

ENES climate data infrastructure

Stephan Kindermann
DKRZ

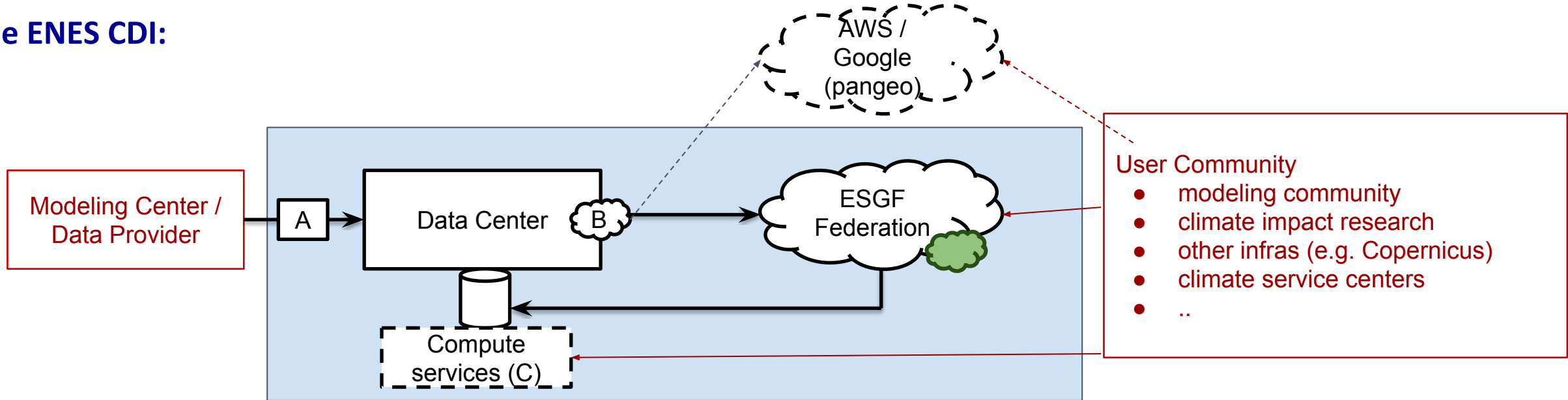
European Network for Earth System Modelling (ENES)



Overview

- **Main objectives, user groups, evolution**
- **Core infrastructural components and services of the data infrastructure**
- **Evolving aspects:**
 - **computation near to the data**
 - **towards cloud based infrastructure(s)**

The ENES CDI:



- User Community
- modeling community
 - climate impact research
 - other infras (e.g. Copernicus)
 - climate service centers
 - ..

- **A: Metadata and Data standards**
- **B: Data centers integrating services in a federation (ESGF based)**
- **C: Enabling data analysis activities near to the data**

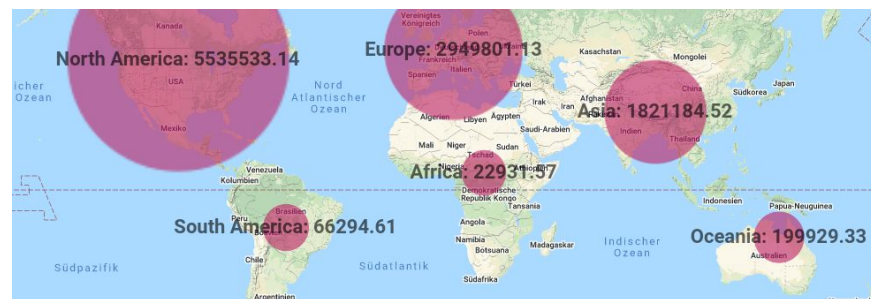
towards cloud deployments:

- a) hosting of core components in commercial clouds 
- b) hosting data and compute in commercial (and institutional) clouds 

Support data distribution for internationally coordinated climate model intercomparison projects (CMIP, currently CMIP6)

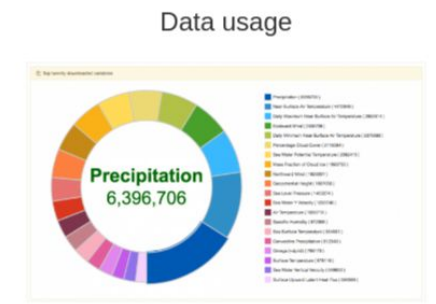
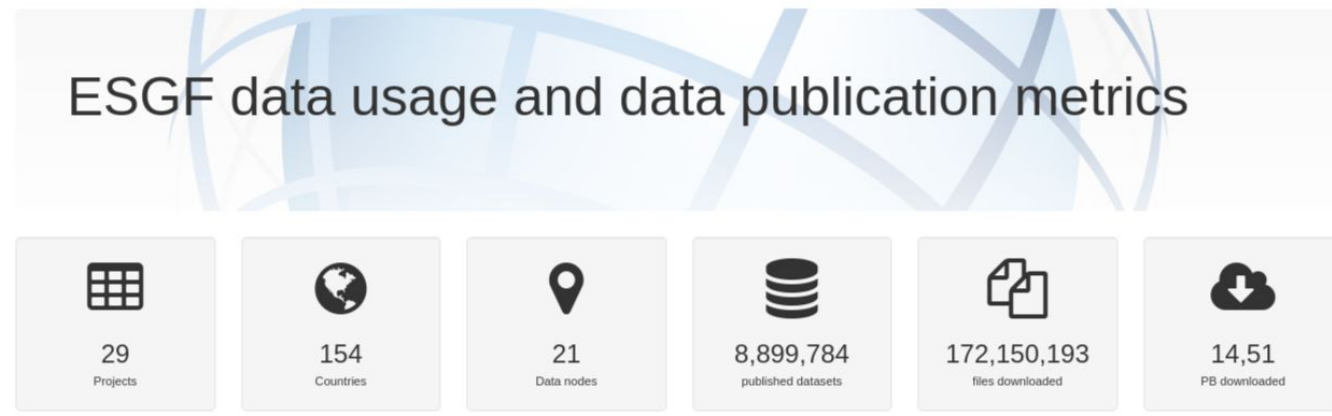
Example: ENES CDI supporting CMIP6:

- 4 ESGF portals (world: 10)
- 11 ESGF data nodes (world: 29)
 - France, Germany, UK (tier1)
 - Ireland, Italy, Norway, Spain, Sweden
- Models: 46 (world: 103)
- Modeling Institutes: 19 (world: 41)



ESGF Data Statistics

- Data Usage Metrics
- Data Publication Metrics
- Geo-downloads
- IS-ENES3 KPIs
- Meta-statistics
- Feedback form

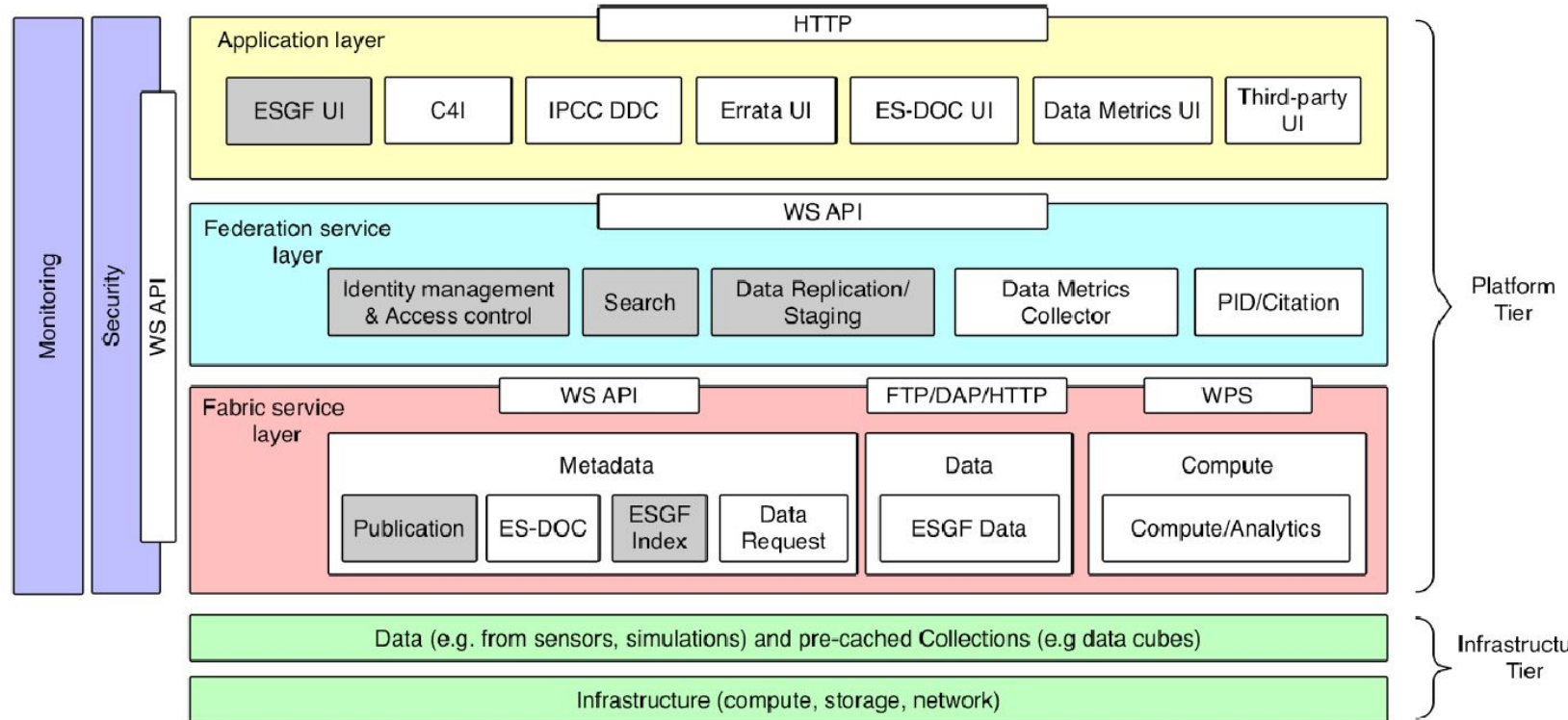
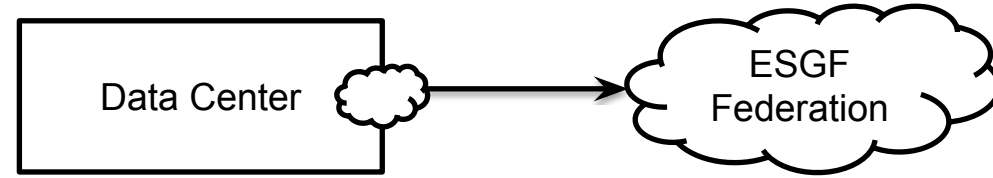


The ESGF SW stack evolved over the past > 10 years as part of an international collaborative effort supported by national/institutional funding

