

Ophidia big data analytics framework

Dr. Sandro Fiore

Director of the CMCC Advanced Scientific Computing Division

on behalf of the Ophidia Team



The Ophidia project

Ophidia (<http://ophidia.cmcc.it>) is a CMCC Foundation research project addressing big data challenges for eScience

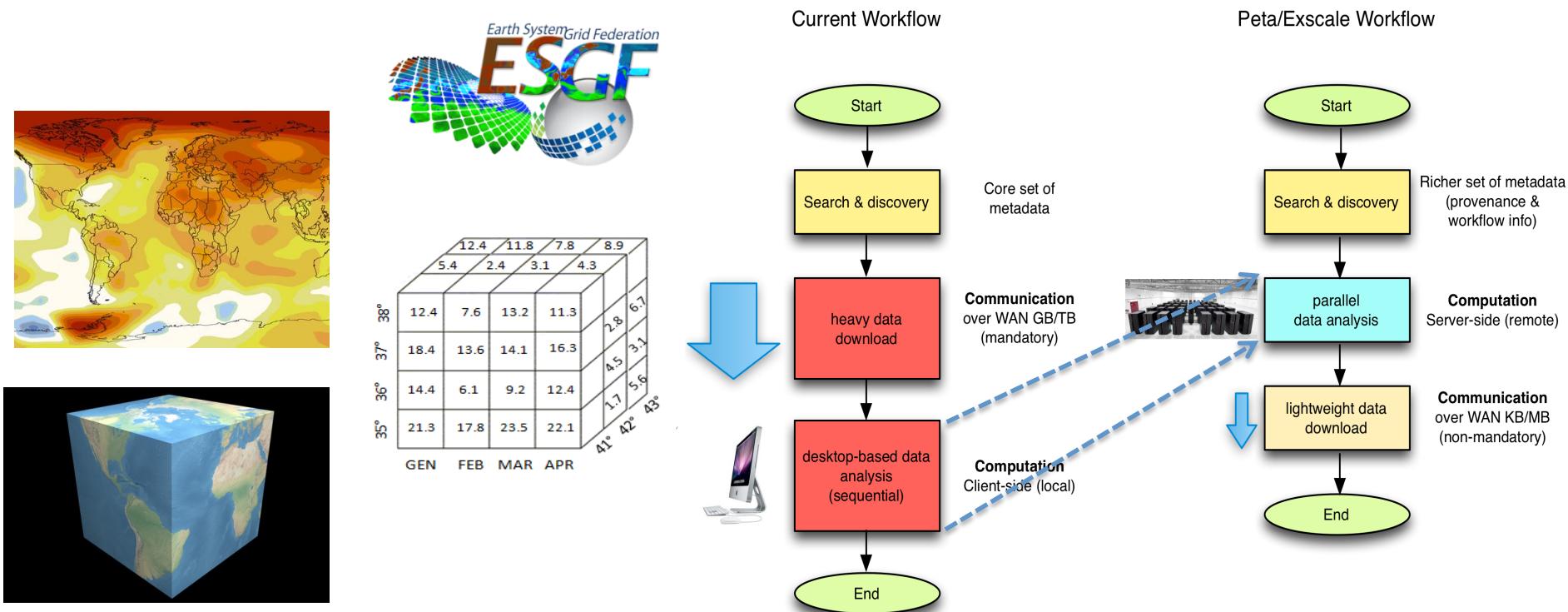
It provides support for declarative and server-side data analysis exploiting high performance computing paradigms and database approaches

Exploits a multidimensional data model providing the data cube abstraction for access and analysis of scientific n-dimensional data



Paradigm shift from client- to server-side

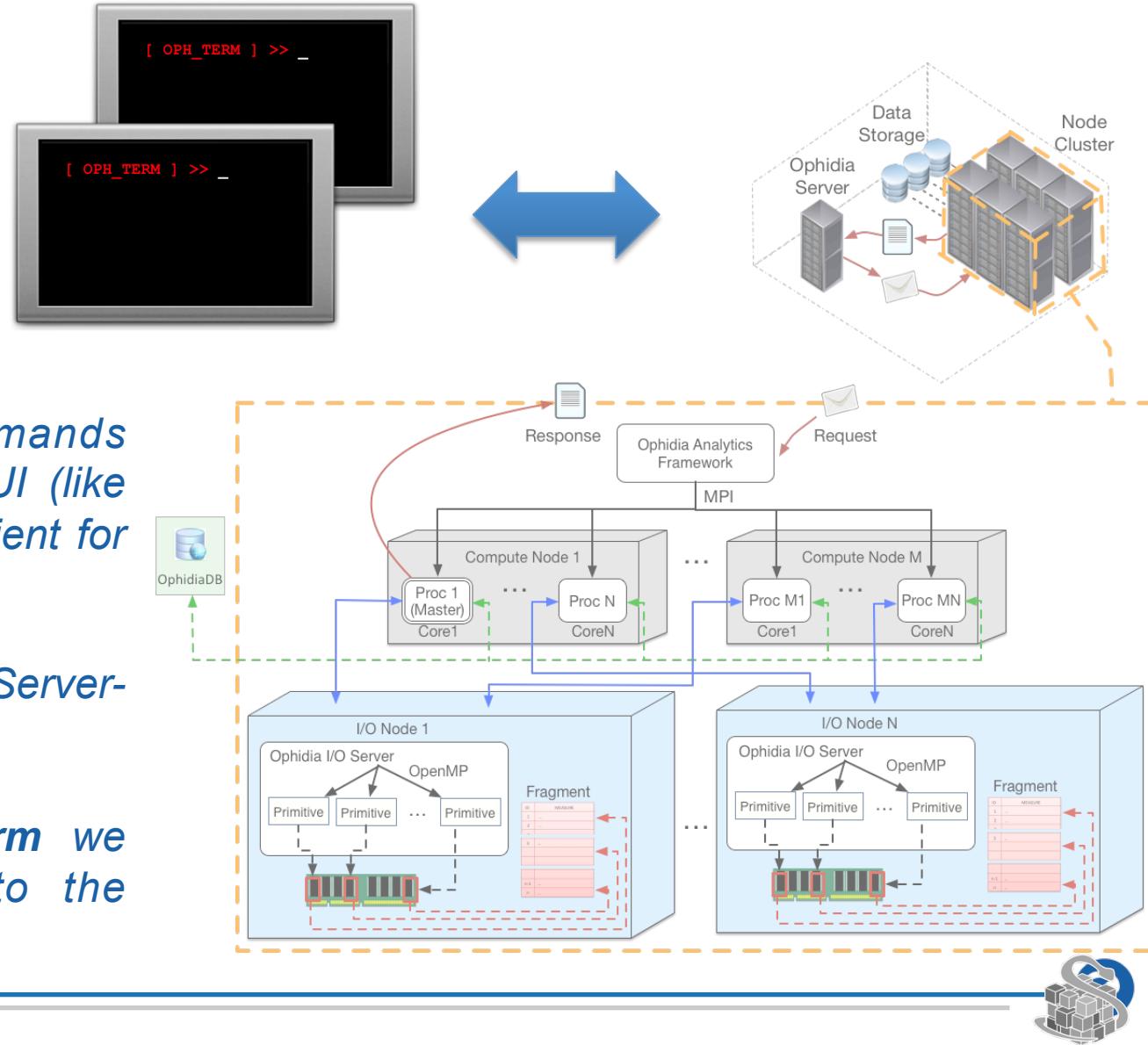
Volume, variety, velocity are key challenges for big data in general and for climate change science in particular. Client-side, sequential and disk-based workflows are three limiting factors for the current scientific data analysis



