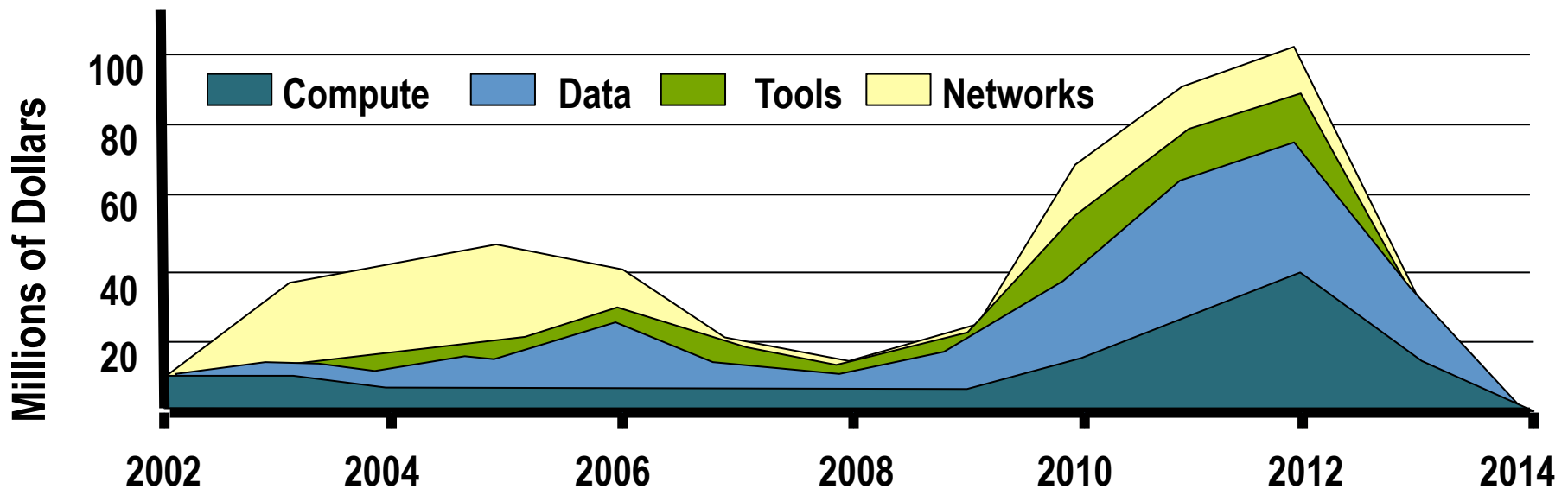




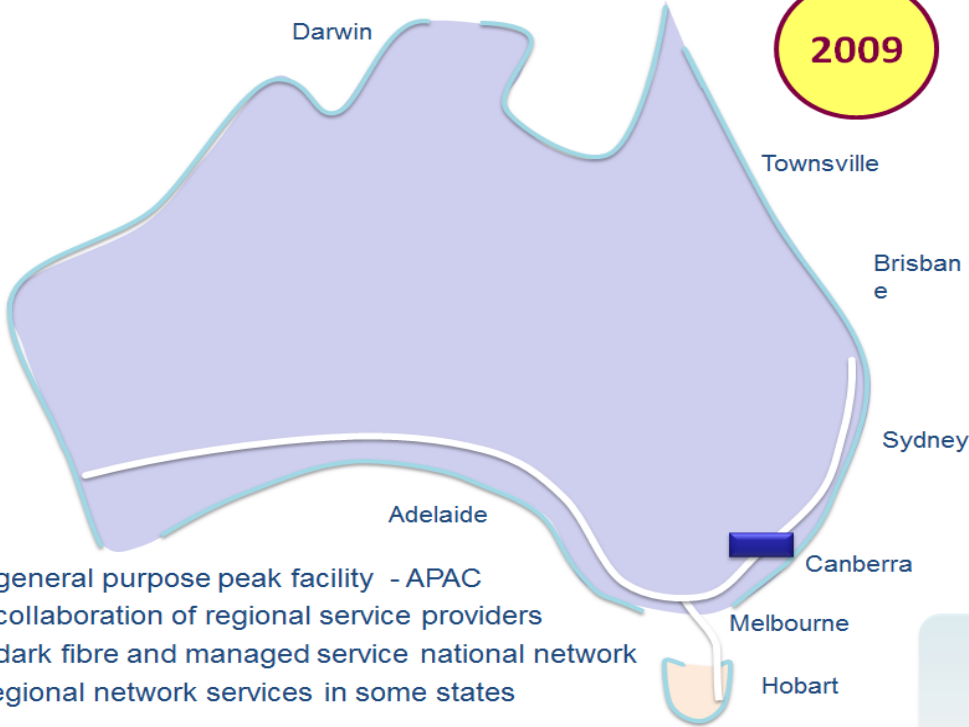
Integrating Big Geoscience Data into the Petascale National Environmental Research Interoperability Platform (NERDIP): Successes and Unforeseen challenges