- IS-ENES3 project is currently coordinating the European contributions (white boxes in picture)
- github hosted and open source

Current challenges:

- Modularization, orchestrated deployment support (docker, ansible, kybernetes)
- Worldwide operating system, best effort funding → sustainability
 - EU: discussed as part of IS-ENES3
 - Worldwide: ESGF XC, WCRP, .. discussions



Beyond data search and access services in support of CMIP6:

- Persistent data identification service (PID service)
- Errata service
- Data citation service (DOI service)

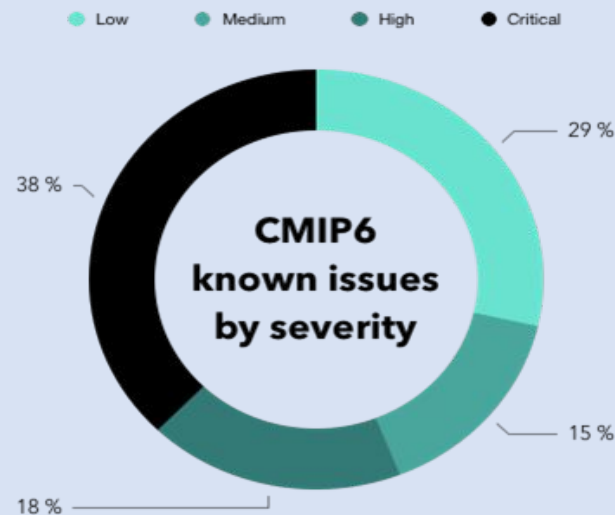
Supporting key aspects
of FAIR data principles

PID service:

- Operational integration based on worldwide distributed message queue
- overall 16.5×10^6 PIDs assigned (including retracted etc.)
- $\ll 0.1\%$ registration problems
- automatic curation scripts in place
- Cooperation with EOSC, EPIC, EUDAT, RDA for sustainability

Errata service:

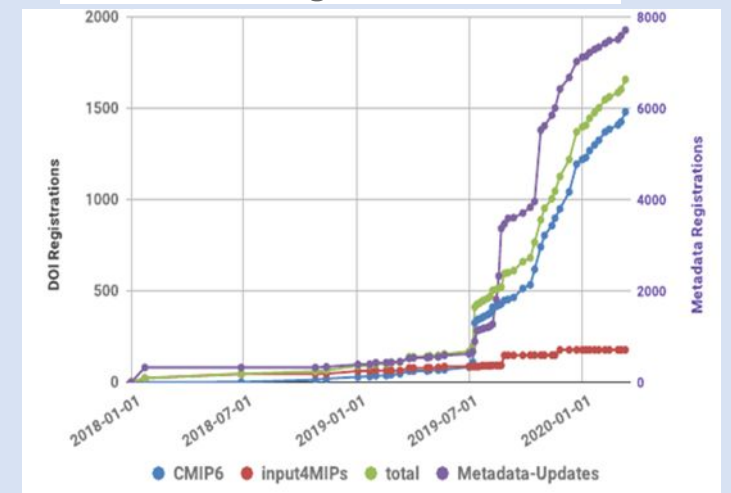
- > 175 issues registered
- > 90.000 data sets affected
3% of CMIP6 archive



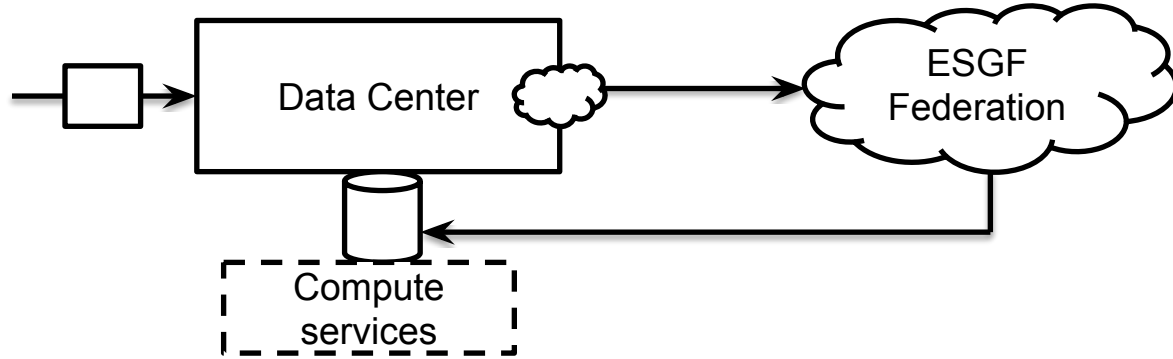
Data citation service:

- > 6000 DOI registrations for CMIP6
- > 8000 citation metadata updates

DataCite DOI Registration for CMIP6

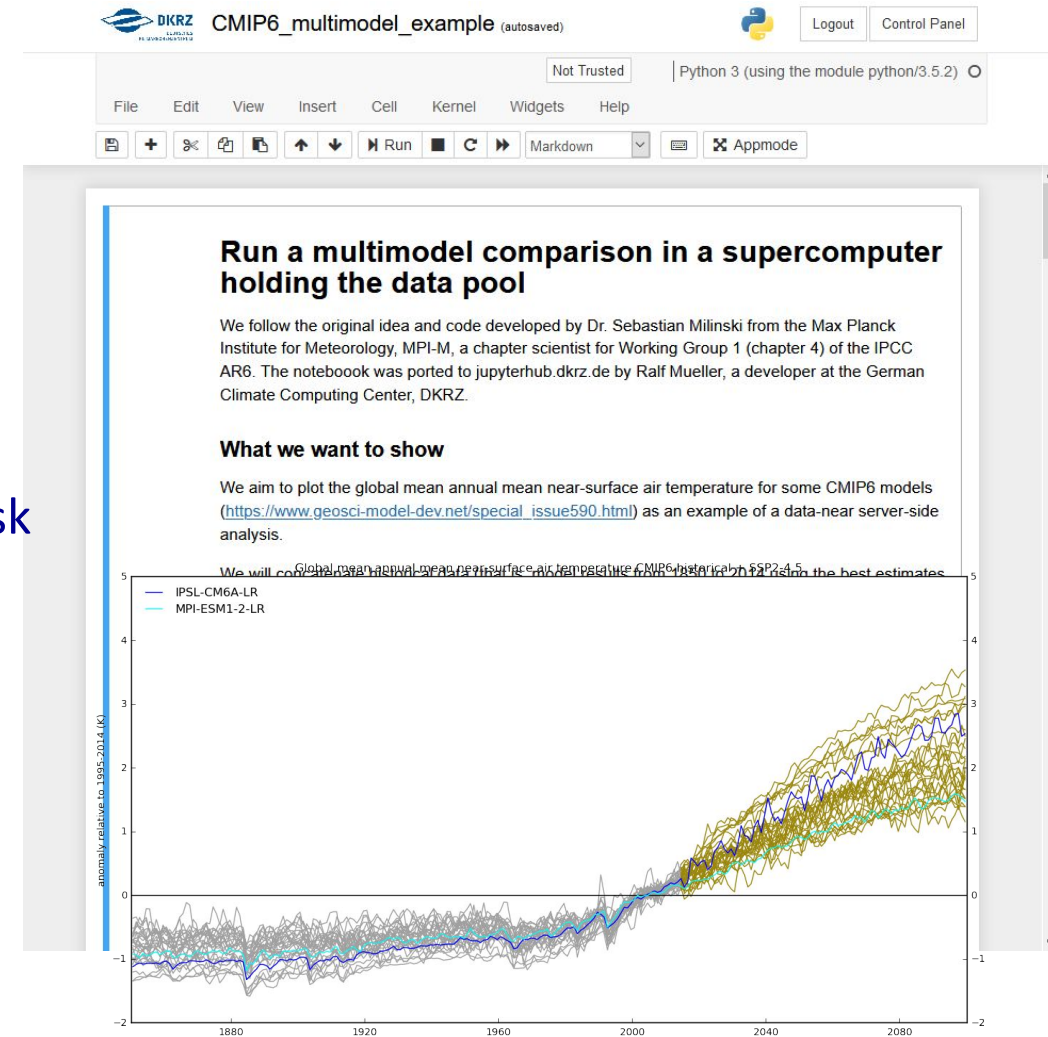


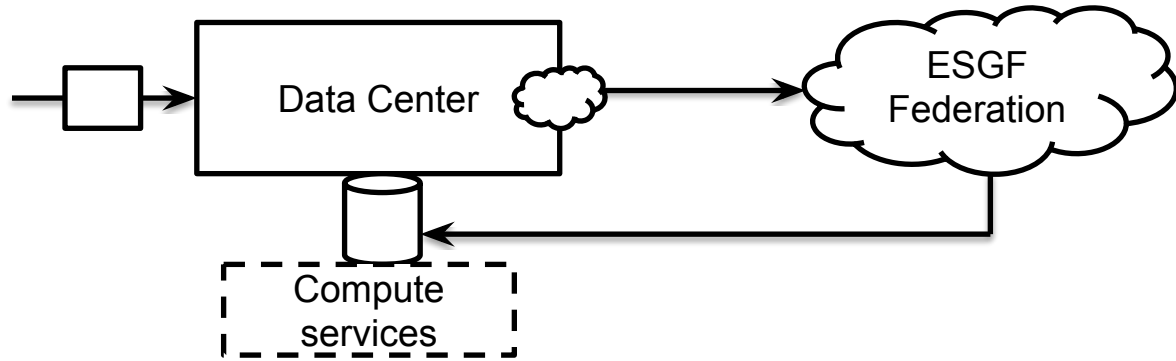
stats: http://bit.ly/CMIP6_Doi_Statistics



