# intake - taking the pain out of data access

Mickaël Lalande (Institut des Géosciences de l'Environnement)

# What is Pangeo?

*"A community platform for Big Data geoscience"*

- Open Community

- Open Source Software
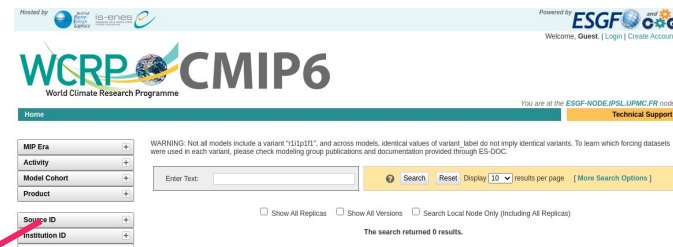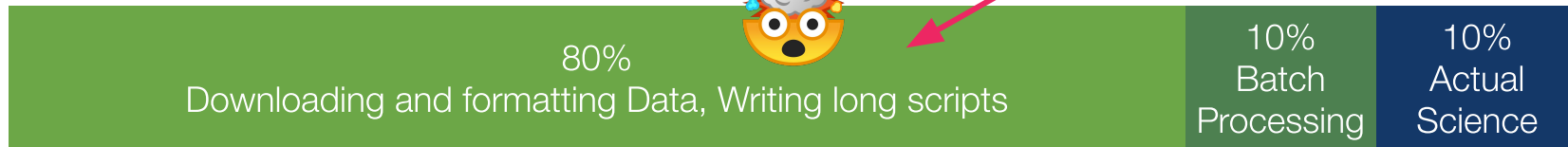
- Open Source Infrastructure

(Aurélie)

# What impacts the velocity of science?
## *Data, Software and Computation*

- Data: time to find, access, clean & format for analysis

- Software: easily available and combinable

- Computation: access and resources

Traditional Analysis Workflow

| 80% Downloading and formatting Data, Writing long scripts | 10% Batch Processing | 10% Actual Science |

Pangeo Analysis Workflow

| 5% Data Preparation | 5% Batch Processing | INTAKE | 90% Actual Science | CliMAF Sharing – Simplifying – Optimizing |

Adopted from https://speakerdeck.com/cgentemann/empowering-transformational-science?slide=4
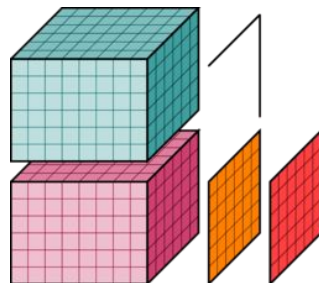
# intake-esm

- Developed by Anderson Banihirwe at NCAR (**@andersy005**)