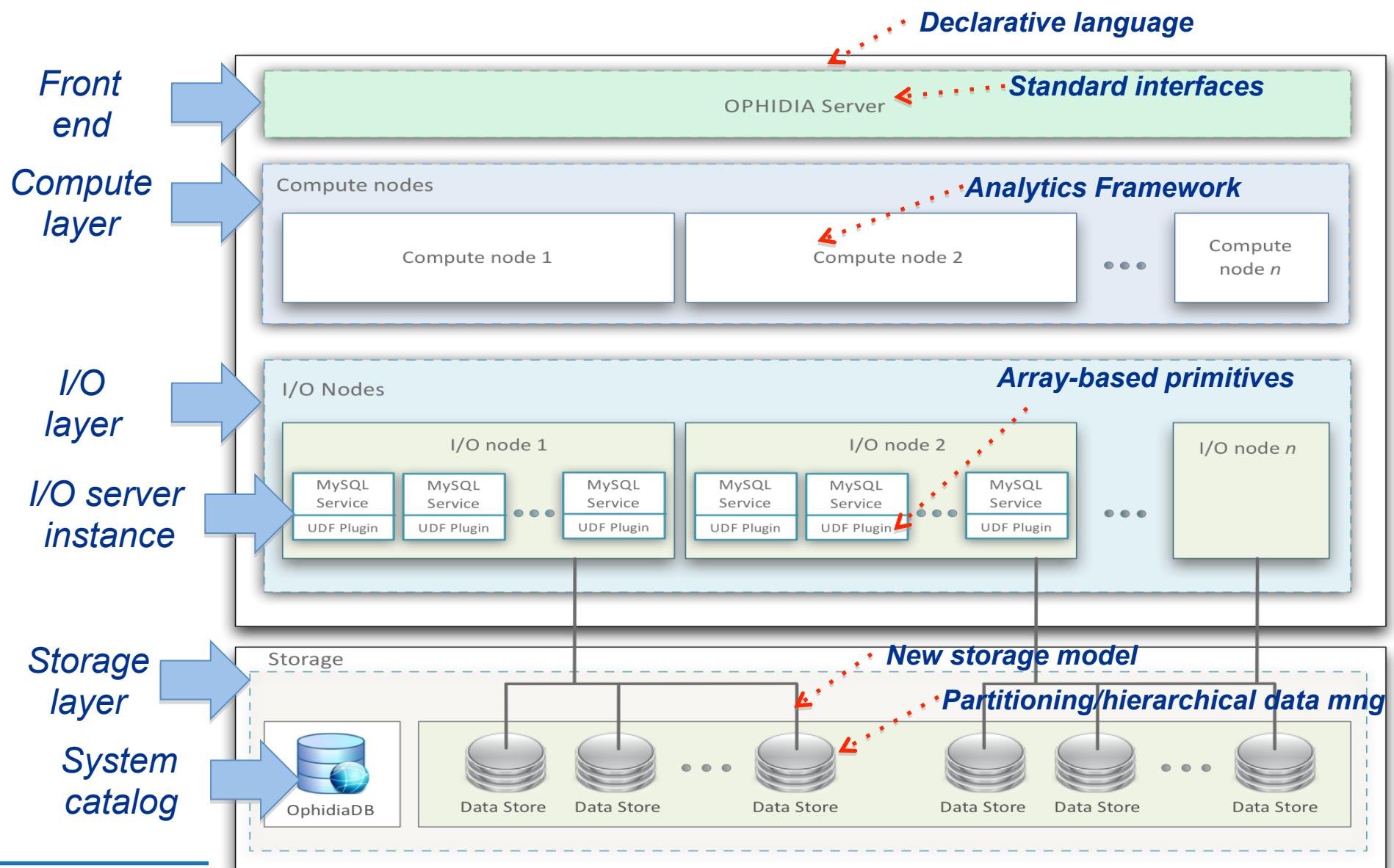
S. Fiore, A. D'Anca, C. Palazzo, I. Foster, D. N. Williams, G. Aloisio, "Ophidia: toward bigdata analytics for eScience", ICCS2013 Conference, Procedia Elsevier, Barcelona, June 5-7, 2013



Ophidia Architecture: end-user view



Ophidia Architecture (sw stack view)