Supporting data near processing

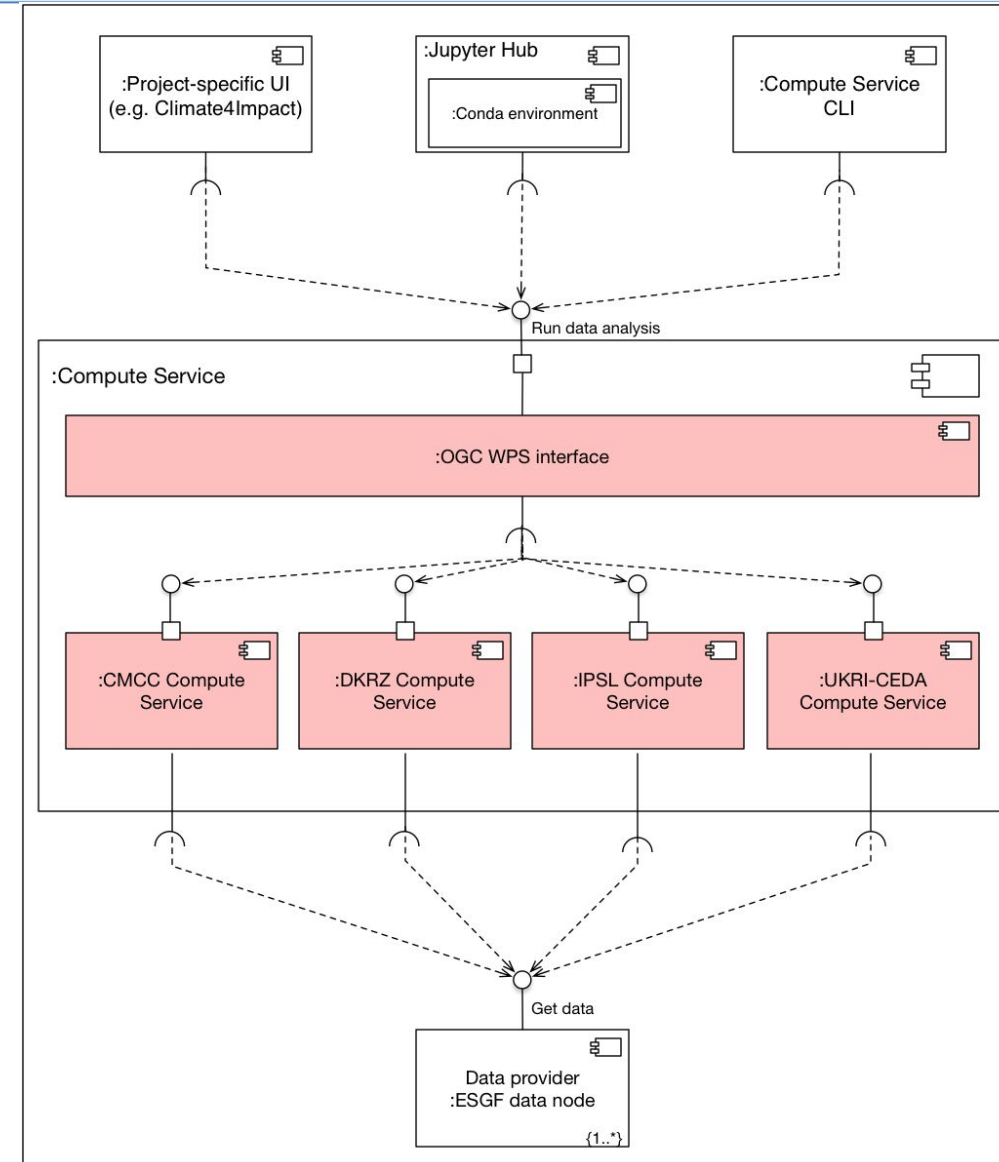
- Prerequisite: Replication of CMIP6 collections into Multi-PByte disk pools (coordinated replication strategy)
- Different processing approaches:
 - Hosting of data near VMs
 - Jupyter-hub installations
 - OGC Web Processing Services (WPS)
 - also for interaction with Copernicus CDS
 - Direct access to virtual workspaces





Supporting the diverse data near processing requirements

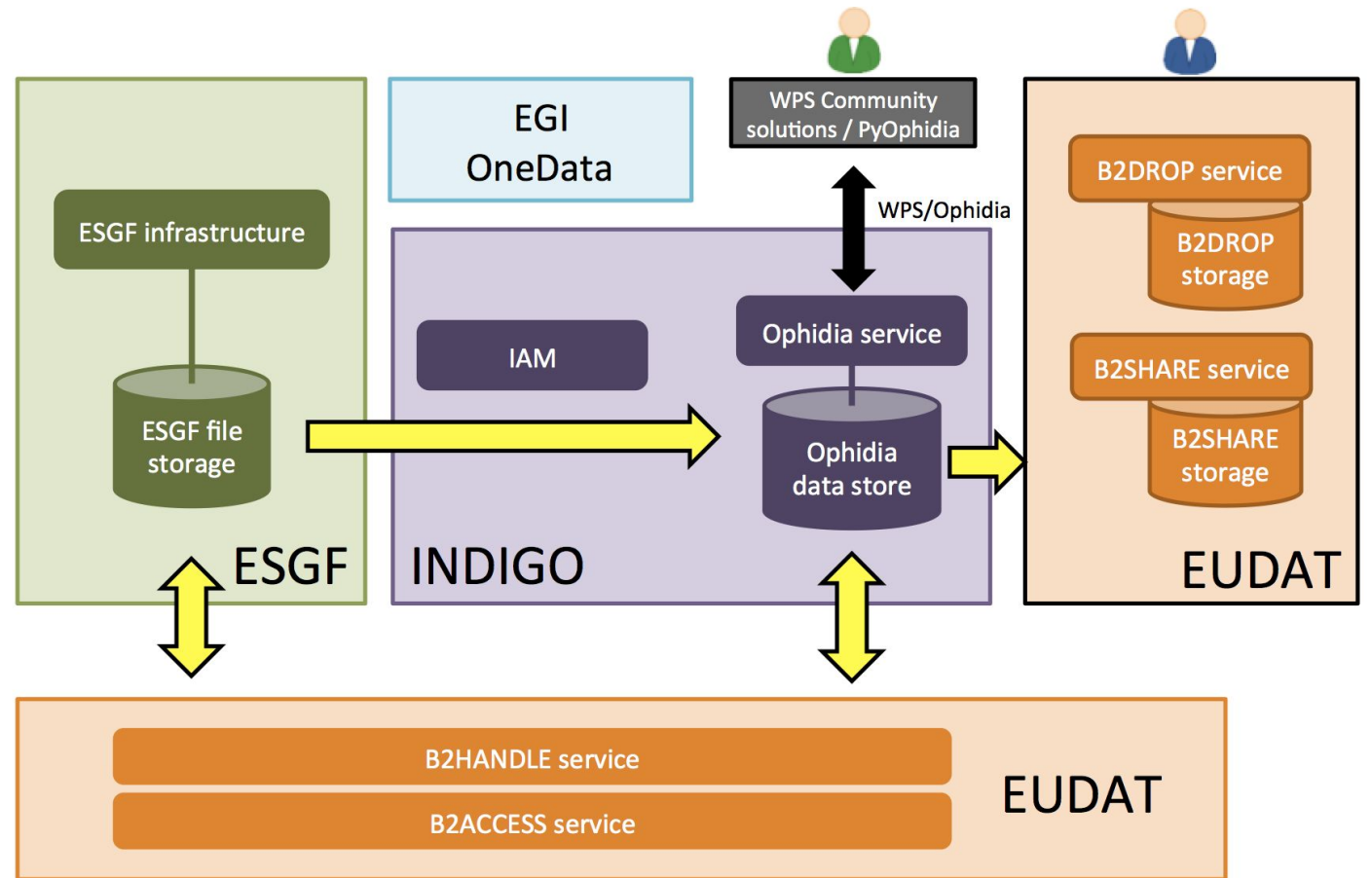
- data center organized data pools
- basic data reduction and transformation services as web processing services
- jupyter-hub installations
- dedicated portals/services for e.g. climate impact community (C4I portal)
- Coordinated service activities at European level (e.g. IS-ENES3, Copernicus)
- different, flexible compute backends (xarray, dask, ophidia, ..)
- ..



