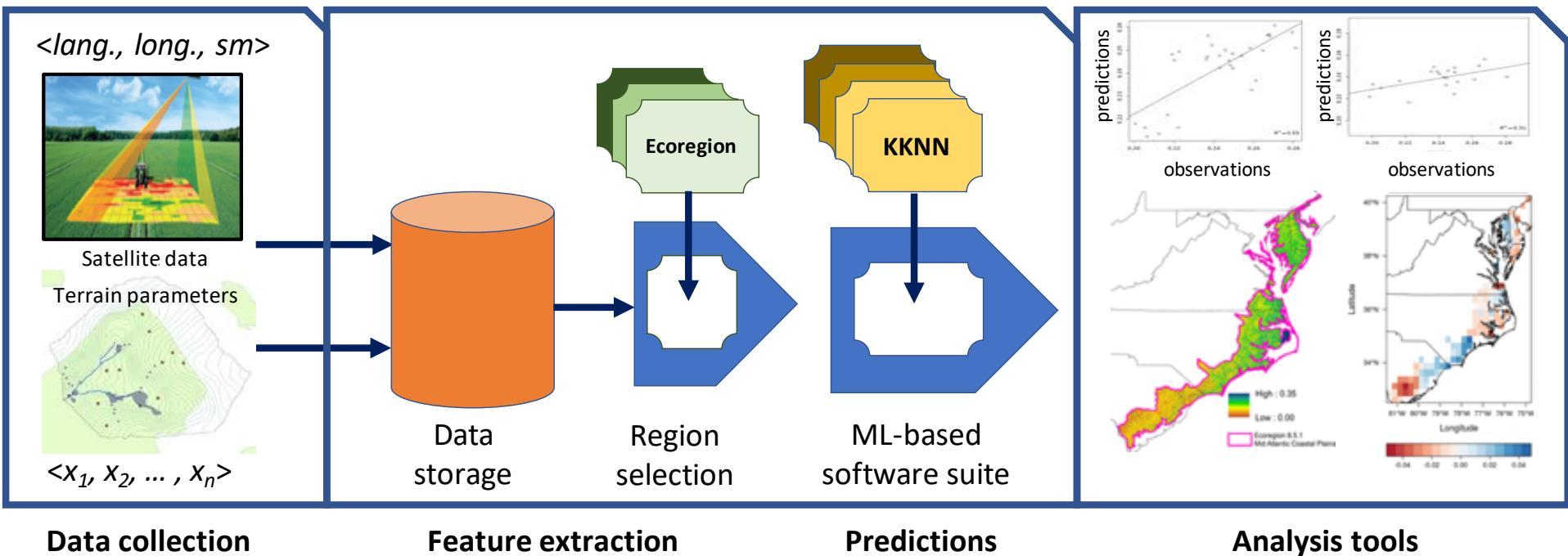
Feature extraction

Predictions

Analysis tools

D. Rorabaugh, M. Guevara, R. Llamas, J. Kitson, R. Vargas, and **M. Taufer**. SOMOSPIE: A Modular SOil MOisture SPatial Inference Engine based on Data Driven Decisions. eScience 2019.

# SOMOSPIE: SOil MOisture SPatial Inference Engine



Data collection

Feature extraction

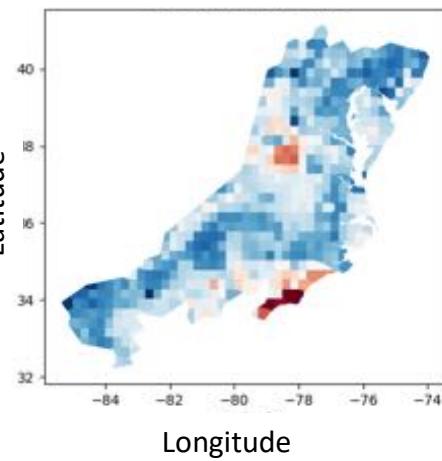
Predictions

Analysis tools

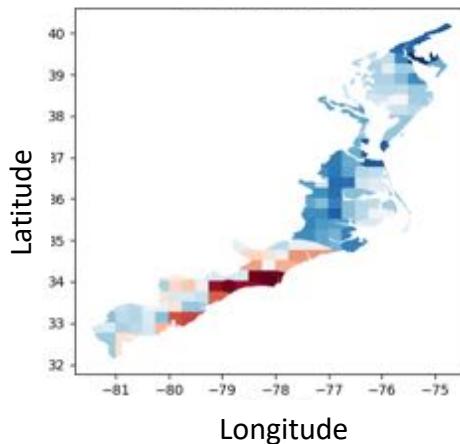
D. Rorabaugh, M. Guevara, R. Llamas, J. Kitson, R. Vargas, and **M. Taufer**. SOMOSPIE: A Modular SOil MOisture SPatial Inference Engine based on Data Driven Decisions. eScience 2019.

# Region selection: format of regions of interest

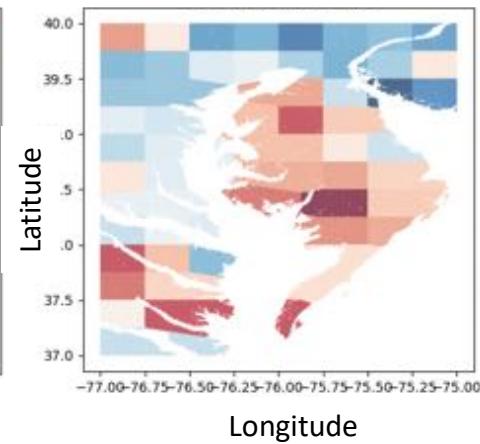
("NEON", "Mid Atlantic")



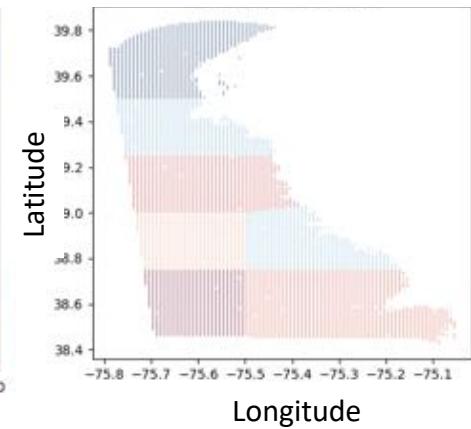
("CEC", "8.5.1")



("BOX", "-77\_-75\_37\_40")



("STATE", "Delaware")



# Algorithmic solutions: ML-based software suite

## Random Forest

- Compute weighted mean of 500 prediction trees

## KKNN:

- Use **local data**
- Compute k and distance kernel using cross validation automatically
- Compute weighted means with the kernel (**many values**)

## Surrogate based model (SBM):

- Use **all sampled data**
- Use regression to generate one single polynomial model (**single polynomial model**)

# Algorithmic solutions: ML-based software suite

## Random Forest

- Compute weighted mean of 500 prediction trees

## KKNN:

- Use **local data**
- Compute k and distance kernel using cross validation automatically
- Compute weighted means with the kernel (**many values**)

## Surrogate based model (SBM):

- Use **all sampled data**
- Use regression to generate one single polynomial model (**single polynomial model**)

