# Earth System Grid Federation Future Architecture, Copernicus, Cloud and ESA

## ClimateData.ca Meeting, 23 November 2021

Philip Kershaw, Technical Manager
Centre for Environmental Data Analysis

Philip Kershaw, Technical Manager
Centre for Environmental Data Analysis

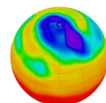Centre for Environmental Data Analysis
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
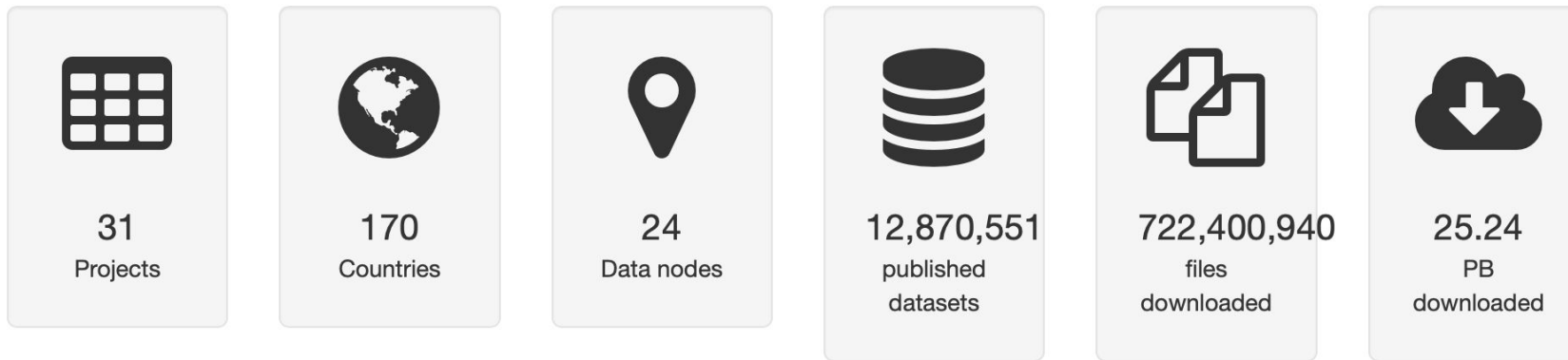NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Earth Observation
NATURAL ENVIRONMENT RESEARCH COUNCIL

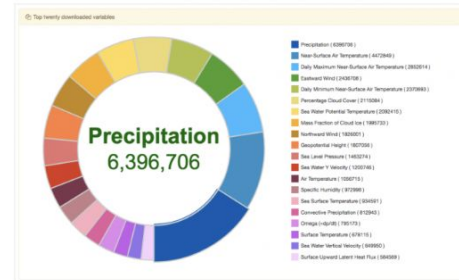# Earth System Grid Federation: a globally distributed data archive for climate data

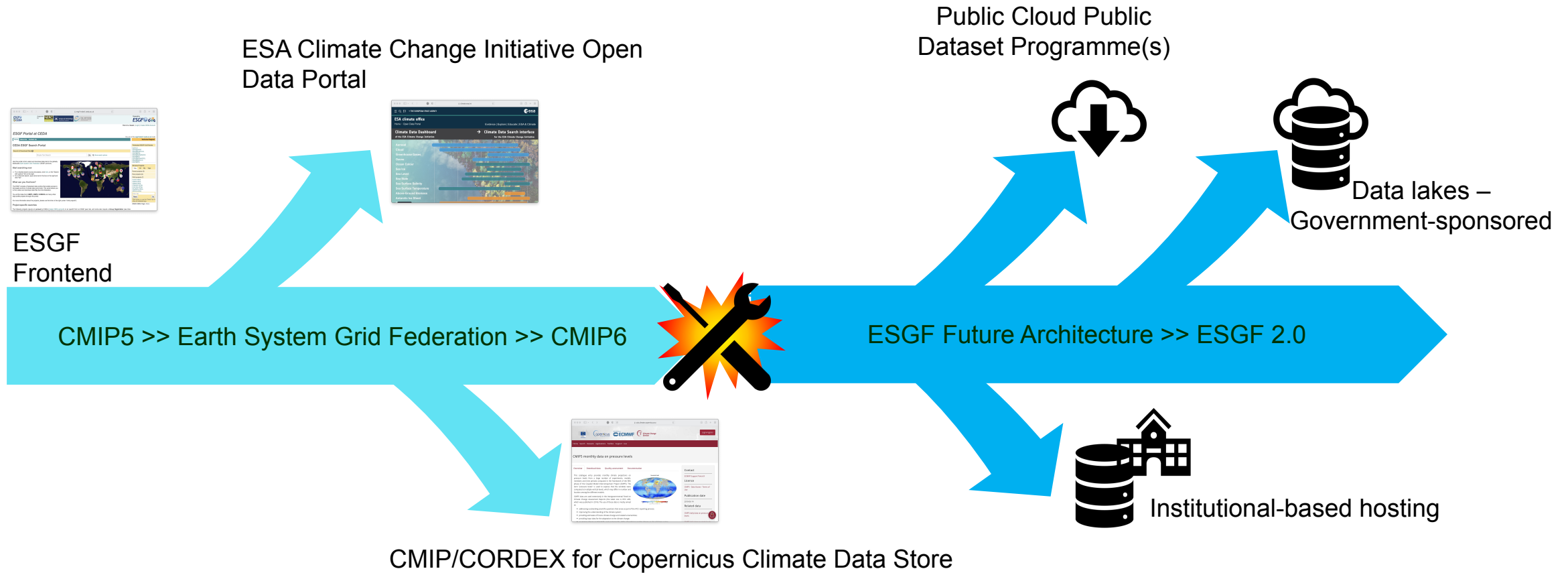| | | | | | |
|---|---|---|---|---|---|
| **31** Projects | **170** Countries | **24** Data nodes | **12,870,551** published datasets | **722,400,940** files downloaded | **25.24** PB downloaded |

