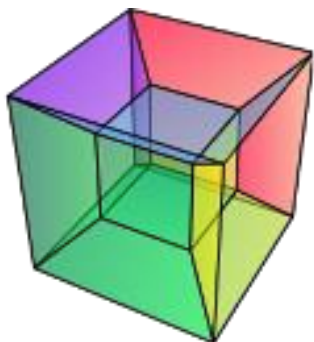


# Sustainable open-source software development with HyperSpy

Eric Prestat<sup>1,2</sup>

<sup>1</sup>School of Materials, University of Manchester, Manchester, UK

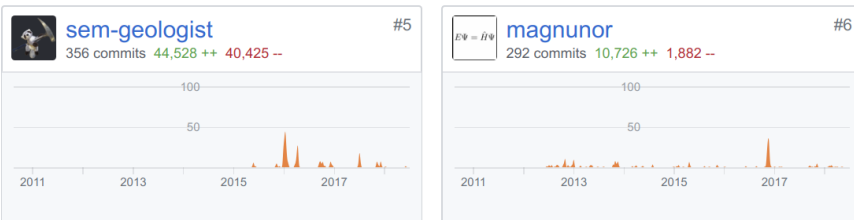
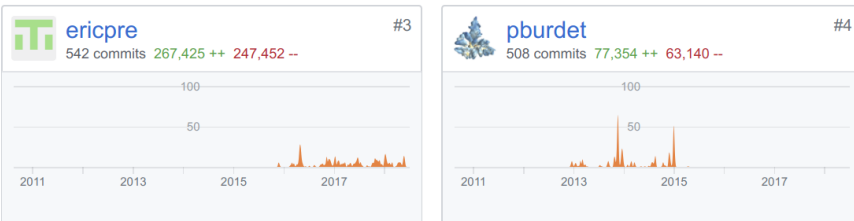
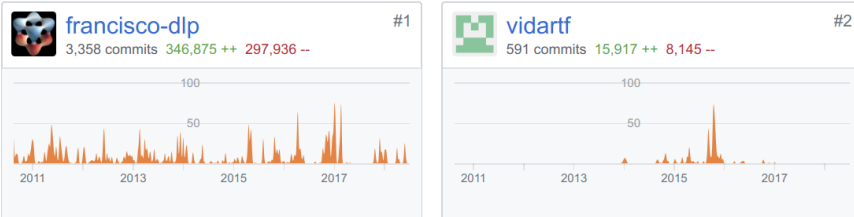
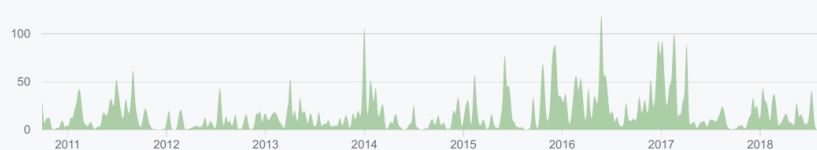
<sup>2</sup>SuperSTEM Laboratory, SciTech Daresbury Campus, Daresbury, UK



- What is HyperSpy?
- What HyperSpy has to offer?
- Where shall we go as an open source scientific community?

Contributions to RELEASE\_next\_minor, excluding merge commits

## Number of commits over time



Francisco de la Peña<sup>1,2,3,4</sup>, Tomas Ostasevicius<sup>1</sup>, Vidar Tonaas Fauske<sup>5</sup>, Pierre Burdet<sup>1</sup>, **Eric Prestat**<sup>6,7</sup>, Petras Jokubauskas<sup>8</sup>, Magnus Nord<sup>5,9,10</sup>, Mike Sarahan<sup>7</sup>, Katherine E. MacArthur<sup>11</sup>, Duncan N. Johnstone<sup>1</sup>, Joshua Taillon<sup>12</sup>, Alberto Eljarrat<sup>13,14</sup>, Vadim Migunov<sup>11</sup>, Jan Caron<sup>11</sup>, Tom Furnival<sup>1</sup>, Stefano Mazzucco<sup>2</sup>, Thomas Aarholt<sup>15,16</sup>, Michael Walls<sup>2</sup>, Florian Winkler<sup>11</sup>, Gaël Donval<sup>4,18</sup>, Robert McLeod<sup>19</sup>, Thomas Slater<sup>6</sup>

<sup>1</sup> Department of Materials Science and Metallurgy, University of Cambridge, Cambridge, U.K.

<sup>2</sup> Laboratoire de Physique des Solides, Paris-Sud University, Orsay, France

<sup>3</sup> Unité Matériaux et Transformations, University Lille1, Lille, France

<sup>4</sup> CEA Grenoble, Grenoble, France

<sup>5</sup> Department of Physics, NTNU, Trondheim, Norway

<sup>6</sup> School of Materials, The University of Manchester, Manchester, U.K.