- Search and load ESM output

- Catalog builds easily from CMORized output

- Query in `pandas.DataFrames`

- Share/archive data sources used for your particular analysis

- Load data with `dask` into `xarray`

# xarray

- Analysis of multi-dimensional data

- Self-describing data

- Efficient: based on `numpy` and `dask`

- Simple: API inspired by `numpy` and `pandas`

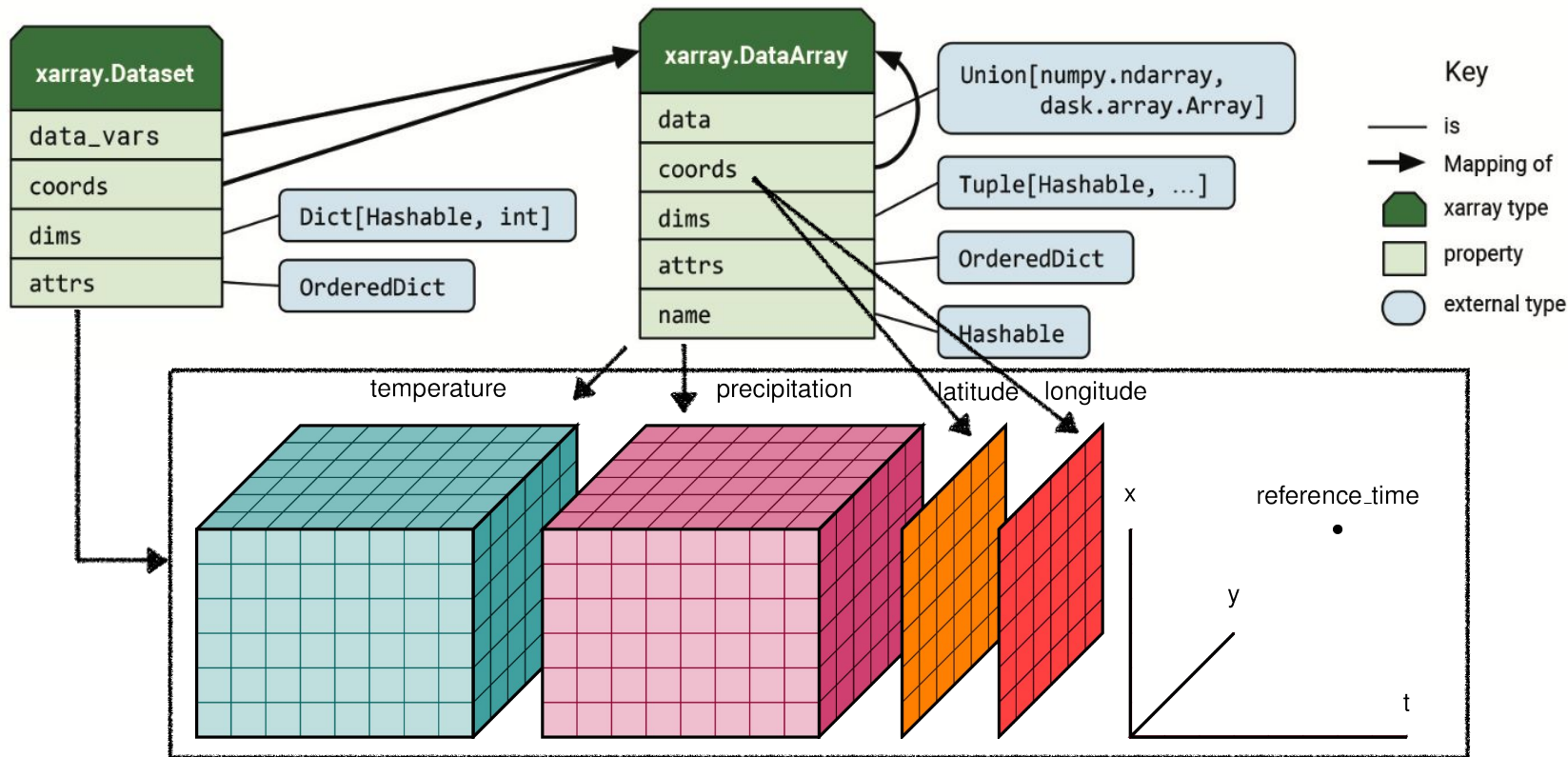- Stephan Hoyer and Joe Hamman ([2017](#)) "Xarray: N-D Labeled Arrays and Datasets in Python"

https://github.com/mickaellalande/MC-Toolkit/tree/master/conda_environment_x
array_xesmf_proplot/xarray