## HYPPO (Hybrid Piecewise Polynomial Modeling):

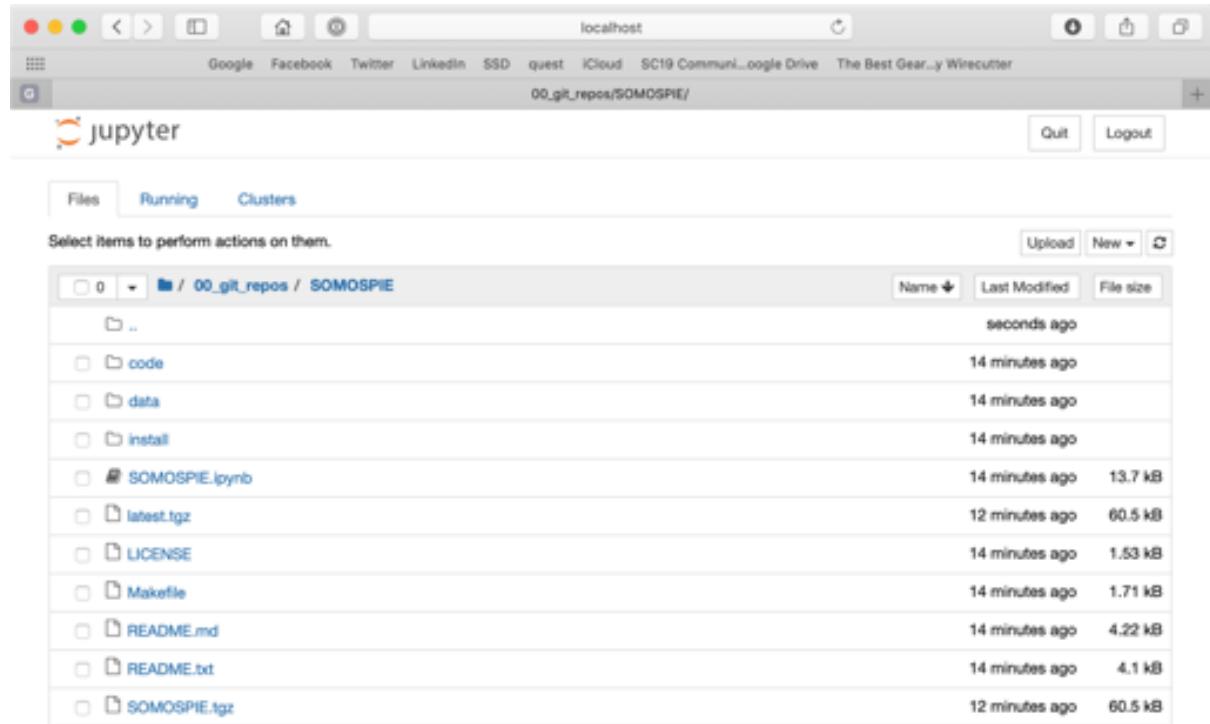
- Use **local data**
- Determine local polynomial degree using cross validation



- Use regression to generate local polynomial model (**many polynomial models**)



# Computational solutions: Jupyter + XSEDE Jetstream



# Computational solutions: Jupyter + XSEDE Jetstream

The screenshot shows a dual-Jupyter Notebook setup running on XSEDE Jetstream. The left window is titled "jupyter" and displays a file tree for a repository named "SOMOSPIE". The right window is titled "jupyter SOMOSPIE (autosaved)" and contains a Jupyter Notebook interface with the following content:

## SOMOSPIE

### SOil MOisture SPatial Inference Engine

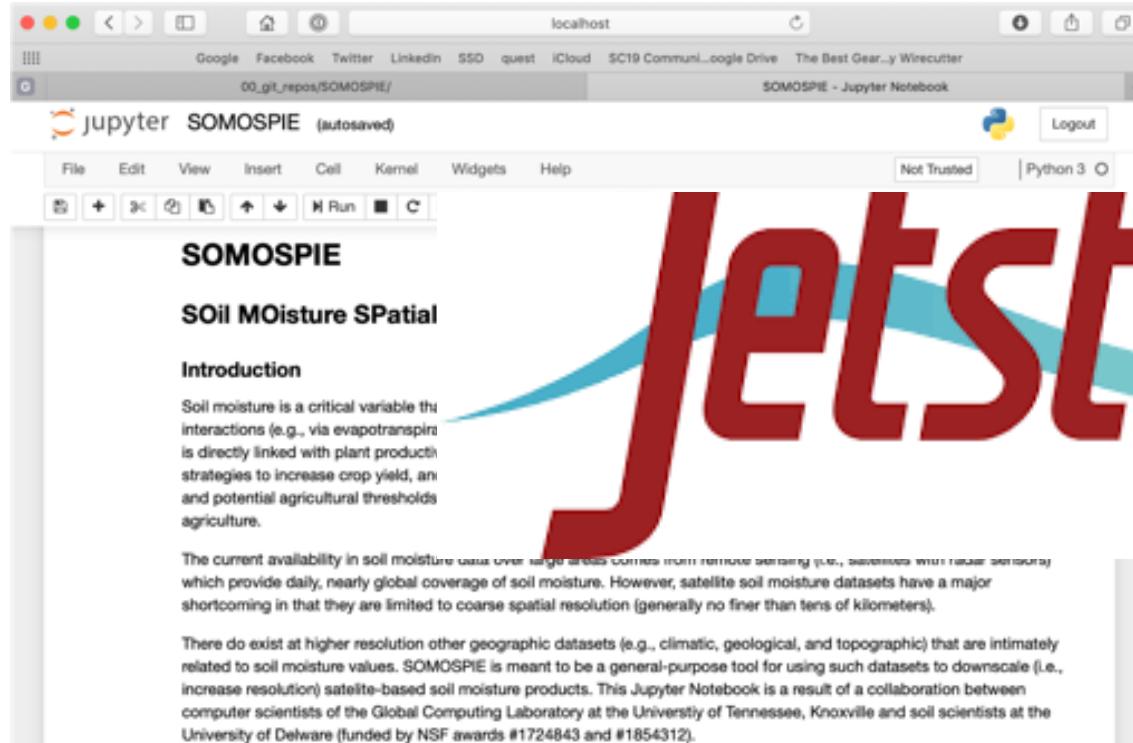
#### Introduction

Soil moisture is a critical variable that links climate dynamics with water and food security. It regulates land-atmosphere interactions (e.g., via evapotranspiration—the loss of water from evaporation and plant transpiration to the atmosphere), and it is directly linked with plant productivity and survival. Information on soil moisture is important to design appropriate irrigation strategies to increase crop yield, and long-term soil moisture coupled with climate information provides insights into trends and potential agricultural thresholds and risks. Thus, information on soil moisture is a key factor to inform and enable precision agriculture.

The current availability in soil moisture data over large areas comes from remote sensing (i.e., satellites with radar sensors) which provide daily, nearly global coverage of soil moisture. However, satellite soil moisture datasets have a major shortcoming in that they are limited to coarse spatial resolution (generally no finer than tens of kilometers).