<sup>7</sup> SuperSTEM, SciTech Daresbury Campus, Warrington, U.K.

<sup>8</sup> Institute of Geochemistry, Mineralogy and Petrology, University of Warsaw, Poland

<sup>9</sup> School of Physics and Astronomy, University of Glasgow, Glasgow, U.K.

<sup>10</sup> EMAT, University of Antwerp, Antwerp, Belgium

<sup>11</sup> Forschungszentrum Jülich GmbH, Jülich, Germany

<sup>12</sup> Material Measurement Laboratory, NIST, Gaithersburg, MD, U.S.A.

<sup>13</sup> Laboratory of Electron NanoScopies, Universitat de Barcelona, Barcelona, Spain

<sup>14</sup> Department of Physics, Humboldt University of Berlin, Newtonstr. 15, 12489 Berlin, Germany

<sup>15</sup> Department of Materials, University of Oxford, Oxford, United Kingdom

<sup>16</sup> Department of Physics, University of Oslo, Oslo, Norway

<sup>18</sup> Institut des Matériaux Jean Rouxel (IMN), Université de Nantes, Nantes, France

<sup>19</sup> Entropy Reduction Algorithmics, Victoria, Canada

+ many more: 35 contributors (~10 is a one off)



## 'Enhanced Data Generated by Electrons'

Following the very successful workshops at Lake Tahoe, Leukerbad, Port Ludlow, Guadeloupe, Grundlsee and Banff, the next international workshop on electron energy loss spectroscopy and imaging will be held in

**Sainte Maxime, France,**

- Francisco's talk on HyperSpy
  - Mostly about open source and scientific python community
  - NOT about the cool stuff that HyperSpy could already do at the time





## 'Enhanced Data Generated by Electrons'

Following the very successful workshops at Lake Tahoe, Leukerbad, Port Ludlow, Guadeloupe, Grundlsee and Banff, the next international workshop on electron energy loss spectroscopy and imaging will be held in

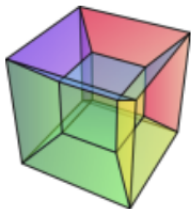
**Sainte Maxime, France,**

- Francisco's talk on HyperSpy
  - Mostly about open source and scientific python community
  - NOT about the cool stuff that HyperSpy could already do at the time



What I did not understand at the time:

- Open source is not only about opening the code but also about community
- *Small* paradigm shift required



# HyperSpy

## multi-dimensional data analysis

---

[Home](#) · [Download](#) · [Documentation](#) · [Demos](#) · [News](#) · [Support](#) · [Citing](#) · [Credits](#)

---

## HyperSpy: multi-dimensional data analysis toolbox

HyperSpy is an open source Python library which provides tools to facilitate the interactive data analysis of multi-dimensional datasets that can be described as multi-dimensional arrays of a given signal (e.g. a 2D array of spectra a.k.a spectrum image).

HyperSpy aims at making it easy and natural to apply analytical procedures that operate on an individual signal to multi-dimensional arrays, as well as providing easy access to analytical tools that exploit the multi-dimensionality of the dataset.

Its modular structure makes it easy to add features to analyze different kinds of signals.

### Highlights

---

- Two families of named and scaled axes: *signal* and *navigation*.
- Visualization tools for multi-dimensional spectra and images.
- Easy access multi-dimensional curve fitting and blind source separation.
- Built on top of NumPy, SciPy, matplotlib and scikit-learn.
- Modular design for easy extensibility.

Originally develop by electron microscopists but designed to analyse any kind of multi-dimensional data



# HyperSpy's core concept

- Signal object
  - Signal dimension
  - Navigation dimension
  - Slicing
  - Generic:
    - Signal1D, Signal2D
  - Specialised:
    - EDS, EELS, etc.
- Syntax
  - Consistent API
  - Try to be not too verbose for the user

```
data = (np.arange(5*10*20*50).reshape(5, 10, 20, 50))
s = hs.signals.Signal1D(data)
s
```

```
<Signal1D, title: , dimensions: (20, 10, 5|50)>
```

```
s.T
```

