

# Computing services for climate data analysis



**Ezequiel Cimadevilla Álvarez**

[ezequiel.cimadevilla@unican.es](mailto:ezequiel.cimadevilla@unican.es)

Meteorology Group, Instituto de Física de Cantabria (IFCA, CSIC-UC), Santander, Spain

This work has been partially supported by:

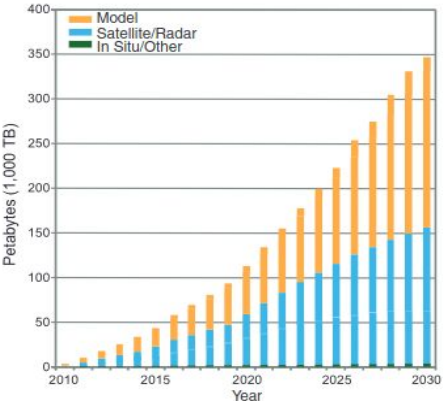
- Grant PID2020-116595RB-I00 funded by MCIN/AEI/10.13039/501100011033.
- Grant PRE2021-097646 funded by MCIN/AEI/10.13039/501100011033.

Project CORdYS (PID2020-116595RB-I00) funded by:



# Introduction

- Climate data, both from climate models and satellite observations, are seeing rapid growth in volumes.
- Data stored historically either on storage media in labs or in department, university, or organizational data centers.
  - “Download and analyze” model.
- Lack of *decision support systems* - Scientific data analysis becomes complex due to the need of deal with complex file formats and large number of files.



Institute	+
Model	-
Model	MPi-ESM-LR (3)
Experiment Family	+
Experiment	-
Experiment	rcp45 (3)
Time Frequency	+
Product	+
Realm	+
Variable	+
Variable Long Name	+
CMIP Table	-
CMIP Table	8hrPlay (3)
CF Standard Name	+
Ensemble	+
Data Node	+

Enter Text:

Display 10 results per page

Show All Replicas  Show All Versions  Search Local Node Only

Search Constraints ■ 8hrPlay ■ rcp45 ■ MPi-ESM-LR

Total Number of Results: 3  
-1-

Add all displayed results to Data Cart Remove all displayed results from Data Cart  
Expert Users: you may display the search URL and return results as XML or return results as JSON

---

1. project=CMIP6,model=MPi-ESM-LR, Max Planck Institute for Meteorology (MPI-M), experiment=RCP4.5,time\_frequency=8hr,modeling\_realm=tmos, ensemble=111p1,version=2011006

Description: MPi-ESM-LR model output prepared for CMIP6 RCP4.5  
Data Node: csf1.dkrz.de  
Version: 20111006  
Total Number of Files (for all variables): 380  
[\[ Show Metadata \]](#) [\[ Hide Files \]](#) [\[ THREDDS Catalog \]](#) [\[ WGET Script \]](#)

Total Number of Files: 380

ps1\_8hrPlay\_MPi-ESM-LR\_rcp45\_r111p1\_2100010100-2100123118.nc  
Checksum: 47c2b008a5103dbdb2c6d9338129ea9a6a8311bc2836d464273c13bed8750  
Size: 107668296  
Tracking Id: 4e248992-88d8-4e2e-b40f-01128007545e  
[\[ More File Metadata \]](#)

