

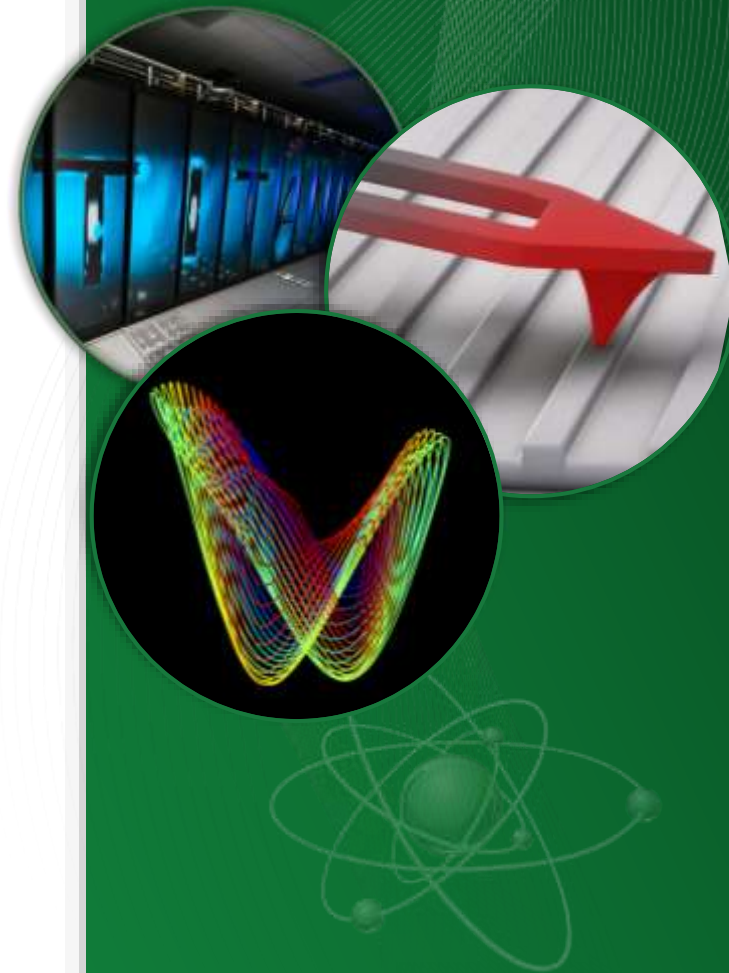
Pycroscopy –

A python package for storing
and analyzing imaging and
spectroscopy data

- Suhas Somnath
- Chris R. Smith
- Stephen Jesse

 **OAK RIDGE** | OAK RIDGE
National Laboratory | LEADERSHIP
COMPUTING FACILITY

 INSTITUTE FOR FUNCTIONAL
IMAGING OF MATERIALS | CENTER FOR
OAK RIDGE NATIONAL LABORATORY | NANOPHASE
MATERIALS SCIENCES



ORNL is managed by UT-Battelle
for the US Department of Energy

Multitude of Instruments



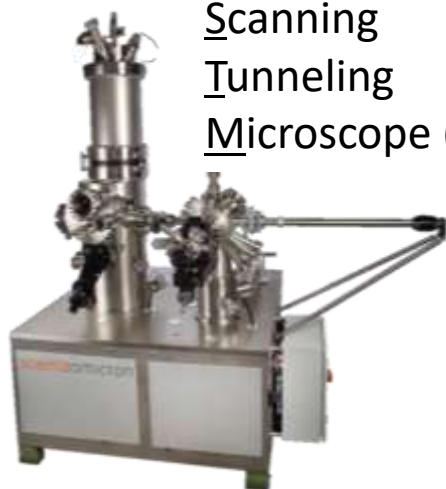
Micro Raman Microscope



Atomic Force
Microscope (AFM)



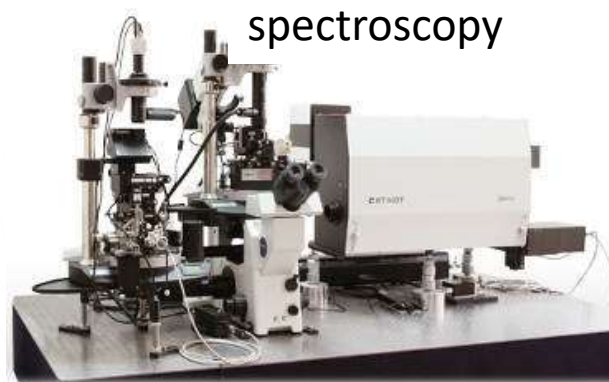
AFM with Infrared
spectroscopy (AFM-IR)



Scanning
Tunneling
Microscope (STM)



Scanning
Transmission
Electron
Microscope (STEM)



AFM with Raman
spectroscopy

What we wanted



Instrument Tier

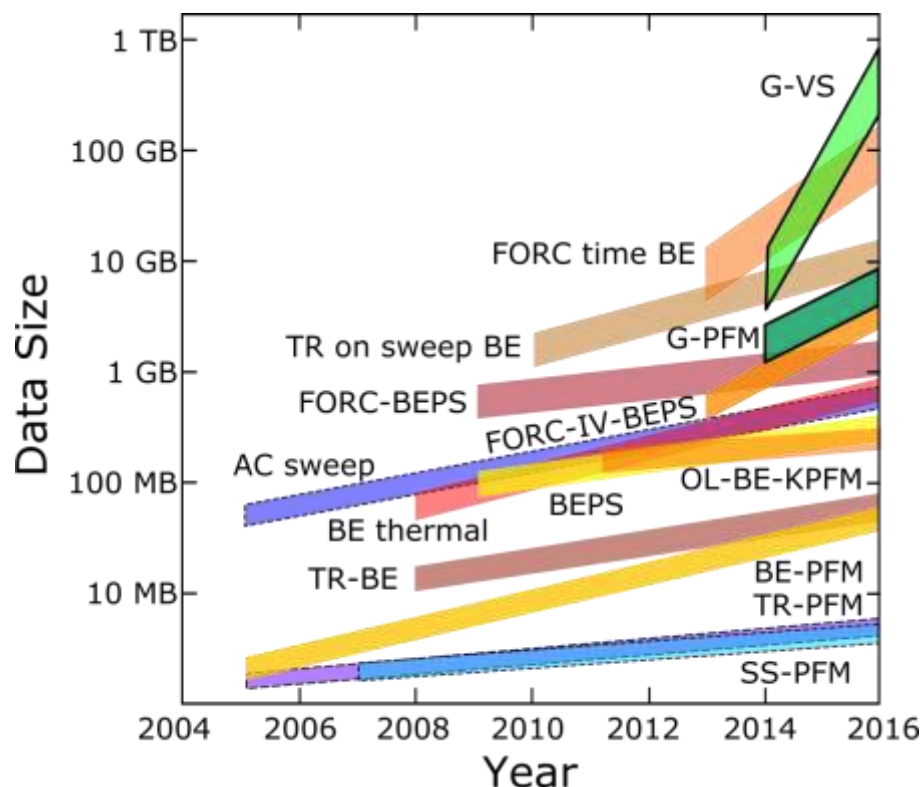
?