Lesley Wyborn and Ben Evans.

- Two main tranches of funding:
 - National Collaborative Research Infrastructure Strategy (NCRIS)
 - \$542M for 2006-2011 (\$75 M for cyberinfrastructure)
 - Super Science Initiative
 - \$901 million for 2009-2013 (\$347M for cyberinfrastructure)
 - Annual Maintenance funding of around \$180M pa since 2014-2015
- All programmes were designed ensure that Australian research continues to be competitive and rank highly on an international scale.

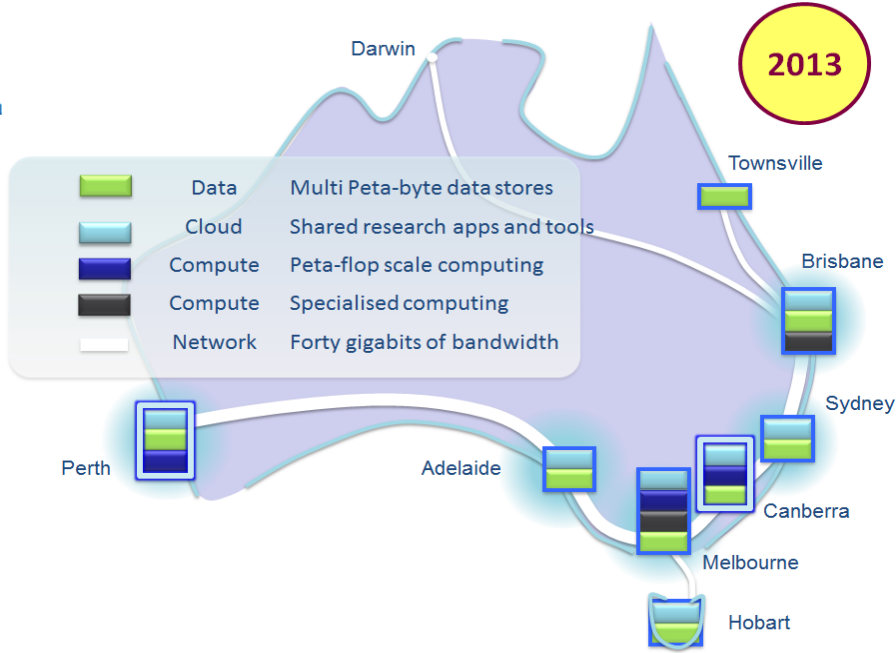


2009