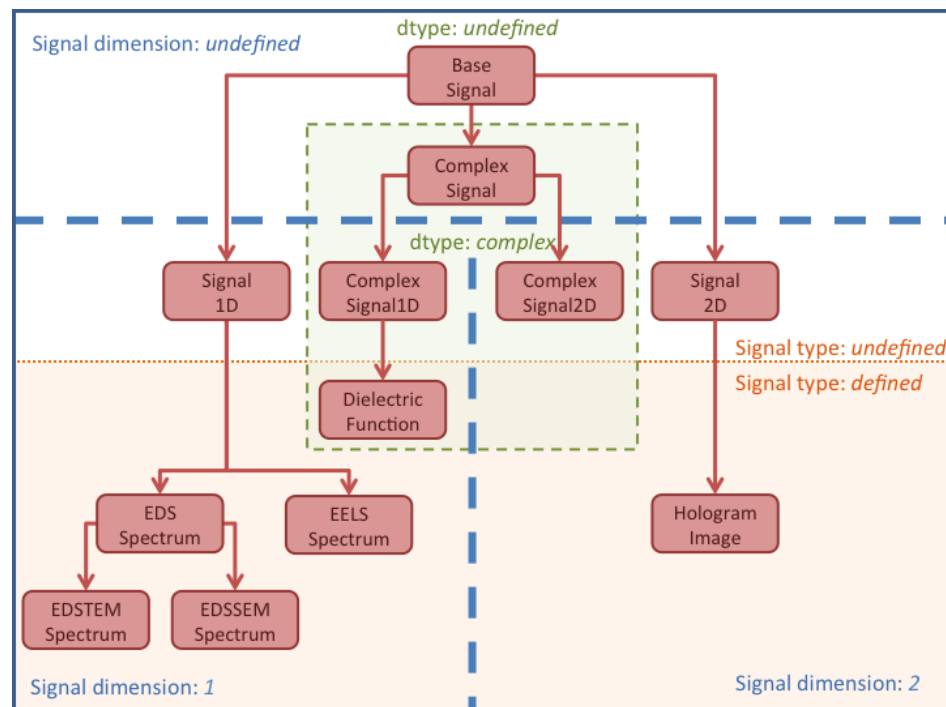
```
<BaseSignal, title: , dimensions: (50|20, 10, 5)>
```

```
s.inav[1:15, ...]
```

```
<Signal1D, title: , dimensions: (14, 10, 5|50)>
```

```
s.isig[:25]
```

```
<Signal1D, title: , dimensions: (20, 10, 5|25)>
```

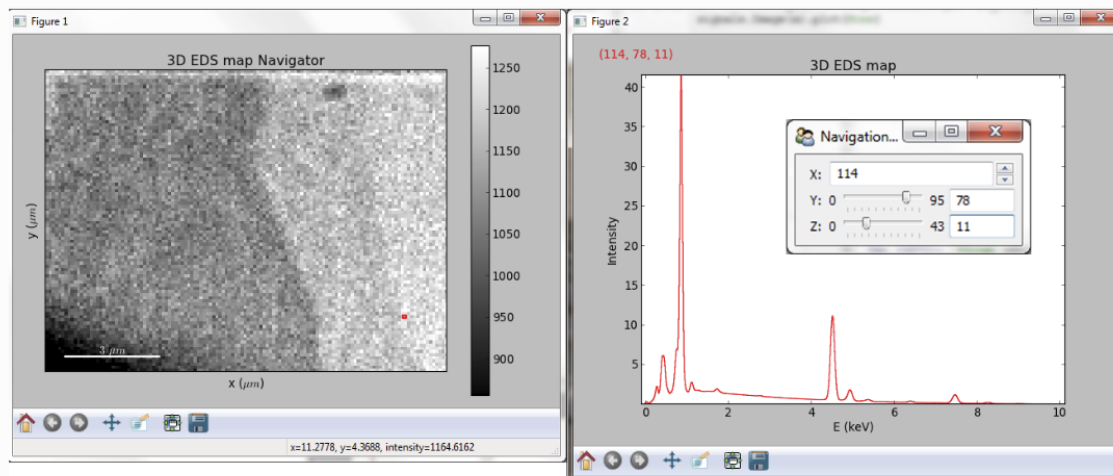


## HyperSpy main features

- Input/Output support almost all format in electron microscopy
- Data visualisation and interactivity rely on matplotlib

matplotlib

```
>>> s = hs.load('Ni_superalloy_0*.rpl', stack=True).as_signal1D(0)
>>> s.plot()
```



Visualisation of a 3D spectrum image with sliders.