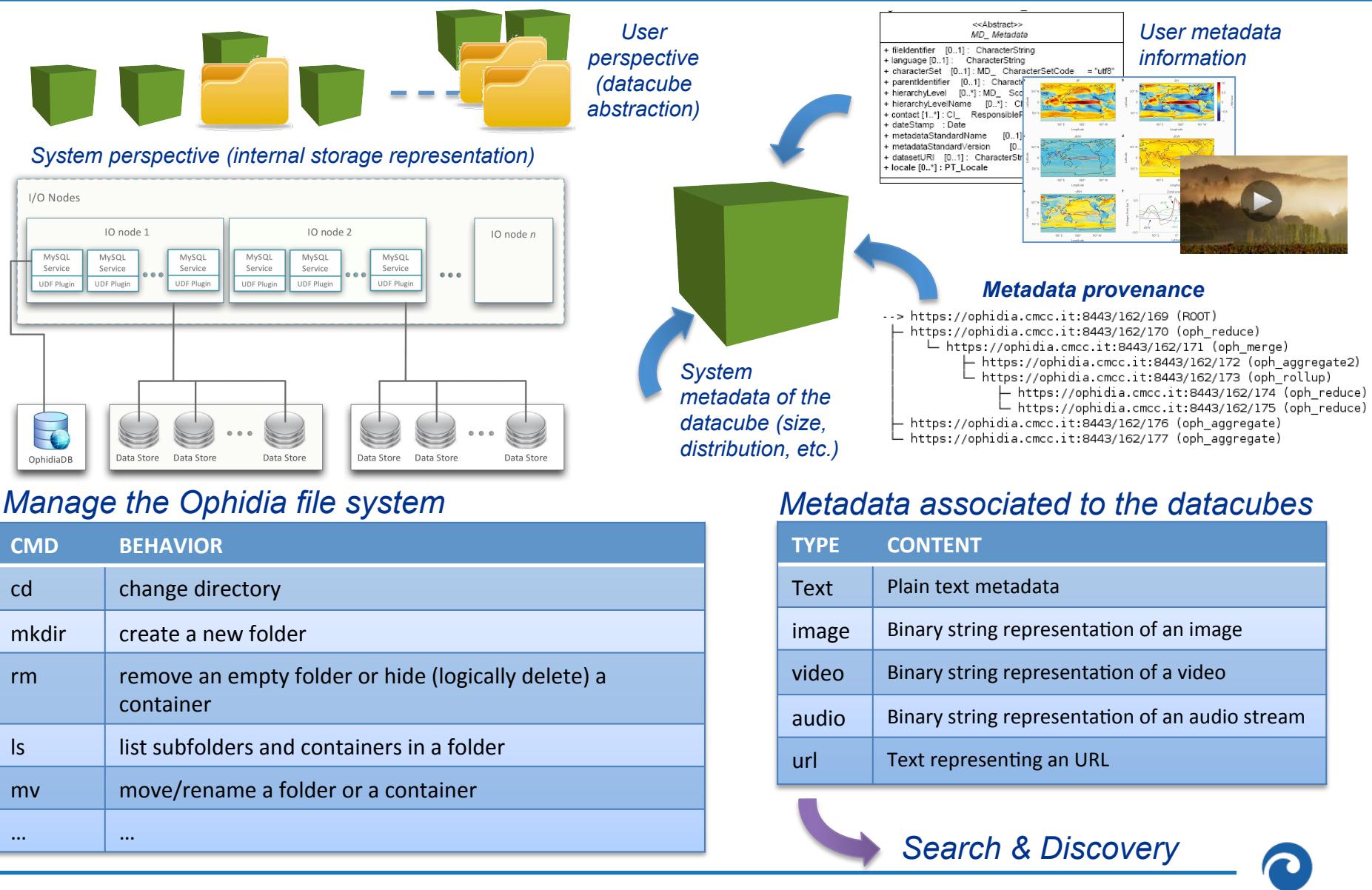


Ophidia in a nutshell

- ✓ *Big data stack for scientific data analysis*
- ✓ **Features:** *time series analysis (array-based analysis), data subsetting (by value/index), data aggregation, model intercomparison, OLAP, etc.*
- ✓ *Use of parallel operators and parallel I/O*
- ✓ **Support for complex workflows / operational chains**
- ✓ *Extensible: simple API to support framework extensions like new operators and array-based primitives*
 - ✓ *currently 50+ operators and 100+ primitives provided*
- ✓ **Multiple interfaces available (WS-I, GSI/VOMS, OGC-WPS).**
- ✓ *Programmatic access via C and Python APIs*
- ✓ *Support for both batch & interactive data analysis*

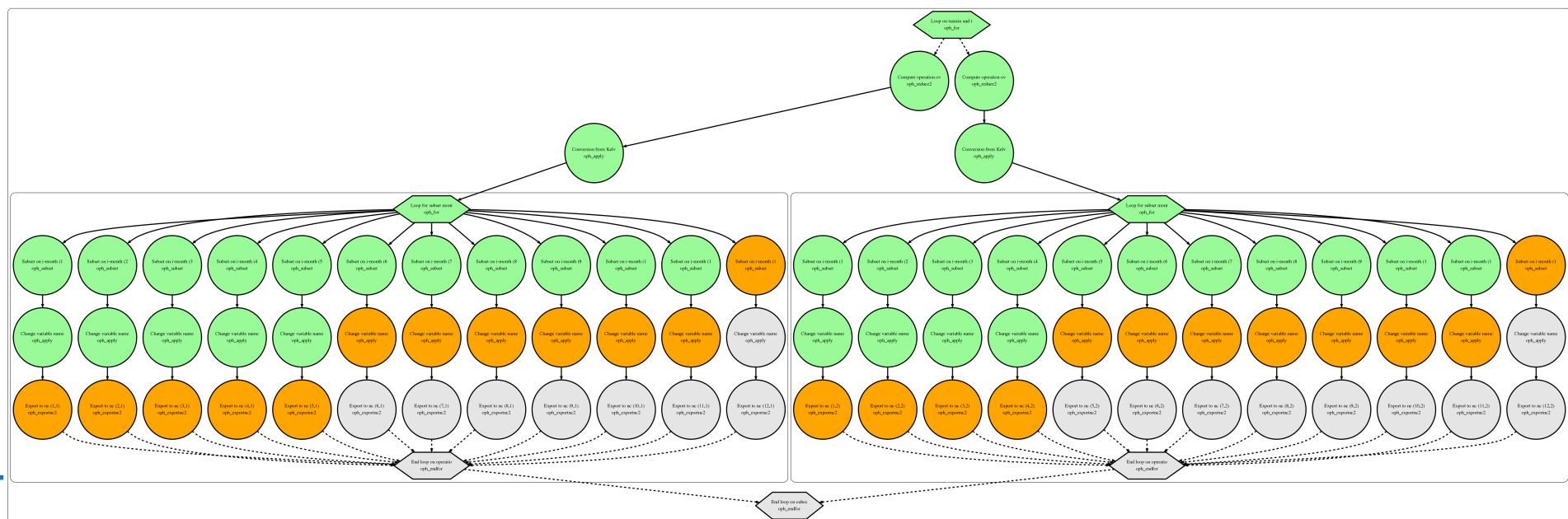
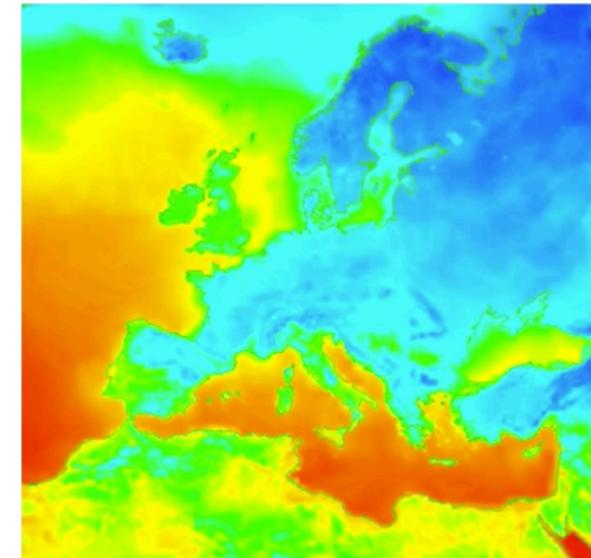


Data abstraction: cube space perspective (OLAP)