# xarray **data types**
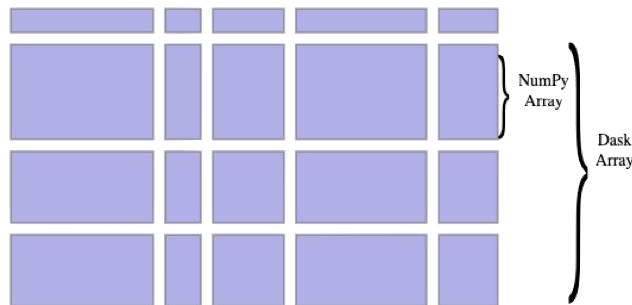
# Extensions to xarray

- scipy : (nearly) all functions callable with `xr.apply_ufunc`
- `dask_jobqueue` : parallelise `dask` across nodes
- `xskillscore :` verification metrics
- `cartopy` : projections of maps
- `geoviews` : dynamic visualisation of geo data
- `regionmask` : spatial aggregation based on shapefiles
- x`esmf` : regridding
- `xgcm` : grid aware operations
- `cmip6_preprocessing` : data cleaning for CMIP6 output
- `climpred` : verification of multi-dim ensemble forecasts
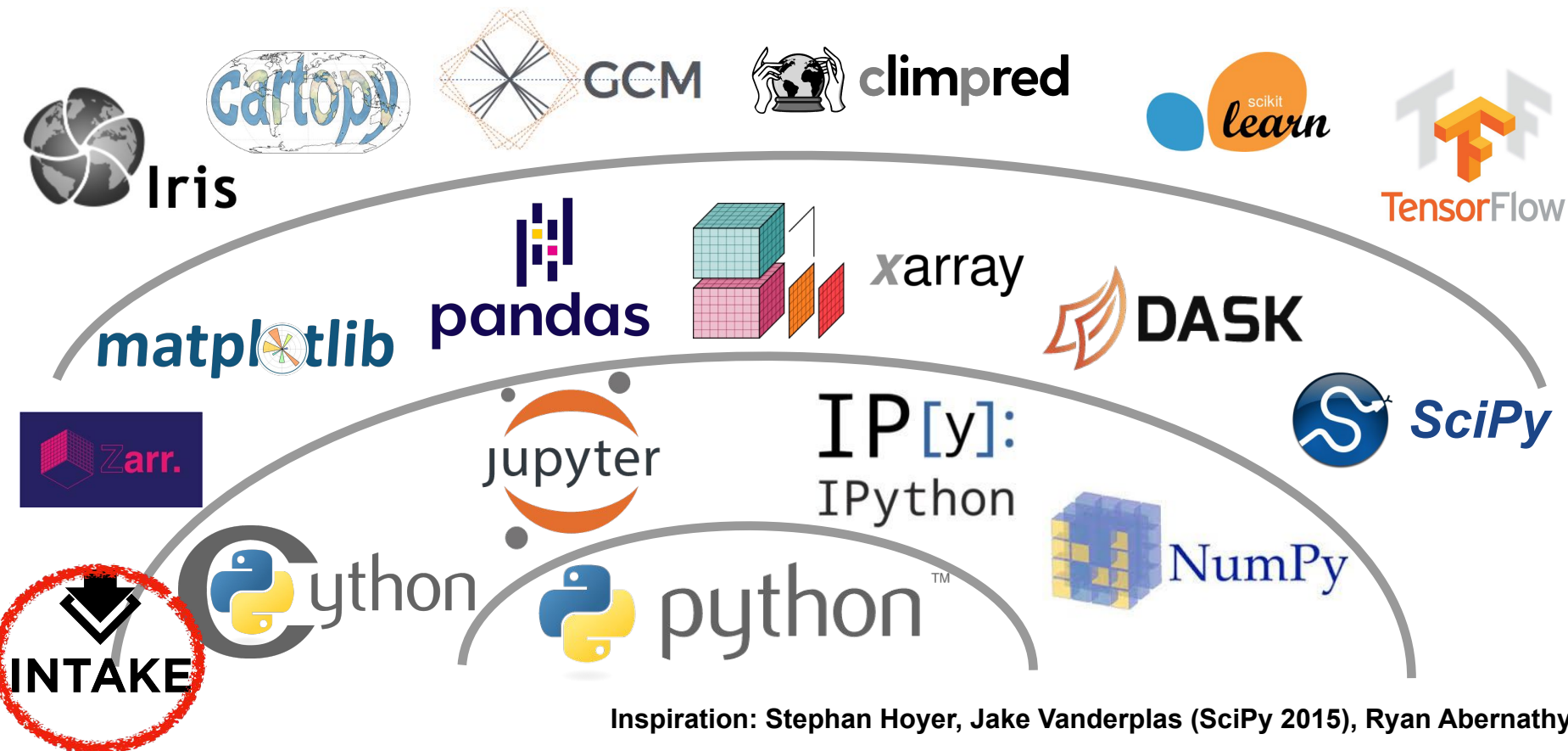- `intake-xarray` : intake for netcdf files
  ... http://xarray.pydata.org/en/stable/related-projects.html