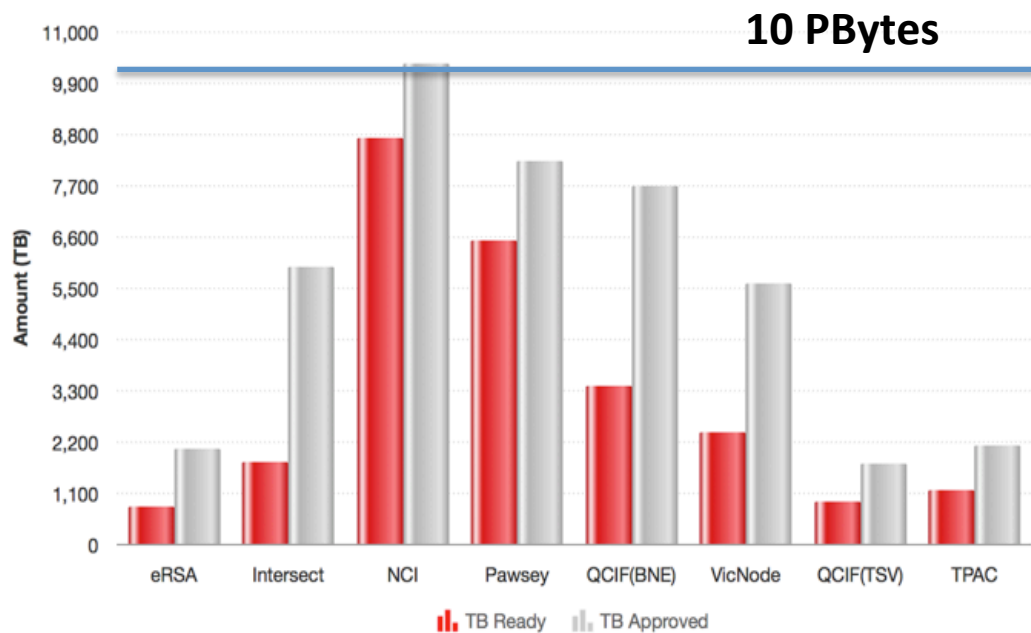
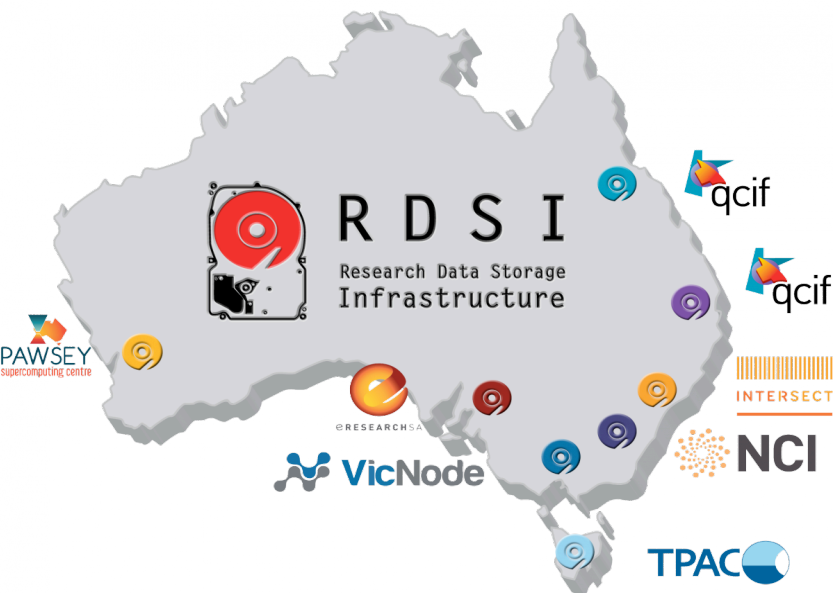


- A general purpose peak facility - APAC
- A collaboration of regional service providers
- A dark fibre and managed service national network
- Regional network services in some states

2013



Progress on Data Ingest as of 16 October, 2015:
~43 Petabytes in 8 distributed nodes

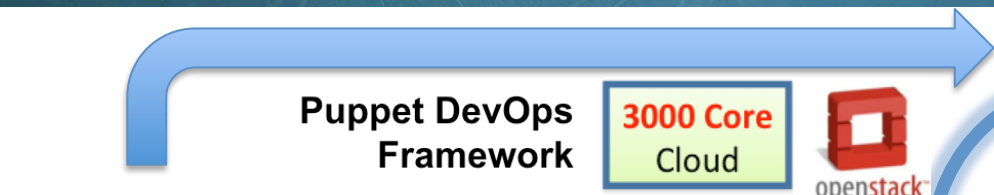


eRSA	Intersect	NCI	Pawsey	QCIF(BNE)	VicNode	QCIF(TSV)	TPAC
785	1743	8709	6504	3406	2404	907	1160
2049	5957	10296	8202	7699	5575	1729	2105