Interactive visualization, analysis,
storage on supercomputers

Growing Data Sizes and Dimensionality

Evolution of Scanning Probe Microscopy Data



- Data sizes have grown from ~ 10 MB to ~ 1 TB in 10 years!
- Dimensionality ranges from 1D spectra to 7D hyperspectral datasets
- Cannot use laptops to analyze data

Instrumentation Software Inadequate for Analysis



- Software provided for controlling instruments typically only comes with basic data analysis capabilities
- Integrating user-developed functionality often impossible



Multitude of File Formats

- Proprietary
- Incompatible

.wdf

.ibw

.cdb

.asc

.dm3

.mat

≠



Disjoint & Unorganized Communities



- Clustering
- Fit spectra ...



- Filter Image
- Register Image ...



- Fit Spectra
- SVD Filtering ...

- FFT Filtering
- SVD Filtering ...



- FFT Filtering
- Classify Images ...



- Register Images
- Clustering



Cannot Share Code Efficiently

- HIGHLY instrument-specific code
- Different programming languages
- Often licensed / costly software like Matlab
- Most popular sharing method = email!
- No centralized repository

Problems Opportunities in Imaging

1. Closed science
 - a. No traceability for data analysis
 - b. Results not (readily) reproducible
2. Multiple, incompatible, proprietary data formats
3. Disorganized and unorganized communities
4. No proper analysis software
5. Growing data volumes, variety, and dimensionality

The Solution



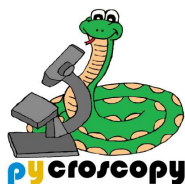
Instrument Tier



Automated, standardized,
modularized data acquisition



Instrument-agnostic, self-describing,
model in HPC-friendly file format



Centralized repository for data
processing, analysis



Interactive visualization + analysis +
storage on supercomputers

Expectation of Data Model / File Format

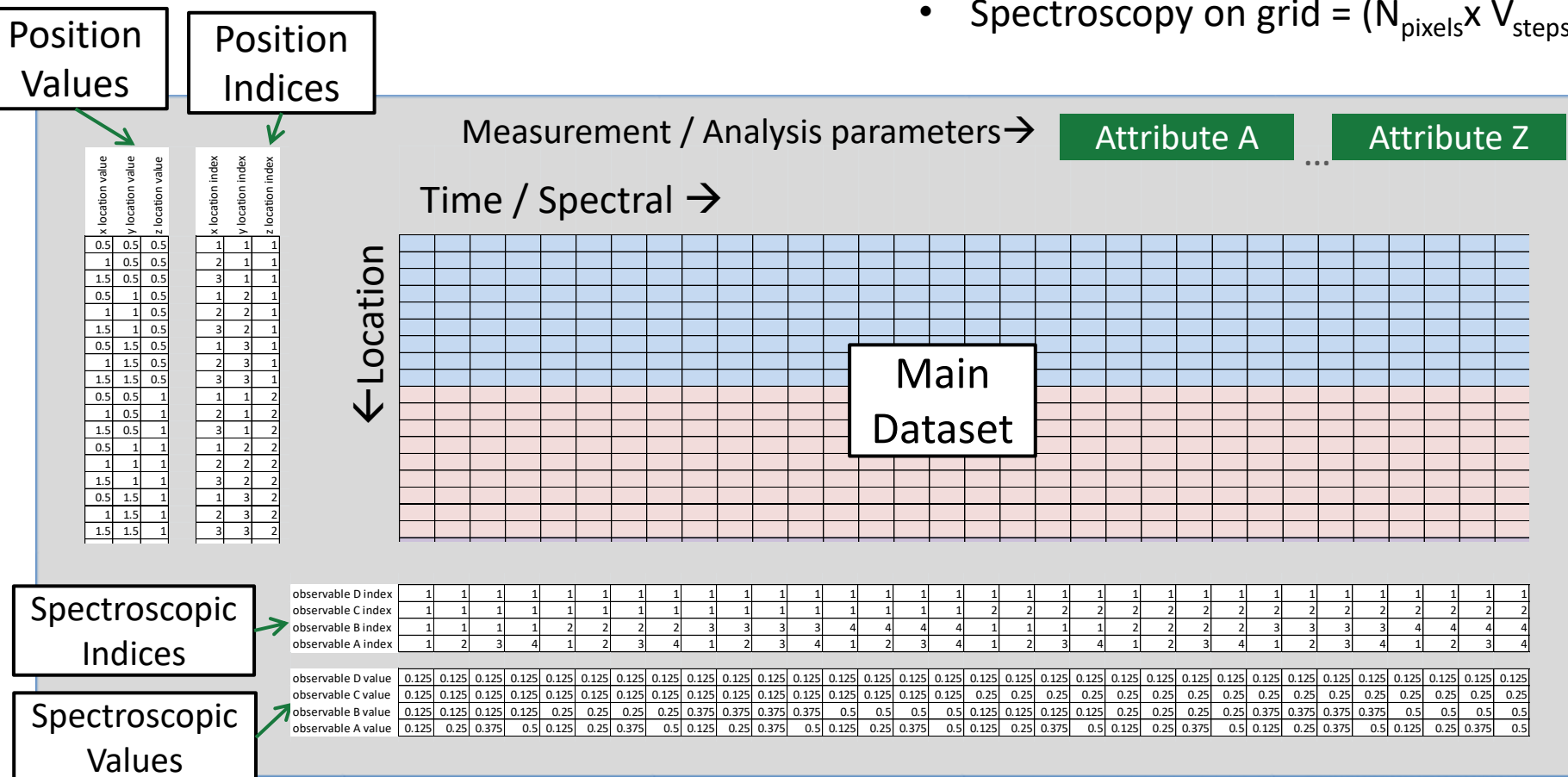
- File format - Established standard in scientific research
- Store multiple datasets of different shapes, dimensionalities, precision and sizes.
- Scale very efficiently from few kilobytes to several terabytes
- Able to read and write data using any programming language including Python, R, Matlab, C/C++, Java, Fortran, Igor Pro, etc.
 - (without requiring installation of modules that are hard to install)
- Store metadata - experimental or analysis parameters
- Highly flexible and poses minimal restrictions on how the data can and should be stored.
- Compatible with cloud and high-performance computing paradigms (support parallel read and write)

Universal Data Model

- Data stored as 2D matrix of (position x spectral values) regardless of dimensionality
- Ancillary datasets explain the data

Example data types:

- 2D images = ($N_{\text{pixels}} \times 1$)
- Single spectra = ($1 \times Z_{\text{steps}}$)
- Spectroscopy on grid = ($N_{\text{pixels}} \times V_{\text{steps}}$)

