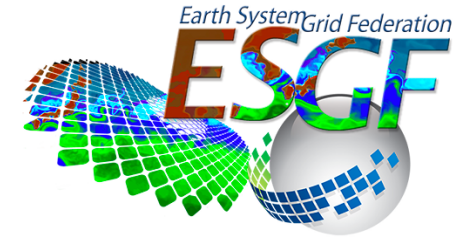


Science and
Technology
Facilities Council

Natural
Environment
Research Council



Evolution and Future Architecture for the Earth System Grid Federation

Philip Kershaw¹

Ghaleb Abdulla², Sasha Ames², Ben Evans³, Tom Landry⁴, Michael Lautenschlager⁵,
Venkatramani Balaji⁶ and Guillaume Levavasseur⁷

1 STFC Rutherford Appleton Laboratory, RAL Space, Didcot, UK

2 LLNL, Livermore, USA

3 NCI, Australian National University, Acton, Australia

4 CRIM, Montréal, Canada

5 DKRZ, Hamburg, Germany

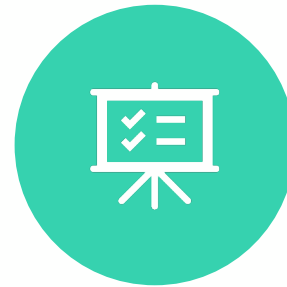
6 Princeton University, Princeton, USA

7 IPSL, Paris, France

Introduction



The Earth System Grid Federation (ESGF) is a globally distributed e-infrastructure for the hosting and dissemination of climate-related data.



ESGF was originally developed to support CMIP5 (5th Coupled Model Intercomparison Project) ...



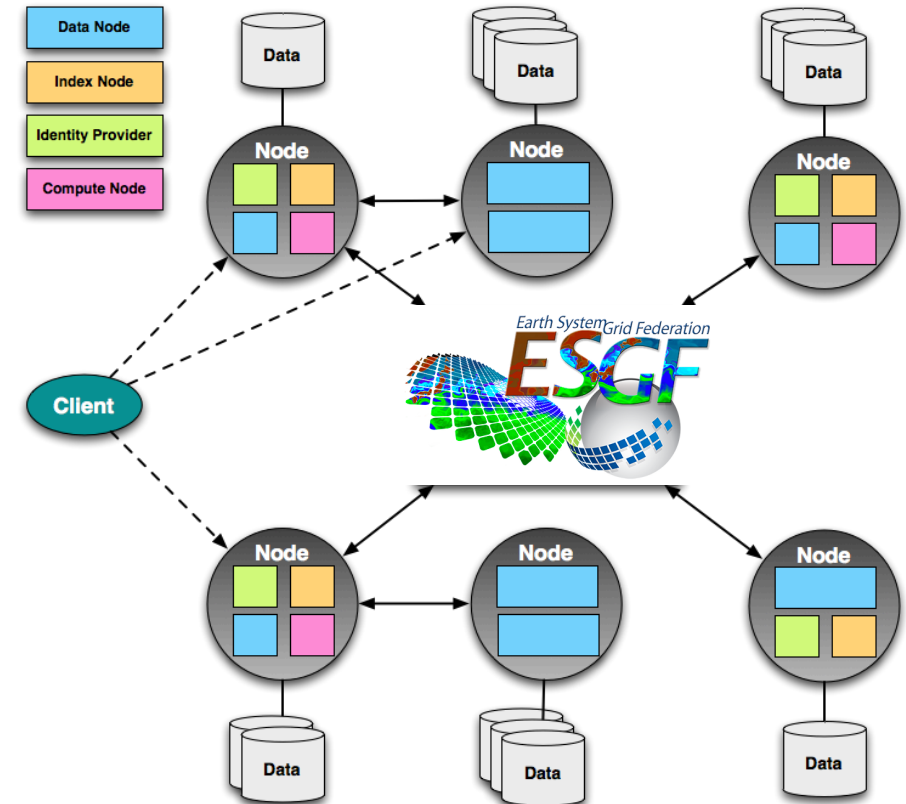
Provide a means for climate research community to access and analyse the data output



For 5th Assessment report made by the IPCC (Intergovernmental Panel on Climate Change).