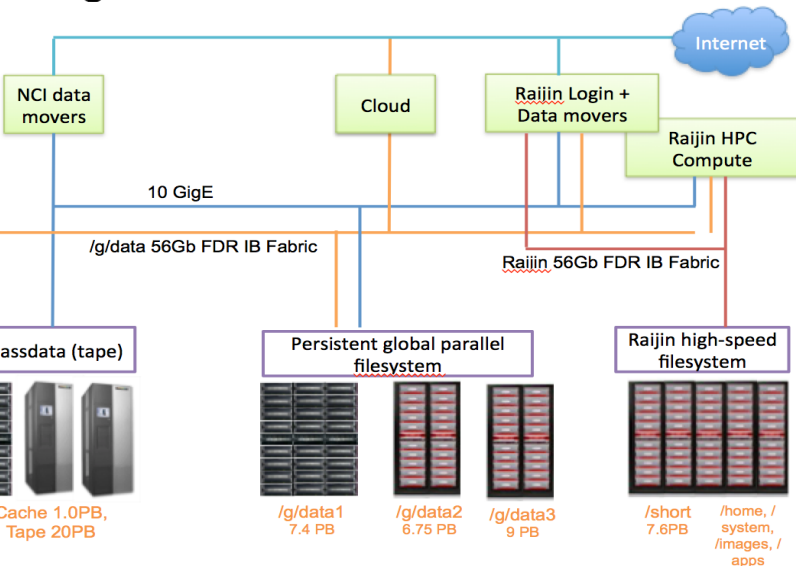
Primary Nodes: eRSA, Intersect, NCI, Pawsey, QCIF(BNE), VicNode, QCIF(TSV)
Additional Nodes: TPAC

Source: <https://www.rds.edu.au/>

Total: TB Ready	25,617	Total: TB Approved	43,611
-----------------	--------	--------------------	--------



Integrated HPC-HPD Environment

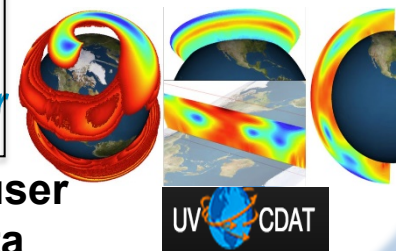


Data Services

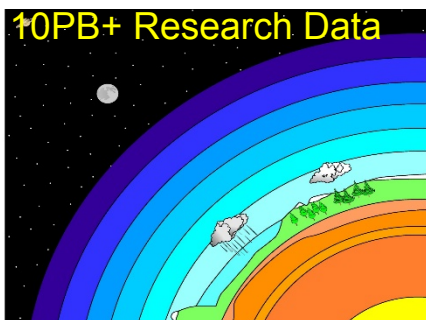
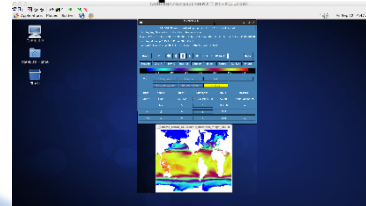
THREDDS **ESGF**
Earth System Grid Federation