### Supported file formats

Format	Read	Write	lazy
Gatan's dm3	Yes	No	Yes
Gatan's dm4	Yes	No	Yes
FEI's emi and ser	Yes	No	Yes
HDF5	Yes	Yes	Yes
Image: jpg	Yes	Yes	Yes
TIFF	Yes	Yes	Yes
MRC	Yes	No	Yes
MRCZ	Yes	Yes	Yes
EMSA/MSA	Yes	Yes	No
NetCDF	Yes	No	No
Ripple	Yes	Yes	Yes
SEMPER unf	Yes	Yes	Yes
Blockfile	Yes	Yes	Yes
DENS heater log	Yes	No	No
Bruker's bcf	Yes	No	Yes
Bruker's spx	Yes	No	No
EMD (NCEM)	Yes	Yes	Yes
EMD (FEI)	Yes	No	Yes
Protochips log	Yes	No	No
EDAX .spc and .spd	Yes	No	Yes

# HyperSpy main features : machine learning

- SVD, BSS, NMF
  - from scikit-learn



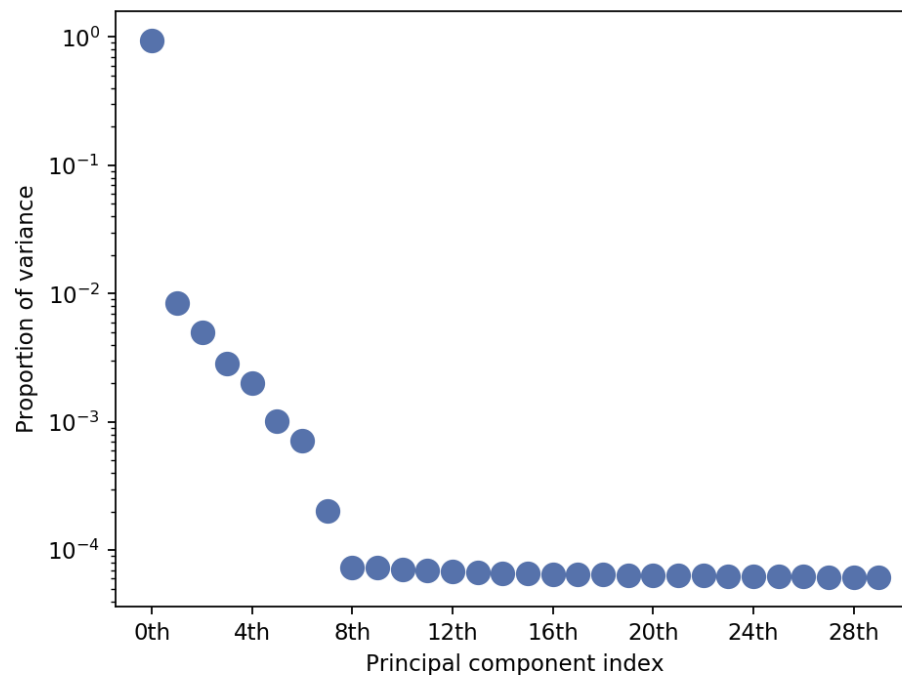
- Robust PCA to deal with outliers
- “online” algorithm for big data

```
s1.decomposition(True)
```

```
s1.plot_explained_variance_ratio()
```

Figure 1

EELS Spectrum Image (low-loss)  
PCA Scree Plot



<matplotlib.axes.\_subplots.AxesSubplot at 0x7fb26d832ef0>

```
s1_denoised = s1.get_decomposition_model(9)
```



## HyperSpy main features: multidimensional curve-fitting

- Any model can be build (1D, 2D)
- Model interactivity

```
m = s.create_model()

m.append(hs.model.components.Gaussian())
m.append(hs.model.components.Polynomial(5))
m.append(hs.model.components.EELSCLEdge('C_K'))
m.append(hs.model.components.ScalableFixedPattern(s.deepcopy(

m.notebook_interaction()
```

