



# **CCIN2P3-CNES : Calcul**

## Des développements

- ❖ Session visu (VNC, Fluxbox) + IDE
  - Quelques soucis de perfs en visu 3D
- ❖ Environnement logiciel et compilateur complet (modules lmod)
  - Gestion des piles logicielles faites manuellement, cronophage
- ❖ Jobs interactifs
- ❖ Jupyter notebook (voir après)

## Du HTC

- ❖ Couche orchestration simple (script batch ou Python)
- ❖ Job arrays souvent, mpi4py parfois
- ❖ Parfois agressif envers le Scheduler PBS. On oriente vers GNU-parallel ou Dask
- ❖ Souvent intensif en IO :
  - Petits fichiers ou random access
  - Audit de code
  - Utilisation d'espace local
- ❖ Simulations paramétriques (dynamique du vol par exemple), Traitement d'image...

## Modélisation, CFD

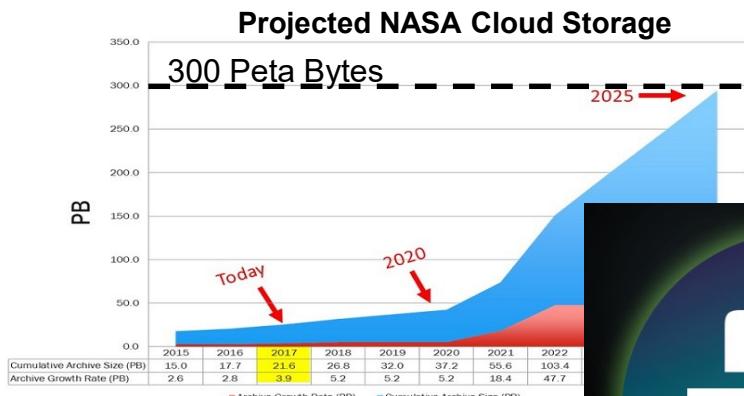
- ❖ Logiciels spécifiques, parfois sous licence
- ❖ MPI/OpenMP : Lancés via les logiciels, peu utilisés sinon.

## Analyse interactive, Big Data

- ❖ Voir Pangeo (après)
- ❖ Tests sur Spark également.

## Chaine de traitement

- ❖ Orchestration lourde :
  - logiciels maison :
    - gèrent jobs uniques ou job array
    - Saturent facilement le Scheduler PBS
    - Utilisent GPFS pour communiquer par fichier
  - Ou scripts complexes utilisant mpi4py, Dask.
- ❖ Généralement plutôt orienté HTC, mais avec des profils de jobs très variés.
- ❖ Campagne de retraitement, nœuds dédiés.



## Mission

To cultivate an ecosystem in which the next generation of open-source analysis tools for the geosciences can be developed, distributed, and sustained.



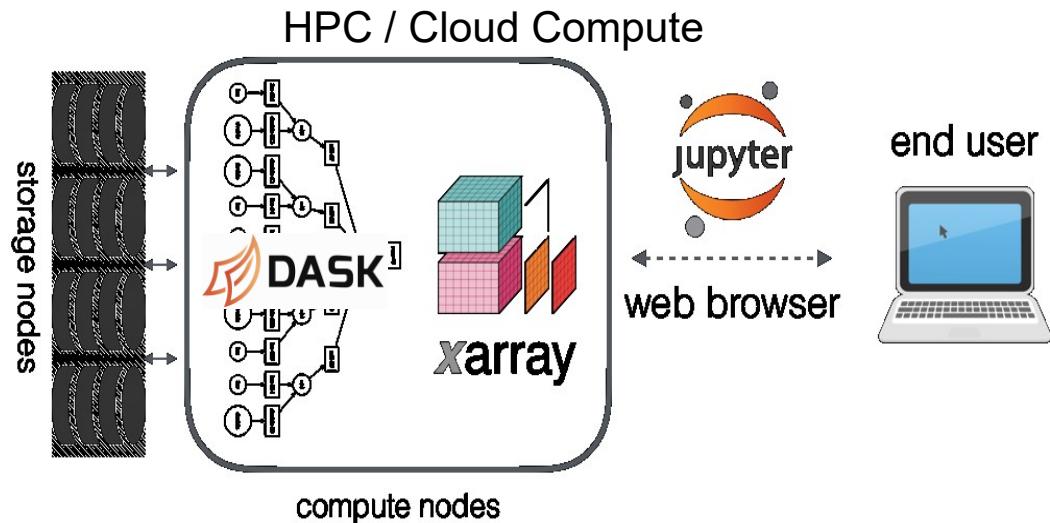
## Problems

- Data volume crisis in (geo)sciences
  - Software multiplication, non reproducibility
- Many copies of the same datasets  
Local vs HPC vs Cloud  
Technology gap: industry vs academia

## Goals/vision

- Cultivate an ecosystem to foster collaboration around the open source Scientific Python ecosystem:
- open and collaborative development
  - Welcoming and inclusive culture
  - Support the development with domain-specific (geo)science and transverse packages
  - Improve scalability of these tools to handle gigabytes to petabytescale datasets

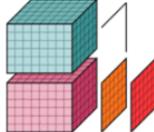
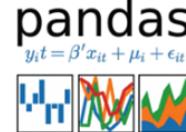
- Set of tools that will facilitate science at all scales
- Platform agnostic
- The core of the Pangeo ecosystem includes:
  - **Xarray** (data-model and toolkit for working with N-dimensional labeled arrays)
  - **Dask** (parallel computing)
  - **Jupyter** (interactive computing)
- Extensible: Series of 3rd party packages that build on top of core libraries
- Flexible: Individual components may be swapped in/out