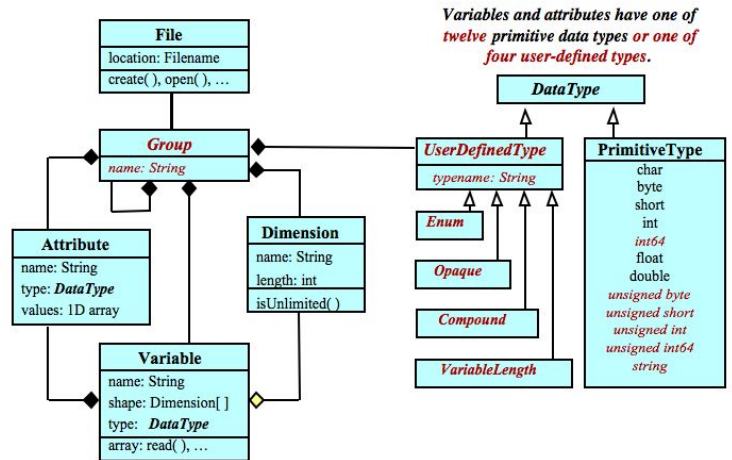
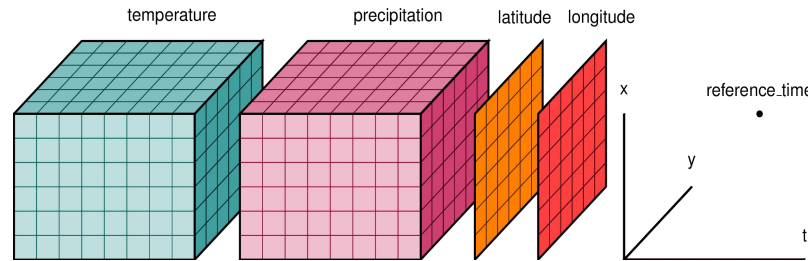
HTTServer  
OPENAP

ps1\_8hrPlay\_MPi-ESM-LR\_rcp45\_r111p1\_2099010100-2099123118.nc  
Checksum: d058aa932c334f5b1c1a7f5c800b619c42e8a94942e84e870bb27359e00f61a  
Size: 107668296  
Tracking Id: e338ca83-9403-46be-ab77-e72b40e033c4  
[\[ More File Metadata \]](#)

HTTServer  
OPENAP

### Multidimensional data

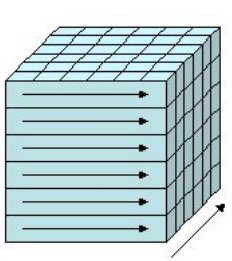
- Multidimensional arrays allow to organize and manipulate data with multiple dimensions or axes.
- Climate multidimensional arrays are both dense and quantitative, rather than sparse and qualitative.
- Different file formats for multidimensional data
  - netCDF
  - HDF5
  - Zarr
- Application data models are built on top of multidimensional arrays:
  - netCDF-java Common Data Model.
  - CF (Climate and Forecast conventions).



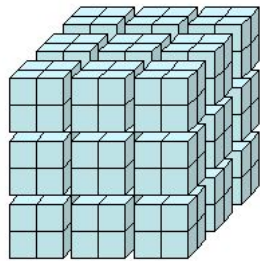
*A file has a top-level unnamed group. Each group may contain one or more named subgroups, user-defined types, variables, dimensions, and attributes. Variables also have attributes. Variables may share dimensions, indicating a common grid. One or more dimensions may be of unlimited length.*

### Multidimensional data

- Dense multidimensional arrays can be physically stored contiguously or chunked:
  - Contiguous - Stored as one long array in the file.
  - Chunked - Stored in chunks of user-defined size.