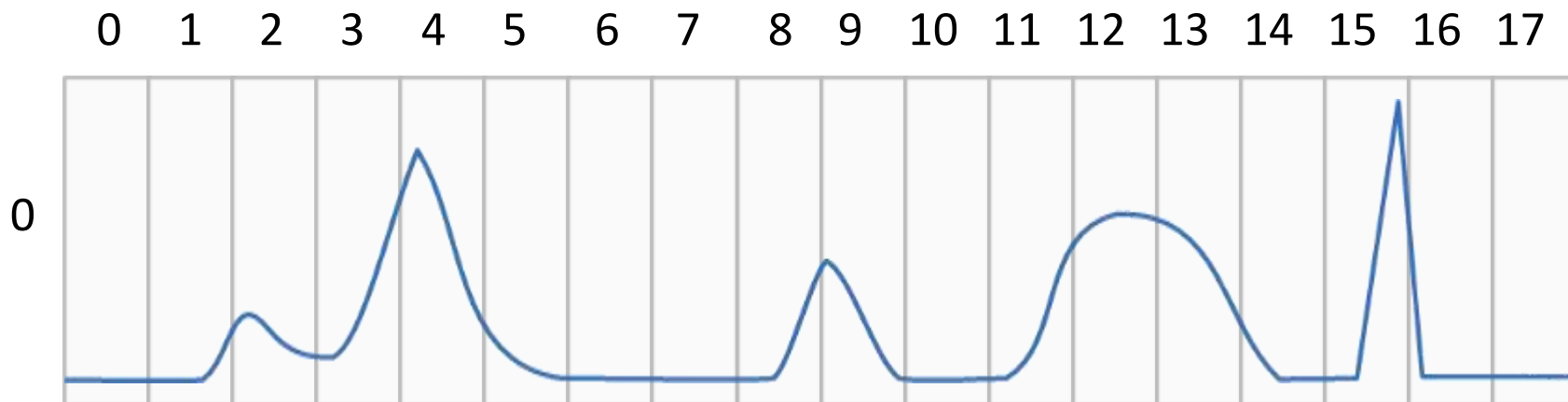


Additional information available at:

Universal Data Model – 1D spectra

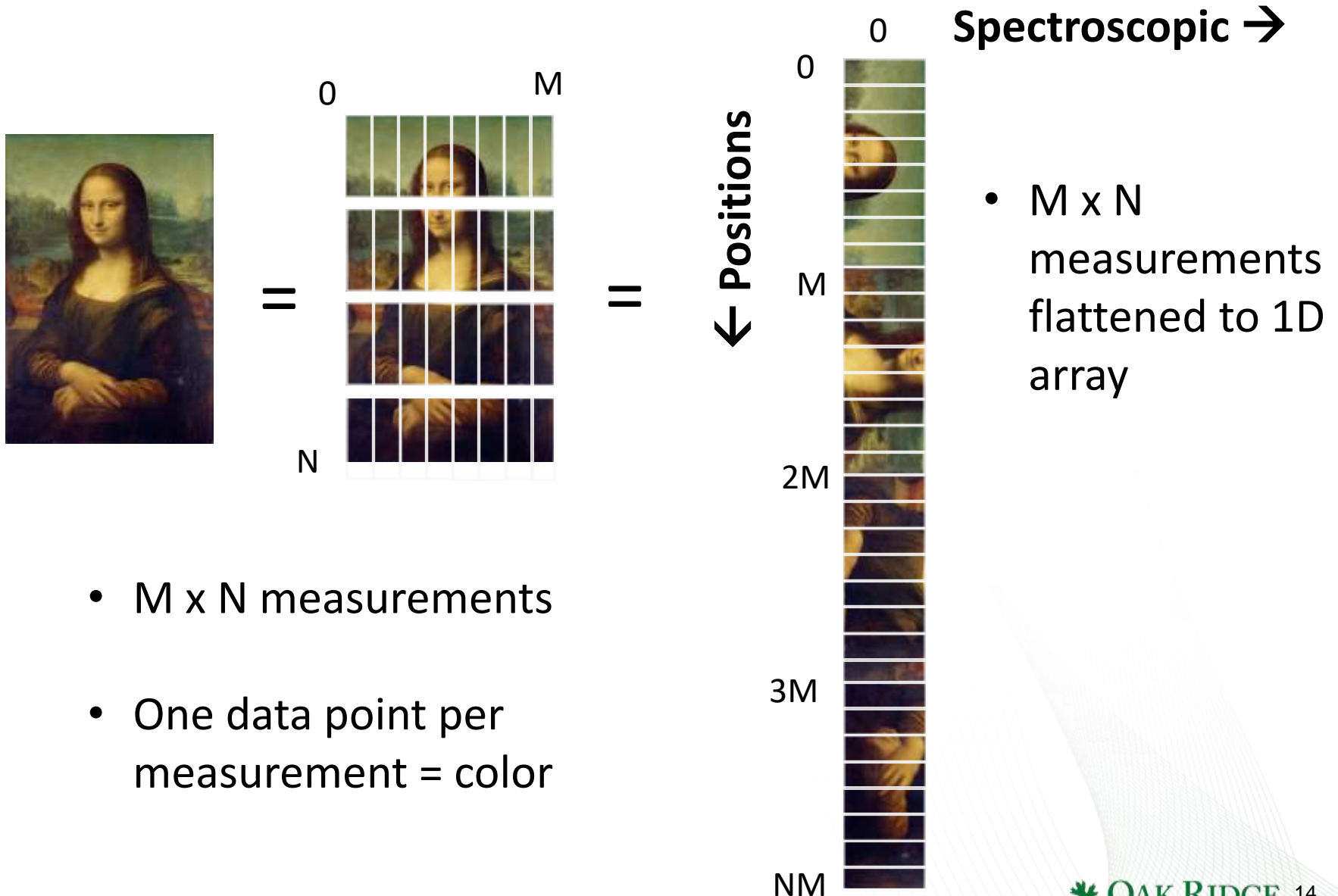
Time / Spectroscopic →

← Positions



- Each measurement = 17 data points in time
- Single measurements = single entry in Positions axis