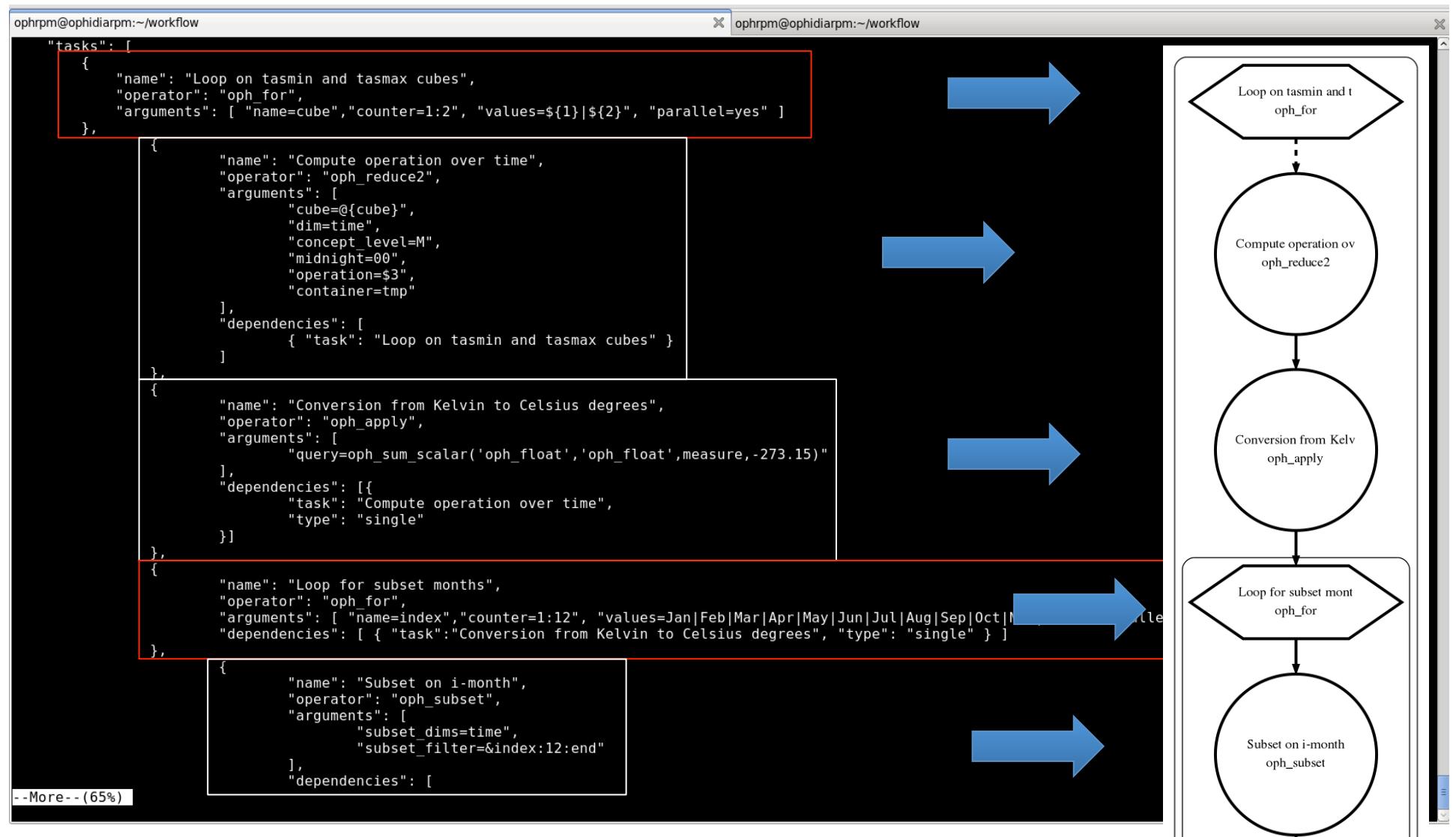


Experiment: climate indicators processing

- ✓ In the CLIPC project, processing chains for data analysis are being implemented with Ophidia to compute **climate indicators**
 - ✓ First set of indicators includes: **TNn**, **TNx**, **TXn**, **TXx**
 - ✓ Input files: 12GBs (TasMin & TasMax)
 - ✓ Workflows have been already implemented
 - ✓ Parallel approach
 - ✓ Inter-parallelism & Intra-parallelism