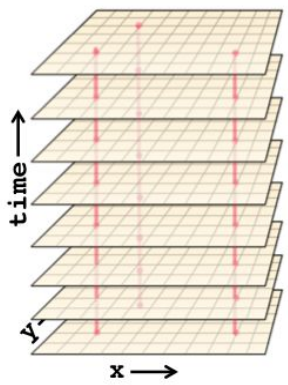


index order

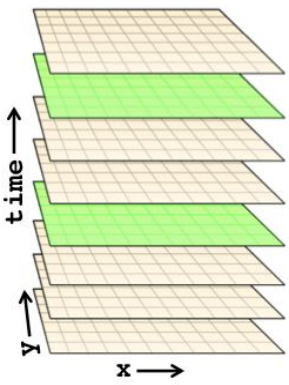


chunked

Time series access



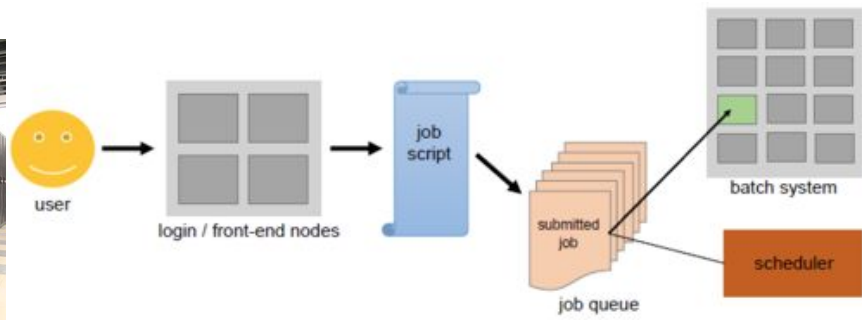
Spatial access



Storage layout, chunk shapes	Read time series (sec)	Read spatial slice (sec)	Performance bias (slowest / fastest)
Contiguous favoring time range	0.013	180	14000
Contiguous favoring spatial slice	200	0.012	17000
Default (all axes equal) chunks, 4673 x 12 x 16	1.4	34	24
36 KB chunks, 92 x 9 x 11	2.4	1.7	1.4
8 KB chunks, 46 x 6 x 8	1.4	1.1	1.2

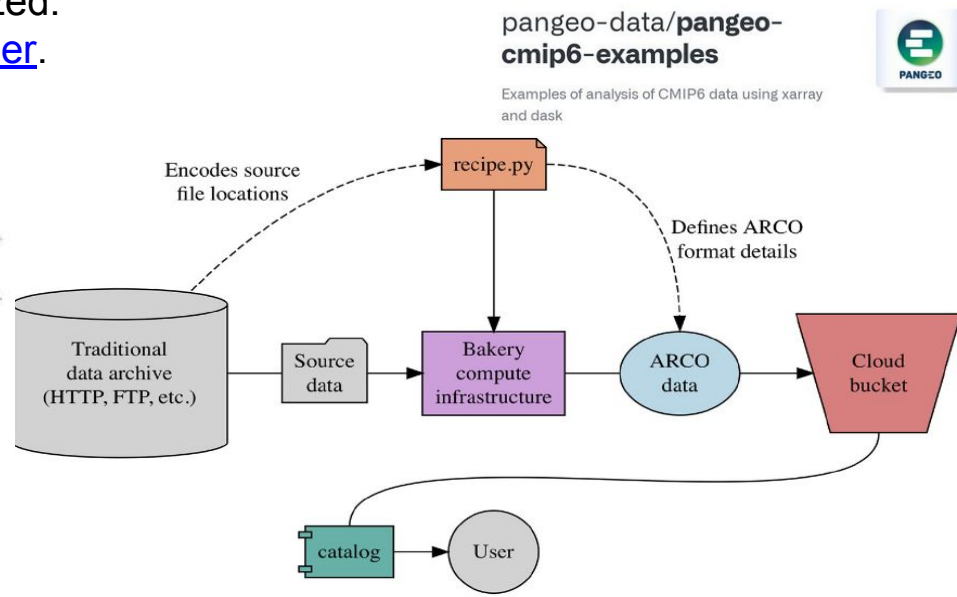
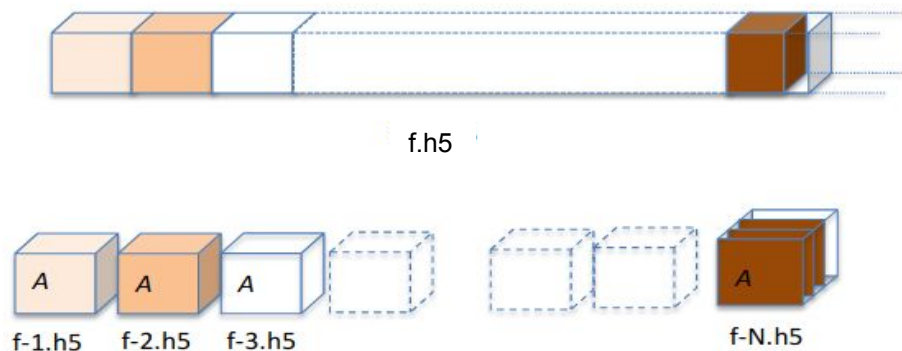
## HPC