Ten years of Operations: History and Evolution

- ESGF has grown to support over
 - 25000 registered users
 - 33 registered nodes at climate research centres spread across 21 countries, sites including DoE, EU IS-ENES collaboration, NASA, NOAA, NCI Australia ...
- Besides the CMIPs, supports a range of other projects such as the Energy Exascale Earth System Model, Obs4MIPS, CORDEX and the European Space Agency's Climate Change Initiative Open Data Portal.
- Over the course of its evolution, ESGF has pioneered technologies and operational practice for distributed systems
 - Federation inherently supports redundancy and large-scale data replication capabilities
 - Search, metadata modelling and capture, identity management
- Important experience gathered over the years about community collaboration for a distributed infrastructure - operational procedures and governance

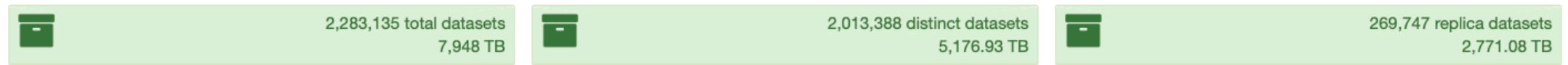


Federation stats: ESGF Dashboard



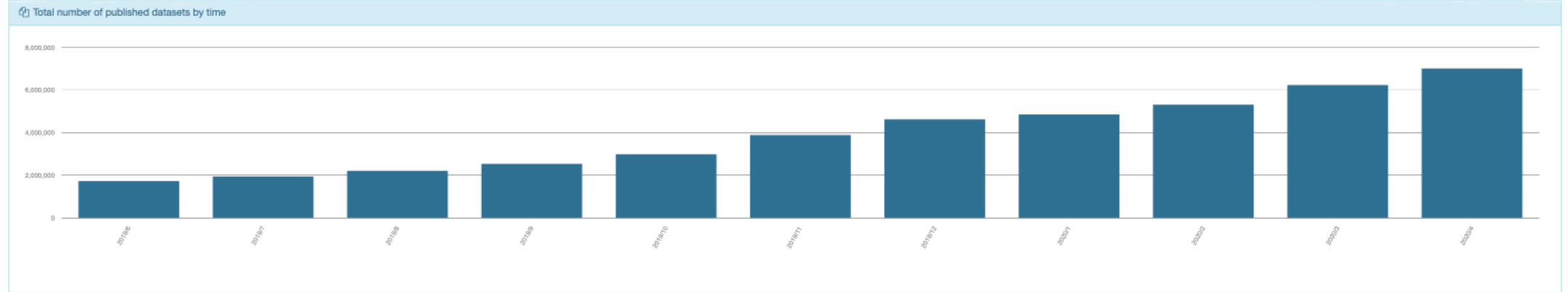
Published data over the federation

Current values

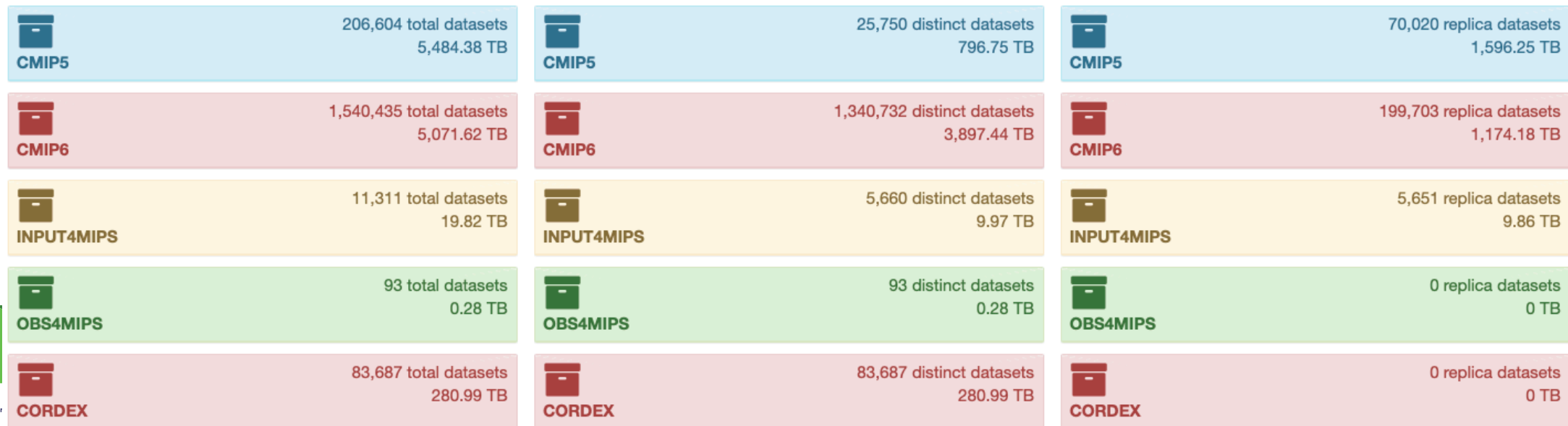


Published data over time (updated every first day of the month)

Total datasets number Total data size [TB] CMIP6 Datasets number CMIP6 Data size [TB] CORDEX Datasets number CORDEX Data size [TB]

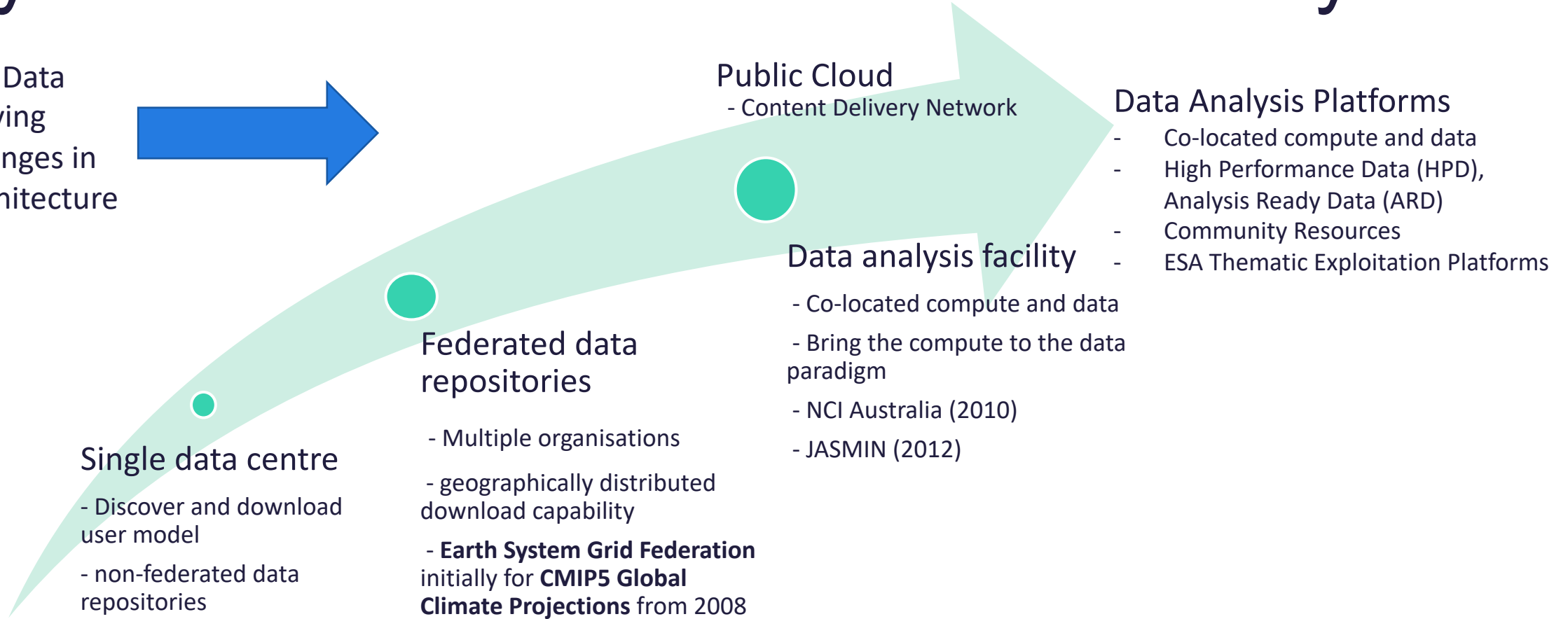


Top projects (updated daily)