ESGF Dashboard: http://esgf-ui.cmcc.it



**ESGF Federation**

**Data usage**

Precipitation 6,396,706

**Data publication**

# ESGF – Application and Evolution



ESA Climate Change Initiative Open Data Portal

Public Cloud Public Dataset Programme(s)

Data lakes – Government-sponsored

ESGF Frontend

CMIP5 >> Earth System Grid Federation >> CMIP6

ESGF Future Architecture >> ESGF 2.0

Institutional-based hosting

CMIP/CORDEX for Copernicus Climate Data Store

# ESA Climate Change Initiative Open Data Portal



2 Phases:

1) Leveraged ESGF
   1) Quick win with search and download
   2) Bespoke search API incompatible with other community standards - OGC CSW
   3) THREDDS Data Server couldn't scale to our needs
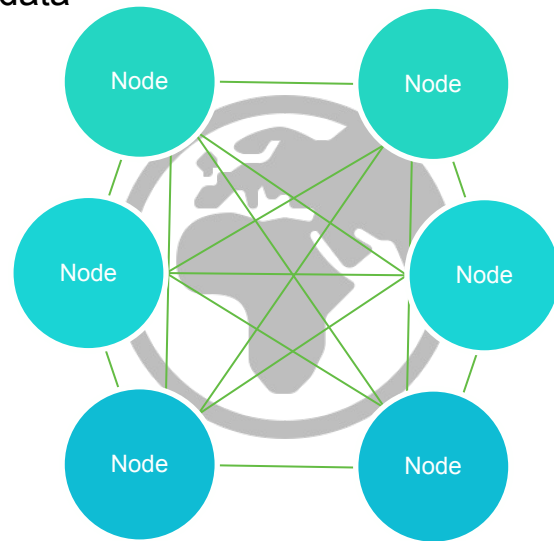
2) Redesigned to address issues
   1) OpenSearch API replaced ESG Search
   2) Scalable data service with Kubernetes
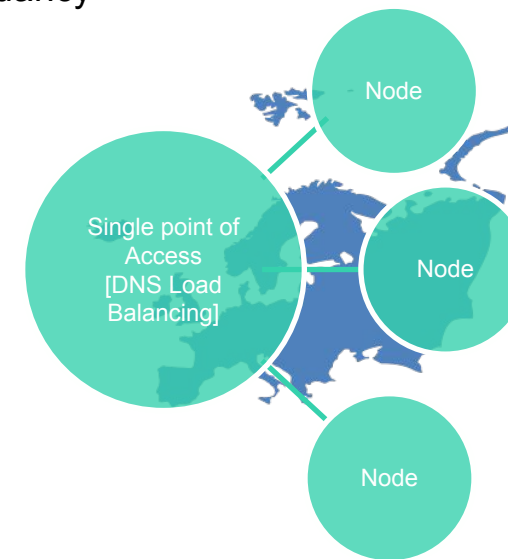   3) Zarr format cache of netCDF data on object store for performance