- Regarding climate model data, the Coupled Model Intercomparison Project (**CMIP**) coordinates the design and distribution of global climate model simulations (20 PB of [CMIP6](#) data).
- Sixth Assessment Report (**AR6**) of the Intergovernmental Panel on Climate Change (**IPCC**).
- Contributions to the **increase in data volume** include the [systematic increase in model resolution and complexity](#) of the experimental protocol and data request.
- The Earth System Grid Federation ([ESGF](#)) is a global infrastructure and network of internationally distributed sites that together work as a federated data archive, supporting the distribution of global climate model simulations.
- Data is downloaded to **local HPC infrastructures** and analyzed using job queues and ad hoc tools for data processing.



## Cloud

- **Pangeo** and the **ESGF** have established a working group to help coordinate efforts related to storing and cataloging CMIP data in the cloud.
- [Google Cloud CMIP6](#) - Derived from the original CMIP6 data files, as distributed via ESGF.
- **ARCO** data - Analysis Ready Cloud Optimized.
- Complex workflow - See [Pangeo Forge paper](#).
  - Involves data duplication.
  - Aggregation of time series.



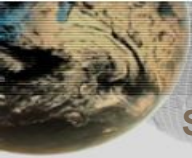
## OLAP - A ~~IT~~ scientific mandate

E.F. Codd et al. - Providing OLAP to User-Analysts: An IT Mandate.

*Most notably lacking has been the ability to consolidate, view, and analyze data according to multiple dimensions, in ways that make sense to one or more specific enterprise analysts at any given point in time. This requirement is called “multidimensional data analysis”.*

Scientific climate data analysis infrastructures can be classified according to:

- Transactional vs Analytical
  - Dealing with files (ESGF) vs dealing with data cubes (Pangeo\*).
  - ETL process to transition from Transactional to Analytical (**physical** vs **virtual** ETL).
- HPC vs cloud
  - Parallel file systems (Lustre, GPFS) vs object stores (Amazon S3, GCS).
  - Institutional job queues vs pay on demand computation services.

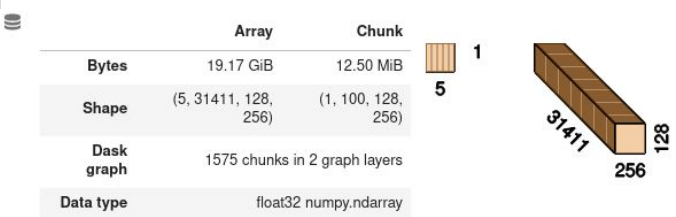


## The ESGF Virtual Aggregation

- “ESGF based Pangeo approach”.
- **Goals**
  - Remote data access and analytical capabilities on top of the ESGF infrastructure.
  - Single entry point catalog (see [here](#)).
- **Requirement** - Use the already existing software stack from ESGF.

```
# CMIP6_ScenarioMIP_CNRM-CERFACS_CNRM-CM6-1_ssp245_day_tas_gr_v20190410_aims3.llnl.gov.ncml
ds = xarray.open_dataset(url)
tas = ds["tas"]
tas = tas.chunk({"variant_label": 1, "time": 100})
tas
```

xarray.DataArray 'tas' (variant\_label: 5, time: 31411, lat: 128, lon: 256)



Coordinates:

lat	(lat)	float64	-88.93 -87.54 ... 87.54 88.93
lon	(lon)	float64	0.0 1.406 2.812 ... 357.2 358.6
height	()	float64	...
time	(time)	datetime64[ns]	2015-01-01T12:00:00 ... 2100-12-...
variant_label	(variant_label)		[S64 b'r2i1p1f2' ... b'r6i1p1f2']

Indexes: (4)  
Attributes: (12)