background

Gaussian

Polynomial

C\_K

ScalableFixedPattern

☒ active

☒ interpolate

min -9

yscale

1.00

max 11

min -9

xscale

1.00

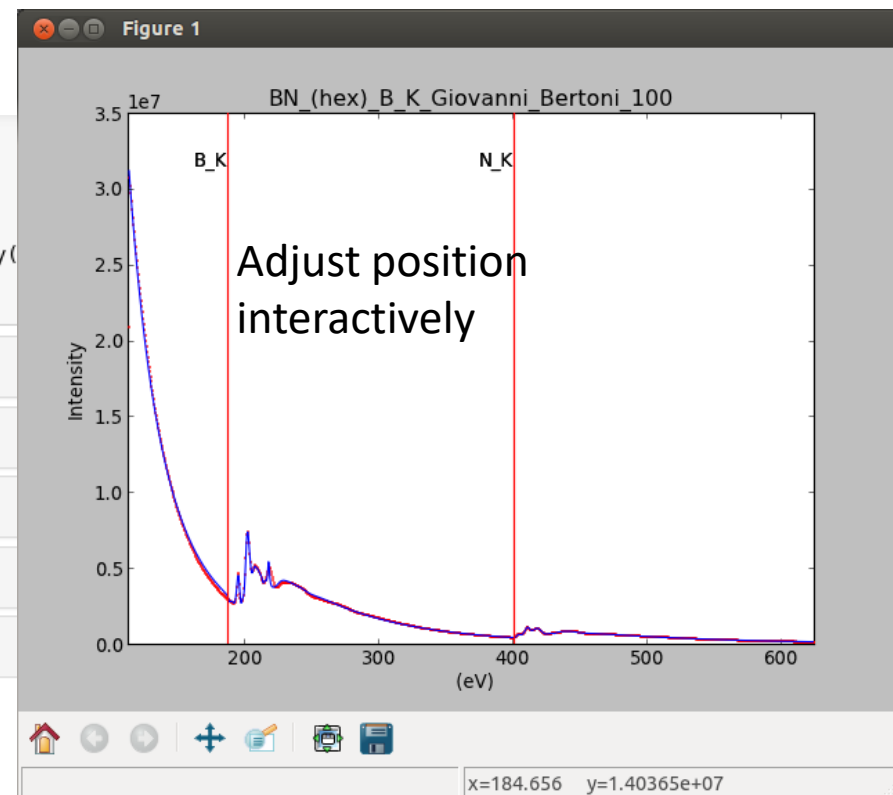
max 11

min 0

shift

0.00

max 499



HyperSpy main features: multidimensional curve-fitting

- Any model can be build (1D, 2D)
- Model interactivity
- Write your own expression



Create a 2D component from a string, which is turned into a function with SymPy

```
>>> g = hs.model.components2D.Expression(
... "k * exp(-((x-x0)**2 / (2 * sx ** 2) + (y-y0)**2 / (2 * sy ** 2)))",
... "Gaussian2d", add_rotation=True, position=("x0", "y0"),
... module="numpy", )
```

- Smart Adaptive Multi-dimensional Fitting (SAMFire)

<https://doi.org/10.1002/9783527808465.EMC2016.6233>



<https://youtu.be/kVlf3bMZcsc>

## HyperSpy main features: out-of-core computation

- Perform operation “lazily”
- Delay the execution until the result is requested and then perform them in a blocked fashion

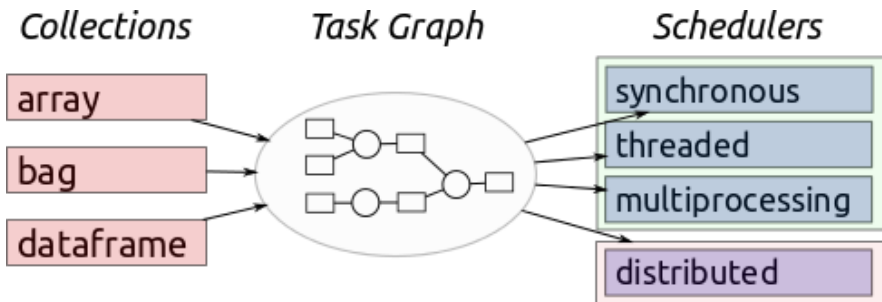


```
In [1]: import dask.array as da
```

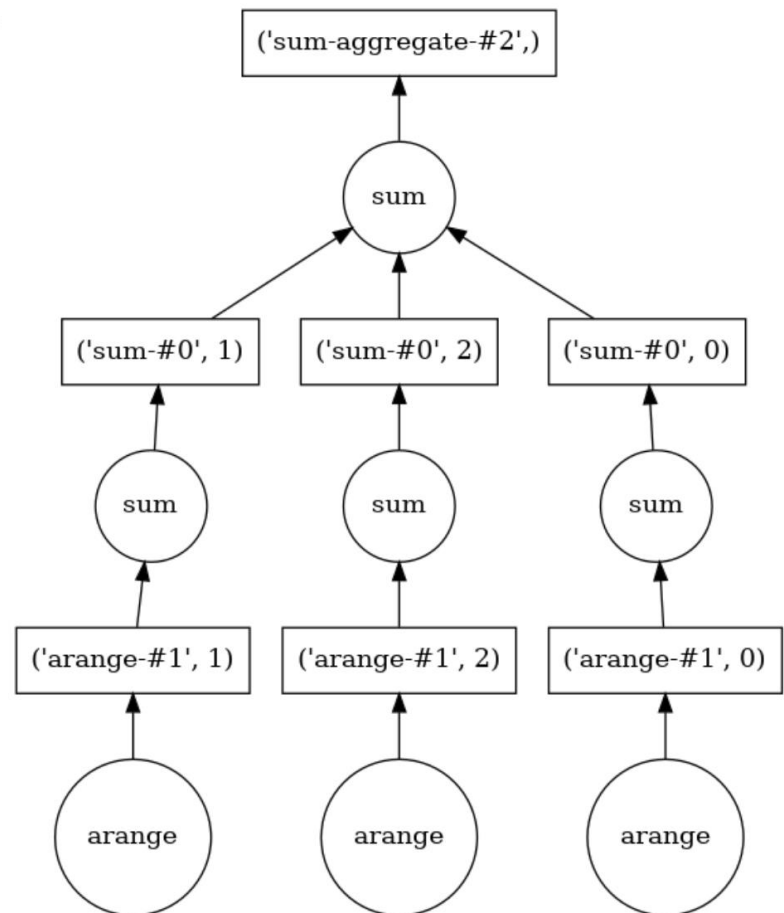
```
In [2]: a = da.arange(30, chunks=(10,)).sum()
```

```
In [3]: a.visualize()
```