# C3S 34[a-f] Projects for the CDS

- Architected a system for delivering resilient CMIP5 and CORDEX data access for the CDS by creative application of federated architecture for Earth System Grid Federation
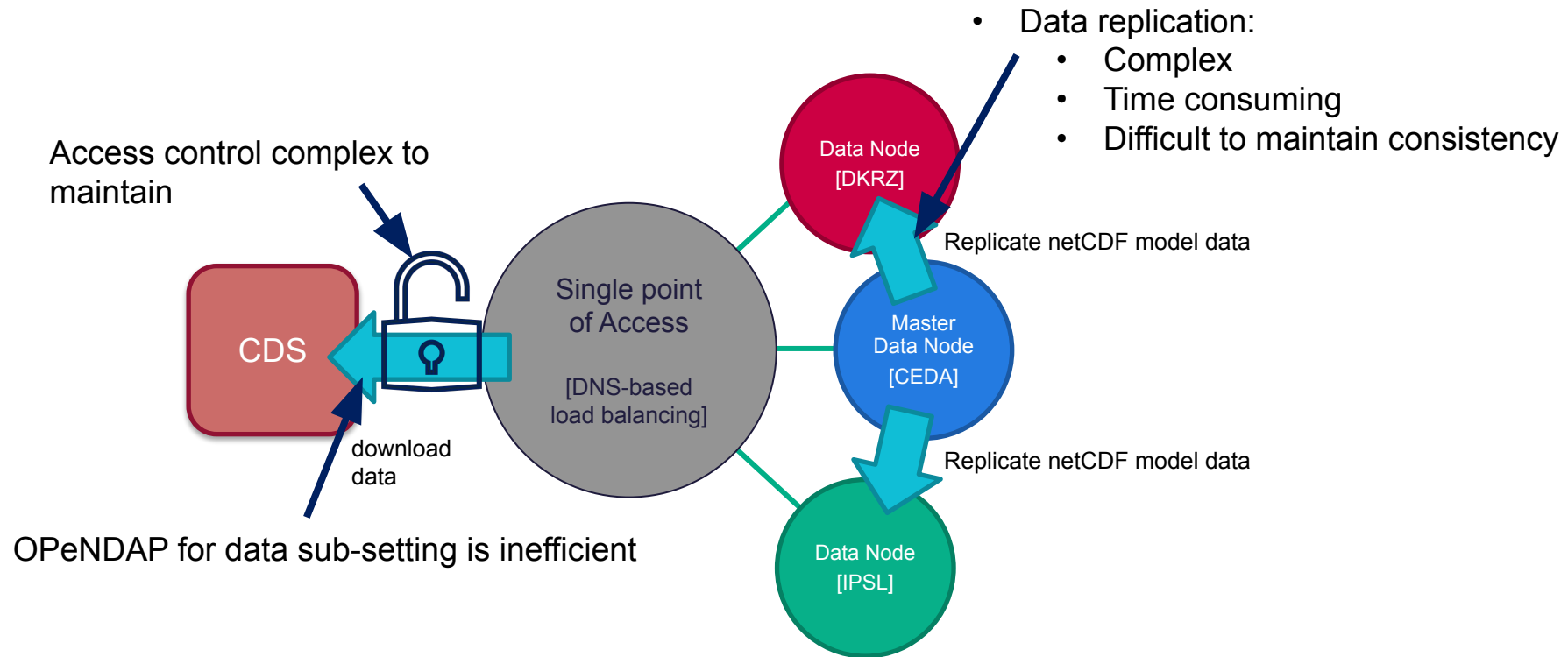
ESGF: an international federation of nodes providing a network of access points to model data
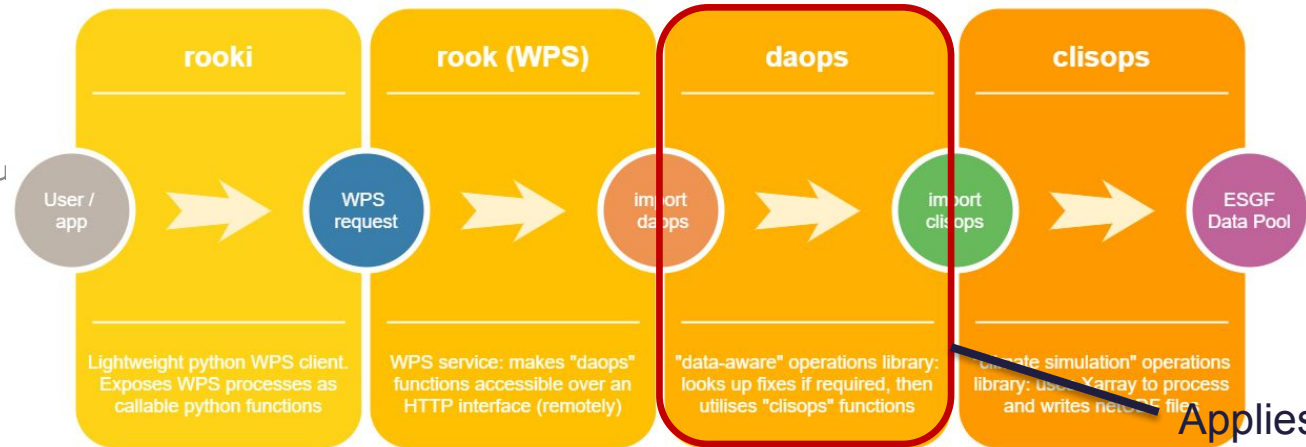
C3S 34a/b system: a single resilient point of access to data delivered through replication and redundancy
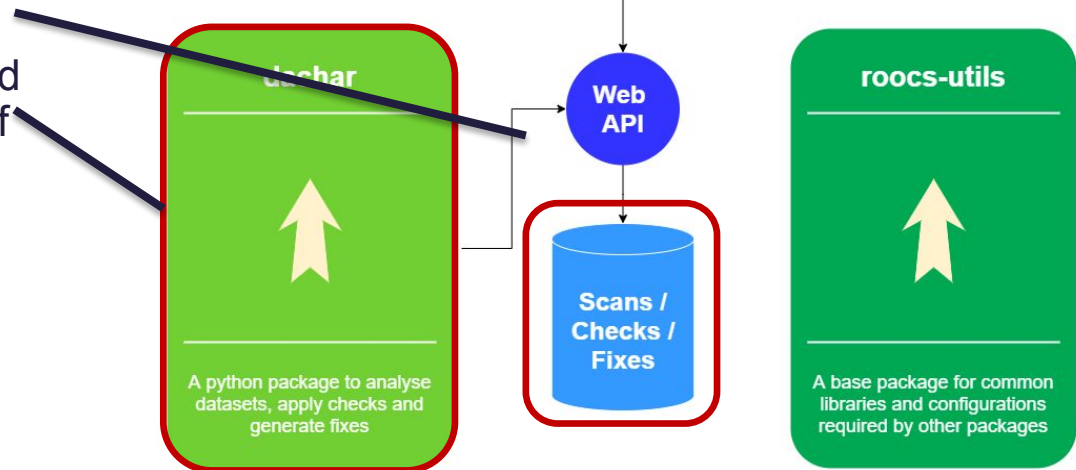
# C3S Resilient CMIP and CORDEX Data Access