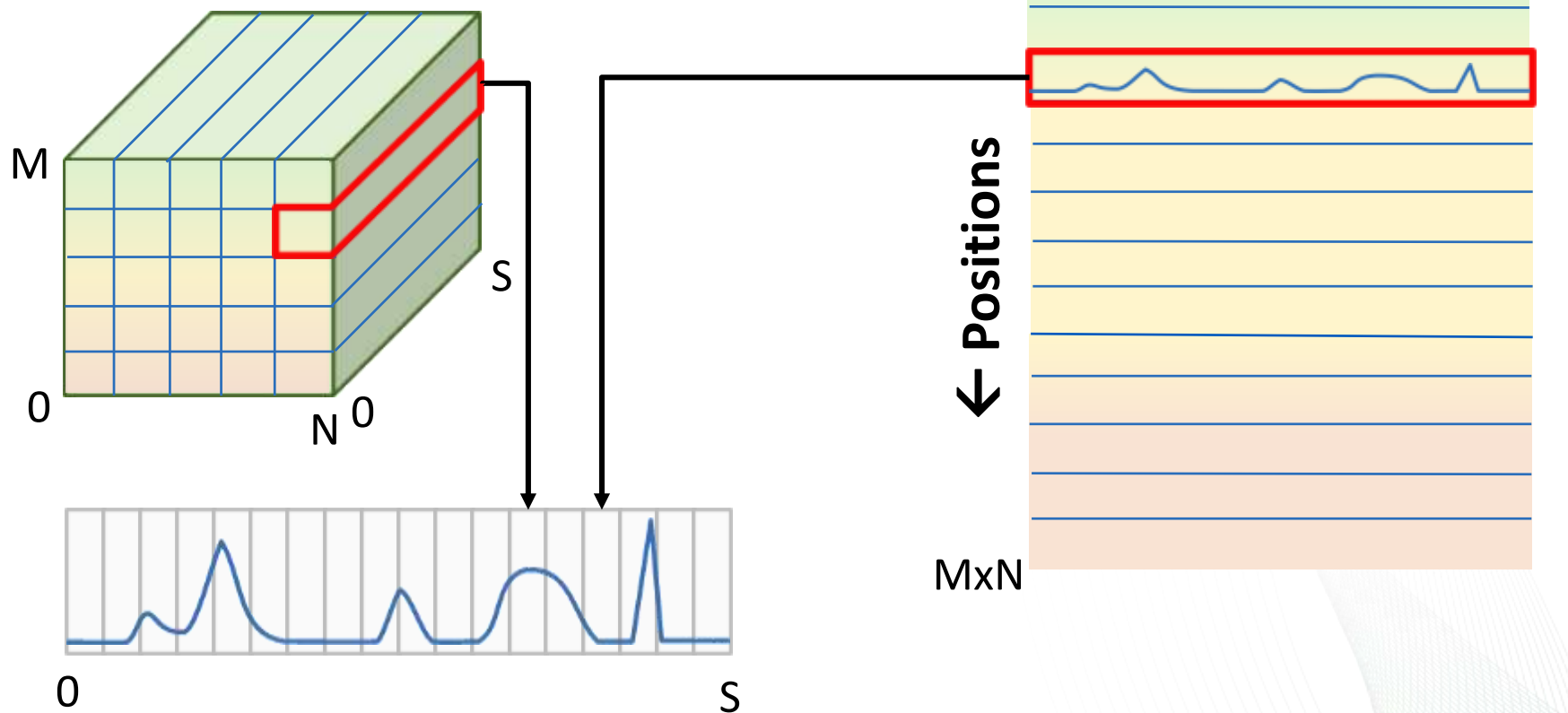
Universal Data Model – 2D Image



- $M \times N$ measurements
- One data point per measurement = color

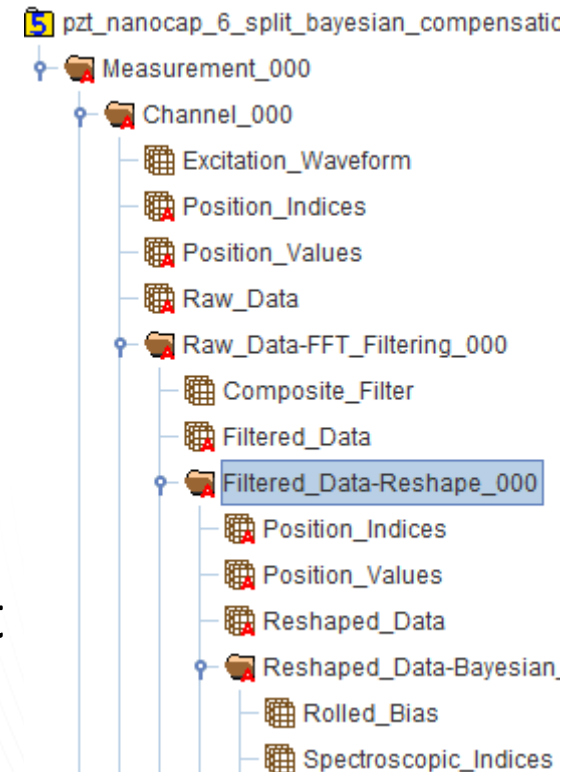
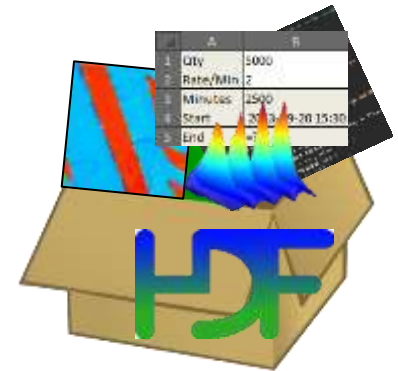
Universal Data Model – Spectra on Grid

- S data points per measurement
- $M \times N$ measurements



Hierarchical Data Format (HDF5)

- A HDF5 file is a smart container
 - Capable of storing multidimensional datasets, Images, text, measurement parameters, etc.
 - Contents organized like traditional folders and files
 - **Groups** - Analogous to file folders
 - **Dataset** – 1 to N dimensional data
 - Integer, floating point, complex numbers etc
 - **Attributes** – {Key : value} pairs useful for describing data and experimental parameters, etc.
- Easily accessible – C, C++, python, Java....
- Tree structure + nomenclature + attributes are **records of workflow** applied to dataset
- Parallel read / write, HPC compatible