Out[3]:



<http://dask.pydata.org>



## Performance optimisation

- Numba
  - JIT compilation



```
@jit_ifnumba()
def _fill_array_with_stream_sum_frames(spectrum_image, stream,
                                      first_frame, last_frame, rebin_energy=1):
    # jit speeds up this function by a factor of ~ 30
    navigation_index = 0
    frame_number = 0
    shape = spectrum_image.shape
    for count_channel in np.nditer(stream):
        # when we reach the end of the frame, reset the navigation index to 0
        if navigation_index == (shape[0] * shape[1]):
            navigation_index = 0
            frame_number += 1
            # break the for loop when we reach the last frame we want to read
            if frame_number == last_frame:
                break
        # if different of '65535', add a count to the corresponding channel
        if count_channel != 65535:
            if first_frame <= frame_number:
                spectrum_image[navigation_index // shape[1],
                              navigation_index % shape[1],
                              count_channel // rebin_energy] += 1
        else:
            navigation_index += 1
```

Particularly interesting when vectorisation is not possible

- Numexpr (CPU cache optimisation for numpy)

```
# algebraic expression:
z = 2*y + 4*x

# transcendental expression:
z = np.sin(x)**2 + np.cos(y)**2
```

NumPy version: 1.15.1

Time for an algebraic expression: 0.343 s / 3.263 GB/s

Time for a transcendental expression: 0.597 s / 1.873 GB/s

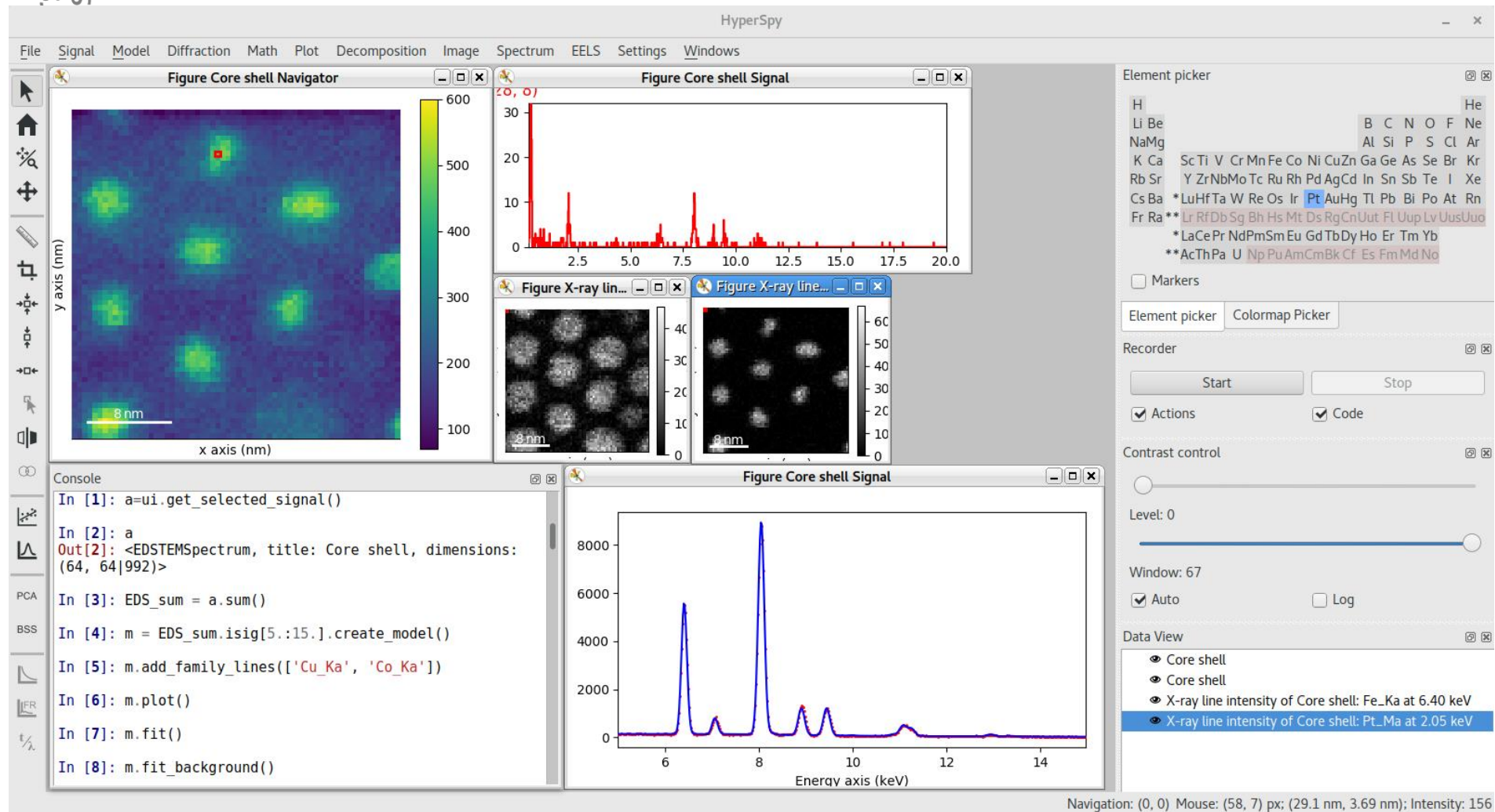
Numexpr version: 2.6.8. Using MKL: True

Time for an algebraic expression: 0.061 s / 18.457 GB/s

Time for a transcendental expression: 0.078 s / 14.399 GB/s

# User interface: HyperSpyUI

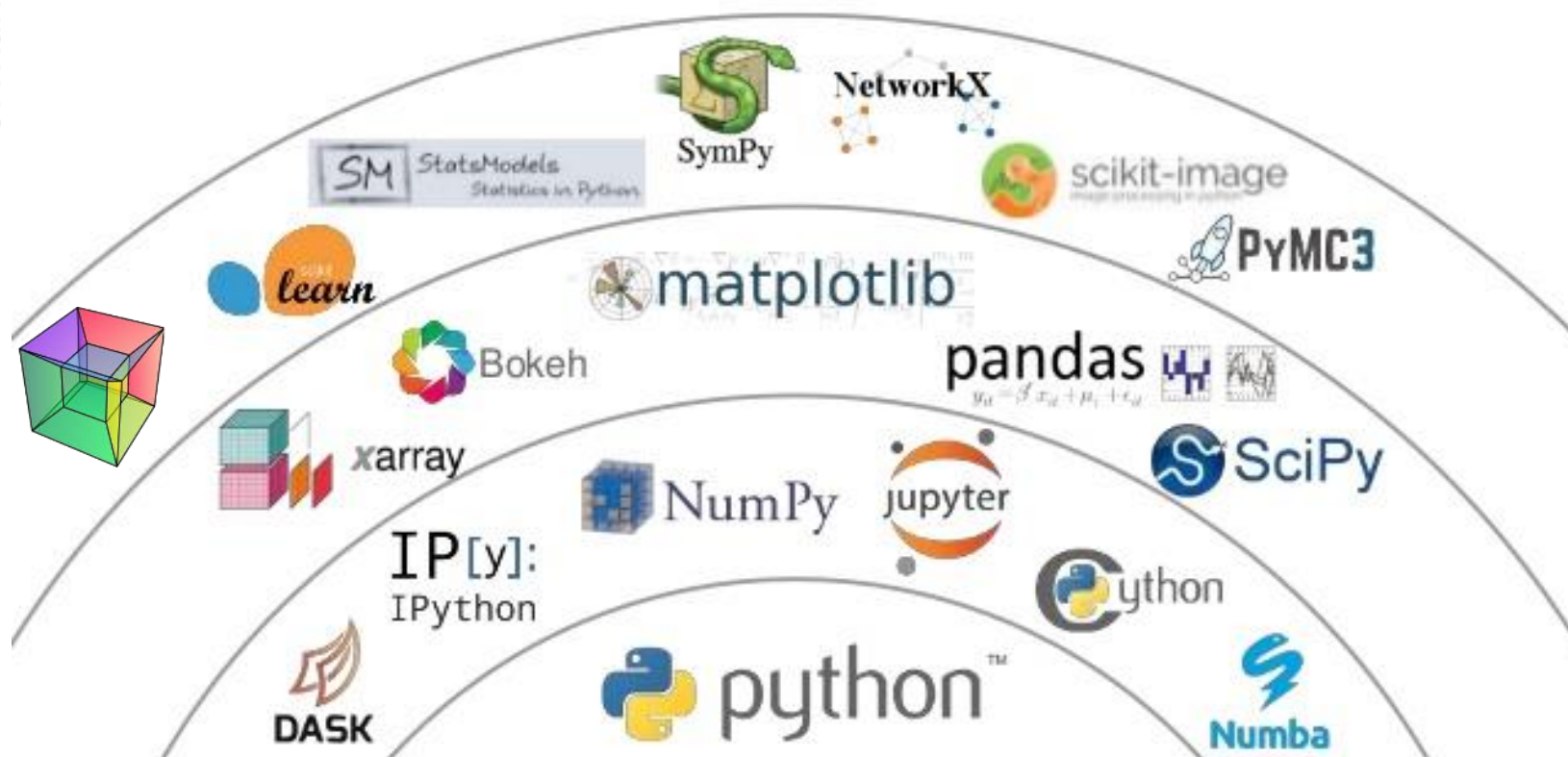
city  
ster