2. [CMIP6.CMIPBCC.BCC-CSM2-MR.historical.r1i1p1f1.3hr.tas.gn](#)  
 Data Node: cmip.bcc.cma.cn  
 Version: 20181127  
 Total Number of Files (for all variables): 22  
 Full Dataset Services: [\[ Show Metadata \]](#) [\[ Hide Files \]](#) [\[ WGET Script \]](#) [\[ LAS \]](#) [\[ Show Citation \]](#) [\[ PID \]](#) [\[ Globus Download \]](#)

Dataset

Total Number of Files: 22

1 [tas\\_3hr\\_BCC-CSM2-MR\\_historical\\_r1i1p1f1\\_gn\\_195001010000-195212312100.nc](#)  
 checksum: b5f270ed53e3ae7cbaa362b8cc1e3961e25750fecbb07ec074bd401ec9d02748  
 size: 1794284104  
 tracking\_id: hdl:21.14100/b880830c-7104-44d5-a02f-59bd431e816d  
[\[ More File Metadata \]](#)

2 [tas\\_3hr\\_BCC-CSM2-MR\\_historical\\_r1i1p1f1\\_gn\\_195301010000-195512312100.nc](#)  
 checksum: 44d89d66384dd5be2d42fa83abeb358a25e7901f7f39625f3a472b8d7c5d592f  
 size: 1794284104  
 tracking\_id: hdl:21.14100/02bcfe96-5445-4454-99f3-448f1617887  
[\[ More File Metadata \]](#)

Files in Dataset

Virtual dataset (NcML file)



Individual netCDF files found in ESGF



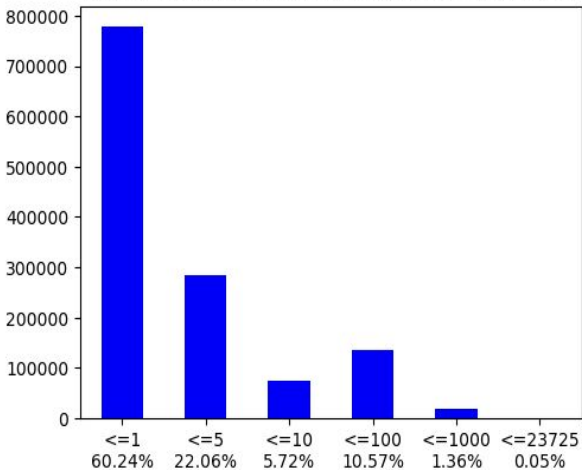


## The ESGF Virtual Aggregation

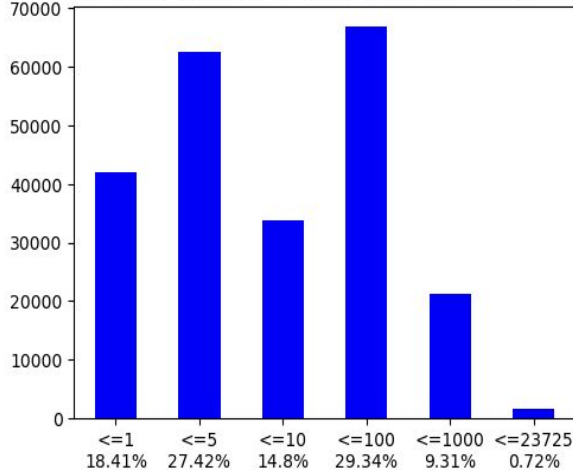
Find the catalog at <https://hub.ipcc.ifca.es/thredds> and Github repo at <https://github.com/zequiug50/eva>.

- Access to 12 PB of ESGF CMIP6 data (file sizes).
- 13,798,484 ESGF netCDF files (out of 30,753,032).
- 26 GB of NcMLs.

