

ClimERI-France Data Analysis Platform

Institut Pierre-Simon Laplace - February 3rd 2022

Levavasseur G. and Ramage K.

Journées ClimERI-France





provides Pierre-Simon Laplace The Institut laboratories with coordination resources and services to help develop major projects and disseminate the results.

The IPSL Computing and Data Centre was born out from several initiatives at the IPSL and its laboratories, wishing to share IT resources and joint projects through a numerical facility for research.

For more than 20 years, the IPSL computing and data centre has been providing Ensemble of Services for (tr. "Pour") Research at the IPSL - ESPRI.

What's "ESPRI"?









Community infrastructure





provide a distributed numerical ESPRI aims to infrastructure IPSL with "FAIR" for research at (Findability, Accessibility, Interoperability and Reusability) data services.



Mission & key figures







- All IPSL-CM production (CMIP5/6 + CORDEX)
- Replica pools:

Project	Thousands of files	Volume		
CMIP6	4 916	1,5 Po		
CMIP5	811	360 To		
CORDEX	458	223 To		
input4MIPs	8	3 То		
C3S-CMIP5	174	25 To		
C3S-CORDEX	582	275 To		
Total	<u>6 949</u>	<u>2,4 Po</u>		





ESPRI centralizes ClimERI-France data services



JON IPSL Kernel Git Settings Help Run Tabs I Launcher C 🚯 Q Filter files by name 0 Notebook Name Last Modified . i miniconda 3 years ago ٦ Ρ R CMIP5_bdd_roots.out 11 hours ago CMIP5_bdd.roots.sh 2 years ago Python 3 Python [conda R [conda panel [7] CMIP6_bdd_roots.out 11 hours ago env:root] env:root] CMIP6_bdd_roots.sh 2 years ago CORDEX_bdd_roots.out 11 hours ago

2 years ago

CORDEX_bdd_roots.sh

Console Python 3

Terminal

R [conda Python [cond env:root] env:root] \$_ Other ≣ Μ \$_ -

Markdown File

Show Contextual Help

Text File

0	File Edit View Run Kernel Git Tabs Settings Help										
	+ 🗈 ± C 🗞		Untitled.ipynb			nb	٠		0.0	-14	
0	Filter files by name Q			+	X			C ⊫⊨ Code ∨	00	git	
	=/			[1]:	: 1	ls /bdd/CMIP6/PMIP					
rel)	Name	Last Modified	1		4	AWIQ		CSIRO@	IPSL@	MRI@	NCC@
40	miniconda	3 years ago			0	CNRM-CERFACS@	a	INM@	MPI-M@	NCAR@	NUIST@
٩	CMIP5_bdd_roots.out 11 hours ago CMIP5_bdd.roots.sh 2 years ago			[3]:		<pre>import dask import intake_esm</pre>					
					1						
≣	CMIP6_bdd_roots.out 11 hours ago				i	import xarray					
	CMIP6_bdd_roots.sh	d_roots.sh 2 years ago			1	import matplotlib					
*	CORDEX_bdd_roots.out	11 hours ago		[]	:						
	CORDEX_bdd_roots.sh	2 years ago									
	Untitled.ipynb	a minute ago									

Towards Python Notebook interface

- Access: https://data.ipsl.fr/jupyter
- **Status**: Pre-production
- **Users**: ~20 "power-users" with extensive usages + Proven with training session.
- **Documentation:** in progress to be available through the new incoming ESPRI website.
- **Python environment**: PANGEO + useful modules.
- Service opening: March 2022

GAIA-Data requirements:

- Data catalogs access (including ClimERI ones)
- Online editors (VSCode, R, Matlab)
- Dashboards (Dask, Tensor)