- Data replication:
  - Complex
  - Time consuming
  - Difficult to maintain consistency

Data Node [DKRZ]

Access control complex to maintain

Single point of Access

[DNS-based load balancing]

CDS

Master Data Node [CEDA]

Replicate netCDF model data

download data

Replicate netCDF model data

OPeNDAP for data sub-setting is inefficient

Data Node [IPSL]

# Sub-setting Services for C3S 34e Project



Credit: Ag Stephens, CEDA

jupyter demo-rooki-subset-by-point

Visit repo    Copy Binder link

File    Edit    View    Insert    Cell    Kernel    Widgets    Help         Not Trusted         Python 3 (ipykernel) O

Markdown    ⇕                    **Memory:** 154.7 MB / 2 GB

Download    GitHub    Binder

# Run subset by (time) point operation

**Rooki** calls climate data operations on the **rook** processing service.

```
In [ ]: import os
os.environ['ROOK_URL'] = 'http://rook.dkrz.de/wps'

from rooki import rooki
```

**parameters of subset operation**

```
In [ ]: rooki.subset?
```

## subset by time interval

```
In [ ]: resp = rooki.subset(
        collection='c3s-cmip6.ScenarioMIP.INM.INM-CM5-0.ssp245.r1i1p1f1.day.
        time='2016-01-01/2016-12-30',
)
resp.ok
```

# ESA Earth Observation Exploitation Platform Common Architecture (EOEPCA)



- Architectural blueprint for the federation of *platforms*
- Interlinked with the work on recent OGC testbeds
- CEDA involved with consultancy role
- Processing and chaining of particular interest
  - ADES and EMS
  - Ability to push customised shrink-wrapped processes to 3rd party WPS instances
- Innovations with ID management: UMA

# ESGF – Application and Evolution

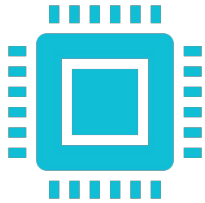# ESGF Future Architecture

**Platforms and systems administration**

Modular, scalable architecture: Containers, Kubernetes

Embrace infrastructure-as-code approach

**Search services**

Modernise, centralise and simplify

Use community standards: STAC

**ID Management and Access Entitlement**

Modernise, centralise and simplify

Use industry standards: OpenID Connect / OAuth 2.0
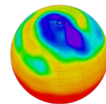
**Progress and Achievements**

- Container and Container+Kubernetes installs available
- **Deployed on AWS (GFDL) and at CEDA**

- Major community engagement on use of STAC for ESM data
- Prototype developed by CEDA
- Integration tests
- CoG and MetaGrid futures??

- OpenID Connect / OAuth 2.0 done
- New Authorisation system with Open Policy Agent
- Authentication integrated with C4I in test
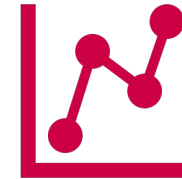
# ESGF Future Architecture

**New modes for Data Access + Storage**

Augment trad. file serving with object store

New models for aggregation and subsetting, retire OPeNDAP

**Compute Services**

Important but no consensus for ESGF-wide standard offering yet

**Metrics Collection**

Leverage advances in industry with standard tooling to exploit - Prometheus and InfluxDb, Grafana

- Factored out TDS
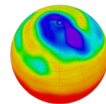- Test CMIP6 data caches on object store at CEDA and DKRZ

- Ported C3S WPS Data Reduction Services for use in ENES CDI
- Used with Climate4Impact
- Reboot of Compute Working Team

- New Metrics system integrated with CMCC

# Future Architecture Node – Phase 1

# Future Architecture Node – Phase 2

# STAC API for ESGF



- Implementation of STAC API

- ElasticSearch backend

- Filter extensions to support faceted search

- Fully featured STAC equivalent API to ESG Search

- Simple frontend created to demonstrate its features

# IS-ENES3 - Data Analytics using Notebooks/icclim

# DestinE and Blueprint Architecture


Various user groups
(New) applications & services
Digital Twins
Open Core Platform
Cloud infrastructures
Data sources & infrastructures
High Performance Computing

Destination Earth (DestinE) - major EU initiative:

- "to develop a very high precision digital model of the Earth (a 'digital twin') to monitor and predict environmental change and human impact to support sustainable development"

# JASMIN



Data Analytics Platform

Data production / processing

High Performance Computing

Data Sources

Cloud Infrastructure

JASMIN

# ESA Digital Twin Earth (DestinE) Precursor - land surface modelling and climate

- Using JULES (Joint UK Land Environment Simulator)
  - the land surface component in the Met Office Unified Model

- Improvements with Data Assimilation
  - LaVEnDAR (The Land Variational Ensemble Data Assimilation fRamework)
  - Feed in satellite observations – SIF and SMAP data

# What could be the future impact of climate change on the soil moisture?

- JULES driven with climate projections from ISIMIP data (Inter-Sectoral Impact Model Intercomparison Project)

# Make a surrogate AI model to JULES

Soil Moisture

- Experimented with Machine Learning (ML) techniques

- Goal: a general-purpose algorithm -

  ```
  time series of daily weather data ☐ time
  series of soil moisture data
  ```

- Successfully applied XGBoost (eXtreme Gradient Boosting) algorithm.