# dask

- Dynamic task scheduling

- Builds upon **multiprocessing**, **threading** and **concurrent**

- out-of-memory computation via chunking

- Scales from laptop to supercomputer

- Intuitive (known) API from **pandas** and **numpy**

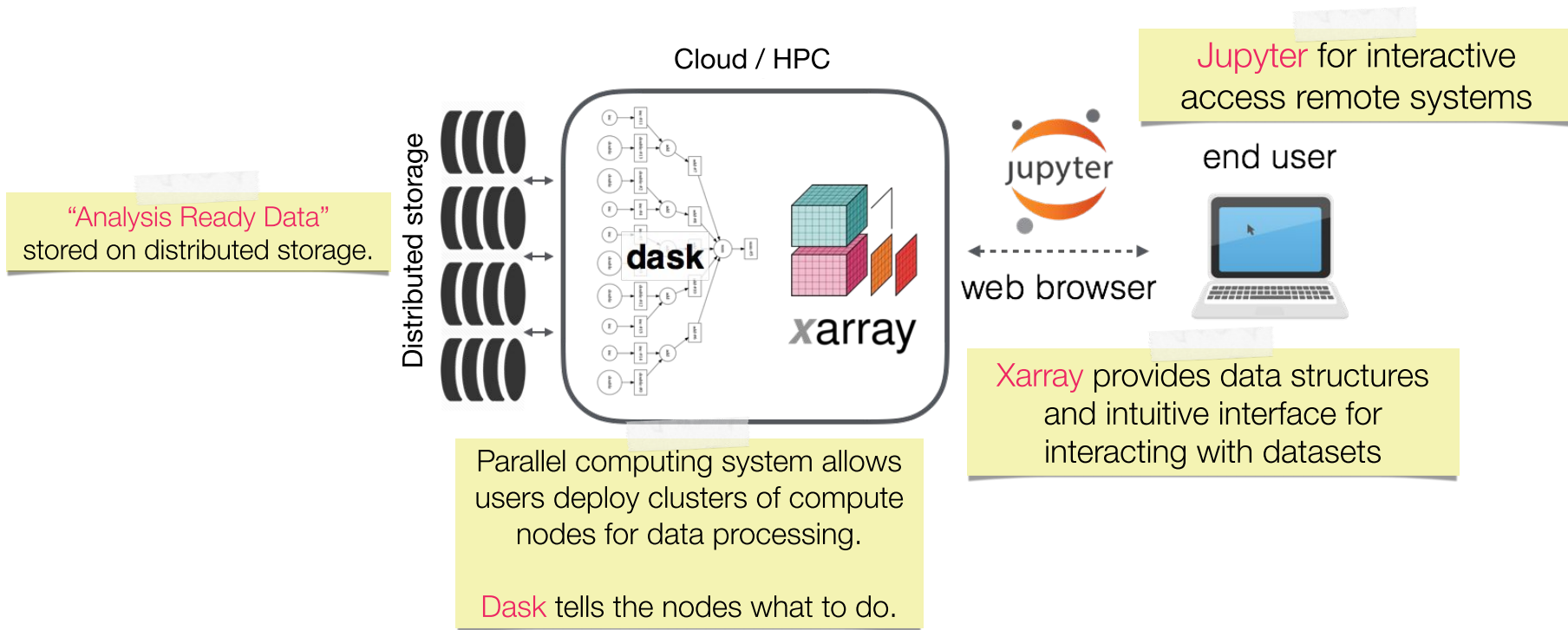- Matthew Rocklin ([2015](#): "Dask: Parallel Computation with Blocked Algorithms and Task Scheduling"
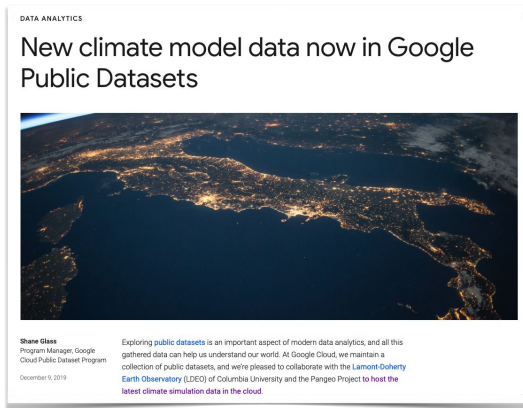
# Pangeo Software Ecosystem