# Python scientific stack: each library find its natural position

## Python's Scientific Stack

Jake Vanderplas PyCon 2017 Keynote



- Strength of the major libraries for the user:
  - work nicely with each other
  - Similarities in the syntax



## Current HyperSpy ecosystem

HyperSpyUI

HyperSpy-ipywidget

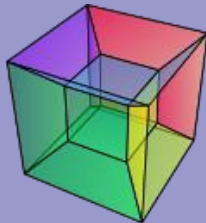
HyperSpy-traitsui

4D STEM Magnetic measurement

**pixSTEM**

Differential contrast imaging

EELS holography  
Big data EDS  
Multi-dimensional IO  
Visualisation and interactivity Machine learning  
Curve-fitting



Strain mapping Orientation mapping

Electron cristallography **pyXem** Scanning Electron diffraction

Atomic position fitting

Quantitative STEM **Atomap** Strain mapping

Future HyperSpy ecosystem and other relevant libraries

EELS

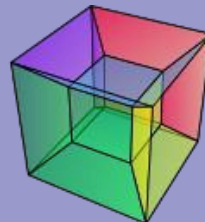
IO

EDS

Holography

HyperSpyUI

Multi-dimensional  
Visualisation and interactivity  
Big data  
Machine learning



Nion Swift

GMS 3

HyperSpy-ipywidget

LiberTEM

tomviz

HyperSpy-traitsui

pycroscopy

4D STEM Magnetic measurement

pixSTEM

Differential contrast imaging

Strain mapping

Orientation mapping

Electron cristallography

pyXem

Scanning Electron diffraction

Atomic position fitting

Atomap

Quantitative STEM

Strain mapping

Other extension

## HyperSpy and its community

- HyperSpy managed to built a *distributed* community of users/contributor
  - Led by its own contributors
  - Does not depend of any academic group/institution
  - Decision based on contributor consensus
- Motivate (at least not discourage) users
  - Pay attention to users feedback
  - From user to contributor: make the learning curve easier
- What are contributors doing?
  - Contribution to user guide, tutorials and online discussion
  - Code writing and/or code review

This is one way to make a library sustainable

