



JASMIN (STFC/Stephen Kill)



Software Development at the Centre for Environmental Data Analysis

RAL Site Software Engineering Community Meeting

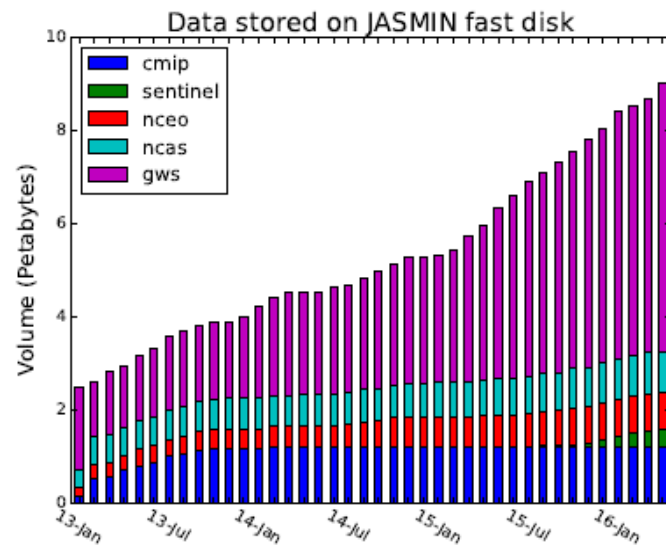
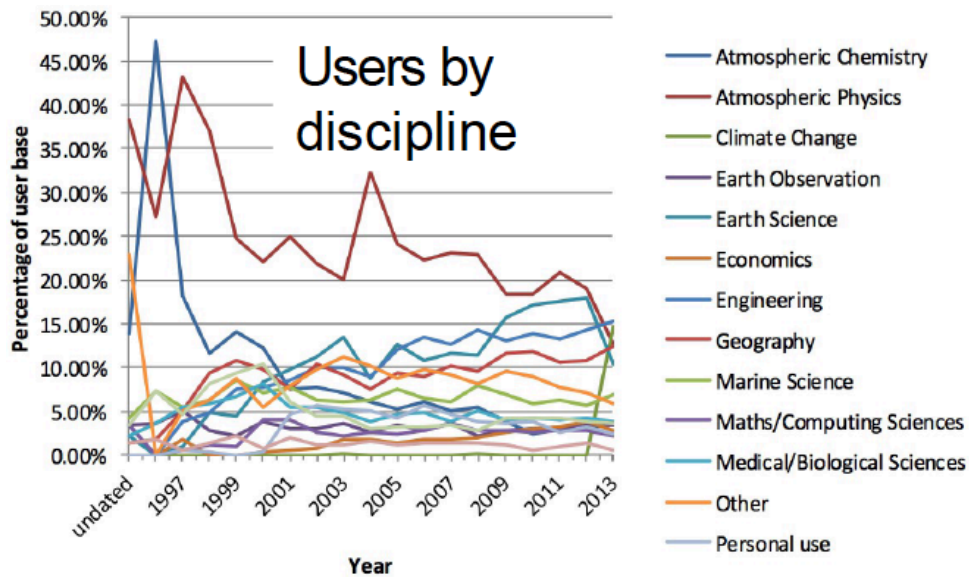
12 February 2018

Philip Kershaw (on behalf of CEDA)