netCDF distribution in the time series aggregation



netCDF distribution in the ensemble aggregation



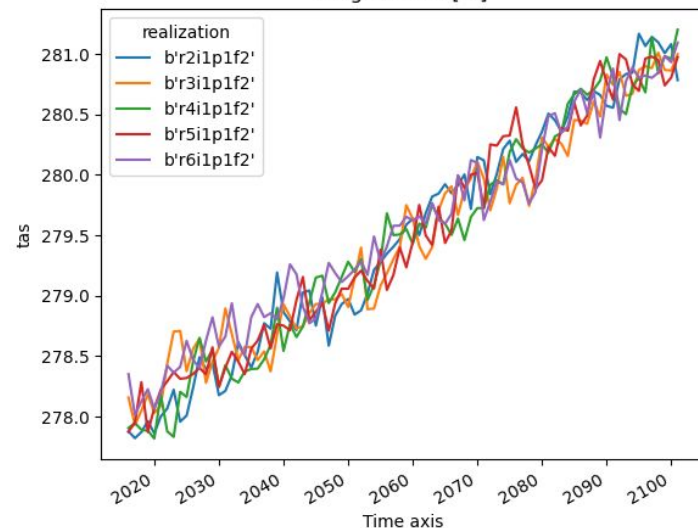
```
%time m = ds["tas"].mean(["lat", "lon"]).compute()
```

CPU times: user 360 ms, sys: 35.9 ms, total: 396 ms  
Wall time: 8min 18s

```
m.resample({"time": "Y").mean().plot.line(x="time")
```

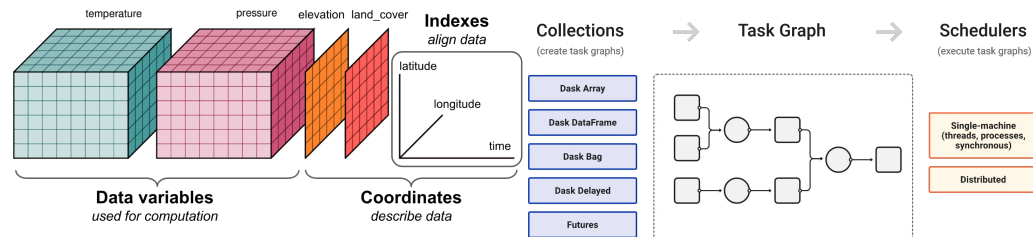
```
[<matplotlib.lines.Line2D at 0x7f901ef33490>,  
<matplotlib.lines.Line2D at 0x7f901ef336d0>,  
<matplotlib.lines.Line2D at 0x7f901ef33700>,  
<matplotlib.lines.Line2D at 0x7f901ef33640>,  
<matplotlib.lines.Line2D at 0x7f901ef33550>]
```

height = 2.0 [m]



## Data hubs

- Simplify access to **data** and **computation**. Move from the “download and analyze” model.
  - Provide remote access to the data using remote data access services.
  - Provide computation services based on web interfaces (**Jupyter/Python**).
- A data hub can be deployed both in HPC and cloud infrastructures.
- Improve **sophistication** and **speed** of information systems.