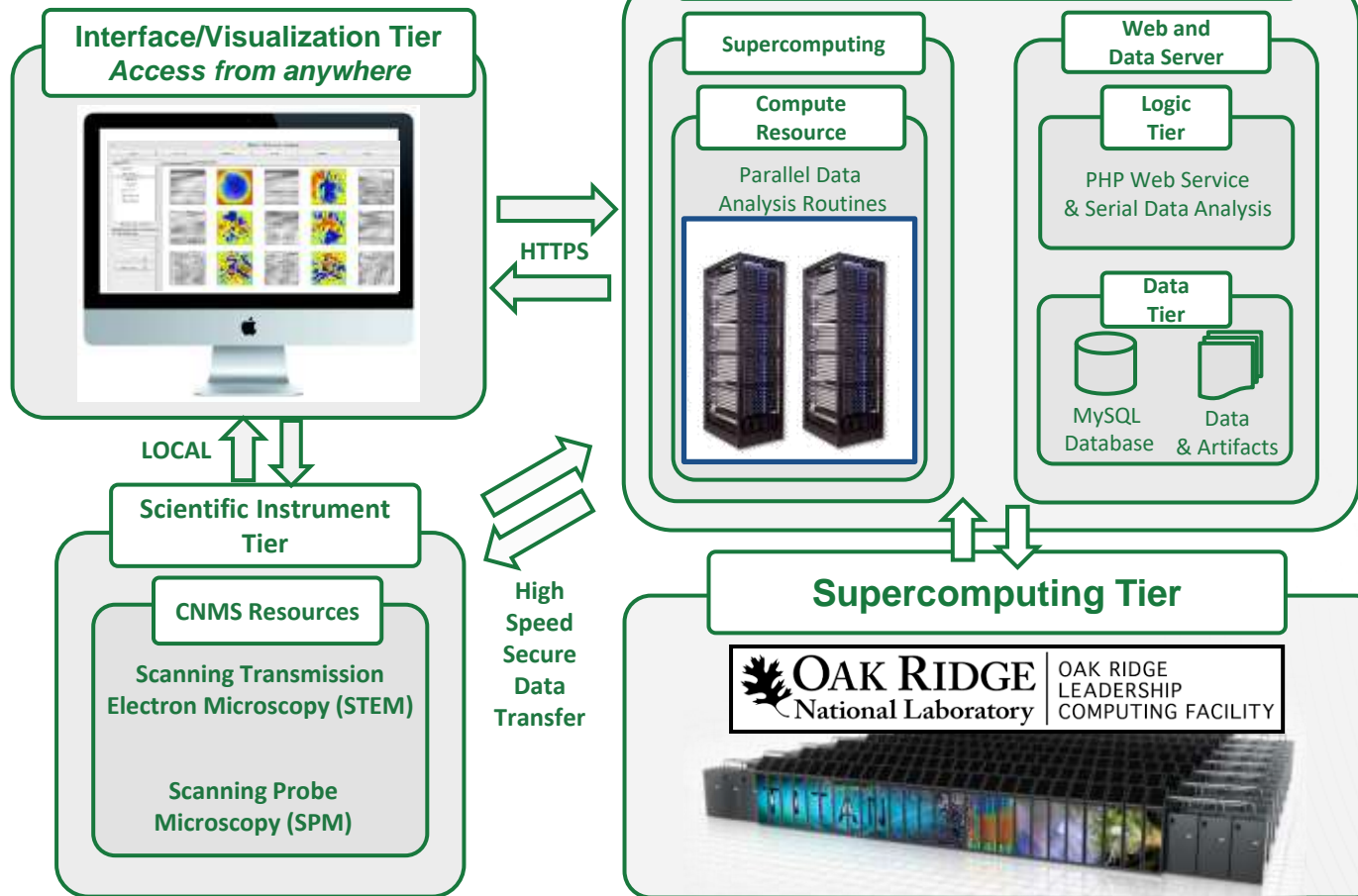


Expectation from Software

- Easy to learn and understand
- Strong support-base
- Established community standard
- Straightforward to implement and maintain
- Optimized libraries for scientific and numeric algorithms
- Access to existing imaging related packages
- Free
- Scalable to multiple CPU cores and HPC

(Purely) Programmer-Driven Solution

Software connecting scientific instruments to supercomputers



- **Successes:**

- Easy to use – Point-click
- Fast – on super-computers

- **Shortcomings:**

- Very long development cycle
- Very expensive
- Brittle (points of failure)
- Scientists had no control!!

Python for Scientific Research

Very easy to learn + code

Numerous, **powerful**
libraries for science



NumPy



SciPy



- Facilitates innovation
- More robust code
- Improved adoption of new methods / standards
- Accelerates scientific progress

Cross-
platform



scalable



Strong user
community



stackoverflow

Established standard for:

- Microscopy
- Microbiology
- Deep learning
- Data science
- Neutron science
- More!

All for a princely sum of ... **\$0!**

pycroscopy



- Python package
- Open source & free
- Written by scientists
- Data centric
- Instrument-independent data model in HDF5
- Instrument-independent analysis algorithms
 - Reusable across scientific domains

Pycroscopy - Organization

pycroscopy

I/O

- Data translators (proprietary formats to HDF5)

Analysis

- Physical model specific
- Fitting to model, etc.
- Physics based Regression

core

- HDF5 file i/o operations
- Base classes, visualization...

Processing

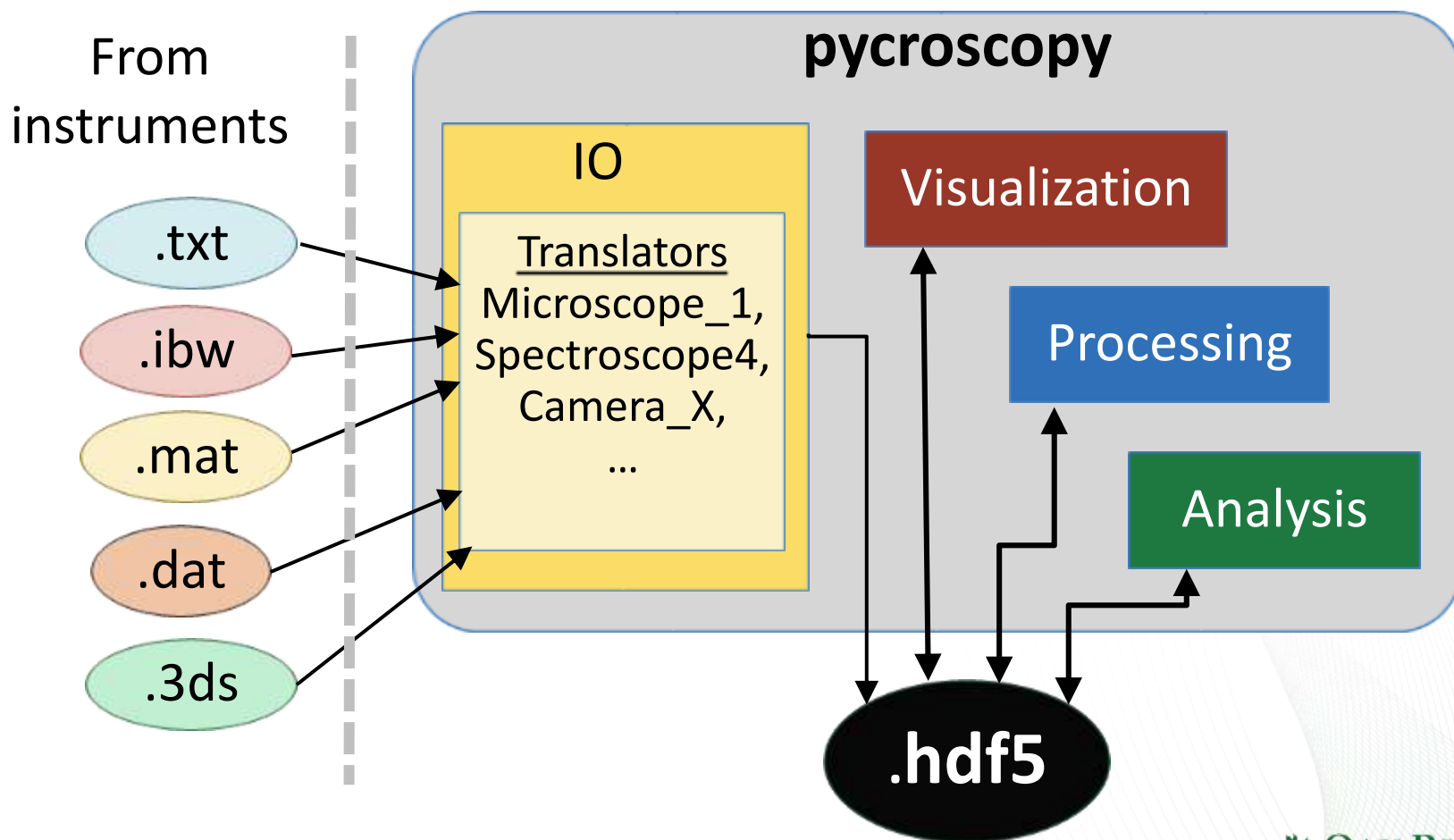
- Physical model agnostic
- Image filtering, registration, etc.