Big Data, cloud and the evolution of systems for data distribution and analysis

Big Data driving changes in architecture



Future Architecture

- An initiative to re-design ESGF
- Drawing from our experience and lessons learnt
 - running ESGF operationally and
 - other large-scale compute/data systems at our facilities
- Also considering other similar systems and initiatives e.g.
 - Interoperable community standards – ISO Catalogues, OpenSearch, FAIR, schema.org, RDA ...
 - ESA's Earth Observation Exploitation Platform Common Architecture (<https://eoepca.github.io>)
 - OGC Testbeds
 - Pangeo (<https://pangeo.io/>), a community and software stack for the geosciences

High-Level Functional Areas to Consider

- 1) User Experience
- 2) Data Repository and Management
- 3) Compute on Data
- 4) Platforms and systems administration

Approach



* Report write-up of the above

Findings

- User Experience
 - Need more integrated search across all data holdings: integrate metadata about models (ES-Doc system) as well as netCDF metadata
- Data Repository and Management
 - Develop new publishing system and adopt new search API standard (later slide)
 - New approach needed over traditional data access mechanisms such as OPeNDAP. Desire for new and efficient means for sub-setting and aggregation
- Platforms and systems administration
 - Strong consensus to build on work to standardize deployment based on Docker and Kubernetes
 - Make architecture more modular with clearly defined interfaces
 - This will better facilitate contributions from the development community
- Compute on Data
 - No consensus to standardize on a single federation-wide reference implementation
 - However, plenty of individual projects leveraging ESGF developing processing services based on OGC Web Processing Service standard
 - These include work tying in with recent OGC Testbeds
 - Individual sites are also deploying services based on Python stacks with Jupyter
- More findings ...

Federation or Cloud for Resilience?



Federation allows scaling, redundancy, sharing of capability



Search and identity management are the linchpins of federation



But federation can also bring complexity:

Search and identity management services are duplicated across sites

There can be a high operational burden

Duplication of services can be confusing for users

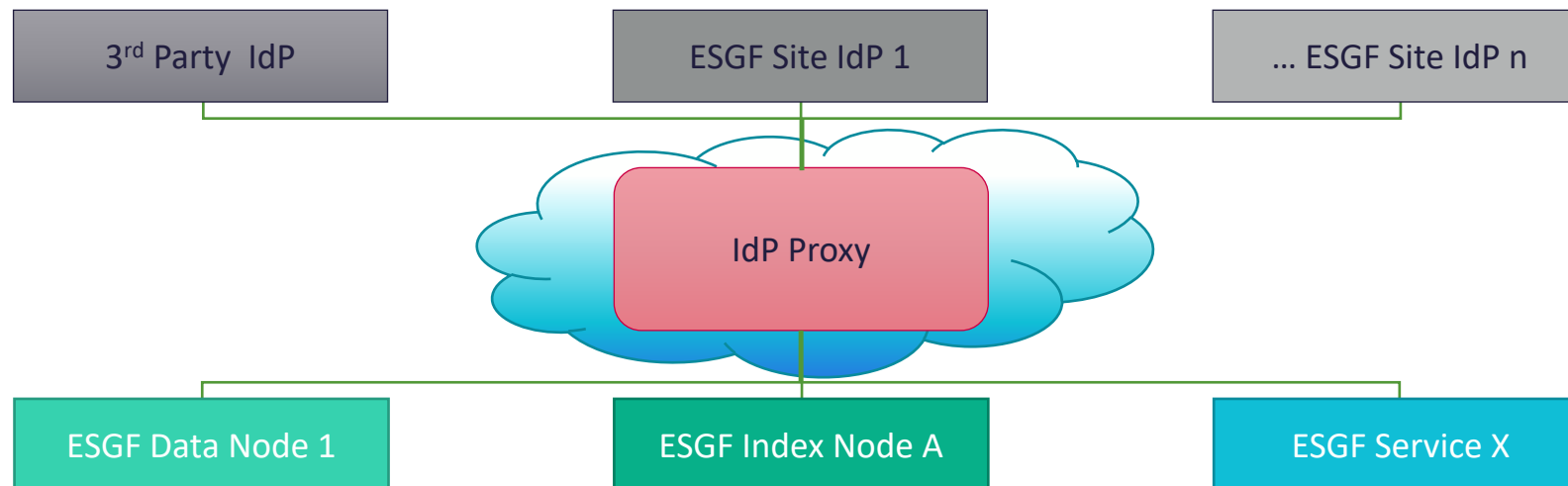


Cloud computing with its impact on deployment practice and hosting architecture can help