OPeNDAP **GeoServer**

Server-side analysis and visualization



VDI: Cloud scale user desktops on data

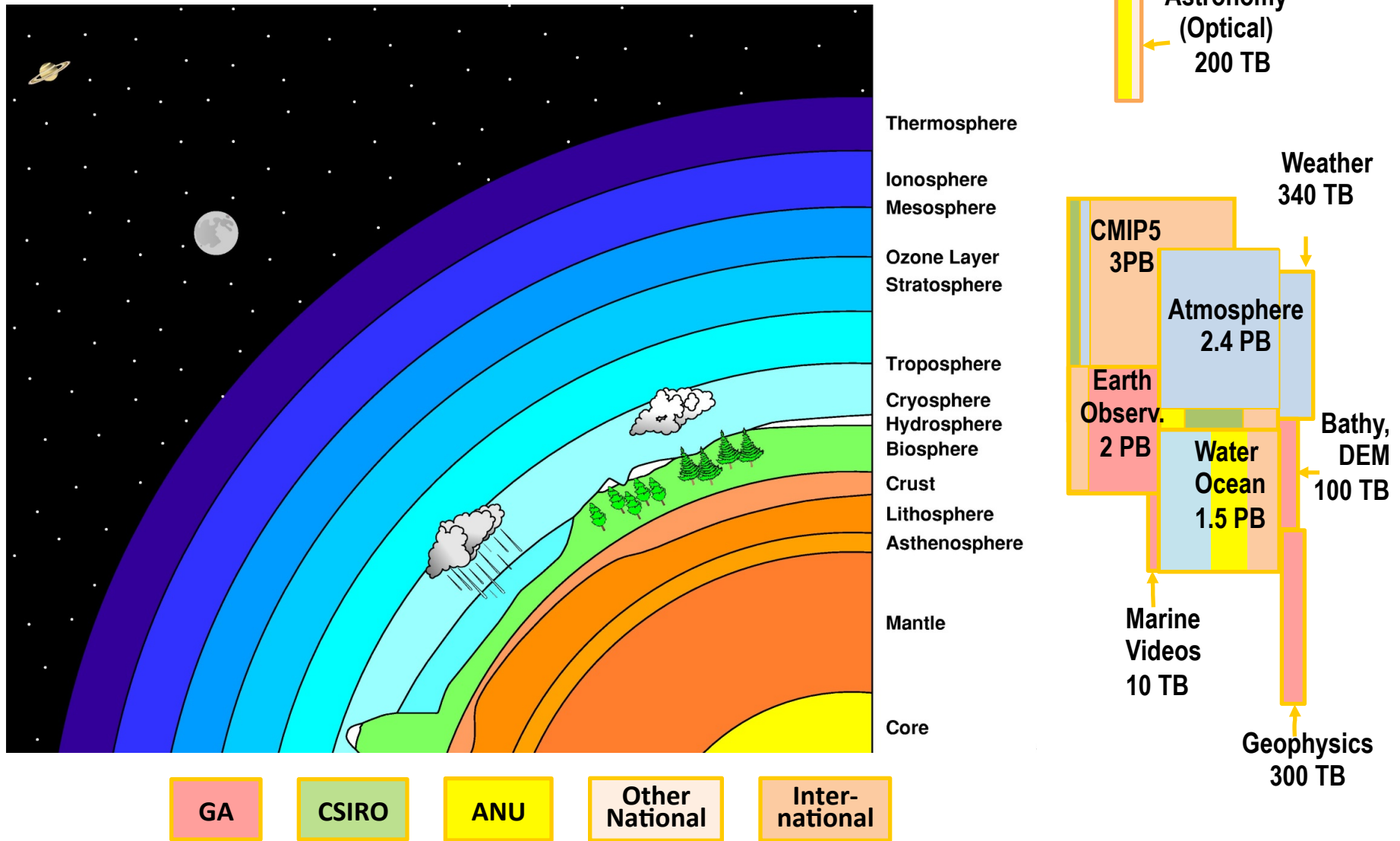


Web-time analytics software



1. Climate/ESS Model Assets and Data Products
2. Earth and Marine Observations and Data Products
3. Geoscience Collections
4. Terrestrial Ecosystems Collections
5. Water Management and Hydrology Collections