- trained on up to 1000 grid cells, representative of the various biomes in continental Africa

- Demonstrated to accurately emulate JULES output at other locations

- The credibility of the model is enhanced by its transparency and explainability

# Digital Twin Precursor on JASMIN: HPC for data production, cloud for analysis

JULES / LAVENDAR Data Assimilation

**Soil Moisture** model outputs netCDF files to regular file system

Data accessed using Jupyter Notebook service

Cluster-as-a-Service deploys ready-made Jupyter service

Batch compute (Lotus)

JASMIN

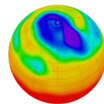netCDF

Group Workspace (GWS)

Cluster-as-a-Service

external cloud tenancy

Managed (Internal) JASMIN

External JASMIN Infrastructure

Move data into object store so that it can be accessed by Jupyter Service on JASMIN cloud

UKRI
Science and Technology Facilities Council
Natural Environment Research Council

Centre for Environmental Data Analysis
SCIENCE AND TECHNOLOGY FACILITIES COUNCIL
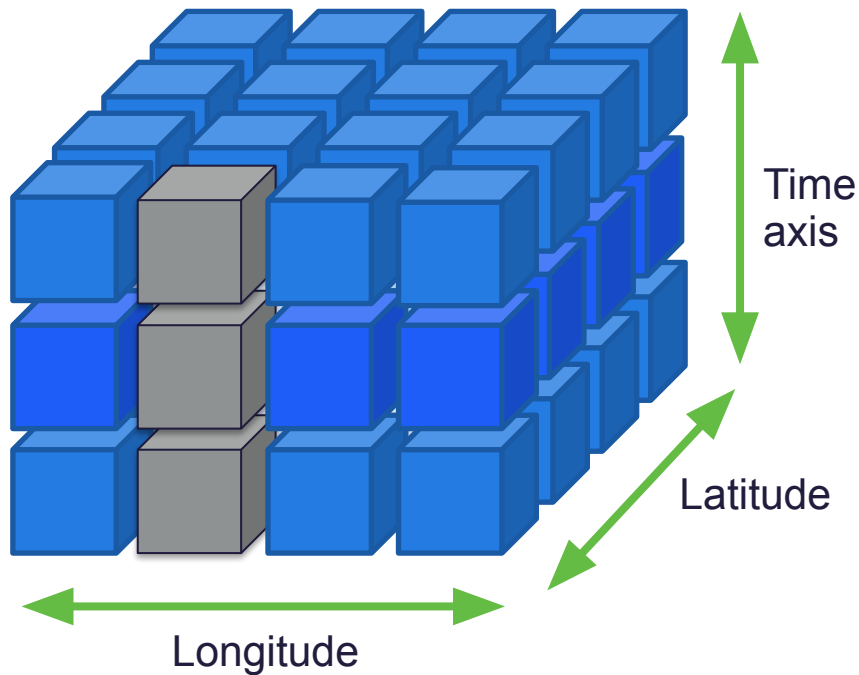NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Atmospheric Science
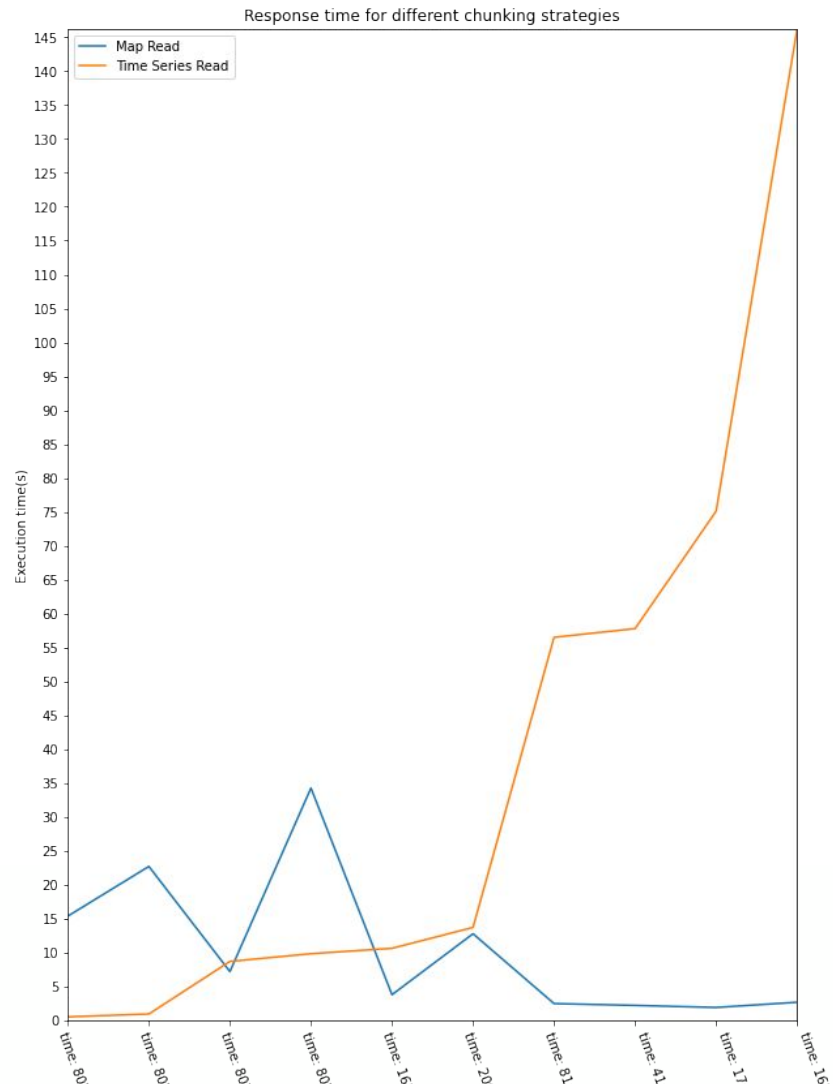NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Earth Observation
NATURAL ENVIRONMENT RESEARCH COUNCIL