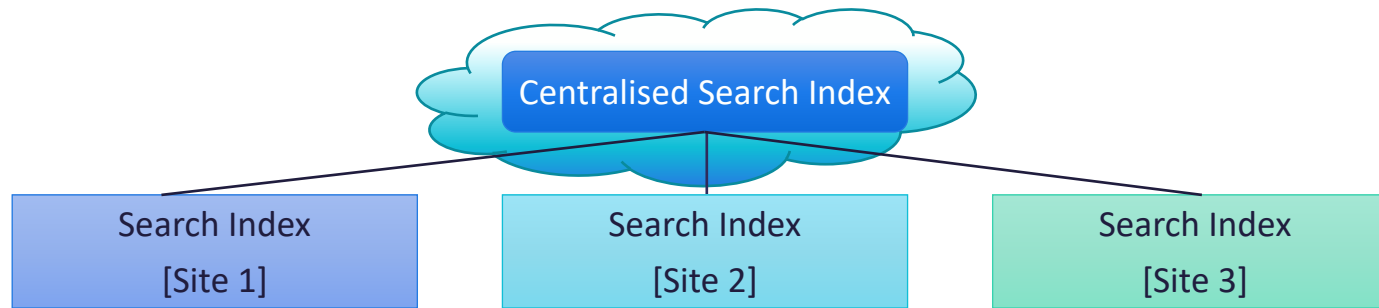
Public cloud with 4 9s resilience capability brings back centralisation as a viable design choice vs. the previously established approach in ESGF of federation for resilience

Cloud and Identity Management

- Existing status quo: Many-to-many relationships: too complicated
- AARC Blueprint: Identity Proxy - Acts as abstraction between ESGF services and IdPs



Cloud and Search: Centralised Search Index



- Centralised cloud-hosted search index harvests content from all providers
- Builds on existing technology in ESGF (Apache Solr) but reduces operational burden for hosting sites and provides a single-entry point for users

Cloud and data access and storage

Opportunities

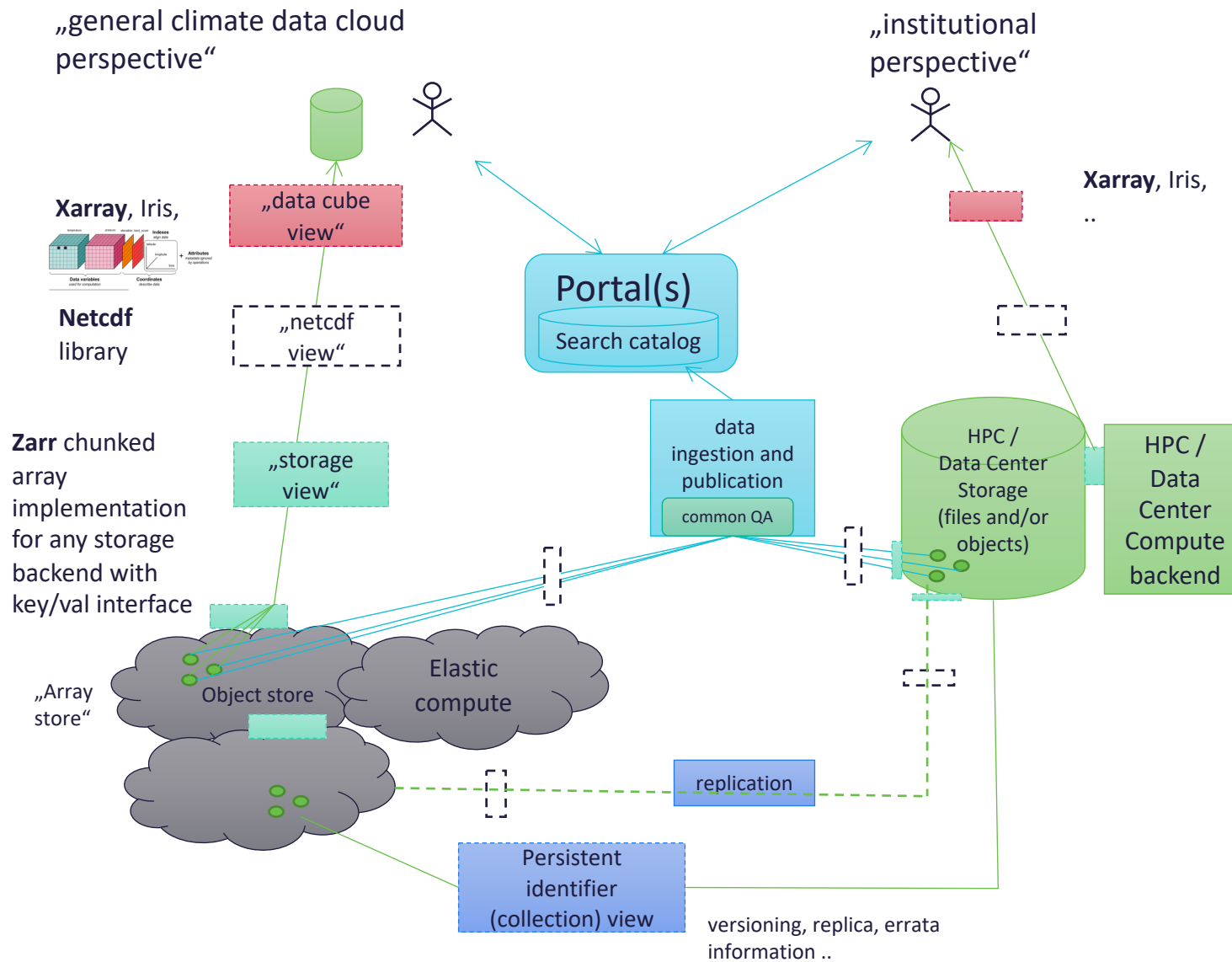
- New paradigms for massively parallel data storage and access, such as object store
 - Xarray and zarr
- Potential point of entry for scientists without access to large-scale computing, analysis, and network resources