There do exist at higher resolution other geographic datasets (e.g., climatic, geological, and topographic) that are intimately related to soil moisture values. SOMOSPIE is meant to be a general-purpose tool for using such datasets to downscale (i.e., increase resolution) satellite-based soil moisture products. This Jupyter Notebook is a result of a collaboration between computer scientists of the Global Computing Laboratory at the University of Tennessee, Knoxville and soil scientists at the University of Delaware (funded by NSF awards #1724843 and #1854312).

# Computational solutions: Jupyter + XSEDE Jetstream



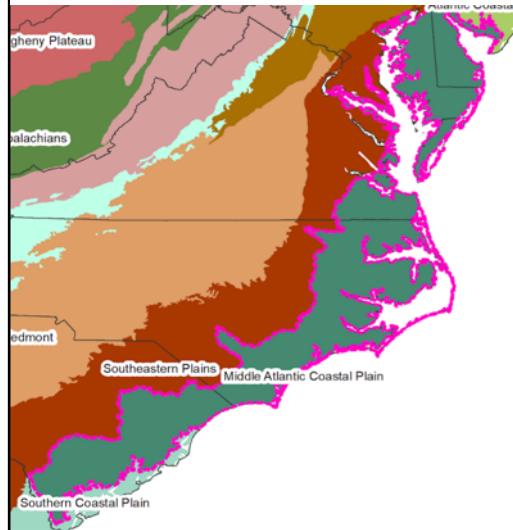
The screenshot shows a Jupyter Notebook interface running on a Mac OS X system. The title bar indicates the notebook is titled "SOMOSPIE" and is "autosaved". The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, Help, Not Trusted, and Python 3. The notebook content starts with a section titled "SOMOSPIE" and "SOil MOisture SPatial". The "Introduction" section discusses soil moisture as a critical variable for crop yield and potential agricultural thresholds. It mentions that current availability over large areas comes from remote sensing (e.g., satellites with radar sensors) which provide daily, nearly global coverage of soil moisture. However, satellite soil moisture datasets have a major shortcoming in that they are limited to coarse spatial resolution (generally no finer than tens of kilometers). The text also notes that there exist higher resolution datasets (e.g., climatic, geological, and topographic) related to soil moisture values. SOMOSPIE is described as a general-purpose tool for using such datasets to downscale (i.e., increase resolution) satellite-based soil moisture products. This Jupyter Notebook is a result of a collaboration between computer scientists at the Global Computing Laboratory at the University of Tennessee, Knoxville and soil scientists at the University of Delaware (funded by NSF awards #1724843 and #1854312).

# Use case I: from 27x27km to 1x1km

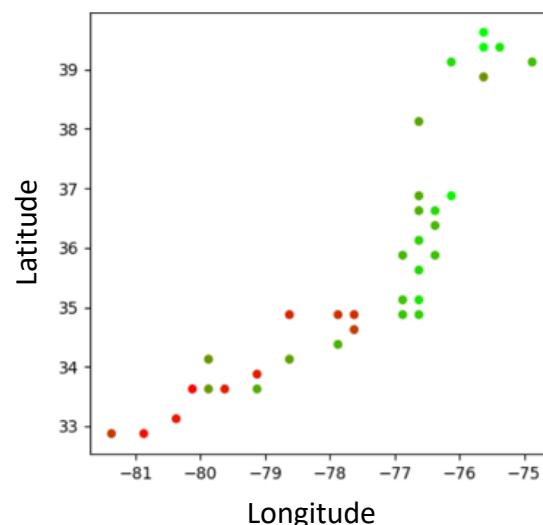
Fine-grained modeling of Mid-Atlantic region in April 2017:

- Terrain parameters: Elevation, Slope, and Wetness Index

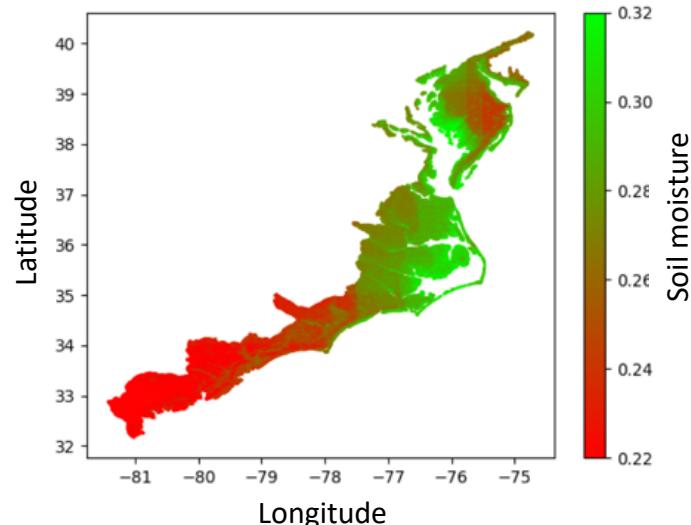
Level III Ecoregions of the  
Continental United States (CELV3)



Original  
satellite data (27x27km)



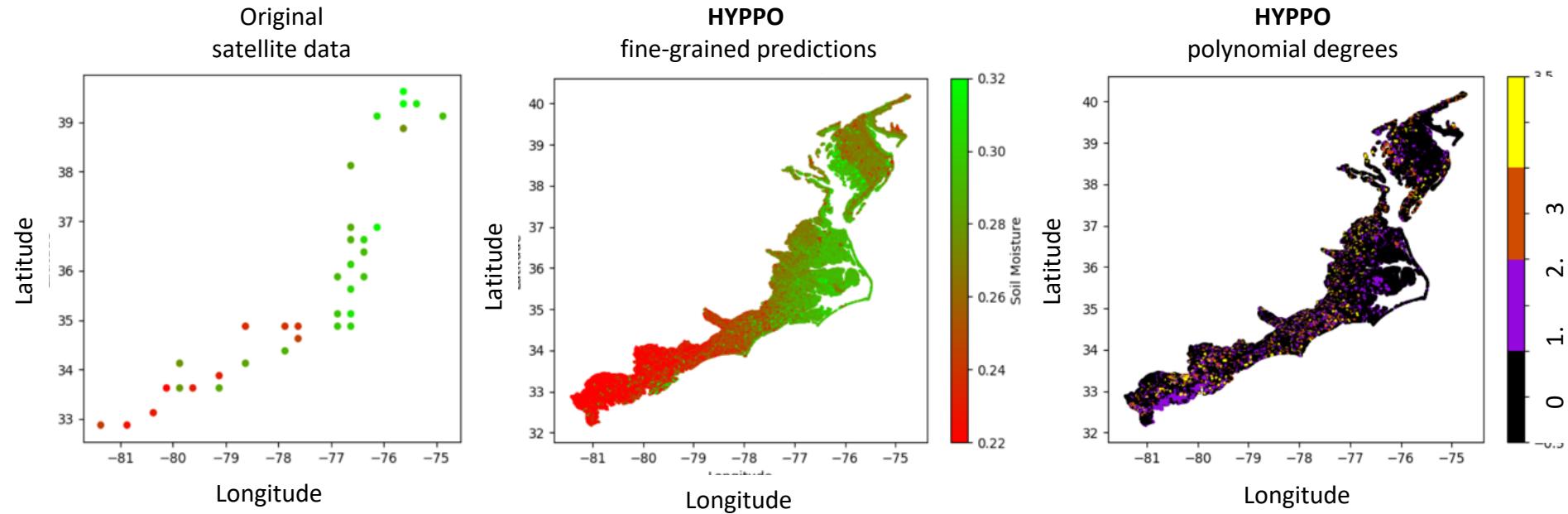
**Random Forest**  
fine-grained predictions 1x1km