Workflow JSON representation



Workflow submission

```
ophrpm@ophidiarpm:~/devel/oph-client/res          ophrpm@ophidiarpm:~/workflow
[37..6380] >> ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max
[JobID]:
http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144

[37..6380] >> view 247
[247] ./Tind_loop.json http://193.204.199.174/ophidia/29/2046 http://193.204.199.174/ophidia/30/2047 max [http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3144]

[Response]:
Workflow Status
-----
OPH_STATUS_COMPLETED

Workflow Progress
-----
+=====+=====+
| NUMBER OF COMPLETED TASKS | TOTAL NUMBER OF TASKS |
+=====+=====+
| 82 | 82 |
+=====+=====+

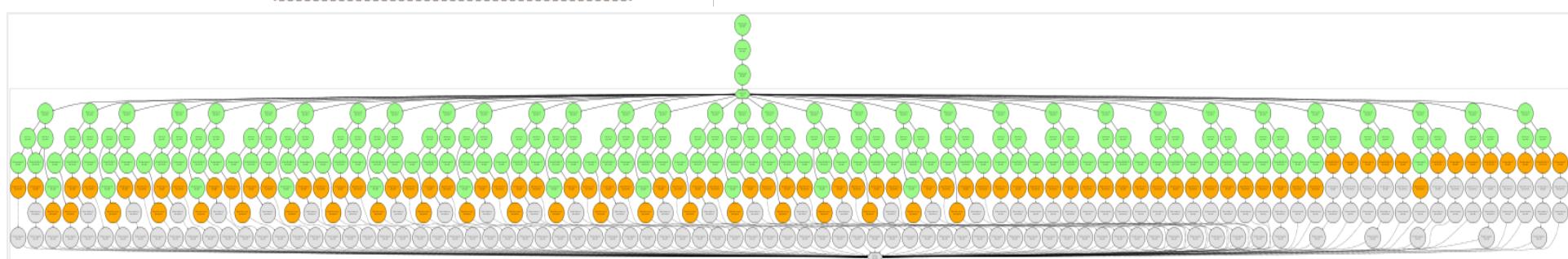
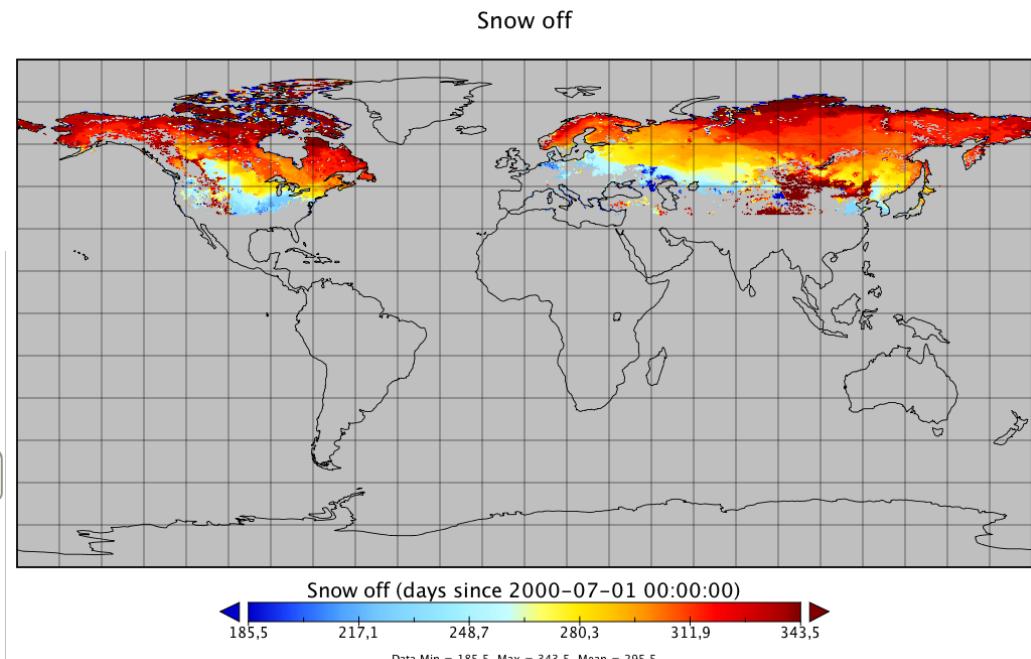
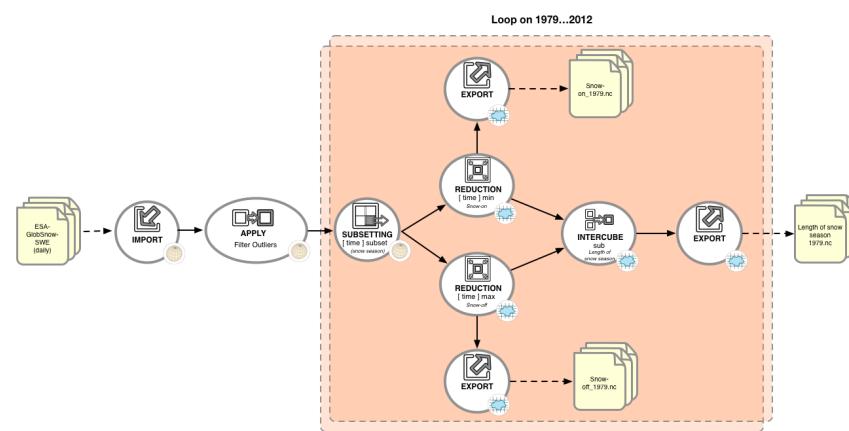
Workflow Task List
-----
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| OPH JOB ID | SESSION CODE | WORKFL | MARKE | PARENT MA | TASK NAME | TYP | EXIT STATUS |  
|           | | OW ID | R ID | RKER ID |           | E |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3145 | 37669923831130223251 | 247 | 3145 | 3144 | Loop on tasmin and tasmax cubes | SIM | OPH_STATUS_PLE |  
| 1449455166146380 | | | | | | | |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3146 | 37669923831130223251 | 247 | 3146 | 3144 | Compute operation over time (1) | SIM | OPH_STATUS_PLE |  
| 1449455166146380 | | | | | | | |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3147 | 37669923831130223251 | 247 | 3147 | 3144 | Compute operation over time (2) | SIM | OPH_STATUS_PLE |  
| 1449455166146380 | | | | | | | |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3148 | 37669923831130223251 | 247 | 3148 | 3144 | Conversion from Kelvin to Celsius degrees (1) | SIM | OPH_STATUS_PLE |  
| 1449455166146380 | | | | | | | |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
| http://193.204.199.174/ophidia/sessions/376699238311302232511449455166146380/experiment?247#3149 | 37669923831130223251 | 247 | 3149 | 3144 | Conversion from Kelvin to Celsius degrees (2) | SIM | OPH_STATUS_PLE |  
| 1449455166146380 | | | | | | | |  
+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+=====+
```



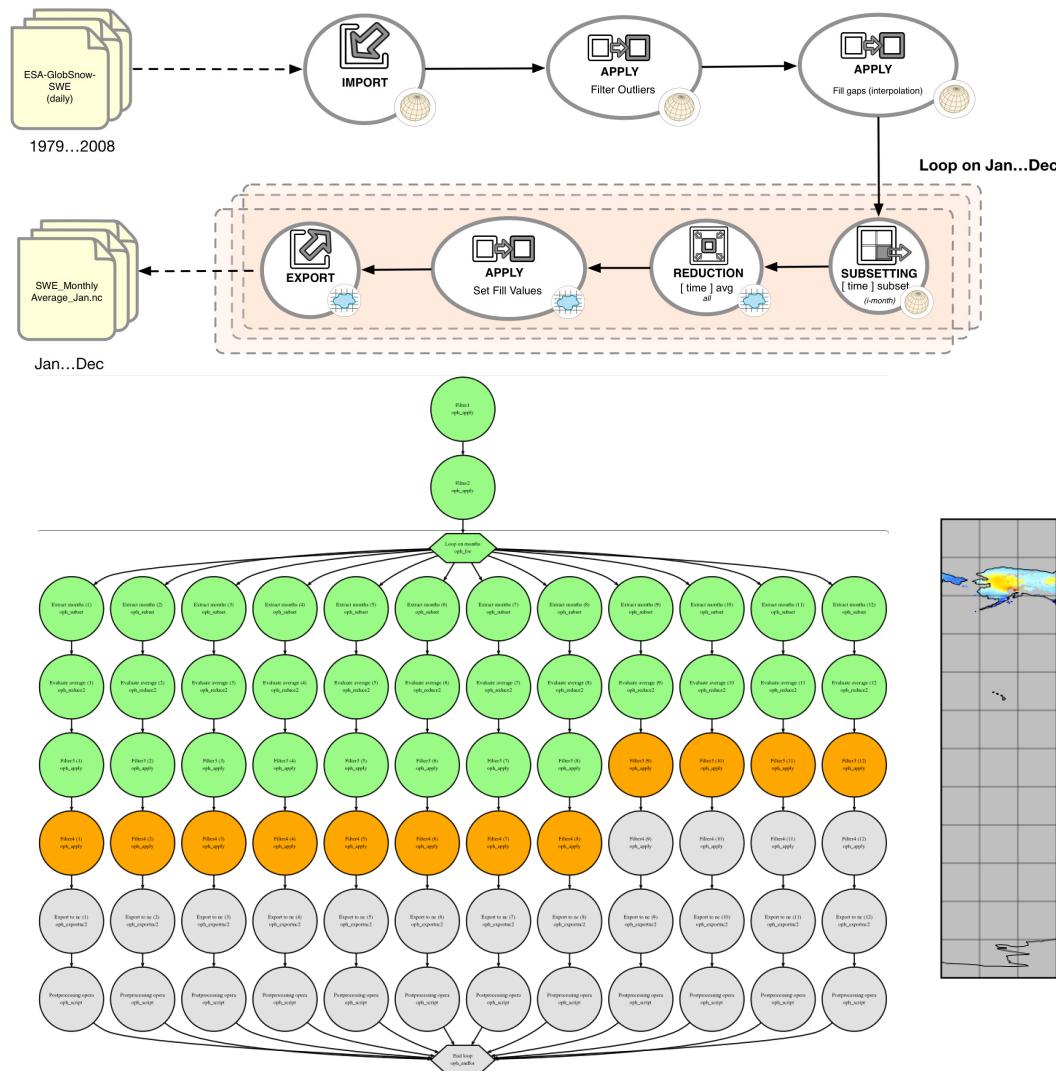
Tier1 indicators using Ophidia

Snow on/off – Length of snow season

- ✓ Dataset time range: 1979-2012
- ✓ 50 GB of input data
- ✓ 434 tasks performed
- ✓ 99 NetCDF output files