# Arrangement of data and efficient access



- Data output from models as netCDF format
- Data in files arranged in spatial dimensions one per time step
- But predominate access pattern for analysis of climate data in the project is time series query (grey blocks)
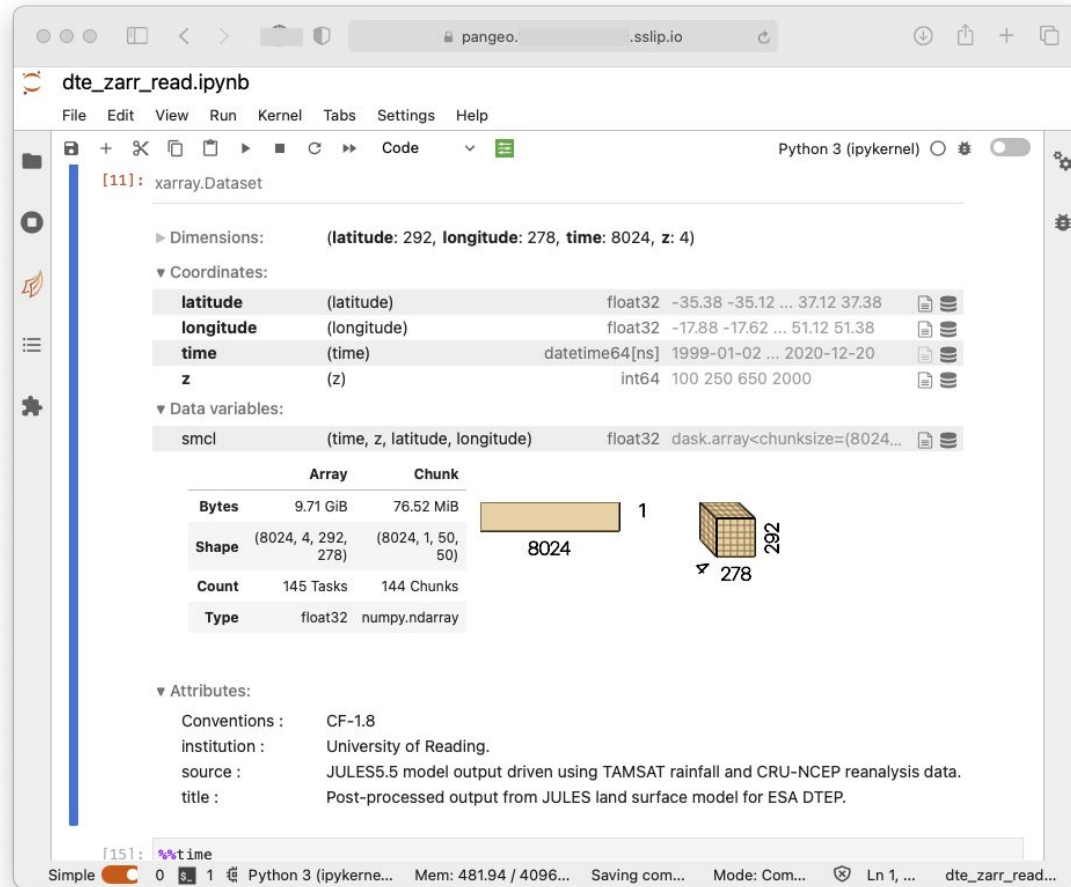
# Object Store: Different storage strategies showed radically different performance
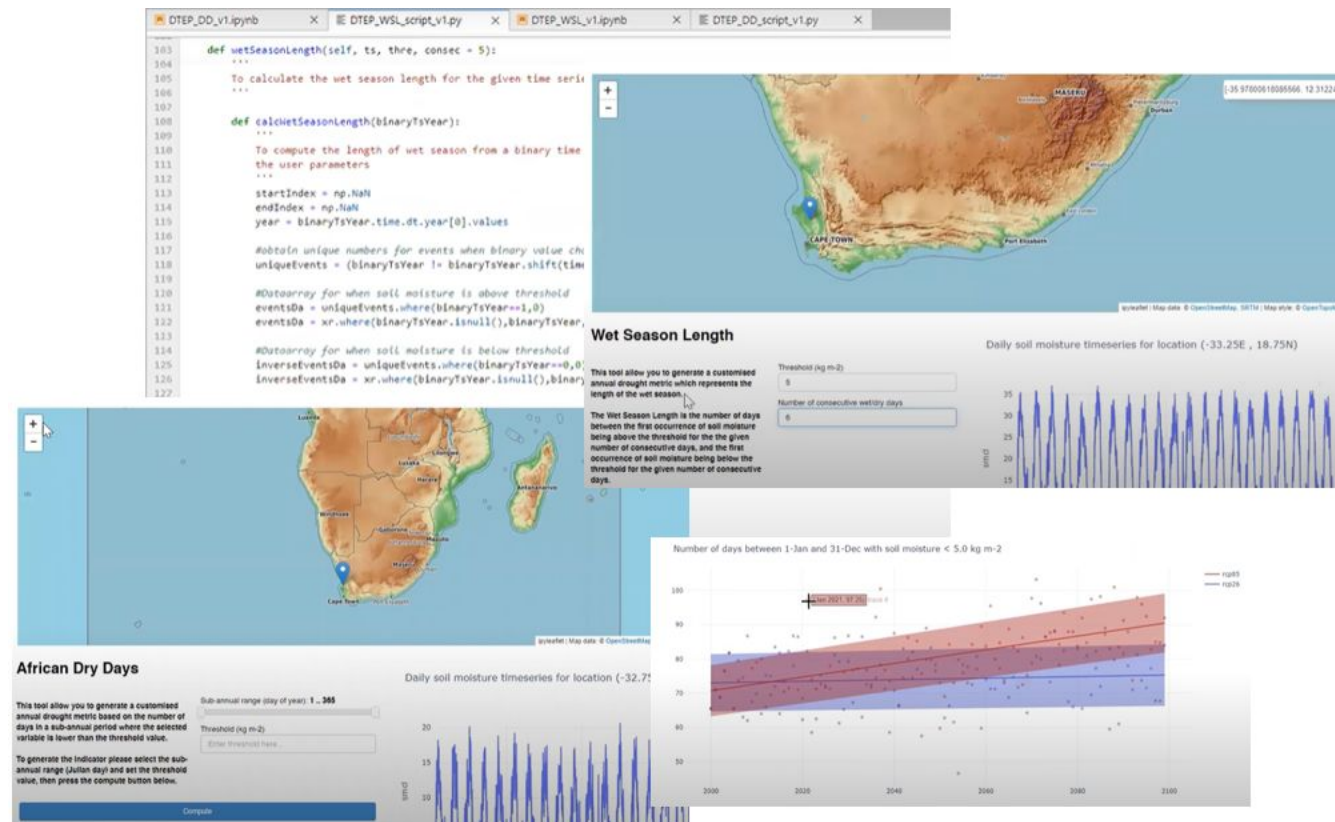