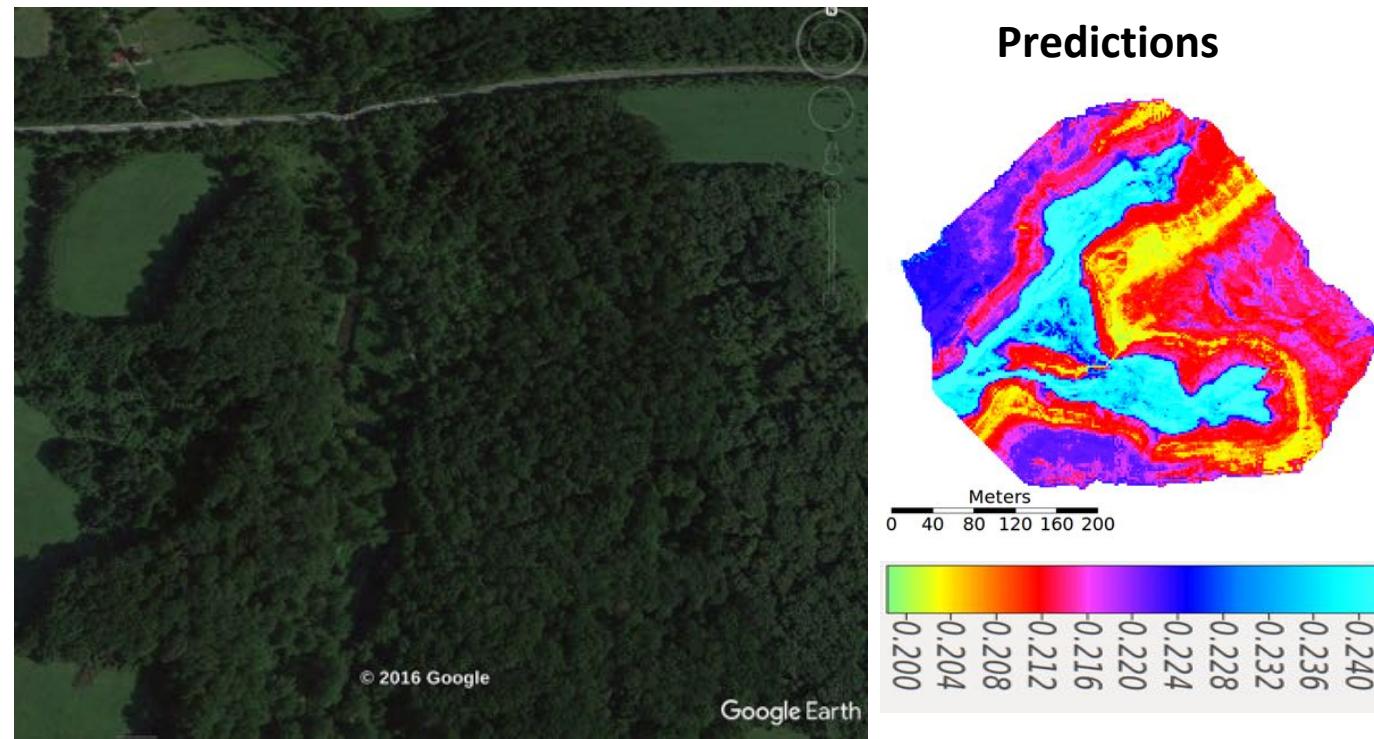
# Use case I: from 27x27km to 1x1km

Fine-grained modeling of Mid-Atlantic region in April 2017:

- Terrain parameters: Elevation, Slope, and Wetness Index

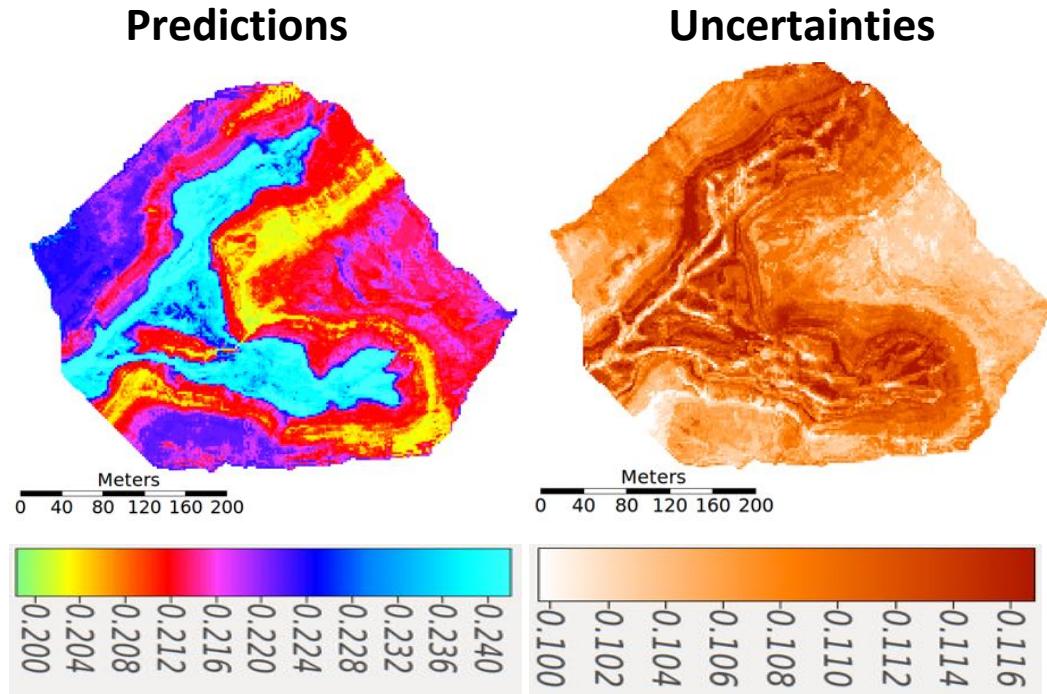
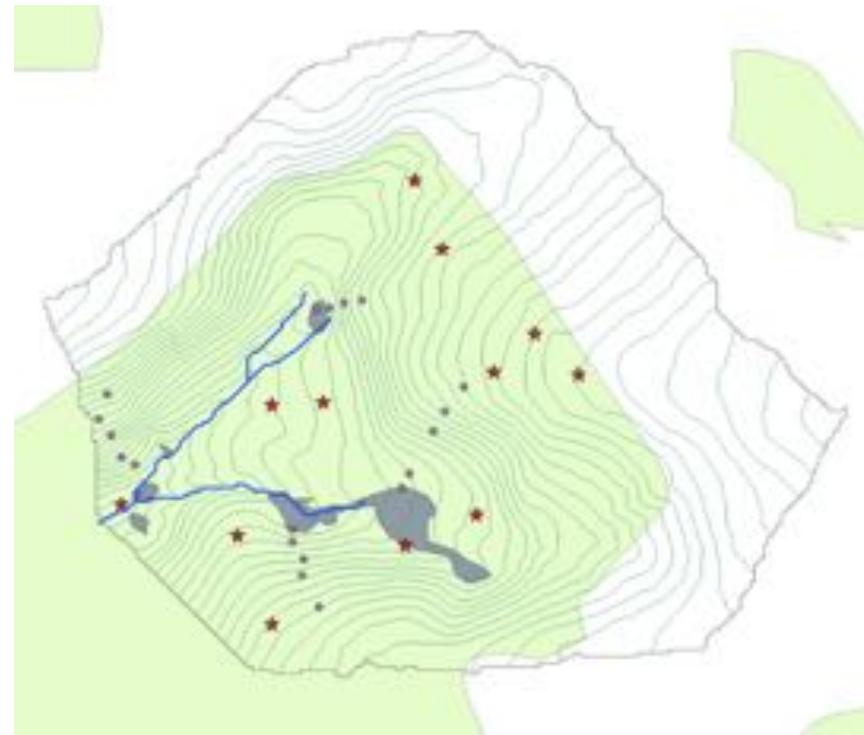