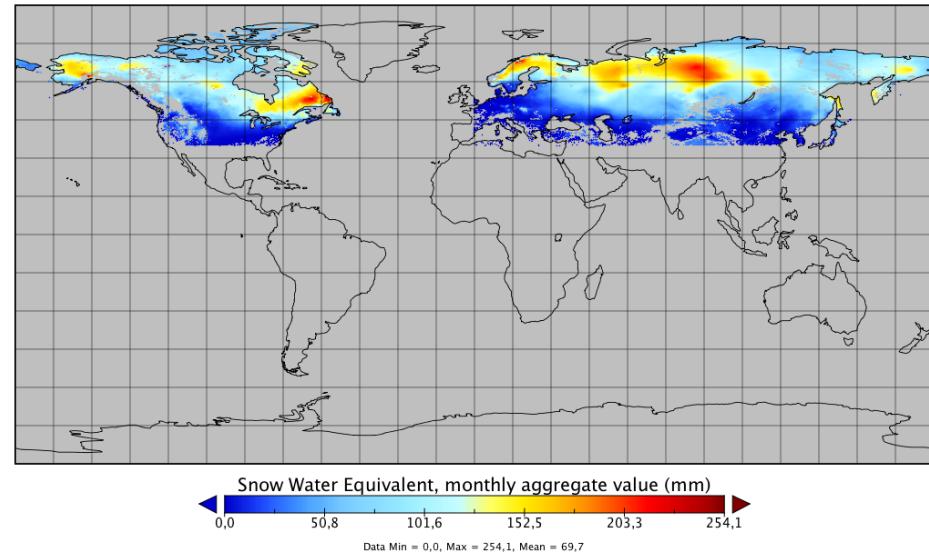
Tier1 indicators using Ophidia



SWE monthly average

- ✓ Dataset time range: 1979-2008
- ✓ 1.7 GB of input data
- ✓ 76 tasks performed
- ✓ 12 NetCDF output files

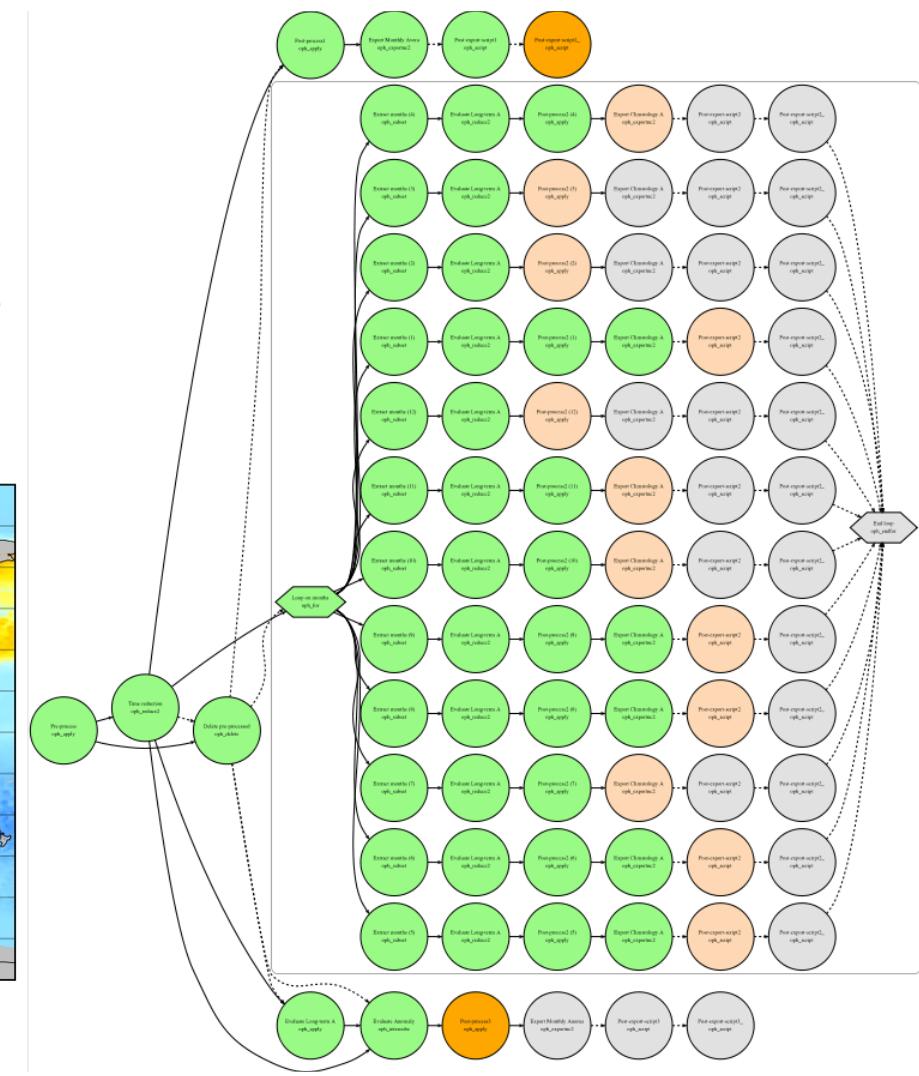
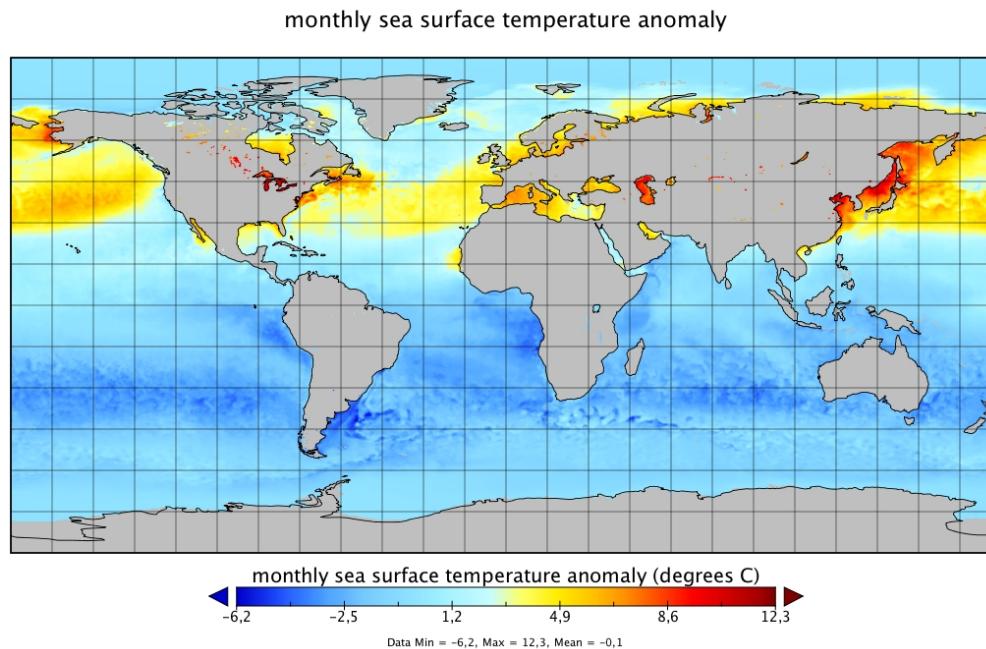
Snow Water Equivalent, monthly aggregate value