Response time for different chunking strategies

- We experimented with different storage chunking arrangements
- 20-year dataset of soil moisture

National Centre for Atmospheric Science
NATURAL ENVIRONMENT RESEARCH COUNCIL

National Centre for Earth Observation
NATURAL ENVIRONMENT RESEARCH COUNCIL

# Using Object Store for re-arrangement of data to suite our access patterns
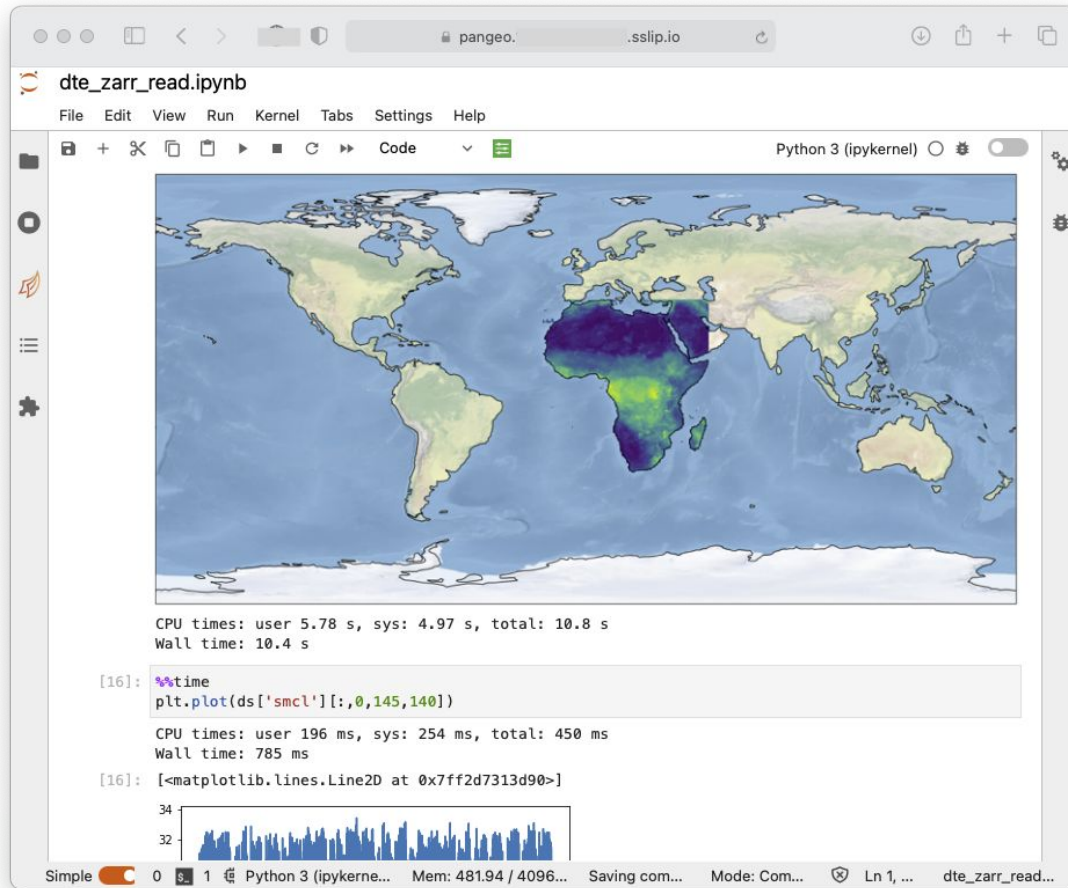


- Using zarr and xarray Python libraries to store and access the data

- Chunked data into a series of strips along the time axis

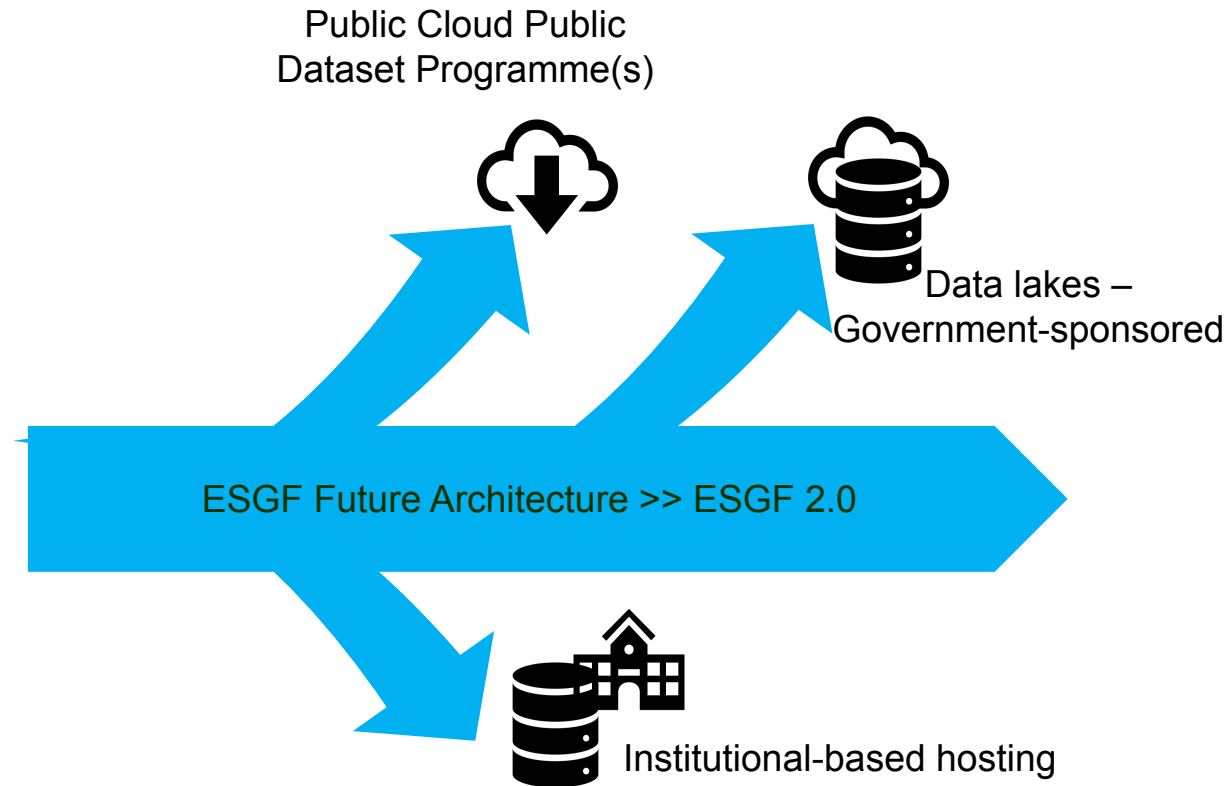# Rechunking of data made possible interactive maps with long time series

# Take home message: object store for analysis-ready cache specific to project needs



- Object store can be efficient for access on cloud

- It is essential to orient data storage to suit predominant access patterns

- Good news – re-writing data into different orientations was fast

# Futures

Public Cloud Public Dataset Programme(s)



Data lakes – Government-sponsored

ESGF Future Architecture >> ESGF 2.0

Institutional-based hosting

# Acknowledgements + Further Info


THE CONSORTIUM

Coordinated by CNRS-IPSL, the IS-ENES3 project gathers 22 partners in 11 countries

ESGF Future Architecture Report: https://doi.org/10.5281/zenodo.3928222

IS-ENES3 website
https://is.enes.org/

@ISENES_RI
@cedanews
@PhilipJKershaw

Contact us at
is-enes@ipsl.fr

Subscribe to the
**IS-ENES3 H2020**
Youtube channel !