## Use case II: Local scale predictions - 1x1m resolution



M. Guevara , M. Taufer, and R. Vargas. Gap-Free Annual Soil Moisture Global across 15km Grids: 1991-2016. Earth System Science Data, 2019.

## Use case II: Local scale predictions - 1x1m resolution



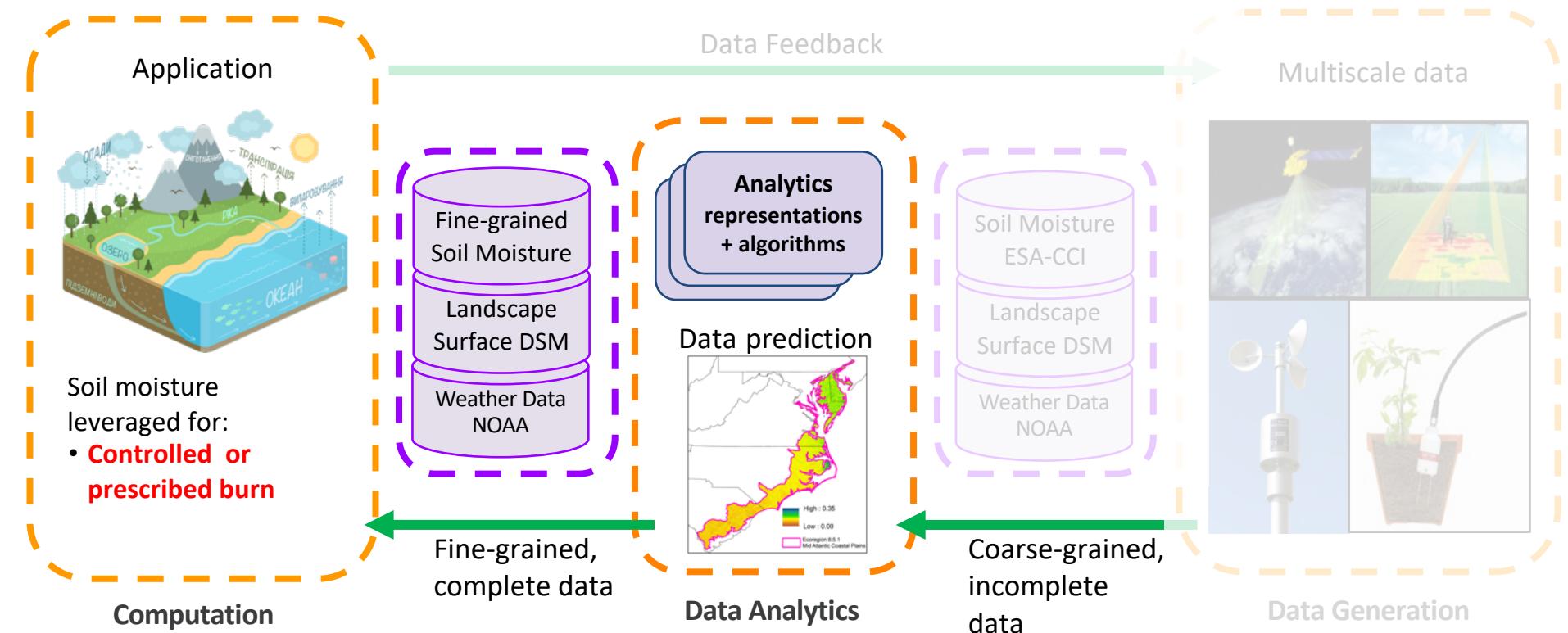
M. Guevara , M. Taufer, and R. Vargas. Gap-Free Annual Soil Moisture Global across 15km Grids: 1991-2016. Earth System Science Data, 2019.

# Combine computing and analytics: integration of soil moisture predictions into controlled (or prescribed) burn

Collaborator: David Icove's group (UTK)

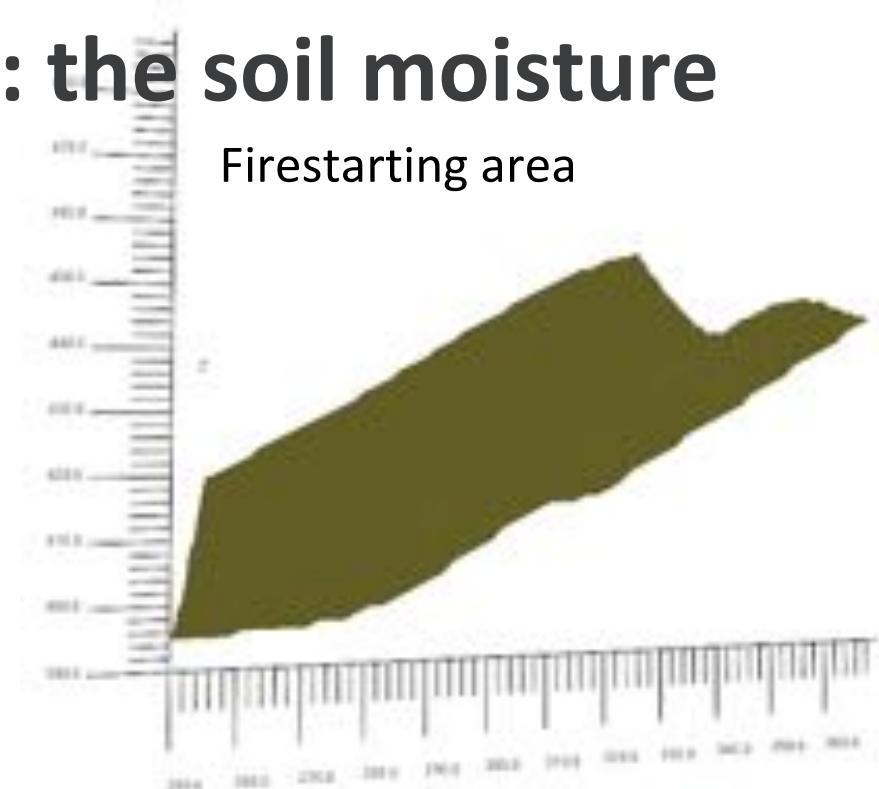
Platform: Tellico cluster (IBM Power9 system) – supported  
by 2019 IBM Shared University Research (SUR) Award