NCAS/NCEO, Centre for Environmental Data Analysis, RAL Space

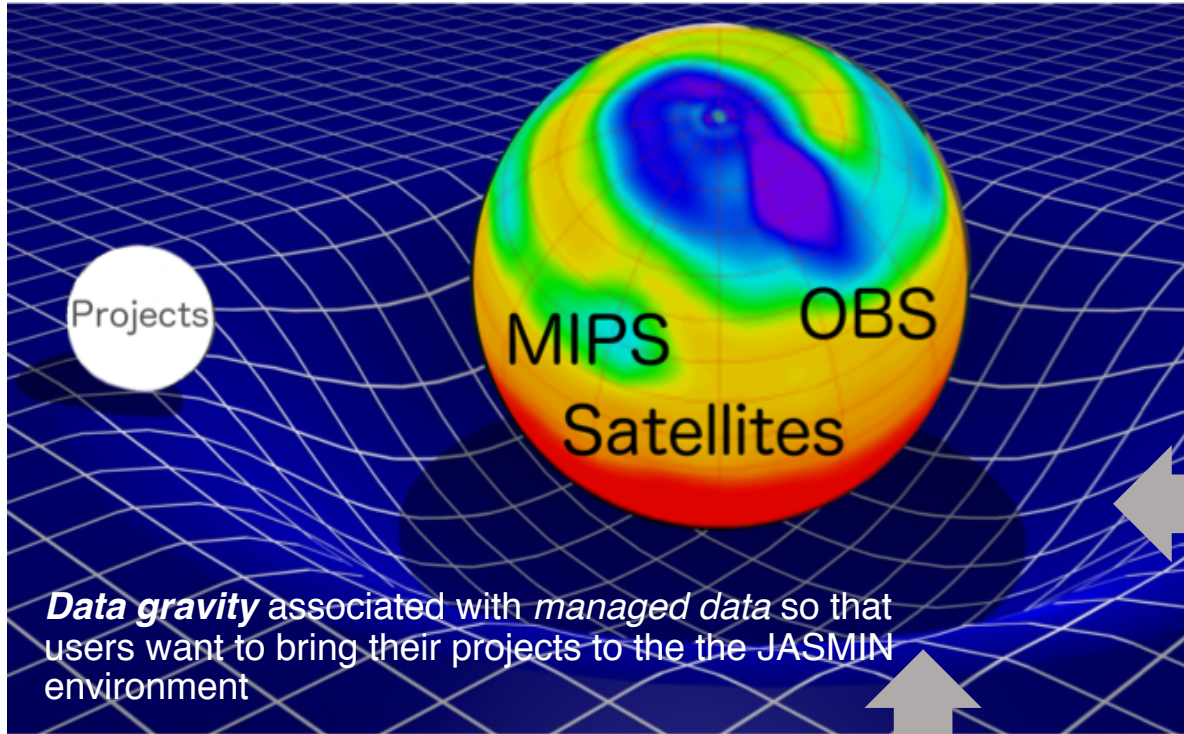
[Thanks and credit to STFC Scientific Computing Department who deploy and operate the JASMIN infrastructure on behalf of CEDA]

Data Growth and Diversification of user community



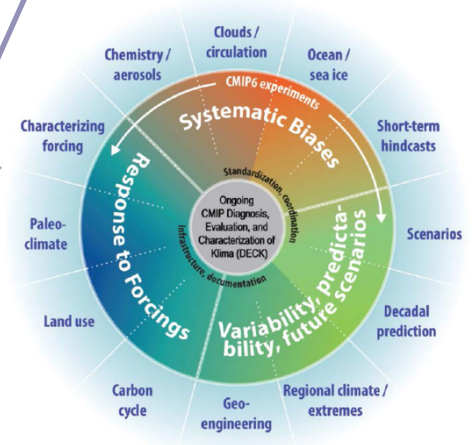
2013-2016 increasing data storage on JASMIN, in Group Workspaces (GWS) and archive

JASMIN as a Data Commons



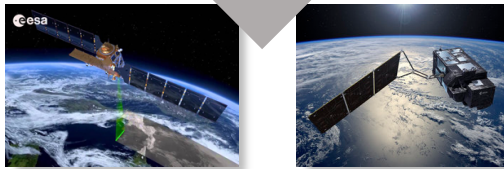
CEDA CREPP system to ingest from MetOffice Hadley Centre