Visualization

- Plotting utilities
- Jupyter widgets

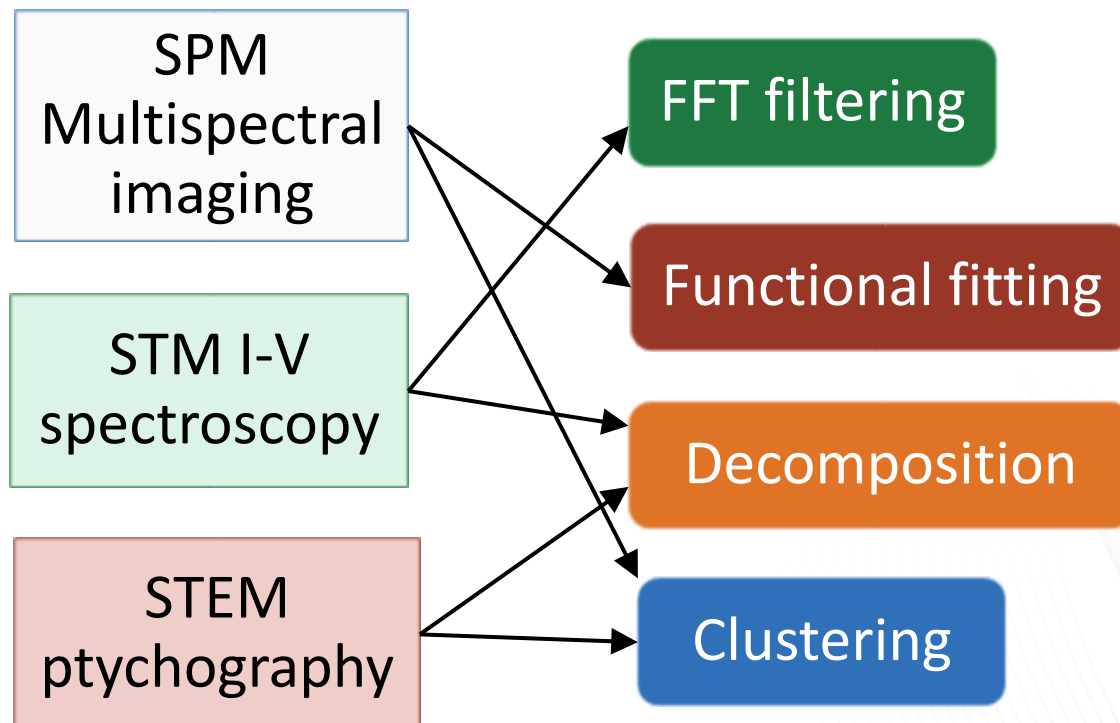
Entering the Pycroscopy Ecosystem

- hdf5 file is the hub for all operations
- Analysis, processing, visualization available after translation to .hdf5



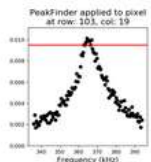
Pycroscopy – Instrument Agnostic Code

- Instrument-agnostic data allows instrument-agnostic code
- Single version of analysis and processing routine
- Brings multiple scientific communities together



Pycroscopy - Well documented

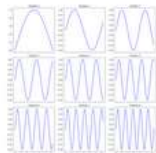
Guides to Pycroscopy



Formalizing Data Processing



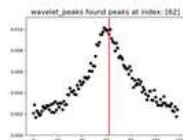
Input / Output / Computing utilities



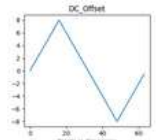
Plotting utilities



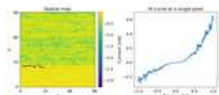
Primer to HDF5 and h5py



Speed up computations with parallel_compute()



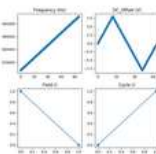
The PycroDataset



Translation and the NumpyTranslator



Utilities for handling data types and transformations



Utilities for reading Pycroscopy HDF5 files

```
reshape_to_Ndims(h5_main, h5_pos=None, h5_spec=None, get_labels=False, verbose=False, sort_dims=False) [source]
```

Reshape the input 2D matrix to be N-dimensions based on the position and spectroscopic datasets.

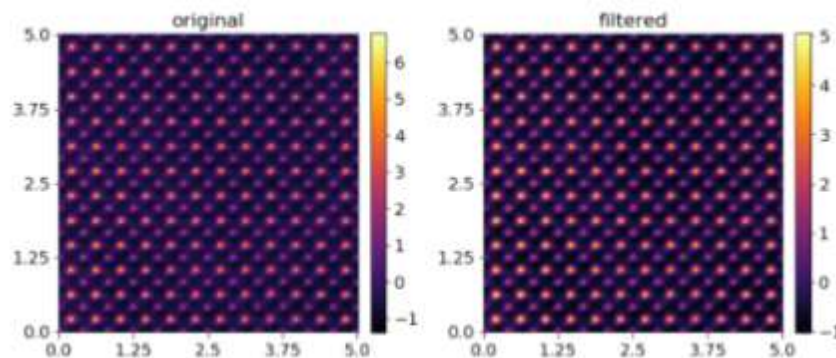
- Parameters:
- `h5_main` (*HDF5 Dataset*) – 2D data to be reshaped
 - `h5_pos` (*HDF5 Dataset, optional*) – Position indices corresponding to rows in `h5_main`
 - `h5_spec` (*HDF5 Dataset, optional*) – Spectroscopic indices corresponding to columns in `h5_main`
 - `get_labels` (*bool, optional*) – Whether or not to return the dimension labels. Default `False`
 - `verbose` (*bool, optional*) – Whether or not to print debugging statements
 - `sort_dims` (*bool*) – If `True`, the data is sorted so that the dimensions are in order from fastest to slowest. If `False`, the data is kept in the original order. If `get_labels` is also `True`, the labels are sorted as well.

- Returns:
- `ds_Nd` (*N-D numpy array*) – N dimensional numpy array arranged as [positions slowest to fastest, spectroscopic slowest to fastest]

To view the filter ('fft2'). Remember necessary to use transform. Also the inverse transform symmetric about result in the inverse times smaller than kept.

```
image_filter
image_filter

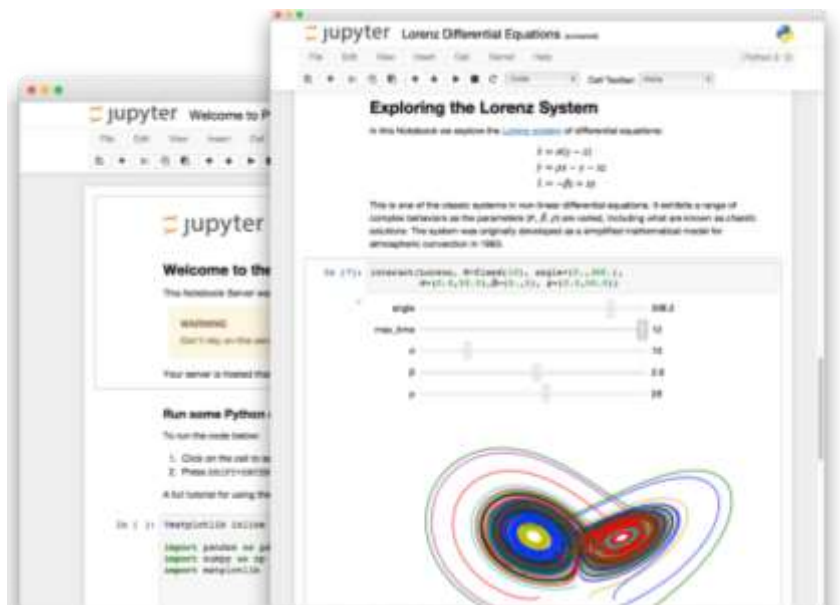
fig, axes = plt.subplots(ncols=2, figsize=(10, 5))
for axis, img, title in zip(axes, [image_raw, image_filtered], ['original', 'filtered']):
    _ = px.plot_utils.plot_map(axis, img, cmap=plt.cm.inferno,
                              x_size=x_edge_length, y_size=y_edge_length, num_ticks=5)
    axis.set_title(title)
fig.tight_layout()
```