The ENES Climate Analytics Service (ECAS), proposed by CMCC & DKRZ in the EU H2020 EOSC-Hub project, supports climate data analysis experiments with a strong focus on data intensive analysis, provenance management, and server-side approaches

It is one of the EOSC-Hub Thematic Services as well as a Compute Service in the IS-ENES3 project

ECAS consists of multiple integrated components from INDIGO-DataCloud, EUDAT, ESGF and EGI, centered around the Ophidia HPDA framework



Moving towards cloud based solutions:

A) data hosting on the cloud

- google cloud hosting of CMIP6 subset ([Pangeo community](#))
 - AWS cloud prototype (US initiated, partners from IS-ENES contributing)
 - institutional pilots
 - data storage based on Zarr → “analysis ready data” → “analytics optimized data store”
- future funding model(s) unclear

B) Individual infrastructural component hosting

- centralized search index
- centralized portal and Identity Proxy
- load balancing components (used e.g. to provide >99% availability services for Copernicus)
- ...

THE CONSORTIUM

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



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement N°824084



Our website
<https://is.enes.org/>



Follow us on Twitter !
@ISENES_RI



Contact us at
is-enes@ipsl.fr



Find our videos on
our channel !

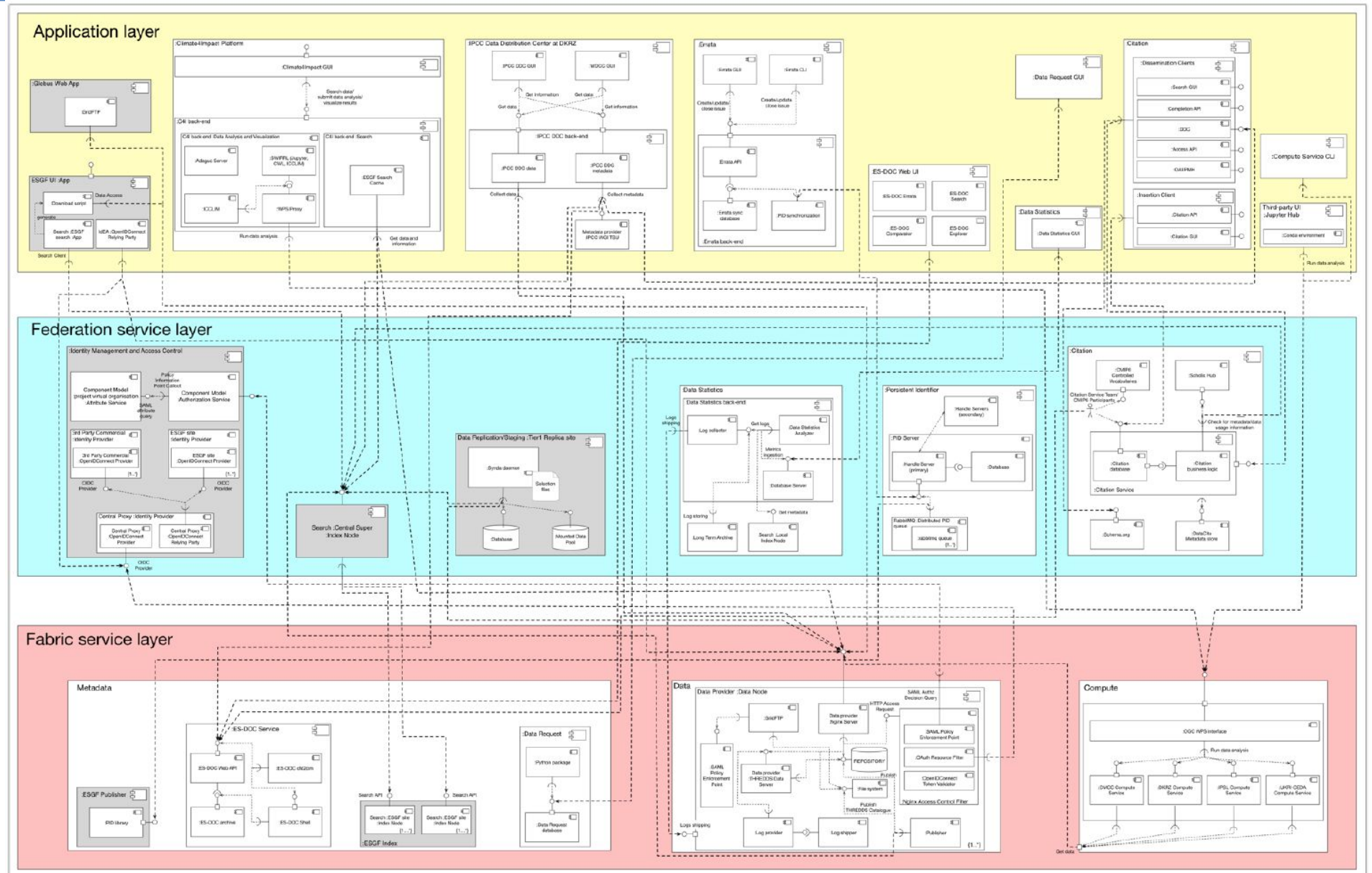


Join our ZENODO
community !

Additional slides

Current architecture:

- complex dependencies