# Soil moisture data for simulating controlled burn

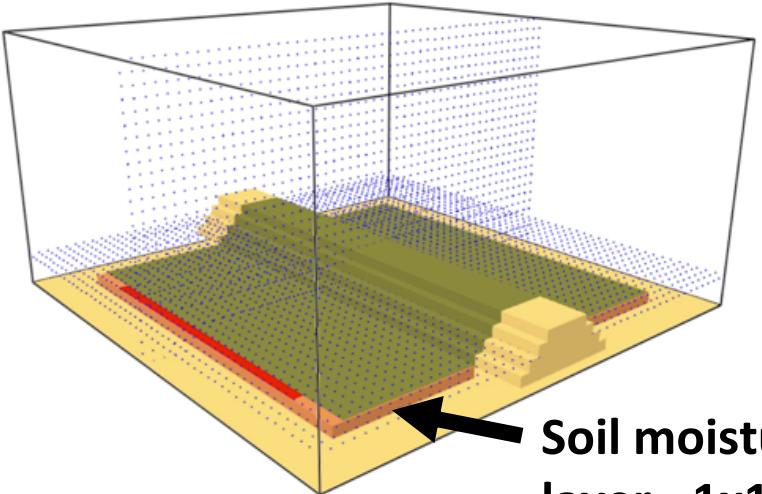


# Elephant in the room: the soil moisture

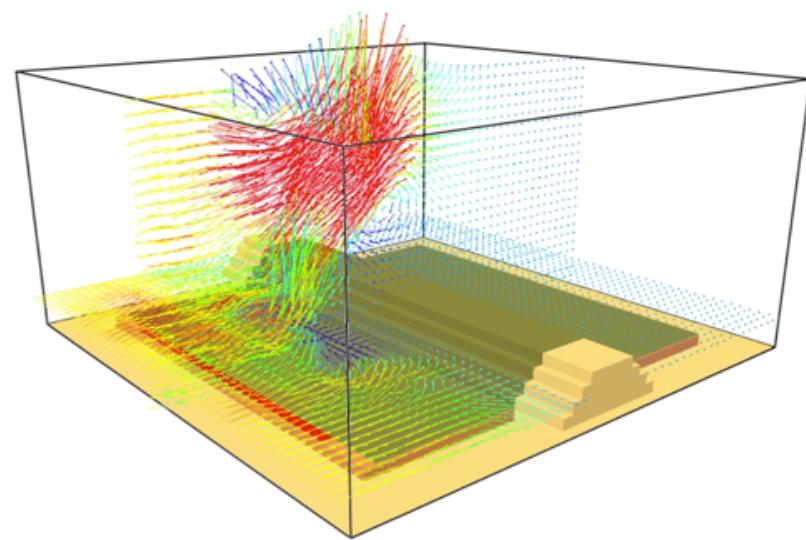
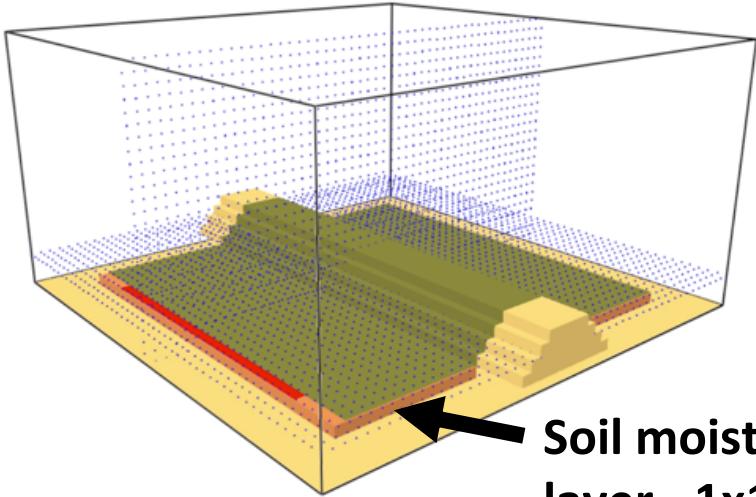
- Simulation of the 2016 Gatlinburg wildfire
- Software:
  - Fire Dynamics Simulator (FDS) - large-eddy simulation (LES) for low-speed flows
- Platform:
  - IBM Power9 cluster at UTK
- Simulation specs:
  - $120\text{m} \times 120\text{m} \times 100\text{m}$  domain
  - 5 frames/sec temp. resolution