CMIP6



Sentinel Earth

Observation Data

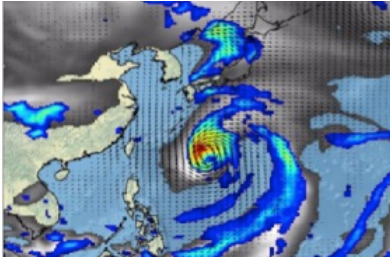


Sentinel missions data rate: ~6PB/year

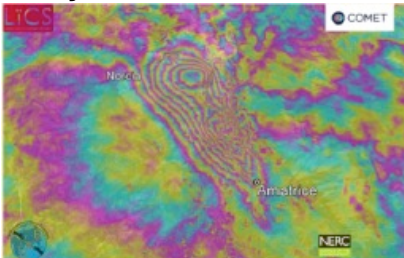
Climate Models

Couple-Model Intercomparison Projects

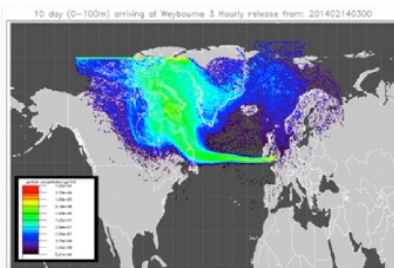
~150 Science projects on JASMIN to date



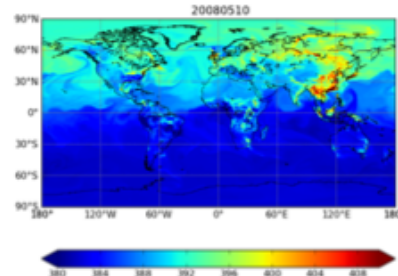
High Res Climate Model analysis



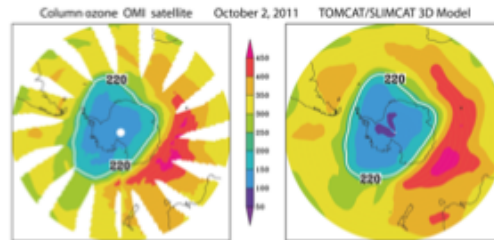
Fault analysis



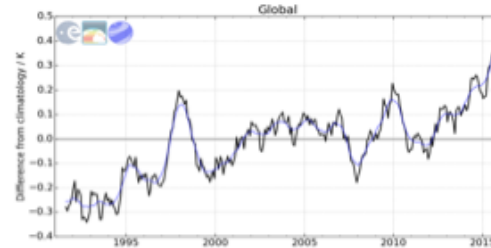
Atmospheric dispersion



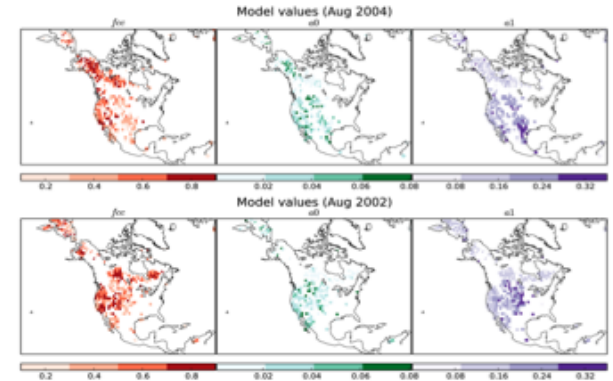
Regional carbon balance on a global scale



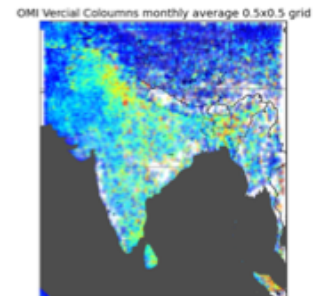
Antarctic Ozone hole: model vs. observations



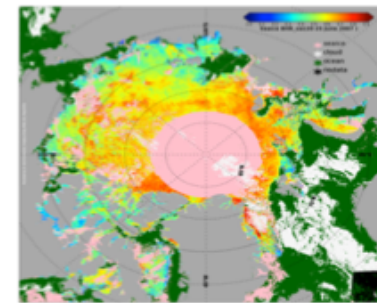
Sea Surface Temperature from satellite observations



Deriving the impact of fire on vegetation from earth observation data



Understanding oxidant chemistry over the Indian subcontinent



Climate variables from European and US instruments/satellites

JASMIN usage: Cloud



ESA Forestry Thematic Exploitation Platform



ESA Climate Change Initiative Open Data Portal



Majic interface to Jules Land-surface model on JASMIN

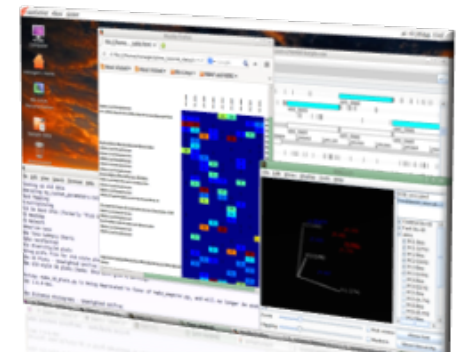


ESA Polar Thematic Exploitation Platform



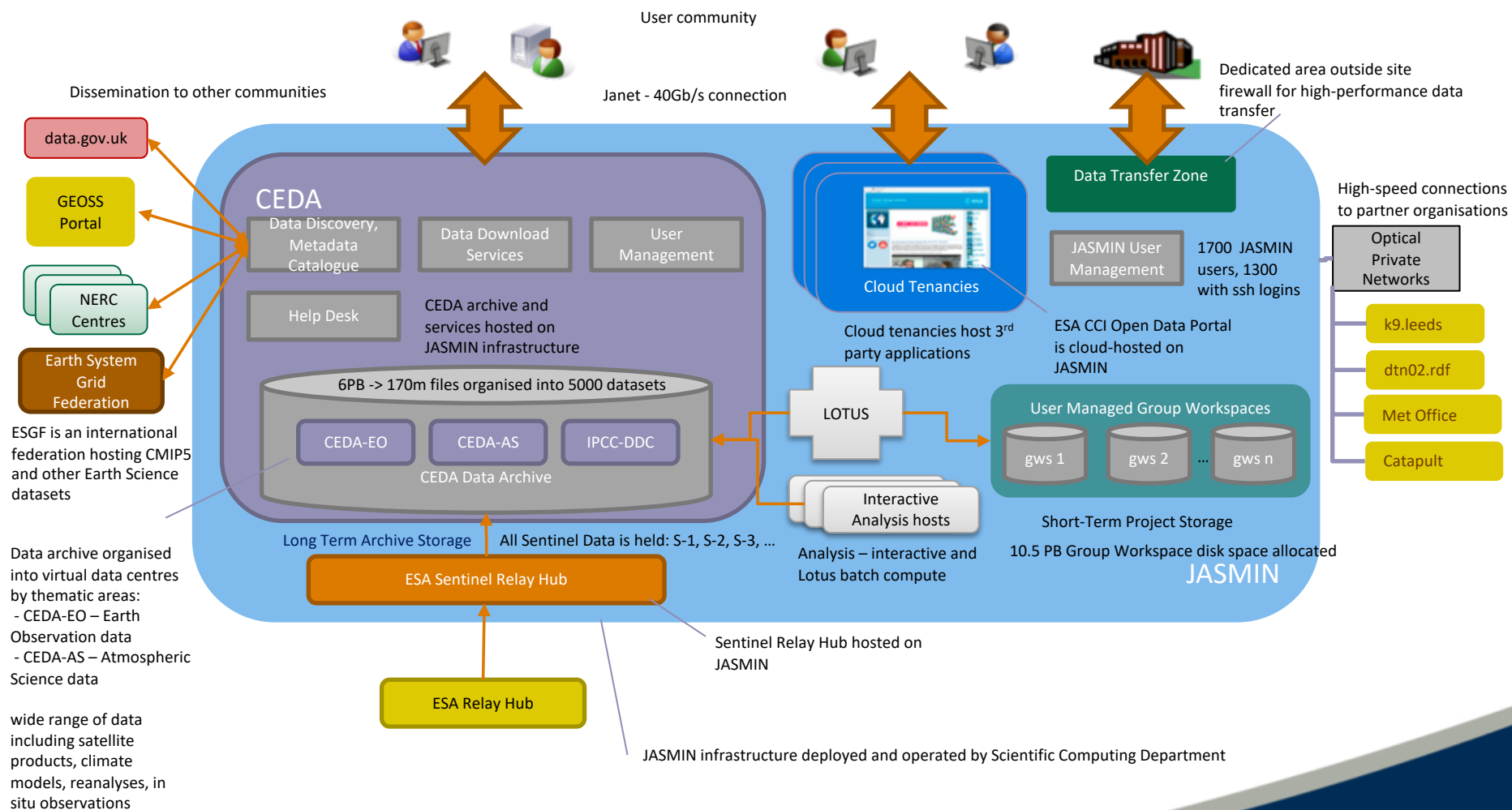
Attendees at ESA Summer school, ESRIN used OPTIRAD Jupyter Notebook environment

– Credit ESA



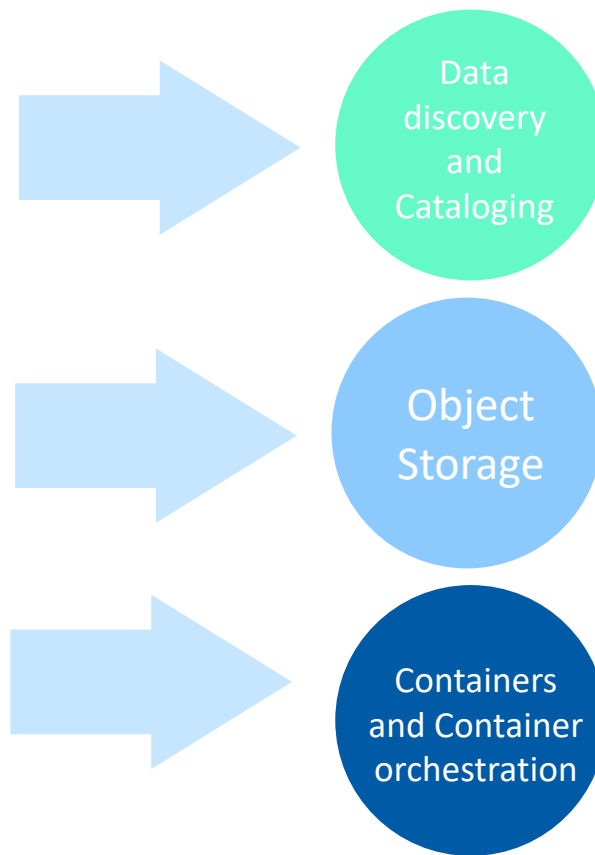
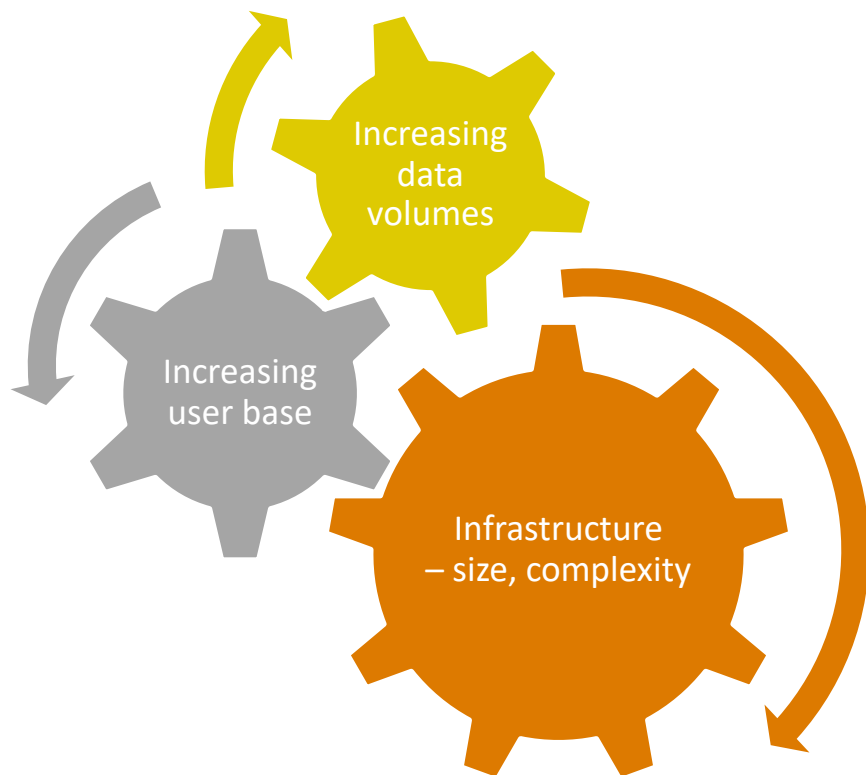
EOS Cloud – Desktop-as-a-Service for Environmental Genomics

JASMIN and CEDA



Challenges and new developments to address them (1)

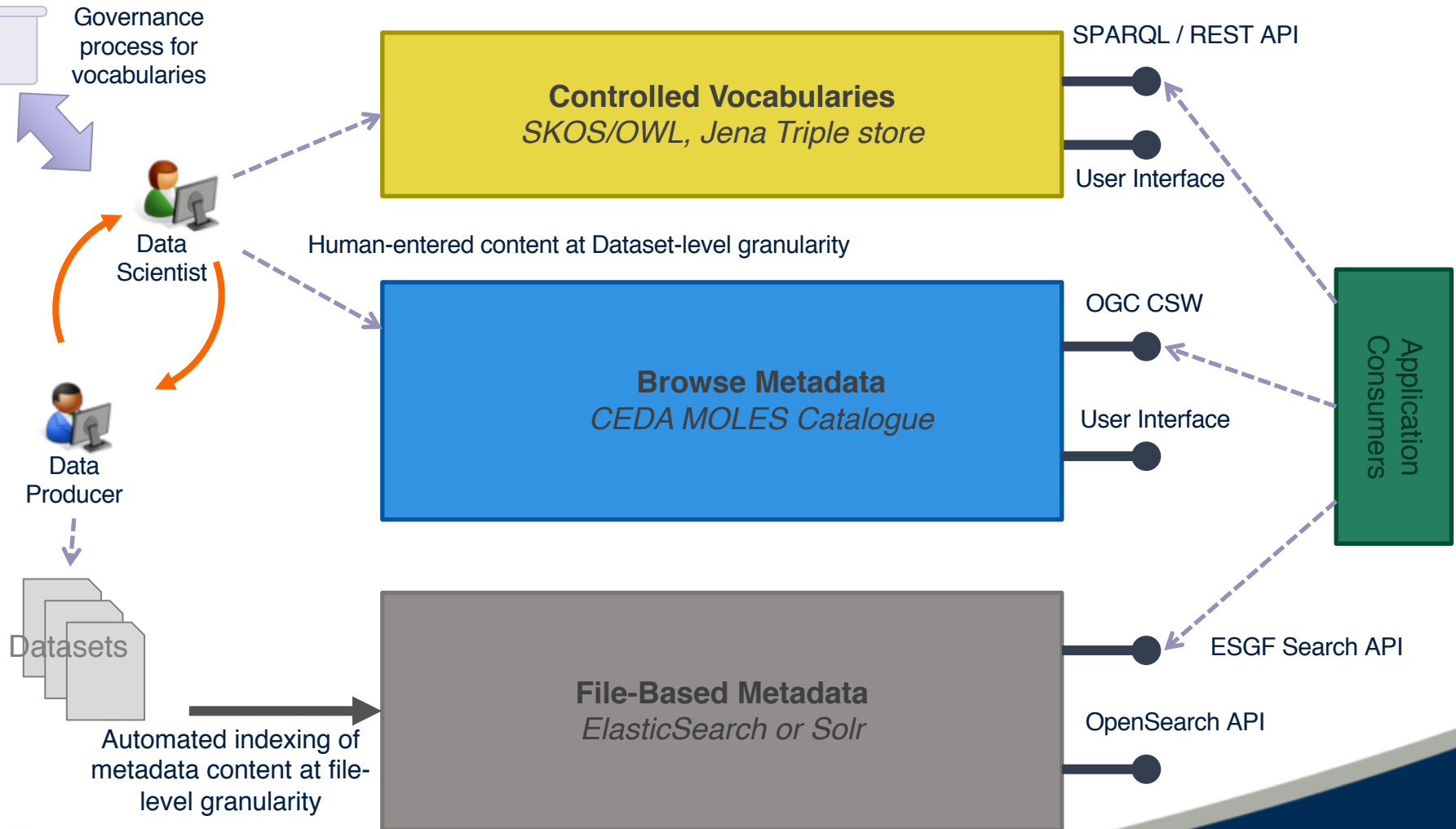
Key technologies to address specific challenges



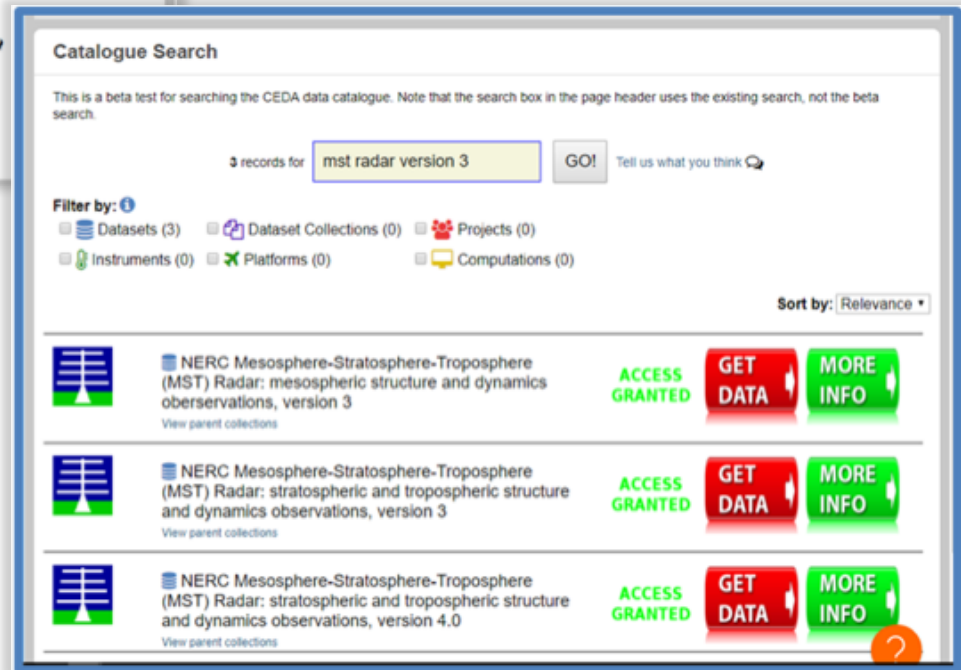
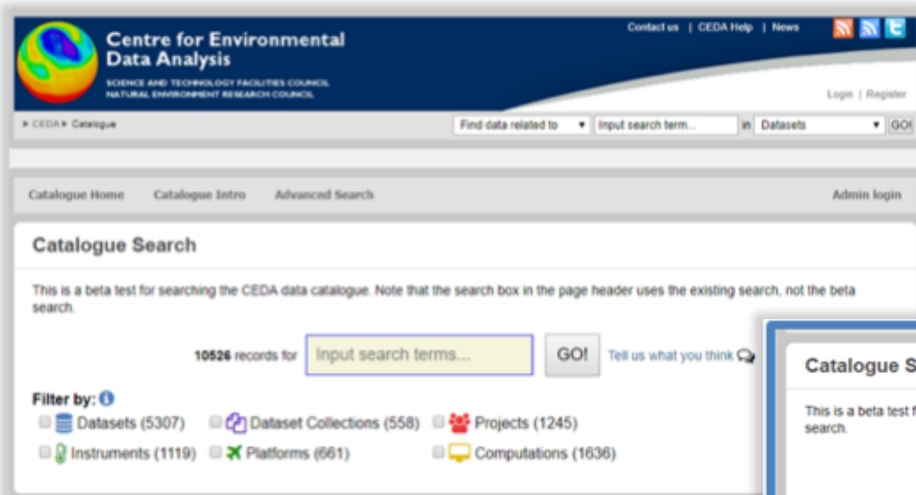
On-prem ↔ Public Cloud portability




Sourcing information for Data discovery



Catalogue search





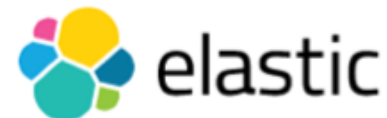
Dataset and File-level metadata ingestion

```
Details Related Datasets Process Variables Tools (3) Docs (2) Comments (0)

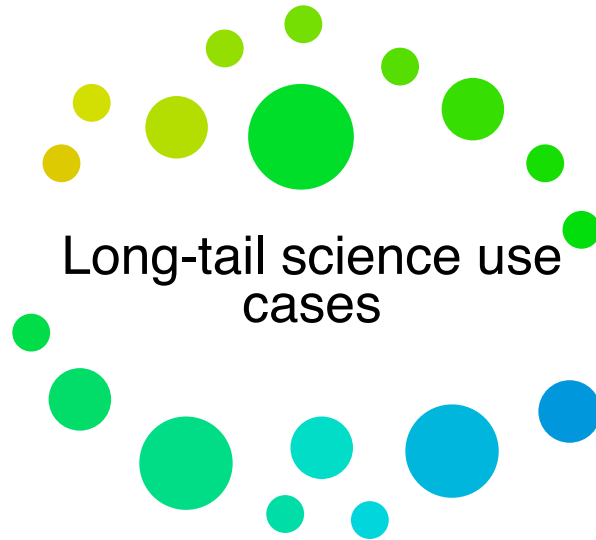
air_pressure (plev) [Pa]
  standard_name: air_pressure
  var_id: plev
  units: Pa
air_pressure_at_sea_level (psl) [Pa]
  standard_name: air_pressure_at_sea_level
  var_id: psl
  units: Pa
air_temperature (ta) [K]
  standard_name: air_temperature
  var_id: ta
  units: K
air_temperature (tas) [K]
  standard_name: air_temperature
  var_id: tas
  units: K
air_temperature (tasmax) [K]
  standard_name: air_temperature
  var_id: tasmax
  units: K
air_temperature (tasmin) [K]
  standard_name: air_temperature
  var_id: tasmin
  units: K
depth (depth) [m]
  standard_name: depth
  var_id: depth
  units: m
depth (lev) [m]
  standard_name: depth
  var_id: lev
  units: m
eastward_wind (ua) [m s-1]
  standard_name: eastward_wind
  var_id: ua
  units: m s-1
eastward_wind (uas) [m s-1]
  standard_name: eastward_wind
  var_id: uas
  units: m s-1
geopotential_height (zg) [m]
```



&



Challenges and new developments to address them (2)

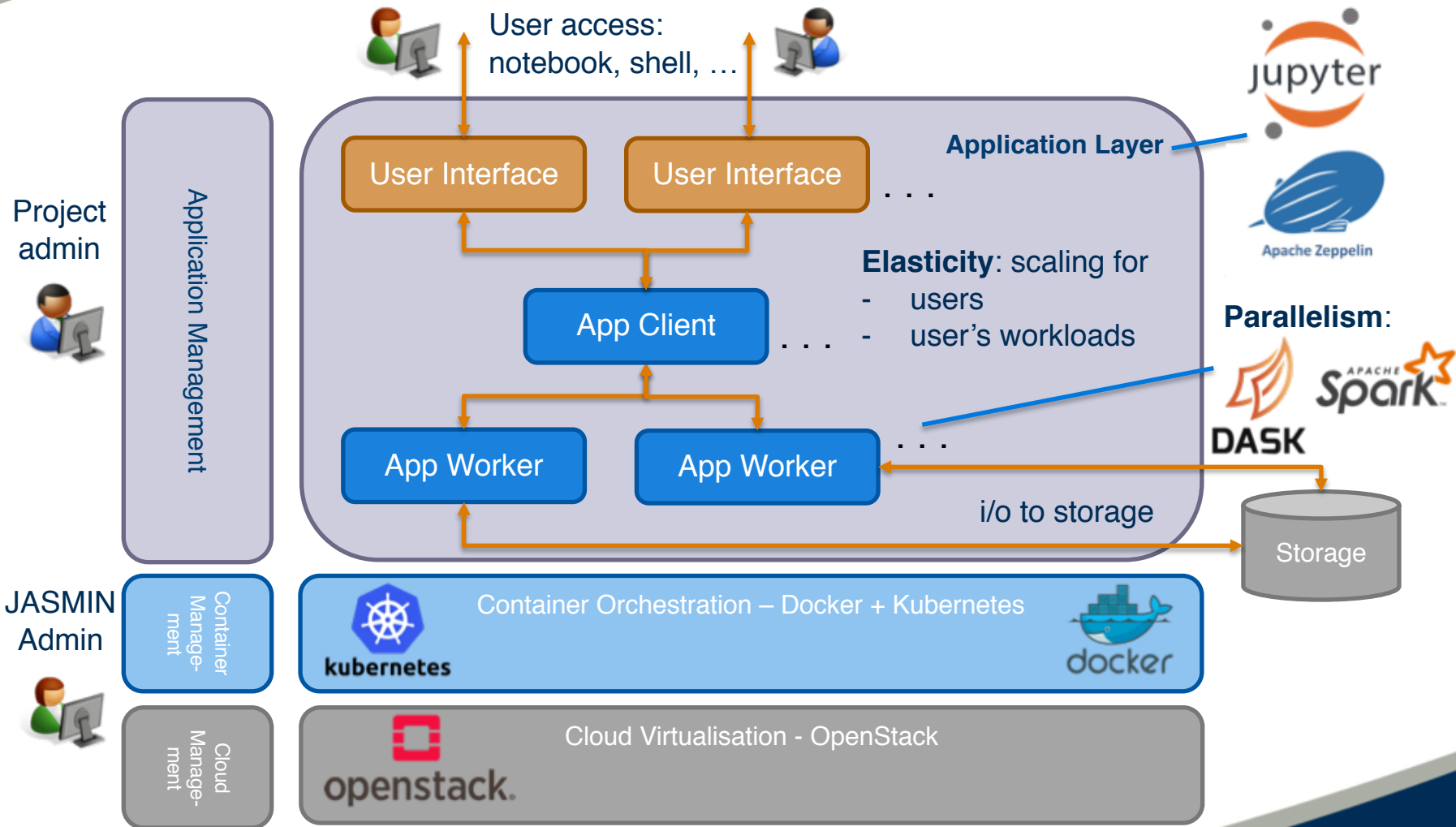


- Need for effective exploitation of parallelism to deliver demonstrable benefit over user's desktop/laptop
- Need for ease of use
- Intuitive user interfaces

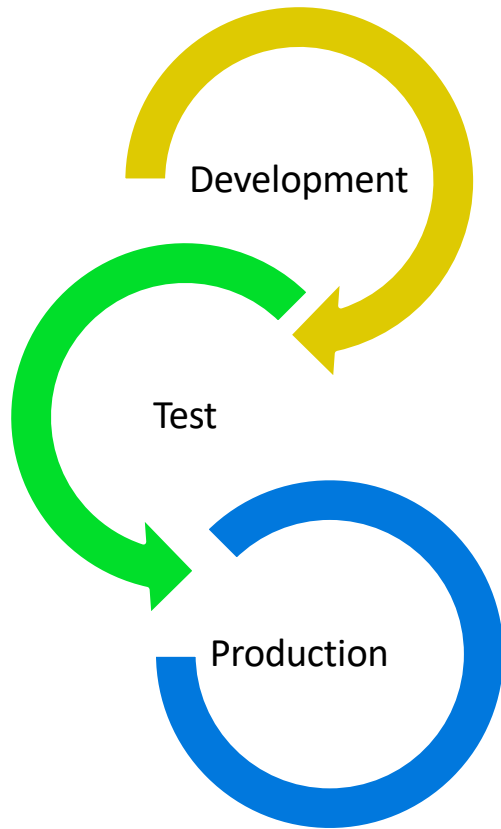
- Software-as-a-Service model
- Desktop app-style user experience
- dynamic provisioning of clusters
- Parallel programming libraries

Data Analytics

Boundaries and responsibilities for administration



Development challenge: size and complexity of code and systems footprint



- CEDA Development team
 - 15 people includes data scientists, ops and managers – approx. 7 f/t s/w development
- Languages
 - Python predominates, all new projects Python 3
 - One Cython project
 - JS - React
 - Some Java
- Development environments: PyCharm, PyDev (Eclipse), CLI + editors
- All projects Open Source by default
 - <https://github.com/cedadev/>
 - Private git for deployment-sensitive content (e.g. Ansible playbooks)
- Training in the user community:
 - Introduction to Scientific Computing Course
- Build, test, integrate, operate
 - Vagrant + Ansible
 - Cloud dev tenancy
 - Standardised on RedHat 6/7
 - Planned: Docker + Kubernetes (OpenShift)
 - Production checkout process and documentation
 - Integrating code tests into Icinga/Nagios operational monitoring



Summary

- JASMIN: data gravity, a data commons for environmental sciences
- Challenges with respect to running at scale:
 - Data volumes
 - Numbers of users
 - Generation and indexing of content for effective discovery and understanding of data for users
 - Effective use of parallelism for long-tail of science users
 - Increasing footprint of code and systems to manage for development and operations
- New infrastructure services to address challenges:
 - Evolution of data discovery and cataloguing systems, AI exploitation?
 - Virtual Research Environments
 - Object store migration
 - Development and operations: Increasing Automation – virtualisation containers and container orchestration





Further Information



CEDA team

- CEDA and JASMIN:
 - <http://www.jasmin.ac.uk/>
 - <http://www.ceda.ac.uk/>
- Github:
 - <https://github.com/cedadev/>
- JASMIN paper
Lawrence, B.N. , V.L. Bennett, J. Churchill, M. Jukes, P. Kershaw, S. Pascoe, S. Pepler, M. Pritchard, and A. Stephens. **Storing and manipulating environmental big data with JASMIN.** *Proceedings of IEEE Big Data 2013*, p68-75, doi:10.1109/BigData.2013.6691556
- philip.kershaw@stfc.ac.uk,
[@PhilipJKershaw](https://twitter.com/PhilipJKershaw)