## Examples of 3<sup>rd</sup> party packages in the Pangeo Ecosystem:

- Data discovery
- Regridding and GIS
- Vector calculus
- Signal processing
- Thermodynamics

# BUILD YOUR OWN PANGEO

Storage Formats			Cloud Optimized COG/Zarr/Parquet/etc.
ND-Arrays	 NumPy		More coming...
Data Models	 xarray		 $y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$
Processing Mode	 Interactive	Batch	
Compute Platform	HPC 	Cloud 	Local 



vs



Mature

Robust

JVM/Python

Query optimized

Collections &  
Dataframes

Python overhead

For big tabular data

Hadoop/Cloud/HPC

Less Mature

Pretty strong

Python only

Science optimized

Collections, DF,  
Arrays, Futures...

Python only

For science data

Hadoop/Cloud/HPC

I don't know much  
about Array DBs!!!

VS

databases  
(Rasdaman,  
SciDB...)

Laptop to cluster

Serverless

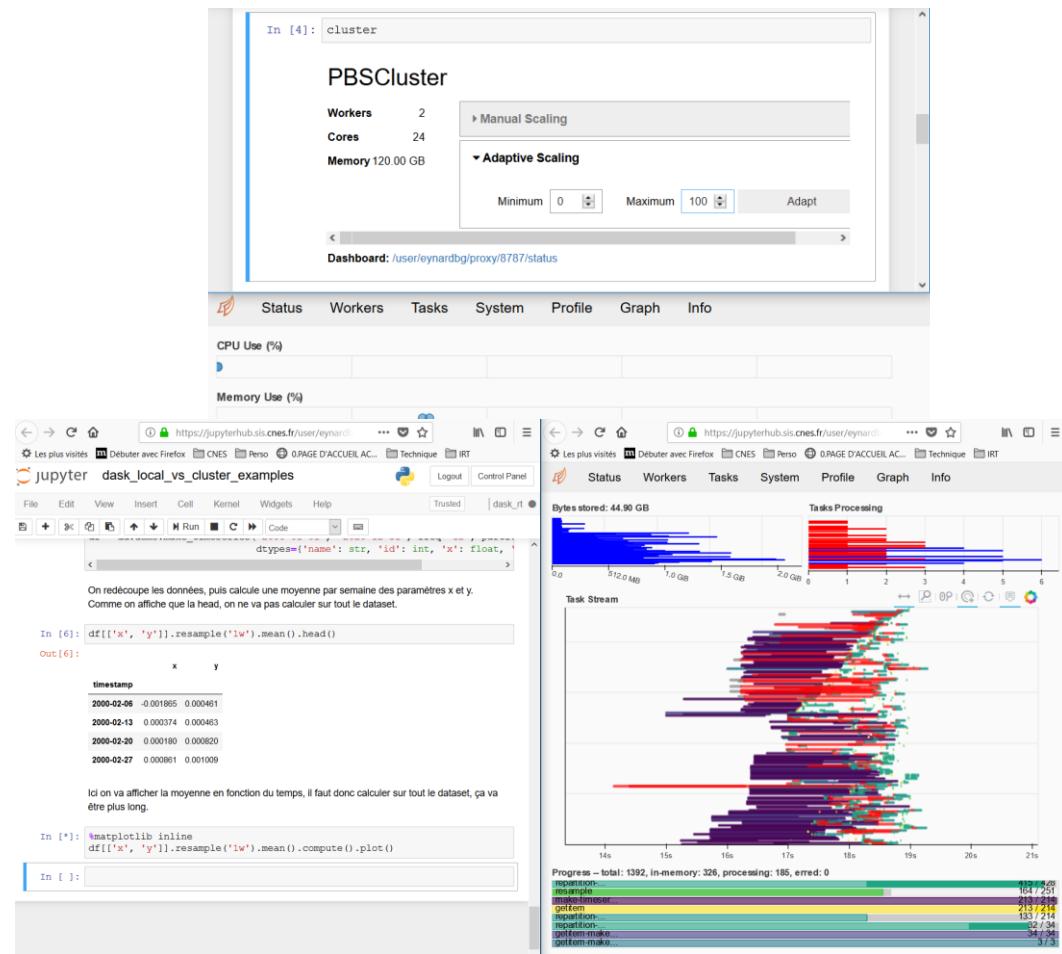
NetCDF/TIFF no ingestion

Scales with Dask

Python only

Can build array db with Pangeo  
(Open data Cube)

- JupyterHub and notebooks for interactive computing
  - Hub on a VM with qsub access
  - Batchspawner, Wrapspawner
- dask.distributed: parallel workers across many HPC nodes
- Xarray for computational toolkit and I/O
- New tool for deploying dask clusters on HPC: **dask-jobqueue**
  - Start a cluster from a notebook
  - Interactive (or not) distributed computing
  - Auto scaling capabilities



- ❖ Actuellement, VM basée sur VMWare, limitée à 4vcpus 16Go
  - Pas orientées traitement
  - Client du cluster HPC possible (annuaire utilisateur, client PBS, montage GPFS en NFS)
- ❖ Mais des réflexions sur Open Stack, Kubernetes. Une étude métier démarre cette année pour HPC 6è génération
- ❖ Calcul dans le cloud : des études, mais rien de réellement fait.