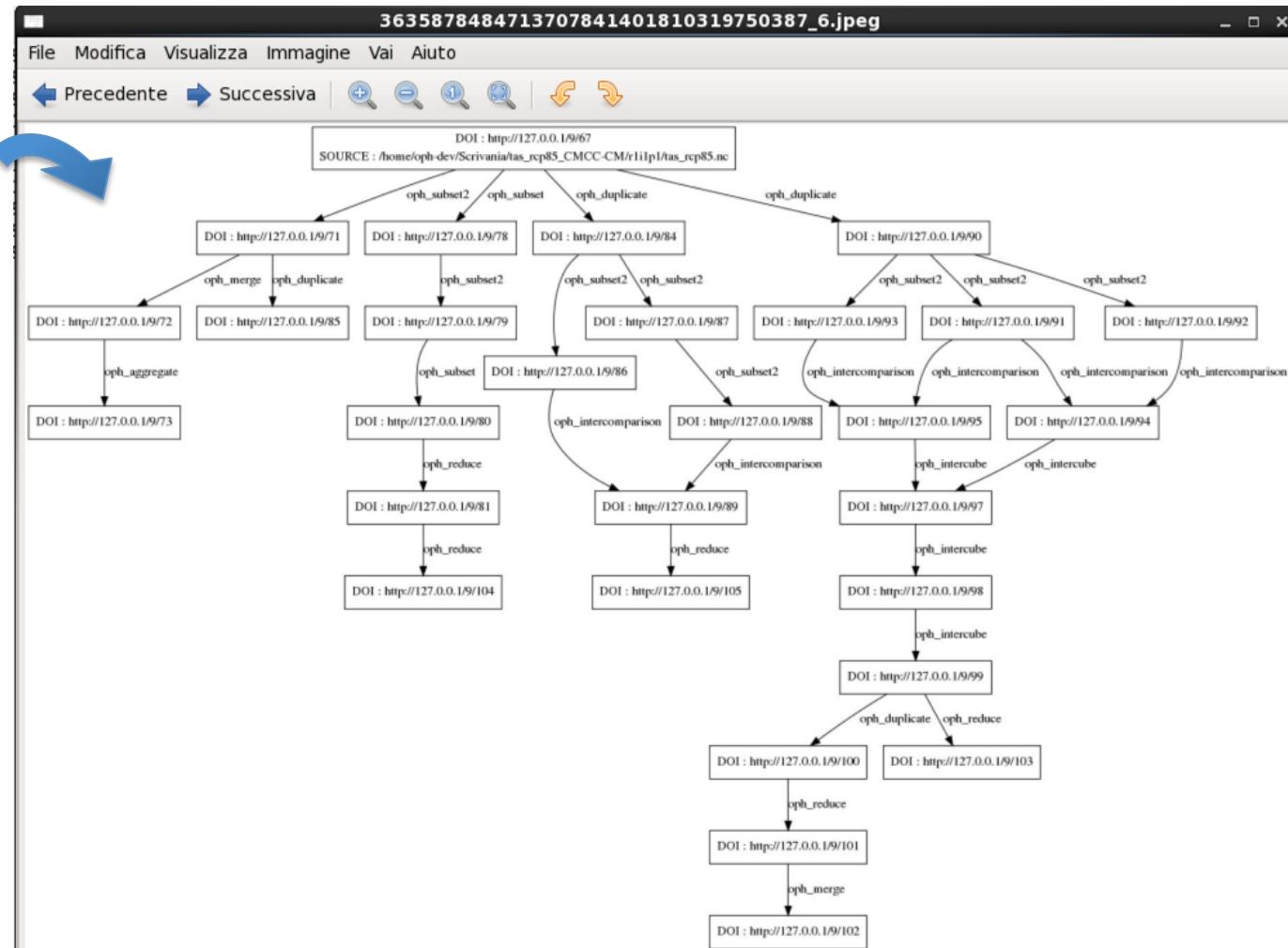
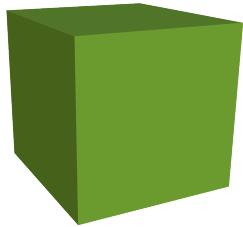
Tier1 indicators using Ophidia

SST mean, anomaly, climatology mean

- ✓ Dataset time range: 1991-2010
- ✓ 350GB of input data
- ✓ 87 tasks performed
- ✓ Expected 12x51MB + 2x12GB of output files



Provenance management (PID-based)



RDA proposal submitted on extending Ophidia with RDA PID recommendation (collaboration with DKRZ)

Ophidia & INDIGO-DataCloud



- An H2020 project approved in January 2015 in the EINFRA-1-2014 call
 - 11.1M€, 30 months (**from April 2015 to September 2017**)
- Who: **26 European partners** in 11 European countries
 - Coordination by the Italian National Institute for Nuclear Physics (INFN)
 - Including developers of distributed software, industrial partners, research institutes, universities, e-infrastructures
- What: **develop an open source Cloud platform** for computing and data ("DataCloud") tailored to science.
- For: **multi-disciplinary scientific communities**
 - E.g. structural biology, earth science, physics, bioinformatics, cultural heritage, astrophysics, life science, climatology
- Where: deployable on **hybrid (public or private) Cloud infrastructures**
 - INDIGO = INtegrating Distributed data Infrastructures for Global Exploitation
- Why: answer to the technological **needs of scientists** seeking to easily exploit distributed Cloud/Grid compute and data resources.