- intervenen partially maintained components



The screenshot shows a JupyterLab window with a file browser on the left and a code editor on the right. The code editor contains the following Python code:

```
[1]: d='/badc/cmip6/data/CMIP6/HighResMIP/MOHC/HadGEM3-GC31-HH/control-1950/r1i1p1f1/day/ta/gn/v20180927'
```

```
[8]: ds = xr.open_mfdataset(f'{d}/*.nc')
```

```
[7]: import xarray as xr
```

```
[14]: import matplotlib.pyplot as plt
import cartopy.crs as ccrs
```

```
[20]: ds.ta[0,0,:]
```

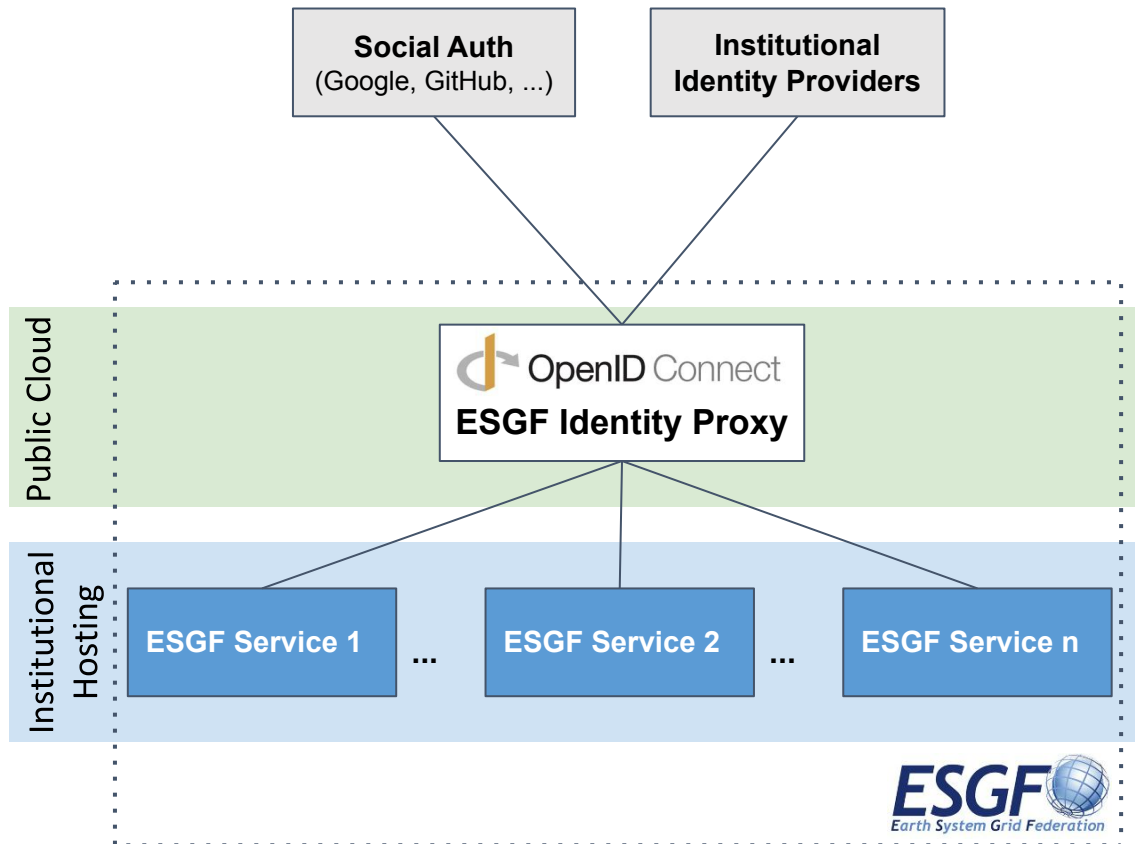
```
[20]: <xarray.DataArray 'ta' (lat: 768, lon: 1024)>
dask.array<shape=(768, 1024), dtype=float32, chunks=(768, 1024)>
Coordinates:
  plev      float64 1e+05
  * lat     (lat) float64 -89.88 -89.65 -89.41 -89.18 ... 89.41 89.65 89.88
  * lon     (lon) float64 0.1758 0.5273 0.8789 1.23 ... 358.8 359.1 359.5 359.8
  time      object 1950-01-01 12:00:00
Attributes:
  standard_name: air_temperature
  long_name:     Air Temperature
  comment:      Air Temperature
  units:        K
  original_name: mo: (stash: m01s30i294, blev: [1000.0, 850.0, 700.0, 500....
  cell_methods: time: mean
  cell_measures: area: areacella
```

```
[26]: %time
fig = plt.figure(figsize=(20, 10))
ax = fig.add_subplot(1,1,1,projection=ccrs.PlateCarree())
ax.set_global()
ax.stock_img()
ax.plot(-0.08, 51.53, 'o', transform=ccrs.PlateCarree())
plt.contourf(ds.lon, ds.lat, ds.ta[0,2,:], 60, transform=ccrs.PlateCarree())
ax.coastlines()
plt.show()
```