Old-fashion discovery

NO M



- \$> ls /bdd/CMIP3
- \$> ls /bdd/CMIP5
- \$> ls /bdd/CMIP6
- \$> ls /bdd/CORDEX
- \$> ls /bdd/obs4MIPs
- \$> ls /bdd/input4MIPs
- \$> ls /bdd/CMIP5-Adjust
- \$> ls /bdd/CORDEX-Adjust
- \$> ls /bdd/C3S-CMIP5
- \$> ls /bdd/C3S-CORDEX

```
$> python3
>>> import intake
>>> cat = intake.open_esm_datastore('/modfs/catalog/CMIP6.json') # Takes up to a minute to load the catalog
>>> cat
<CMIP6 catalog with 11352 dataset(s) from 6868559 asset(s)>
>>> cat.df.columns
Index(['path', 'project', 'activity id', 'institution id', 'source id', ..., 'latest'], dtype='object')
>>> cat.unique(columns='table id')
{'table_id': {'count': 39, 'values': ['Emon', 'Amon', 'CFmon', 'day', ..., 'Eday']}}
>>> subcat = cat.search(experiment id='ssp585',
                    institution id=['IPSL','CNRM-CERFACS'],
                    variable id=['pr', 'tas'],
                    table id='Amon',
                    latest=True)
>>> subcat
<CMIP6 catalog with 4 dataset(s) from 40 asset(s)>
>>> dsets = subcat.to dataset dict()
>>> dsets.keys()
dict keys(['ScenarioMIP.IPSL.IPSL-CM6A-LR.ssp585.Amon.gr', 'ScenarioMIP.CNRM-CERFACS.CNRM-CM6-1.ssp585.Amon.gr',
...])
>>> dsets['ScenarioMIP.IPSL.IPSL-CM6A-LR.ssp585.Amon.gr']
<xarray.Dataset>
                  (time: 3432, lat: 143, lon: 144, axis nbounds: 2, member id: 7)
Dimensions:
[...]
Data variables:
[...]
                  (member id, time, lat, lon) float32 dask.array<chunksize=(1, 1032, 143, 144), meta=np.ndarray>
    \mathtt{pr}
                  (member id, time, lat, lon) float32 dask.array<chunksize=(1, 1032, 143, 144), meta=np.ndarray>
    tas
Attributes: (12/50)
[...]
    tracking id:
                             hdl:21.14100/f09be139-ce58-4719-8fe9-2769aad503c...
```

Improved local data discovery

Discovery through *"intake-esm" catalogs*

Access: /modfs/catalog JupyterHub Supported projects: CMIP5/6 CORDEX (other coming soon). **Quality control**: CV checked

Update:

Weekly

Service opening: February 2022











Journées ClimERI-France



REQUEST

Search criteria called facets are used to select files which to download. They can be set on command line or using a template.



ESGF NODES

Jenn

SDT retrives the certificates and builds the HTTP requests to Solr corresponding to the search criteria.

SYNDA

Transfer

A command line tool to discover and download from the Earth files System Grid Federation (ESGF) archive.

FILESYSTEM

ESGF files are downloaded using the HTTP or GridFTP protocol and managed on local filesystem the following the Data Reference Syntax.

SDT DATABASE

A SQLite database records each downloaded file and dataset. A complete dataset triggers a "dataset_complete" event, which informs the SDP module to start the pipeline.

New paradigm implemented to perform **parallel downloads** asynchronously. Once each workers have their download tasks given out by the scheduler, they are able to asynchronously carry out their duties without having to wait for each other.







This year, the ESPRI team was awarded by the Crystal collective prize from CNRS.



#Inktober2021 01/10/2021 • crystal

Thank you for your support and feedbacks !

INSTITUT PIERRE-SIMON LAPLACE

ESPRI awarded







Towards **STAC catalogs** with common search engine between observational and modelling data sets.

IFEE

- **Extending** a multi-thematic storage at IDRIS
- Extending and automating **DOI** services to codes.

Network mounting point between **CNRS-CERFACS and ESPRI** through a **10Gbps** RENATER link and private VPN, allowing direct access to CNRM-CERFACS climate simulations.

2022 and beyond with GAIA-Data







The Copernicus **Climate Data Store** (CDS) supplies access to a **subset** of **global** and regional climate projections from CMIP5/6 and CORDEX exercices (~300TB).

This service relies on **dedicated** and **load-balanced** servers deployed by **ENES** partners which provides a single resilient point of access to data delivered through (~500 replication and **redundancy** downloads/day or ~300GB/day).

ESPRI will **coordinate** the **operational** maintenance of the infrastructure for the next 4 years.





Thank you for your attention

Institut Pierre-Simon Laplace