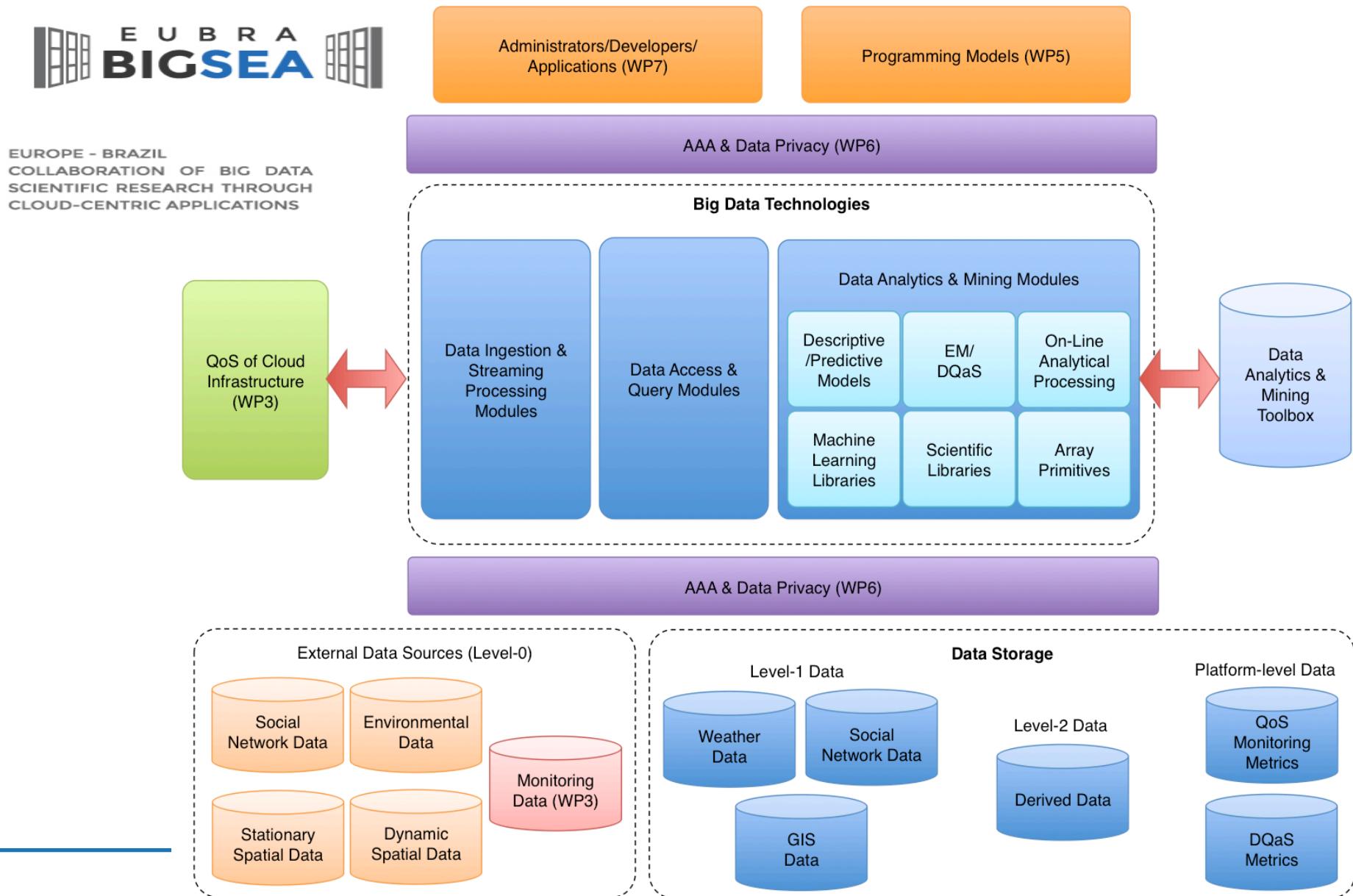
Testbed environment for running multi-model experiments (precipitation trend analysis)



More on this topic in the afternoon talk at 5.30pm



BIGSEA: cloud and QoS based vertical and horizontal elasticity for big data systems



Resources (<http://ophidia.cmcc.it>)

The screenshot shows the Ophidia website as it would appear in a web browser. At the top, the browser's menu bar includes 'Safari', 'File', 'Edit', 'View', 'History', 'Bookmarks', 'Window', 'Help', and various system status icons like signal strength and battery level. The address bar shows the URL 'ophidia.cmcc.it'. Below the browser window, the Ophidia homepage features a large header with the 'Ophidia' logo (an orange stylized 'S' over blue 3D blocks) and the word 'Ophidia'. A navigation bar with links to 'Home', 'Overview', 'Download', 'Documentation', 'News', 'Tutorials', and 'About Us' is positioned above a large banner image of a cloudy sky. The banner also contains the Ophidia logo. Below the banner, the text 'High Performance Data Mining & Analytics for eScience' is displayed. The page is divided into four main sections: 'PARALLEL' (Parallel computing approach for data analytics), 'SCIENTIFIC' (Analytics framework for scientific data management), 'EXTENSIBLE' (API available to enable end-users extensions), and 'SERVER-SIDE' (Remote data processing based on standard interfaces). Each section has a corresponding icon and a 'Learn more' button.

High Performance Data Mining & Analytics for eScience

PARALLEL
Parallel computing approach for data analytics
[Learn more](#)

SCIENTIFIC
Analytics framework for scientific data management
[Learn more](#)

EXTENSIBLE
API available to enable end-users extensions
[Learn more](#)

SERVER-SIDE
Remote data processing based on standard interfaces
[Learn more](#)

Ophidia is a CMCC Foundation research project addressing big data challenges for eScience. It provides support for data-intensive analysis exploiting advanced parallel computing techniques and smart data distribution methods. It exploits an array-based storage model and a hierarchical storage organisation to partition and distribute multidimensional scientific datasets over multiple nodes. The Ophidia analytics framework can be exploited in different scientific domains (e.g. Climate Change, Earth Sciences, Life Sciences) and with very heterogeneous sets of data.

Resources (II)

The image displays three separate screenshots illustrating various Ophidia resources:

- Top Left (Safari Browser):** Shows the Ophidia documentation page at ophidia.cmcc.it/documentation/users/operators/. The page includes a navigation bar with links like File, Edit, View, History, Bookmarks, Window, and a search bar.
- Top Right (Python Package Page):** Shows the PyOphidia package page on PyPI. The URL is <https://pypi.org/project/PyOphidia/>. The page title is "PyOphidia 1.2.1". It describes PyOphidia as a "Python bindings for the Ophidia Data Analytics Platform". The package is licensed under GPLv3. A sidebar on the right shows a "Not Logged In" section with links for Login, Register, Lost Login?, Use OpenID, and Login with Google.
- Bottom (YouTube Channel):** Shows the Ophidia YouTube channel page. The URL is <https://www.youtube.com/channel/UCtPjyfXWzJLcIwzgkVQDwgg>. The channel has 1,000 subscribers and 1,000 views. It features four video thumbnails related to the Data Analytics Terminal:
 - Data Analytics Terminal : using aliases**: 1 year ago, 10 views.
 - Data Analytics Terminal : using environment variables**: 1 year ago, 15 views.
 - Data Analytics Terminal : switching between sessions**: 1 year ago, 24 views.
 - Data Analytics Terminal : autocompletion feature**: 1 year ago, 22 views.



Thanks



<http://ophidia.cmcc.it>



@OphidiaBigData



www.youtube.com/user/OphidiaBigData