As individual is it worth contributing to HyperSpy?

- HyperSpy acknowledgement through zenodo DOI
  - One DOI for each release (important for reproducibility)
  - New contributors will get acknowledged
  - Google scholar not *yet* fully compatible with zenodo DOI (work in progress)
- Contribution to HyperSpy (or any other library) can be useful for career development
  - Github profil can be used as linkedin, etc.
  - Recognition by the microscopy community
- As PI/group leader
  - PI are not acknowledged but they can benefit a lot from the expertise gained through HyperSpy
  - Fairly useful “training” for post-doc/student

## HyperSpy from the developer perspective

- Bundle installer for windows
  - Based on winpython (fairly well tested by the community)
  - Build on appveyor (one click on github to make a release)
  - Independent from HyperSpy release
  - Easy installer (NSIS 2)
  - Can package any python library easily (with C code)
- Bundle installer for mac and linux
  - Use the same tool that build anaconda and miniconda
  - One current issue to fix... coming anytime soon?
- Offer a toolbox
  - Generic multidimensional processing
  - Fitting, big data, IO, etc
  - `map` function (apply any processing on the signal dimension)

Concluding remarks: what would be ideal for the user?

- Good intercompatibility between softwares/libraries; integrated workflow
  - Similar syntax, data format, etc.
  - Users should be able to switch easily from one library to another
- Run the same code in different environment
  - Jupyter notebook or qtconsole
  - HyperSpyUI, Nion Swift, GMS, etc.
  - Supercomputer (High-performance computing)
- Code implemented at the right position of the python stack
  - Usually implementation improve significantly during review
  - Code exposed to larger users base -> better feedback

Thanks for your attention