The output of the code is a contour plot showing air temperature over the Arctic region. The plot uses a PlateCarree projection and shows a color scale from blue (colder) to red (warmer). A red dot is plotted at approximately 51.53°N, -0.08°W. The plot shows a clear temperature gradient across the Arctic region.

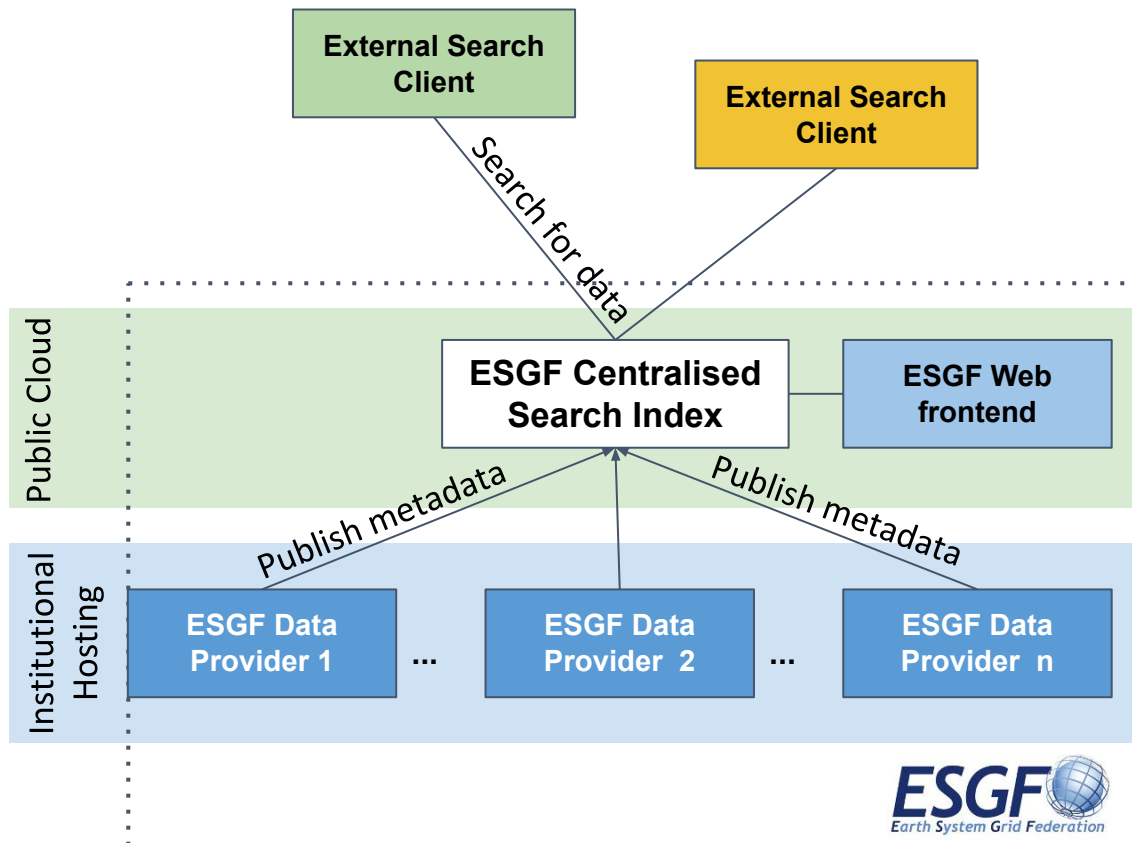
- Jupyter Notebooks enable users to perform analysis of large data holdings at centres like DKRZ and CEDA (JASMIN facility)
- Avoids the need for “download and process at home” model
- Also provides potential as a training tool
- Notebooks could be a shared resource for the community demonstrating different analyses of data

Identity Proxy Model



- Identity management and access control is needed in ESGF to secure access to restricted datasets and other resources
- The current system uses a many-to-many relationship between identity services and those services that are dependent on them
- The European AARC2 Project provides an alternative model
- External identity providers are proxied through a centralised Identity Proxy
- The Identity Proxy enables the presentation of a common interface to ESGF services
- It simplifies the interface between ESGF services and external identity providers

Centralised Search Service



- ESGF Search services provide users with the ability to find the model data they are interested in accessing
- The existing architecture for ESGF uses a system of distributed search services at hosting institutions in the federation
- This is flexible but complex to maintain
- A new proposal is for a centralised search service (index) for ESGF
- This would be hosted on public cloud to make it resilient
- Access is simplified for applications because they can go to the one overall service for search queries