Challenges

- Public cloud not cost effective for *long-term, large-scale* data storage and access
- Most on-premise data centres still have POSIX file systems (Though some are experimenting with object store – JASMIN)
- Concerns about Xarray/zarr for data archiving – consistency, integrity and metadata retention

Cloud vs. institutional Perspective

(Diagram courtesy of Stephan Kindermann, DKRZ)



Role of standards

- Critical for hosting centres in ESGF
 - These increasingly support a broader range of domains and communities in the Earth sciences.
 - Interoperability with other similar systems, minimise duplication and maximise re-use
- Data
 - CF-netCDF continues as established standard for this community
 - Challenges around storage and access with Xarray and zarr: how to maintain data integrity and consistency
 - New standards evolving – with end-user centric perspective e.g. https://portal.ogc.org/files/?artifact_id=91644#PartDAPA
- Identity management
 - Adopt standards and use off-the-shelf solutions where possible - OIDC (OpenID Connect) and Keycloak
- Search – ESGF uses its own standard based on Apache Solr API, these standards are being explored as alternatives:
 - OpenSearch (<http://ceos.org/ourwork/workinggroups/wgiss/access/opensearch/>)
 - STAC (<https://stacspect.org>)
 - ESM Collection Specification (<https://github.com/NCAR/esm-collection-spec>)

Conclusions and next steps

- Cloud enables major architectural changes to be made to simplify
 - Centralised search and identity services
- Storage and data access is at a cross-roads:
 - POSIX access continues to dominate for on-premise
 - Object storage on cloud provides new possibilities
 - No community-wide consensus on a single solution yet
- A roadmap for re-development of ESGF has been established
 - Incremental releases with initial new baseline version ready for the summer

Thank you



- This work has been carried out through IS-ENES3, a project funded by the European Union's Horizon 2020 research and innovation programme under grant agreement No 824084
- Free access to computing platforms for multi-model climate data analyses for CMIP6 and CORDEX ! –
 - <https://portal.enes.org/data/data-metadata-service/analysis-platforms>
- More information at EGU sessions ...
 - CL5.7 Climate Services – Underpinning, 05 May, 10:45–12:30, EGU 2020-19121: <https://meetingorganizer.copernicus.org/EGU2020/session/36737>
 - CL2.6 Detecting and attributing climate change: trends, extreme events, and impacts, 07 May, 08:30–10:15, EGU 2020-19340: <https://meetingorganizer.copernicus.org/EGU2020/session/36768>