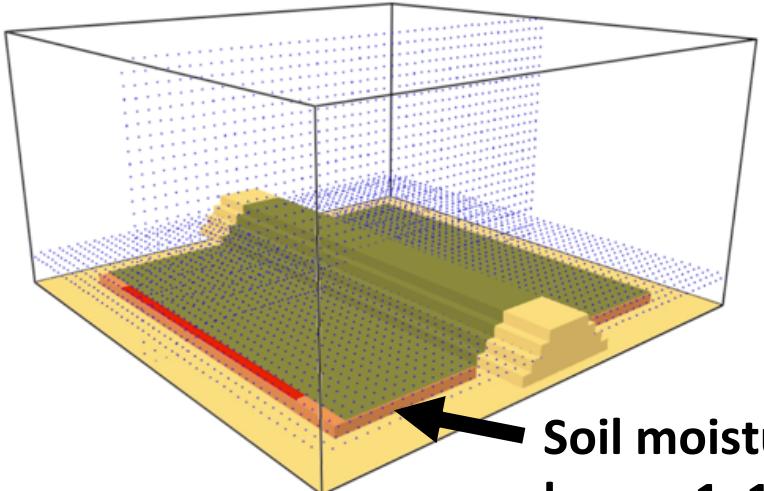


**Soil moisture is missing in FDS**

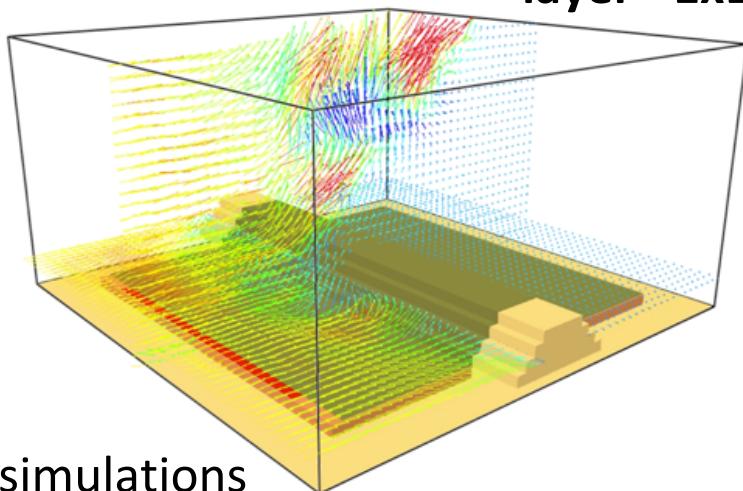
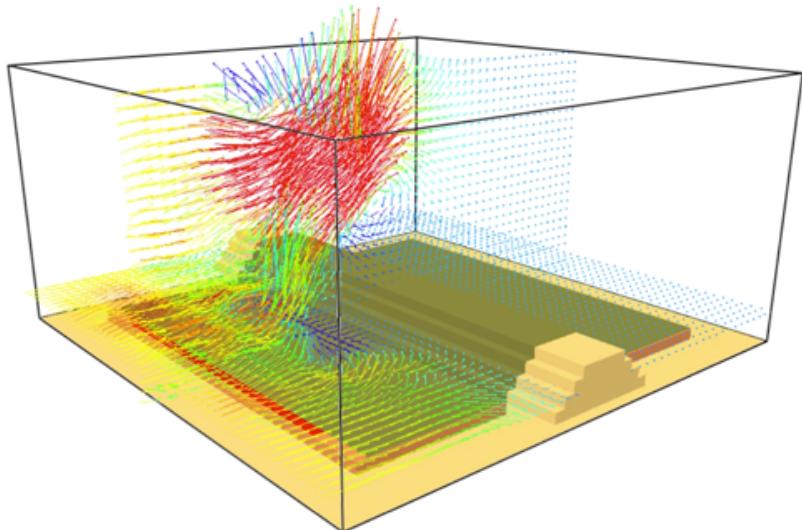


**Soil moisture  
layer - 1x1m**

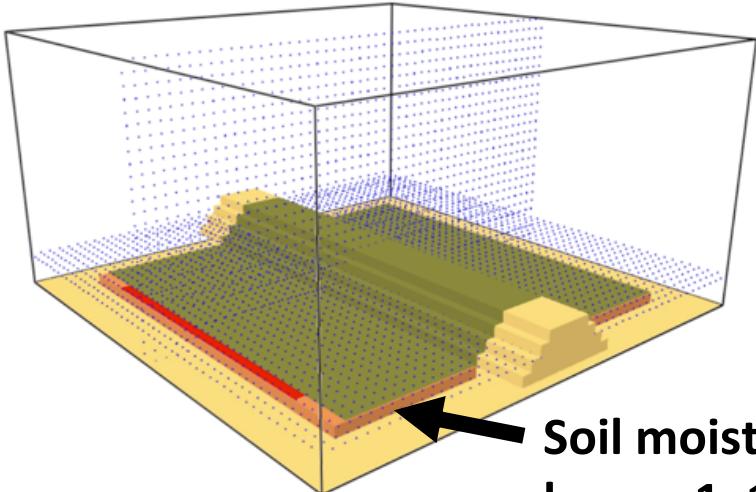




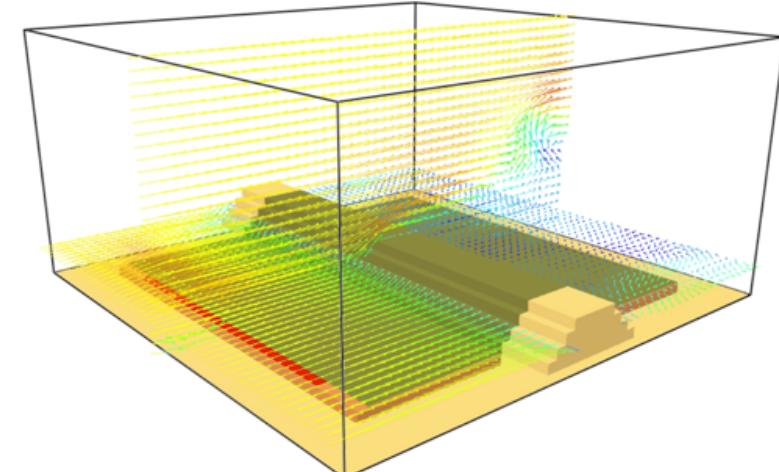
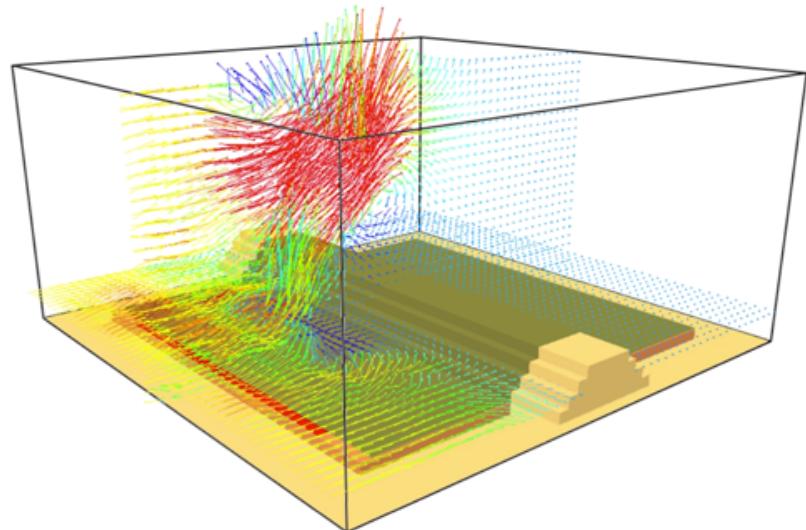
**Soil moisture  
layer - 1x1m**



FDS simulations



**Soil moisture  
layer - 1x1m**



FDS simulations

**Build trust in results through reproducibility,  
replicability, and transparency**

**Collaborator: Victoria Stodden's Group (UIUC)**

**Platform: NSF XSEDE Jetstream**

**NSF OAC 1941443 EAGER: Reproducibility and  
Cyberinfrastructure for Computational and Data-Enabled  
Science**

# Leveraging other NSF projects: Whole Tale

- Building an **open platform for computational reproducibility**
  - Create and publish **executable research objects ("Tales"**)
- Simplify process of creating & verifying reproducible computational artifacts for scientific discovery

Easy-to-access **cloud-based** computing environments



**Transparent** access to research **data**



**Export** and **publish** executable research **objects**



## WHOLE TALE Dashboard

New to WholeTale? Learn more about using the Dashboard...

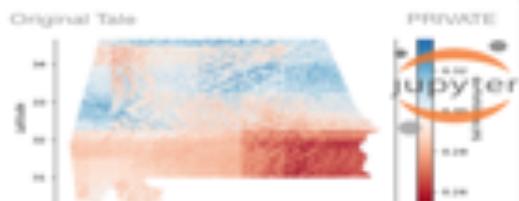
## Tales

Search tales...