Jupyter Notebooks



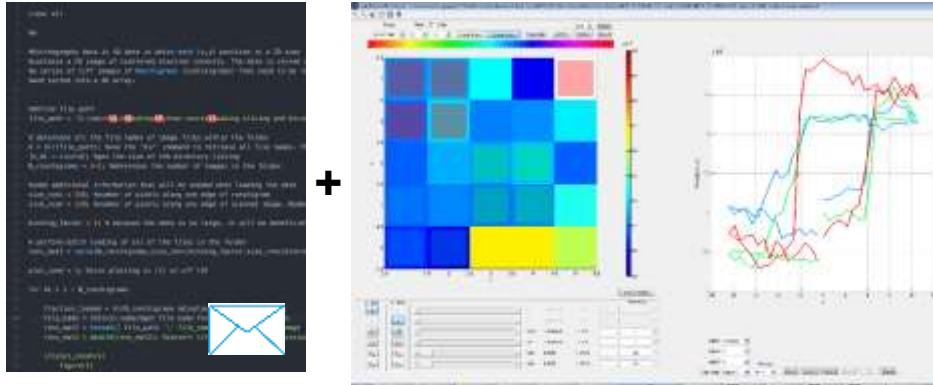
Jupyter Notebook



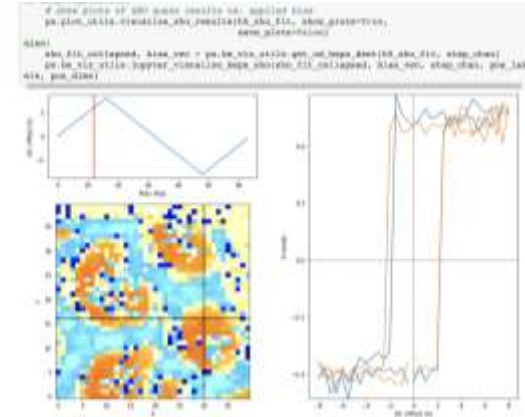
- Interactive documents
- Exploratory programming
- Code
- Text
- Images
- Interactive – slice through data, pan, move, rotate ...

Pycroscopy - Supporting User Research

Before 2016



Since 2016



Scripts + complicated, Matlab GUI

Set of simple Jupyter notebooks

Written by dedicated software engineer

Written by material scientists

Not customizable

Completely customizable.

2-3 hours of training before use

Notebooks include instructions. NO training required!

Deployed only on two offline workstations due to licensing restrictions = queue

Each user gets VMs with jupyter notebook server

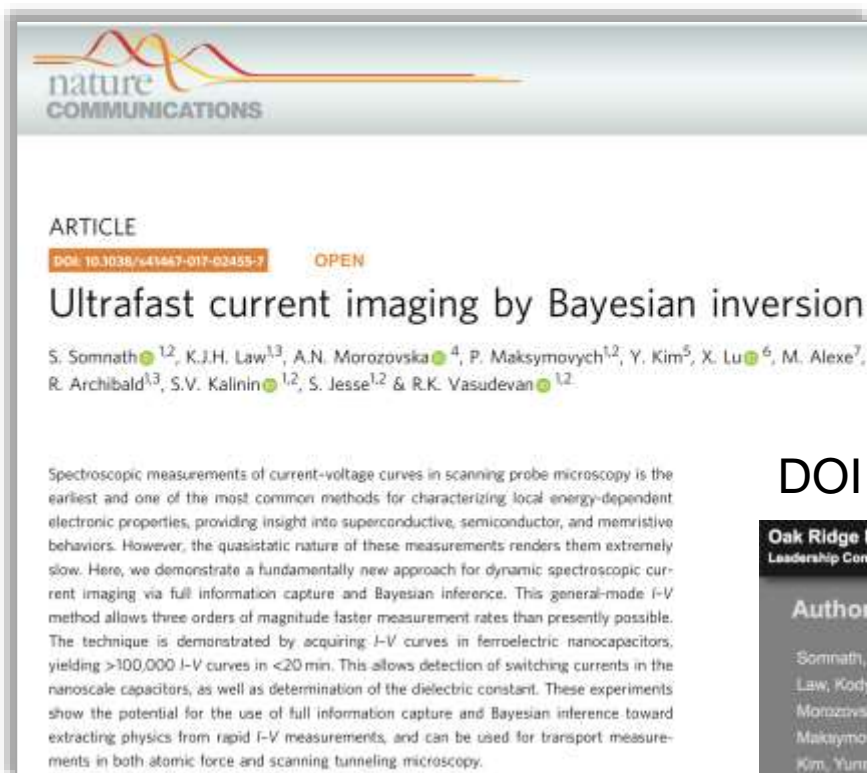
Will remain on off-line desktops

In the process of switching to computations on clusters

Truly Achieving Open Science, Reproducibility

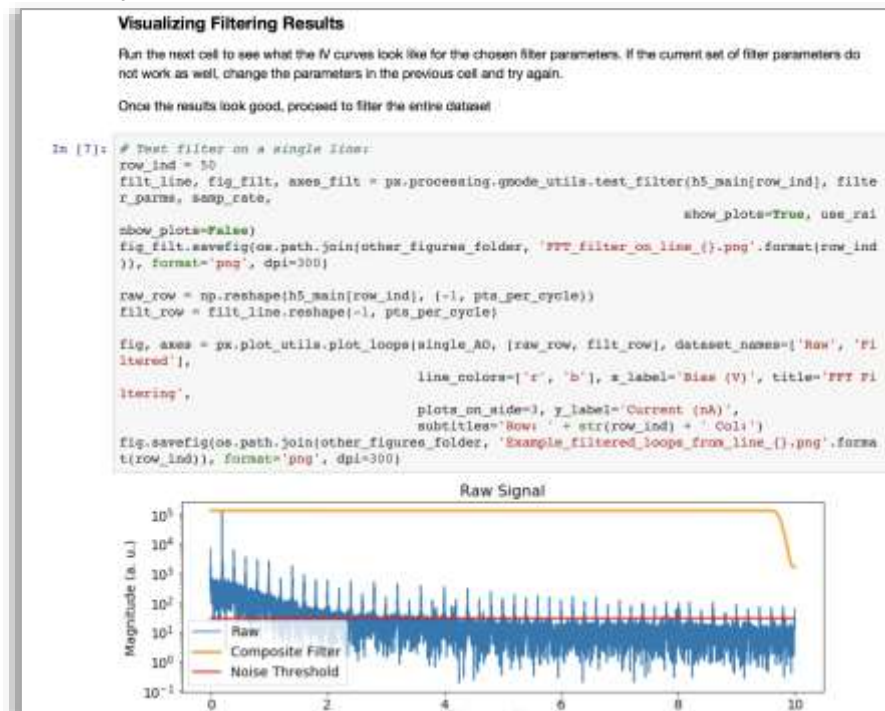
Aim – ALL scientific journal papers accompanied with:

- Jupyter notebook that shows all analysis (raw data → figures).
- Data with DOI number



Suhas Somnath, somnaths@ornl.gov

Jupyter notebook associated with paper



DOI associated with data (raw → paper figures)

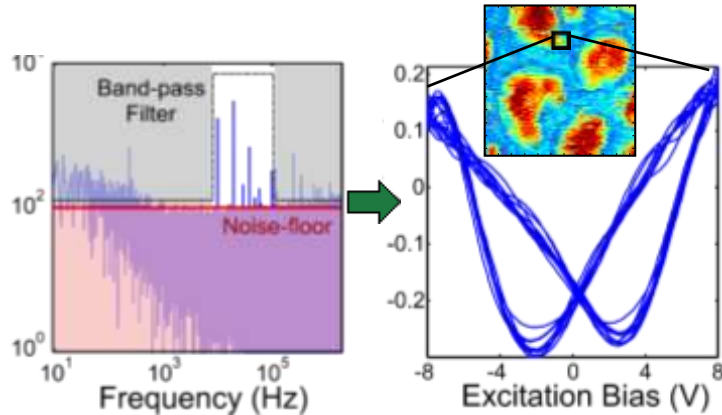
Oak Ridge National Laboratory Leadership Computing Facility 10.13139/OLCF/1410993 Download

Authors

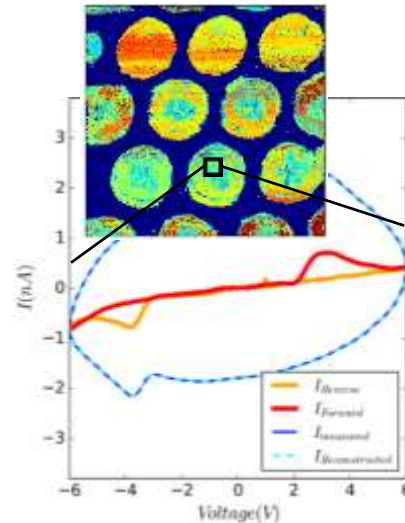
Somnath, Suhas	somnaths@ornl.gov
Law, Kody	lawkj@ornl.gov
Morozovska, Anna	anna.n.morozovska@gmail.com
Maksymovych, Petro	maksymovychp@ornl.gov
Kim, Yurmeek	yloms43@gmail.com
Lu, Xieoli	xliu@xidian.edu.cn
Alexe, Marin	M.Alexe@warwick.ac.uk

Pycroscopy - Scientific Advancements

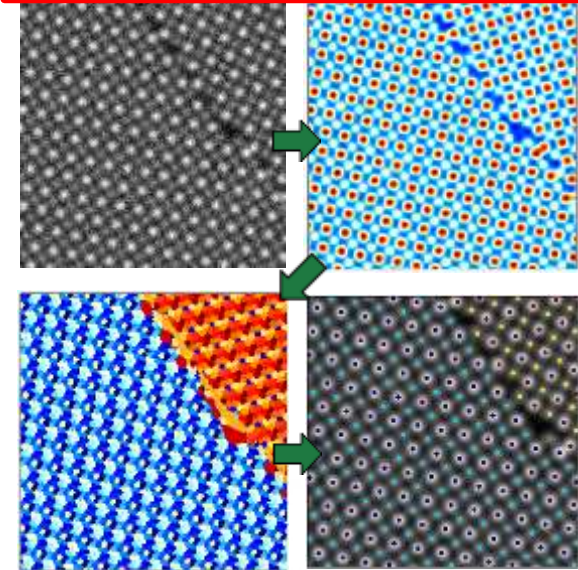
3,500x faster imaging via adaptive signal filtering, linear unmixing of signals



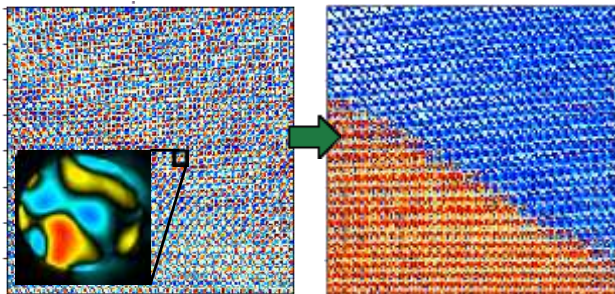
200x faster spectroscopy via Bayesian inference



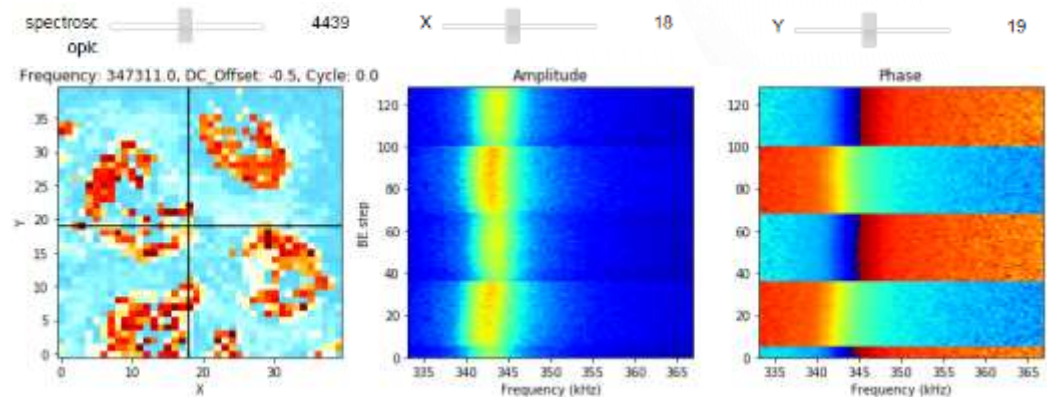
Separating uncorrelated data from correlated data to clean images



Identifying invisible patterns using multivariate analysis



Simplified navigation multidimensional data - users



Pycroscopy Progression

Scaling up Computing:



2016

Single core



2017

Multi-core
Single CPU



2018

Multi-core
Multi CPU



2019

JupyterHub
On HPC

Emphasis always on ease-of-development instead of raw performance

New Scientific Domains:

Atomic
Force
Microscopy

2016

+ Scanning
Transmission
Electron
Microscopy

2017

+
Mass
Spectrometry
+ Bio-
Chemistry

2018

+
Neutron
Science

2019

Thank you

Questions?