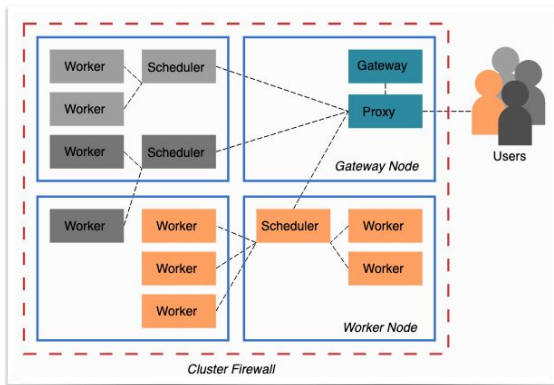Inspiration: Stephan Hoyer, Jake Vanderplas (SciPy 2015), Ryan Abernathy

# HPC Architecture



Cloud / HPC

Distributed storage

Jupyter for interactive access remote systems

end user

web browser

"Analysis Ready Data" stored on distributed storage.

Parallel computing system allows users deploy clusters of compute nodes for data processing.

Dask tells the nodes what to do.

Xarray provides data structures and intuitive interface for interacting with datasets

# Pangeo in the cloud

- Server-side computing

➡️ Play with cloud data yourself: launch binder

- Science in a GitHub repo:
  ‣ http://gallery.pangeo.io/
  ‣ Data in the cloud
  ‣ reproducible with binder

# GMST historical+obs with intake

**INTAKE**

- Demo in Jupyter

- Pangeo-binder:
  https://github.com/mickaellalande/intake_CMIP6/tree/pangeo-notebook

- Intake experimental on CICLAD!
  (voir avec Guillaume Levasseur)

# References

- Papers:

  - ‣ Rocklin, M. (2015). Dask: Parallel Computation with Blocked algorithms and Task Scheduling. 126–132. doi: [10.25080/Majora-7b98e3ed-013](10.25080/Majora-7b98e3ed-013)

  - ‣ Hoyer, S., & Hamman, J. (2017). xarray: N-D labeled Arrays and Datasets in Python. Journal of Open Research Software, 5(1). doi: [10/gdqdmw](10/gdqdmw)

  - ‣ Emanuel, K. (2020). The Relevance of Theory for Contemporary Research in Atmospheres, Oceans, and Climate. *AGU Advances*, *1*(2), e2019AV000129. doi: [10/gg3dzt](10/gg3dzt)

  - ‣ [https://authorea.com/users/372628/articles/490577-cloud-native-repositories-for-big-scientific-data](https://authorea.com/users/372628/articles/490577-cloud-native-repositories-for-big-scientific-data)

- Pictures:

  - ‣ xarray website, dask website, MPIM, DKRZ, pangeo

- Tutorials:

  - ‣ xarray: https://xarray-contrib.github.io/xarray-tutorial/scipy-tutorial/00_overview.html

  - ‣ dask: https://tutorial.dask.org/03_array.html

  - ‣ pangeo: [http://gallery.pangeo.io/](http://gallery.pangeo.io/)

- Similar talks: Empowering Transformational Science - [https://speakerdeck.com/cgentemann/empowering-transformational-science?slide=19](https://speakerdeck.com/cgentemann/empowering-transformational-science?slide=19)