1. File browser ++

2. Run cells in .ipynb file using ctrl+Enter

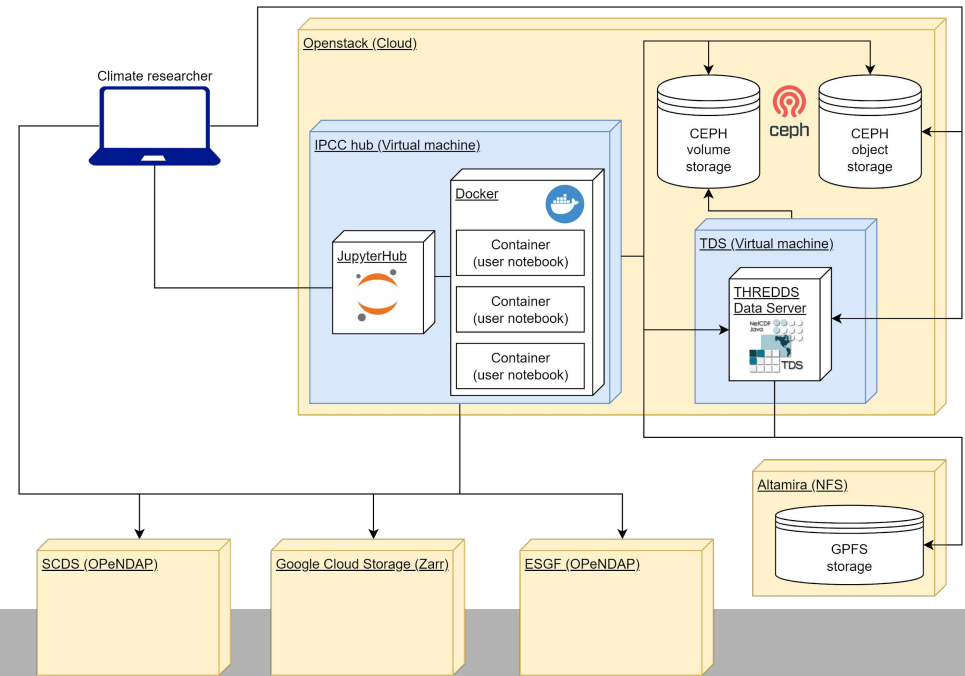
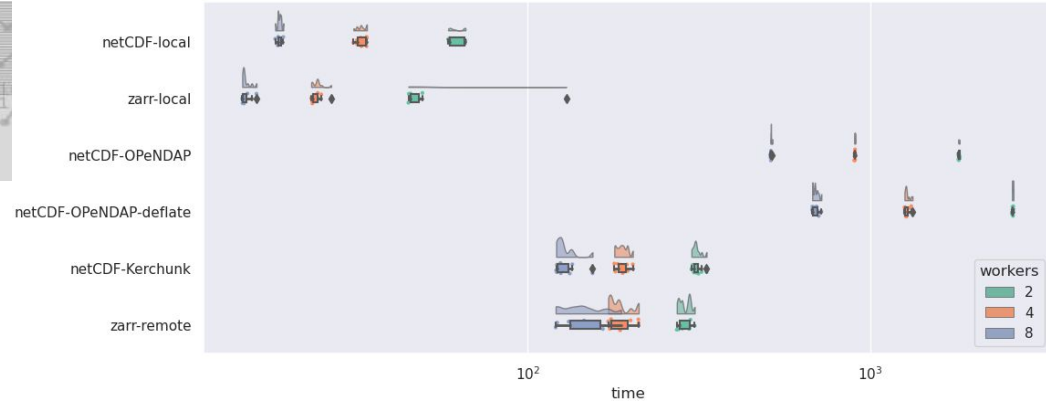
3. Run single line or highlighted text using keyboard shortcut

4. Run code in console using shift+Enter

5. Inspect variables or data frames in console without cluttering notebook output

## Data hubs - Remote data access

- Remote data access allows users to access data remotely without having to download the entire dataset.
- May offer interoperability between different data sources and formats:
  - xarray - netCDF, zarr, ...
- On the one hand, users from different locations can access datasets from their own infrastructures.
  - Far more close to FAIR principles than data downloads.
- On the other hand, lower data transfer throughput due to physical distance.



## Conclusions

- Climate data analysis still based on “file download and analyze” model.
  - New advances to transition from *transactional* to *analytical* systems.
  - New infrastructures appearing as alternatives (cloud vs HPC).
- ETLs needed to offer *data analysis ready* capabilities.
  - Physical ETL - Duplication of data with a new schema.
  - Virtual ETL - Different view of the original data without duplication.
- Data hubs are important tools for acceleration of climate scientific research.
  - Data hubs can be deployed either on HPC or cloud infrastructures, although differences in ecosystems between infrastructures need to be taken into account.

# Computing services for climate data analysis



**Ezequiel Cimadevilla Álvarez**

[ezequiel.cimadevilla@unican.es](mailto:ezequiel.cimadevilla@unican.es)

Meteorology Group, Instituto de Física de Cantabria (IFCA, CSIC-UC), Santander, Spain

This work has been partially supported by:

- Grant PID2020-116595RB-I00 funded by MCIN/AEI/10.13039/501100011033.
- Grant PRE2021-097646 funded by MCIN/AEI/10.13039/501100011033.

Project CORdYS (PID2020-116595RB-I00) funded by:

