



Data repository management in the environmental sciences in the UK

Sarah Callaghan
sarah.callaghan@stfc.ac.uk
@sorca_ni

Geoscience Digital Data Resource and Repository Service
(GeoDaRRS) Workshop
Boulder, CO, USA, 7-9 August 2018



NERC Environmental Data Centres

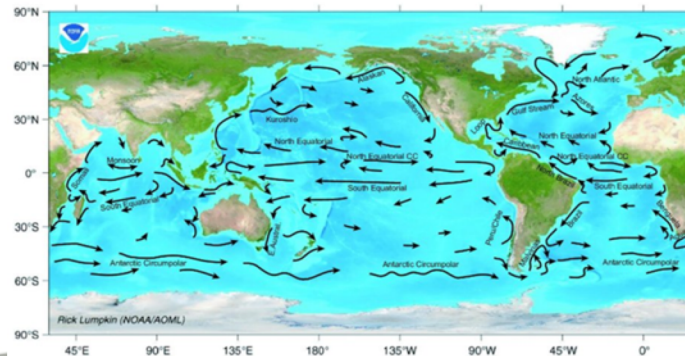
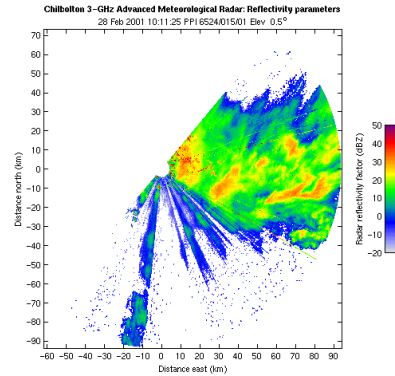
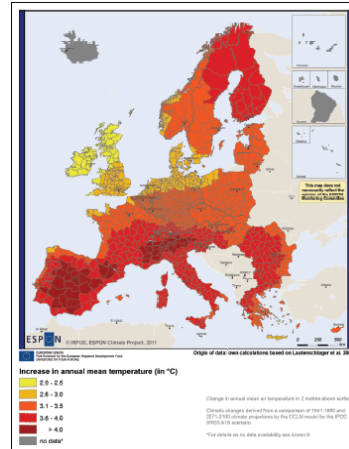
- National Geoscience Data Centre (NGDC)
 - Geoscience/ Sub-surface
- **Centre for Environmental Data Analysis (CEDA)**
 - **Atmosphere/ Solar System/ Earth Observation**
- Environmental Information Data Centre (EIDC)
 - Terrestrial / Freshwater
- Polar Data Centre (PDC)
 - Polar - regional
- British Oceanographic Data Centre (BODC)
 - Oceanographic / Marine





Data heterogeneity

1. Time series, some still being updated e.g. meteorological measurements
2. Large 4D synthesised datasets, e.g. Climate, Oceanographic, Hydrological and Numerical Weather Prediction model data generated on a supercomputer
3. 2D scans e.g. satellite data, weather radar data
4. 2D snapshots, e.g. cloud camera
5. Traces through a changing medium, e.g. radiosonde launches, aircraft flights, ocean salinity and temperature
6. Datasets consisting of data from multiple instruments as part of the same measurement campaign
7. Physical samples, e.g. fossils, ice cores

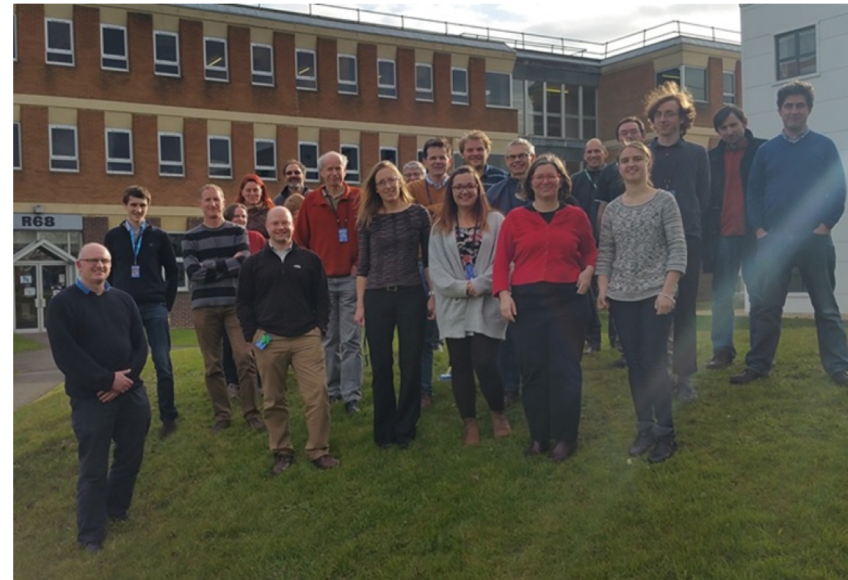


Data Analysis



CEDA exists to enable the efficient delivery of environmental science programmes through effective data and information services. This is achieved by:

- operation of efficient data curation services
- facilitating environmental science by running data and information services
- development of new data service technologies
- close contact with the research community
- the development of a national capability in data management expertise
- contributing to and learning from the international community



Team of ~29; mixture of Data Scientists and Software Engineers

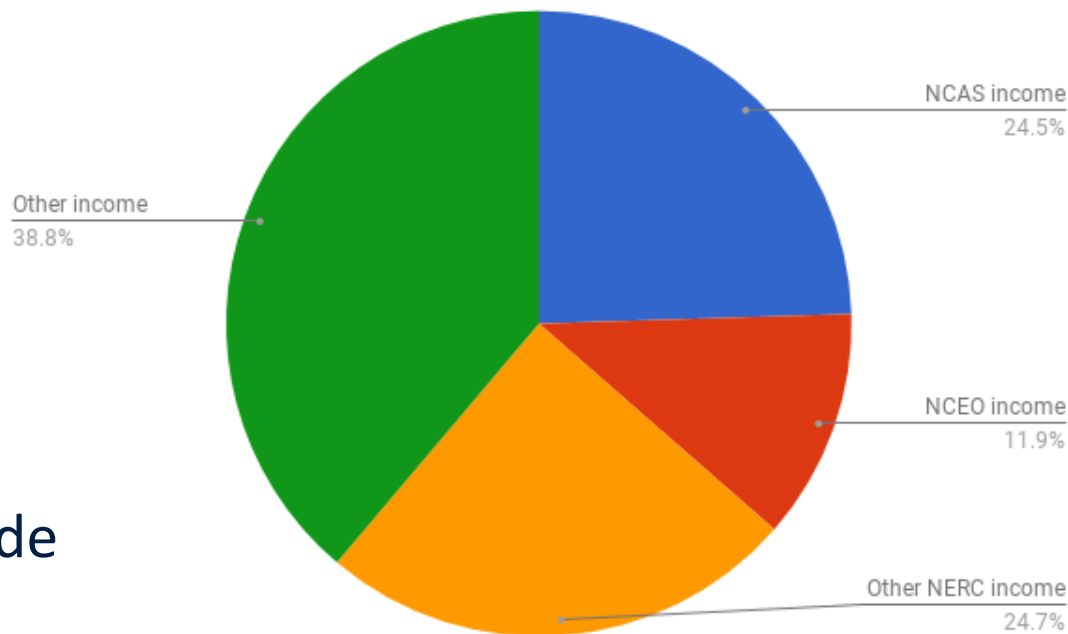
Diverse expertise in:

- Earth observation
- Climate modelling
- Aircraft measurements
- Data standards
- ... and much more!



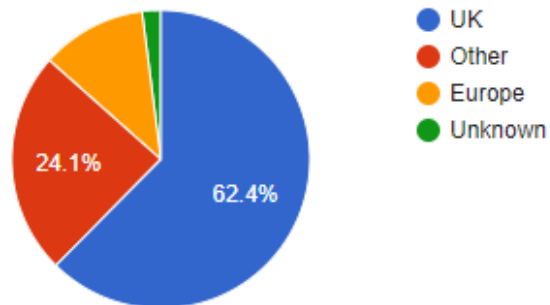
What we do

- CEDA Data Centres
 - NCAS
 - NCEO
 - IPCC-DDC
 - UKSSDC
- JASMIN “super data cluster”
- Projects: funders include EC, ESA, UKSA, Defra, BEIS, Met Office

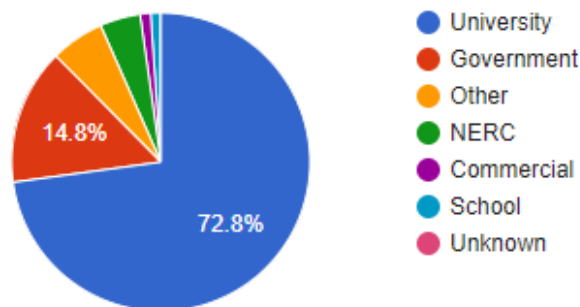


The screenshot shows the CEDA Archive website interface. The browser address bar is 'flight-finder.ceda.ac.uk'. The page title is 'Dataset'. The main heading is 'GBS 20.7GHz slant path radio propagation measurements, Sparsholt site'. Below the heading, there is a 'Download' button and a green 'Access Granted' button. To the right, a table shows metadata: Update Frequency: Not Planned; Status: Completed; Online Status: ONLINE; Publication State: Citable; Publication Date: 2003-09-20; DOI Publication Date: 2011-04-01; Download Stats: last 12 months. Below this is an 'Abstract' section with text about the GBS dataset. To the right of the abstract is a 'Coverage' section with 'Temporal Range' (Start time: 2003-10-08T00:00:00, End time: 2005-03-31T00:00:00) and 'Geographic Extent' (a map showing the Sparsholt site location near Salisbury and Winchester). The left sidebar contains search filters like 'EUFAR' and 'Geog'.

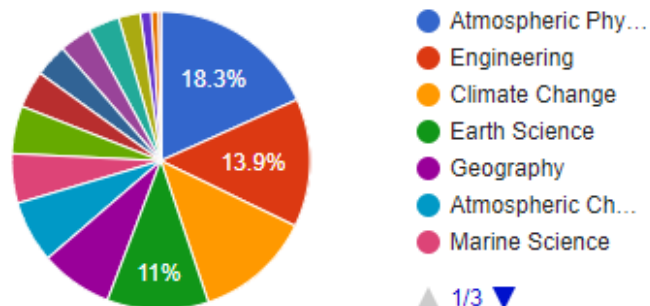
Number of users by Area



Number of users by Institute type

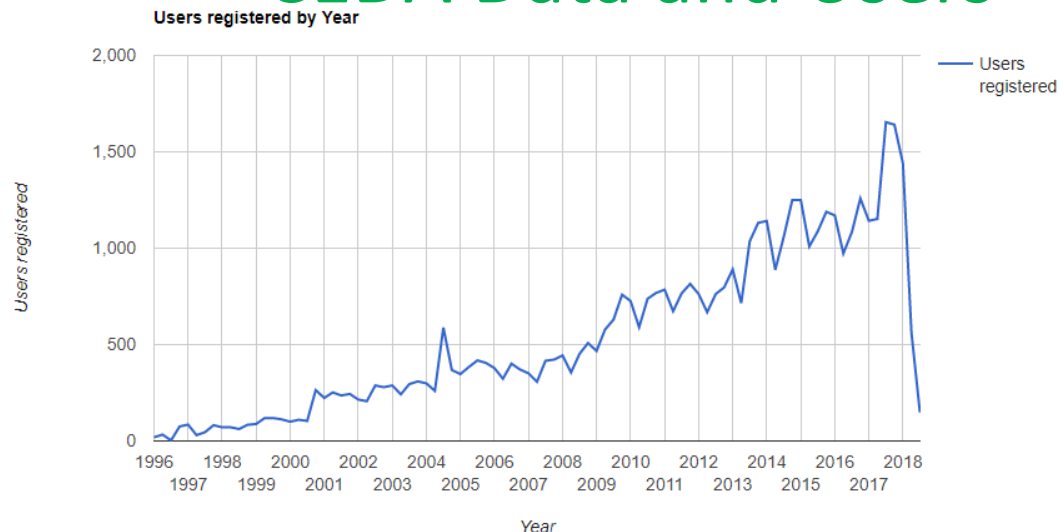


Number of users by Field



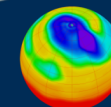
▲ 1/3 ▼

CEDA Data and Users



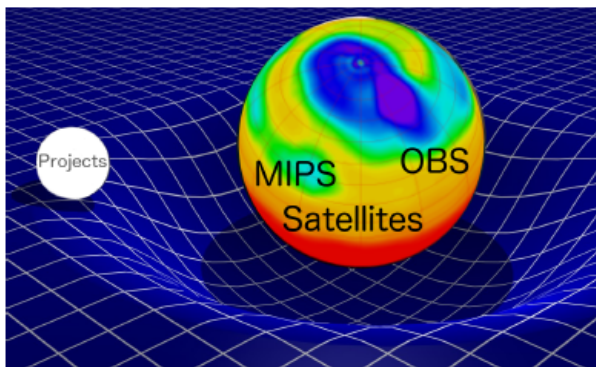
Data Type	Data Volume (Petabytes)
Earth Observation	4
Atmospheric Science	2
Total	6 PB

- ~ 550 datasets
- ~ 150 million files
- > 44,000 registered users

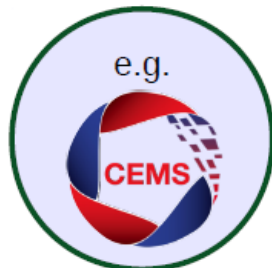




JASMIN – The Data Commons

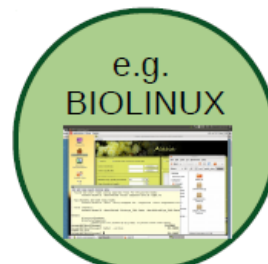


- ▶ Provide a state-of-the art storage and computational environment
- ▶ Provide and populate a managed data environment with key datasets (the “archive”).
- ▶ Encourage and facilitate the bringing of data and/or computation alongside/to the archive!
- ▶ Provide **FLEXIBLE** methods of exploiting the computational environment.



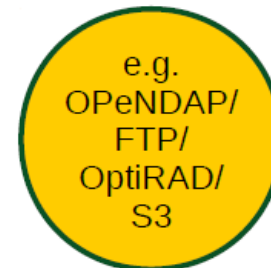
**Platform
as a
Service**

We provide you the “Platform”; you can LOGIN and exploit the batch cluster.



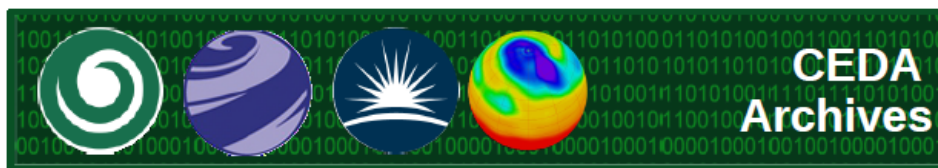
**Infrastructure
as a
Service**

We provide you with a cloud on which you INSTALL your own computing.



**Software
as a
Service**

We provide you with REMOTE access to data VIA web and other interfaces.



JASMIN – Data Intensive Computer

Storage, Compute and Network Fabric
Batch Compute, Private Cloud, Disk, Tape





Logical View

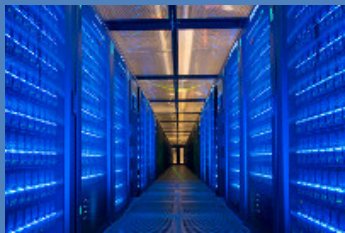


CEDA Archive Services

Data Centres, Curation, DB systems
User management, External Helpdesk

Analysis Environment

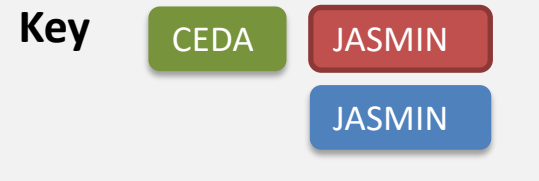
Compute:
Science VMs + User Management, PaaS and IaaS Cloud,
Group Workspaces: Fast Disk & Elastic Tape
External Helpdesk



JASMIN Compute and Storage

Managed Compute (Lotus and GWSs, Tape Store + Data Transfer Zone),
Community Cloud + Internal Helpdesk

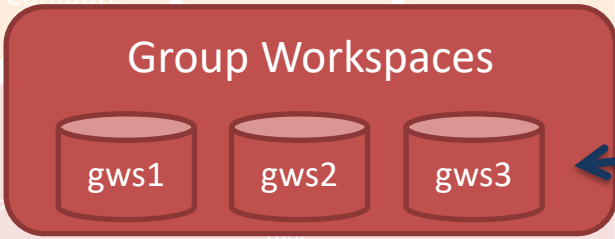




Long Term Archive Storage



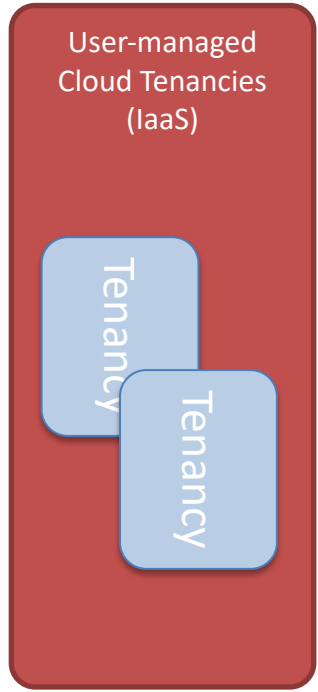
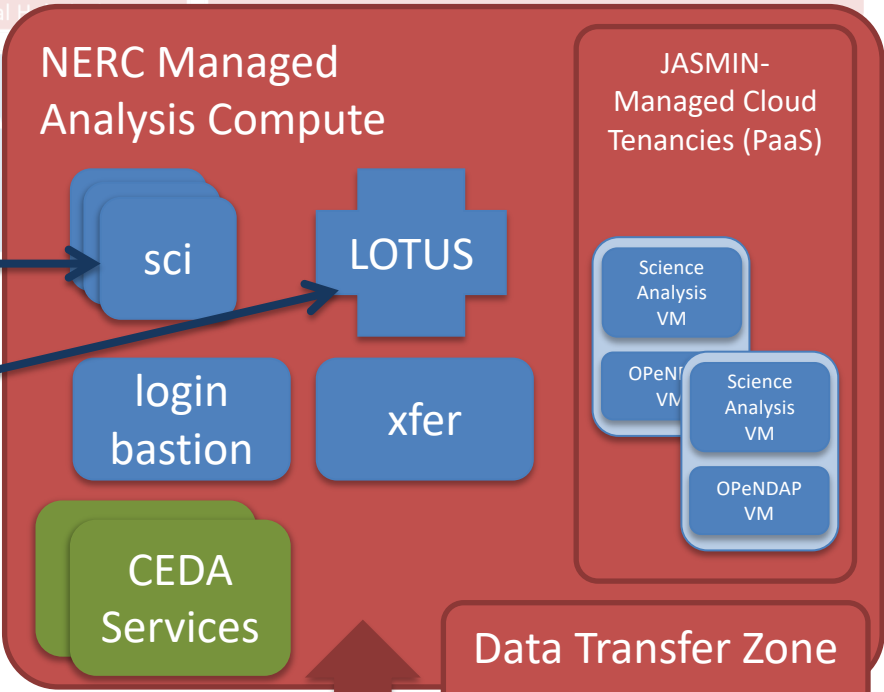
Short-Term Project Storage



Phase 4: 17PB



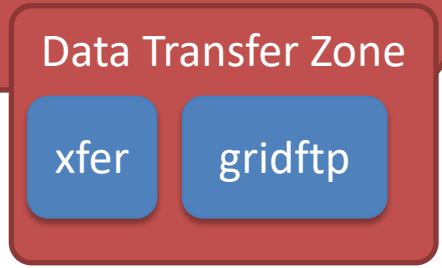
Phase 4: 26PB



Phase 4: 11500 cores – all types of compute

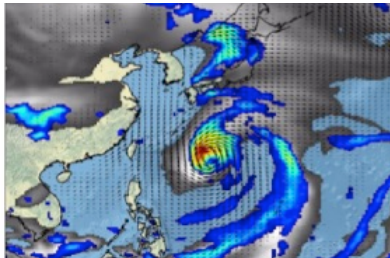
Interactive Compute

Batch Compute

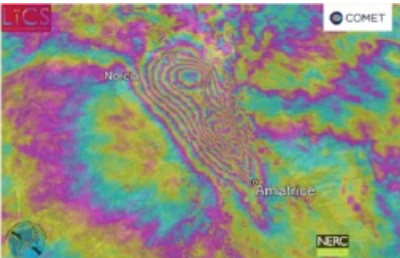


Functional View

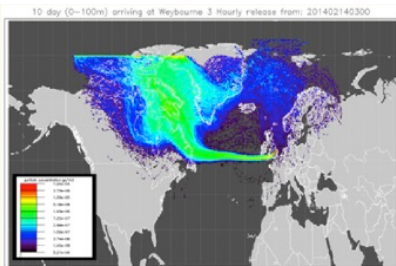
~150 Science projects on JASMIN to date



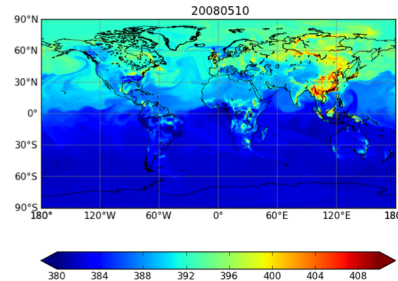
High Res Climate Model analysis



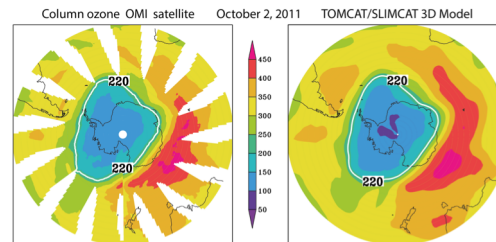
Fault analysis



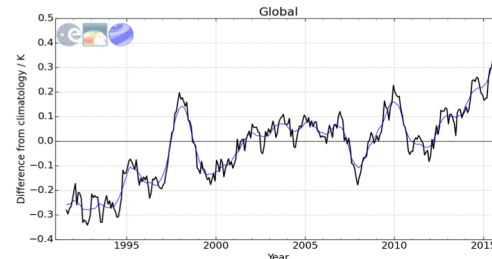
Atmospheric dispersion



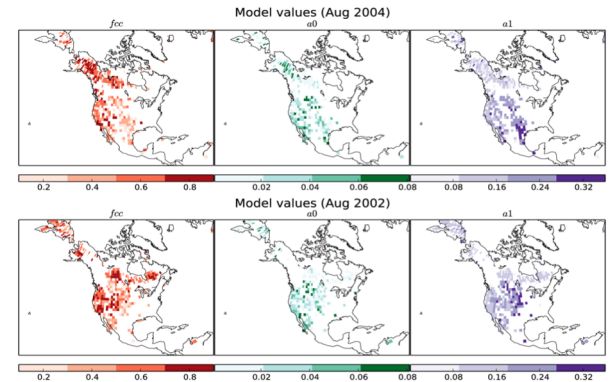
Regional carbon balance on a global scale



Antarctic Ozone hole: model vs. observations

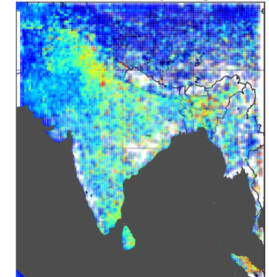


Sea Surface Temperature from satellite observations

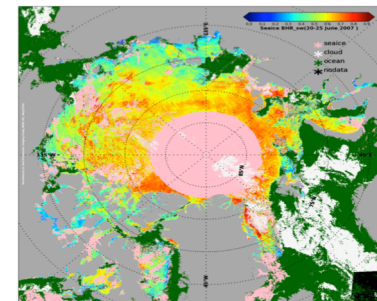


Deriving the impact of fire on vegetation from earth observation data

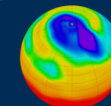
OMI Vertical Columns monthly average 0.5x0.5 grid



Understanding oxidant chemistry over the Indian subcontinent

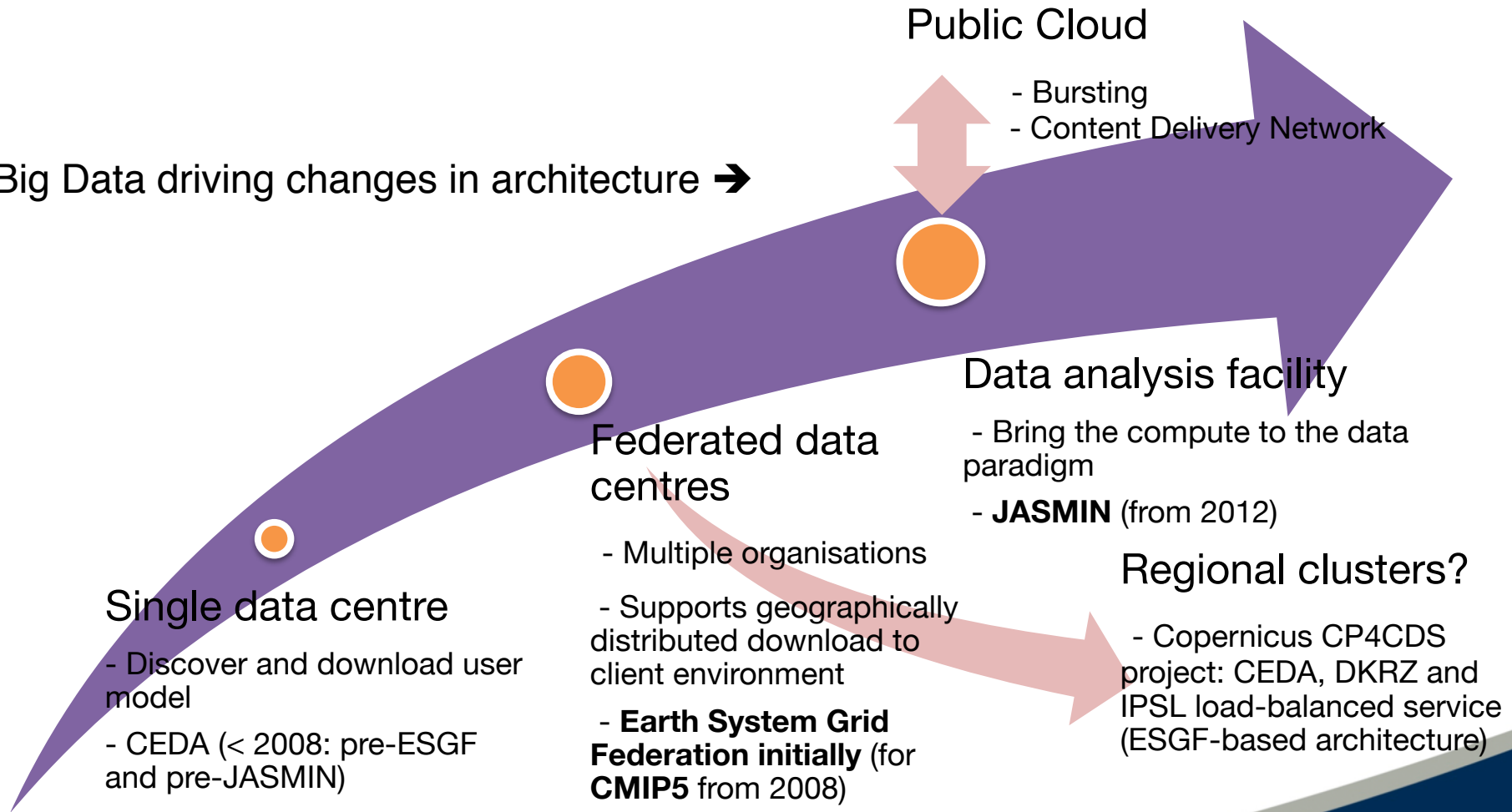


Climate variables from European and US instruments/satellites



Evolution in models for data distribution and analysis

Big Data driving changes in architecture →



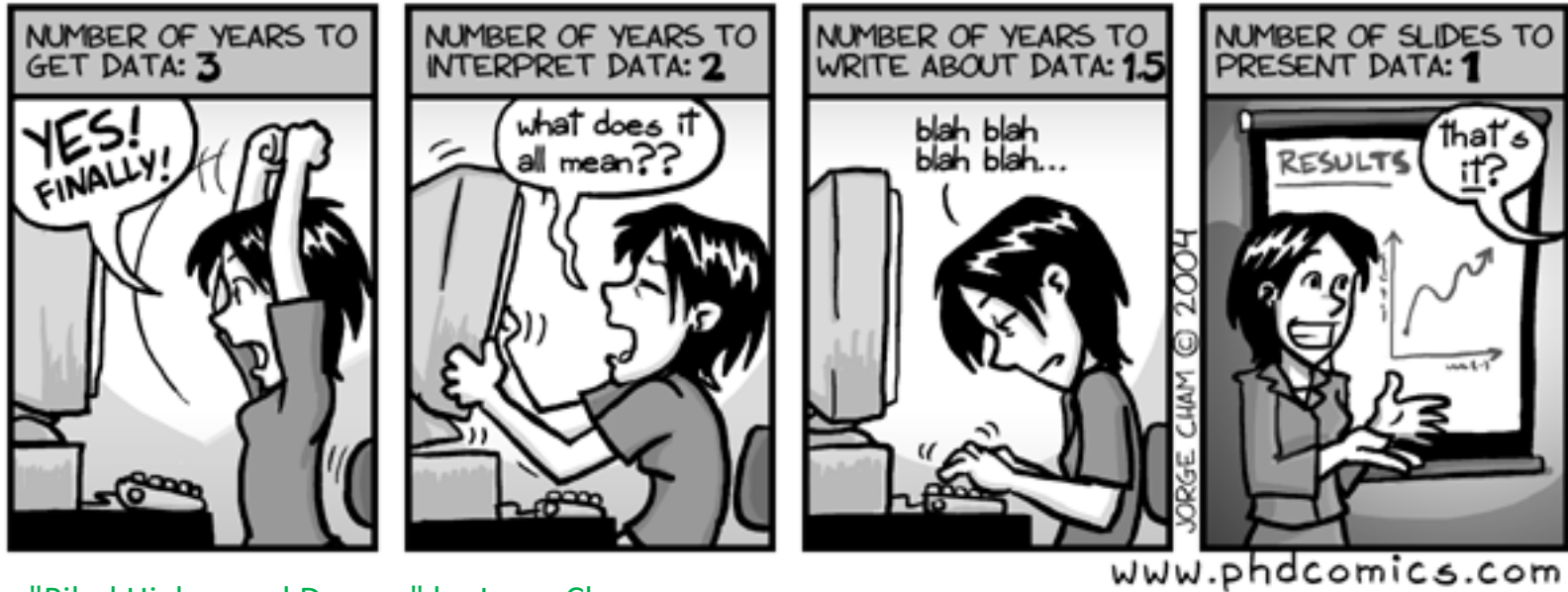


JASMIN Challenges

- Data intensive computing facility
 - Storage is a major cost
 - Phase 4 introduces more cost effective storage technologies for future scaling: object storage and scale-out file system
- Balancing "bleeding edge" technology with demands of operational service
 - Research infrastructure, not 99.999% available
 - Small team (capital-heavy)
- Diverse user community
 - Wide range of workflow requirements
 - Wide range of user skill level
 - Training, user education, engagement

Creating a dataset is hard work!

DATA: BY THE NUMBERS



"Piled Higher and Deeper" by Jorge Cham
www.phdcomics.com

Documenting a dataset so that it is usable and understandable by others is extra work!

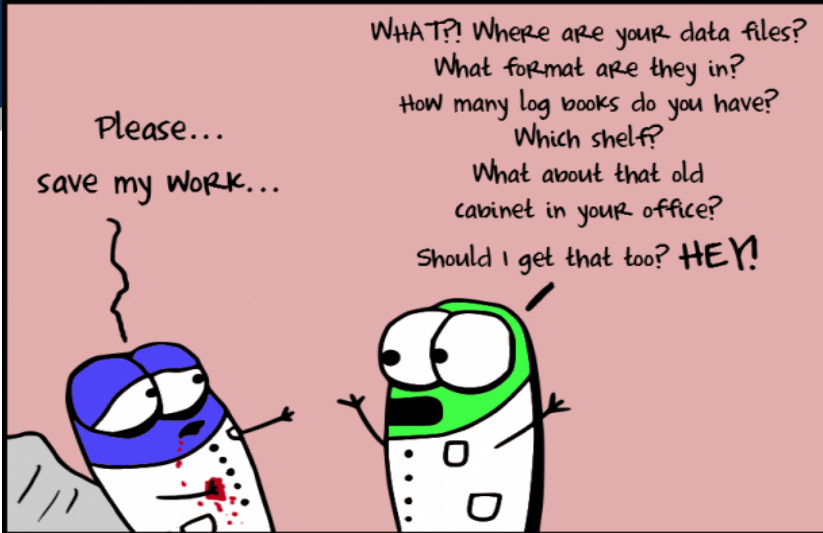
Don't just create a bit bucket

“When required to make the data available by my program manager, my collaborators, and ultimately by law, I will grudgingly do so by placing the raw data on an FTP site, named with UUIDs like 4e283d36-61c4-11df-9a26-edddf420622d. I will under no circumstances make any attempt to provide analysis source code, documentation for formats, or any metadata with the raw data. When requested (and ONLY when requested), I will provide an Excel spreadsheet linking the names to data sets with published results. This spreadsheet will likely be wrong -- but since no one will be able to analyze the data, that won't matter.”

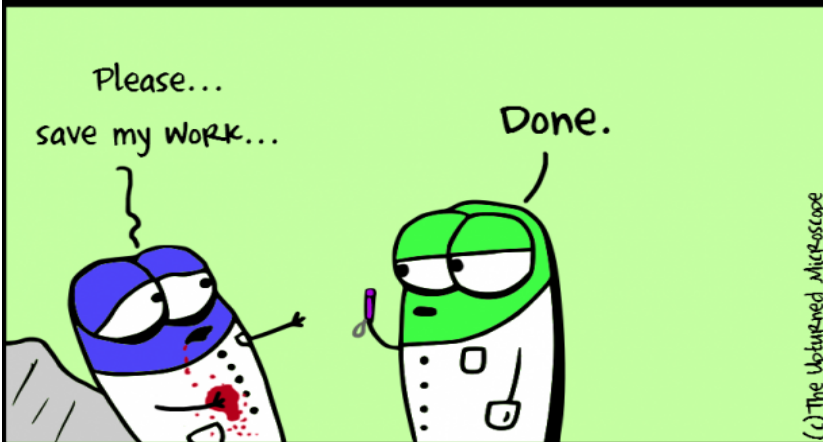
- <http://ivory.idyll.org/blog/data-management.html>



REAL scientist



MOVIE scientist



<http://theupturnedmicroscope.com/comic/real-vs-movie-scientist-3/>

Work with your users

- Build things in a modular fashion which can be easily reused for different communities
- Start small and build up
- Minimise data friction – make it easy for users to submit/download/use data
- If there's no users – what's the point of the repository?



(some of) What we've learned

- Train your depositors in good data management
- Use off the shelf solutions where possible
- Not all data should be kept (data value checklist)
- It's easier to create accurate metadata at the beginning of the project
- Providing secure workspaces makes researchers *more* likely to share their data
- Researchers need incentives to put their data in a repository (especially if there is extra work needed to do so)
- Standard tools and formats make life easier for everyone



THE LAVA IS ENTERING THE SEA, AND
NEW RIFTS ARE OPENING TO THE NORTH!

GET A GIS SURVEY TEAM IN
THE AIR! WE NEED TO REVISE
OUR COASTLINE SHAPEFILES!



I WANT TO MAKE A DISASTER MOVIE
THAT JUST SHOWS SCIENTISTS RUSHING
TO UPDATE ALL THEIR DATA SETS.

“Really, they’ d be rushing around collecting revisions to go into the next scheduled quarterly public data update, not publishing them immediately, but you have to embellish things a little for Hollywood.”

<https://xkcd.com/2029/>

Thanks!

Any Questions?

sarah.callaghan@stfc.ac.uk
@sorca_ni