All Tales

My Tales



**MACHINE LEARNING**  
**SOMOSPIE: A modular SOIL MOisture S...**

by Danny Ronabaugh, Mario Guivara, Ricardo Llamanas, Joy Kifson, Rodrigo Vargas, Michela Taufer

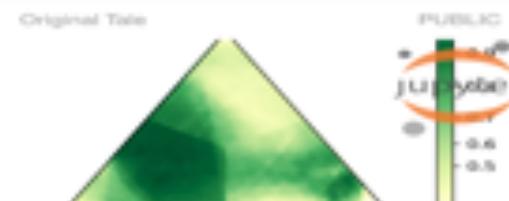
Run Tale



**MATERIALS SCIENCE**  
**Accelerated discovery of metallic i...**

By Logan Ward

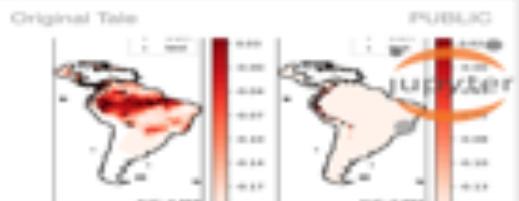
Run Tale



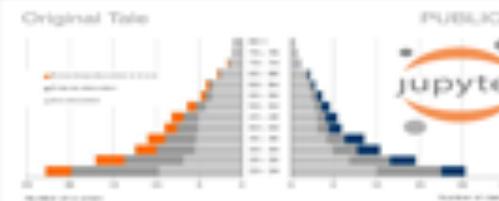
**MATERIALS SCIENCE**  
**Predicting the Properties of Inorga...**

By Logan Ward

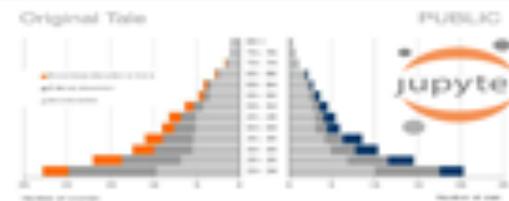
Run Tale



**ECOLOGY**  
**Replication of a classical ecologic...**



**SCIENCE**  
**Concept mapping**



**SCIENCE**  
**Species co-occurrence distribution**

[Return to Dashboard](#)

## SOMOSPIE: A modular SOIL MOisture S...

Danny Rorabaugh, Mario Guevara, Ricardo Llamas, Joy Kilon, Rodrigo Vargas, Michela Taufer

[Run](#)[Close](#)[Dataset](#) [Files](#) [Metadata](#)

Title: SOMOSPIE: A modular SOIL MOisture SPatial Inference Engine based on data driven decisions

**Authors:**

First Name: Danny Last Name: Rorabaugh ORCID: <https://orcid.org/0000-0000-0000-0000>

First Name: Mario Last Name: Guevara ORCID: <https://orcid.org/0000-0000-0000-0000>

First Name: Ricardo Last Name: Llamas ORCID: <https://orcid.org/0000-0000-0000-0000>

First Name: Joy Last Name: Kilon ORCID: <https://orcid.org/0000-0000-0000-0000>

First Name: Rodrigo Last Name: Vargas ORCID: <https://orcid.org/0000-0000-0000-0000>

First Name: Michela Last Name: Taufer ORCID: <https://orcid.org/0000-0000-0000-0000>

Add another author

Category: Machine Learning

Environment: Jupyter Classic

Datasets used: No citable data

License: Creative Commons Attribution 4.0 International

# Capturing metadata



File Edit View Insert Cell Kernel Widgets Help

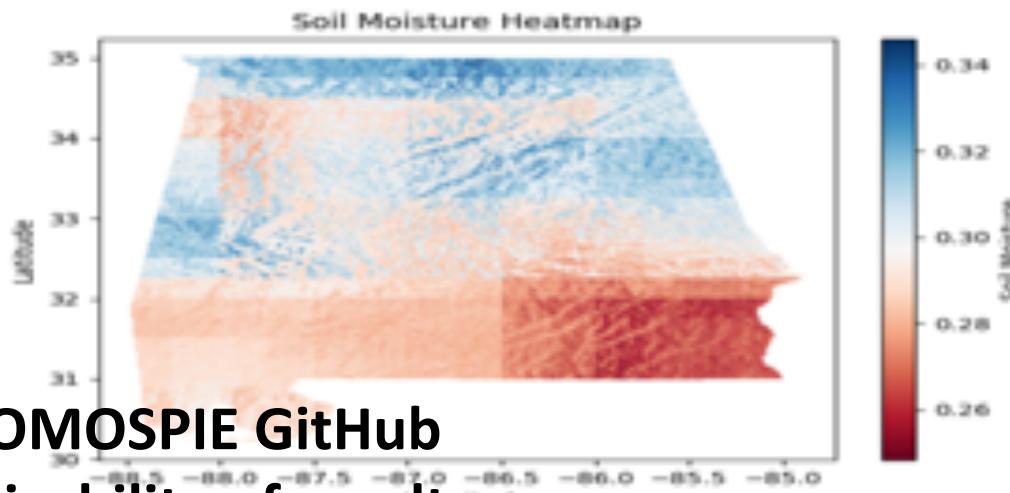
Not Trusted

Python 3

In [22]:

```
try:
    visualiser.display_output()
except ValueError:
    print("You must select an output folder in the widget above before running this cell.")

{0: {"START": "...", "CODE": ".../code/", "DATA": ".../data/", "OUTPUT": ".../out/", "COV_FILE": "topo_predictors/stack.tif", "SM_FILE": "NSA_OCI/2017_ESA_monthly.rds", "EVAL_FIELIST": "...", "MARE_T_R": 1, "USE_PCA": 0, "RAND_SEED": 0, "USE_VIS": 1, "BUFFER": 0, "SUPER": 0, "MIN_TEST_POINTS": 11, "VALIDATE": 1.0, "REG_LIST": [{"STATE": "Alabama"}], "MONTHS": [1], "MODICT": {"RF": {}, "UHMODEL": {}}, "YEAR": [2017], "COV_LAYERS": ["Aspect"], "TRAIN_DIR": ".../2017/L", "EVAL_DIR": ".../2017/m"}, 1: {1: {"seed": 65050}}, 2: 1580838903.2455065, 3: 22.56052279472351},
4: {"Alabama": '/home/jowyan/work/workspace/out/job_2020_02_04_17_54_41/2017/month1=0,30_65050/Alabama'}}
```



Plug into SOMOSPIE GitHub  
Enable replicability of results

Method completion time = t=8.003055095672607



Search or jump to...

Pull requests Issues Marketplace Explore



TauferLab / SOMOSPIE

[Unwatch](#) 4[Star](#) 0[Fork](#) 0[Code](#)[Issues 0](#)[Pull requests 0](#)[Actions](#)[Projects 0](#)[Wiki](#)[Security](#)[Insights](#)[Settings](#)[Releases](#)[Tags](#)[Edit release](#)[Delete](#)[Latest release](#)[v1.2](#)[ca24063](#)[Verified](#)[Compare ▾](#)

## SOMOSPIE version 1.2

leobardovalera released this 13 hours ago - 2 commits to master since this release

[v1.2](#)[Update README.md](#)

<https://github.com/TauferLab/SOMOSPIE/releases/latest>

[Assets 3](#)[SOMOSPIE.tgz](#)

59.1 KB

[Source code \(zip\)](#)[Source code \(tar.gz\)](#)