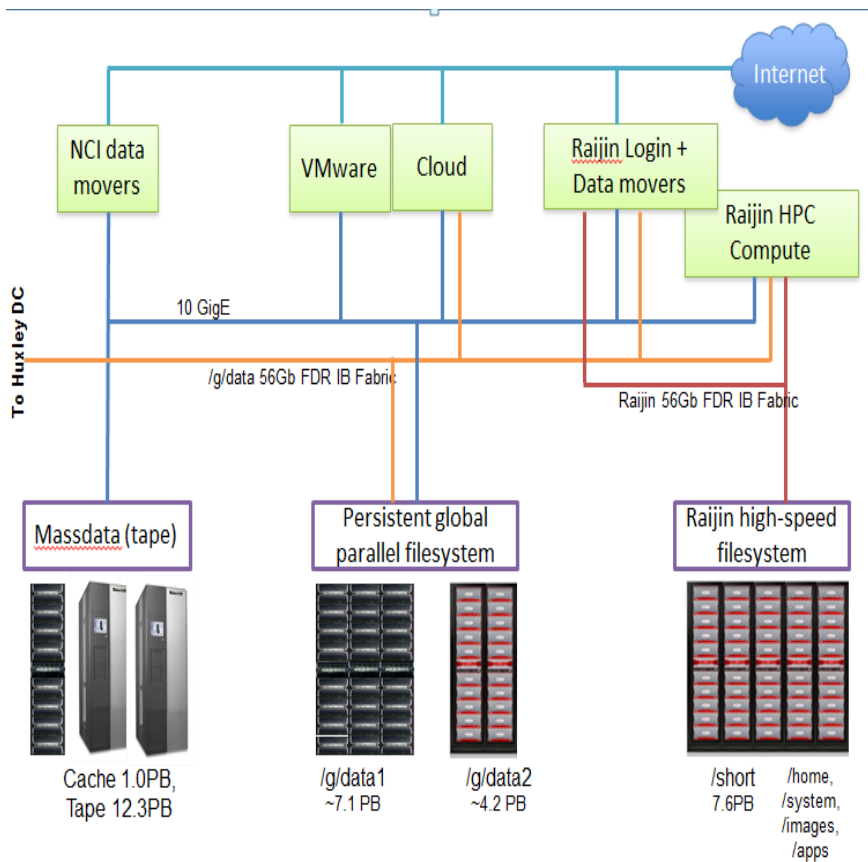
Data Collections	Approx. Capacity
CMIP5, CORDEX	2 Pbytes
ACCESS products	3.3 Pbytes
LANDSAT, MODIS, VIIRS, AVHRR, INSAR, MERIS	2 Pbytes
Digital Elevation, Bathymetry, Onshore Geophysics	400 Tbytes
Seasonal Climate	600 Tbytes
Bureau of Meteorology Observations	400 Tbytes
Bureau of Meteorology Ocean-Marine	220 Tbytes
Terrestrial Ecosystem	290 Tbytes
Reanalysis products	175 Tbytes



- Combined and integrated, the NCI collections are too large to move
 - bandwidth limits the capacity to move them easily
 - the data transfers are too slow, complicated and too expensive
 - even if our data can be moved, few can afford to store 10 PB on spinning disk
- We need to change our focus to:
 - moving users to the data (for sophisticated analysis)
 - moving processing to data
 - having online applications to process the data in-situ
 - Improving the sophistication of users – with our help
- We called for a new form of system design where:
 - storage and various types of computation are co-located
 - systems are programmed and operated to allow users to interactively invoke different forms of analysis in-situ over integrated large-scale data collections



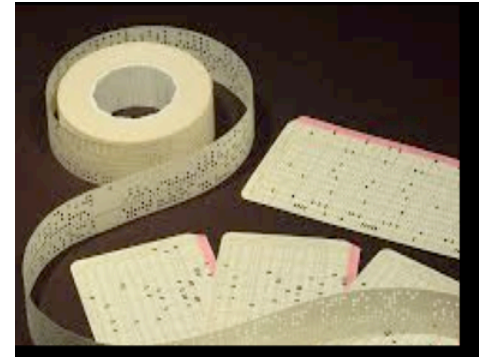
- Work at NCI has also highlighted the need for balanced systems to enable Data-intensive Science including:
 - Interconnecting processes and high throughput to reduce inefficiencies
 - The need to really care about placement of data resources
 - Better communications between the nodes
 - I/O capability to match the computational power
 - Close coupling of cluster, cloud and storage



NCI's Integrated High Performance Environment

1. Volume: data at rest
2. Velocity: data in motion (streaming)
3. Variety: many types, forms and structures (or no structures)
4. Veracity: trustworthiness, provenance, lineage, quality
5. Validity: data that is correct
6. Visualization: data in patterns
7. Vulnerability: data at risk
8. Value: data that is meaningful
9. V?????
10. V?????

- Big Data is a relative term where the volume, velocity and variety of data exceed an organisations storage or compute capacity for accurate and timely decision making
- We define High Performance Data (HPD) as data that is carefully prepared, standardised and structured so that it can be used in Data-Intensive Science on HPC (Evans et al., 2015)
- To get on top of the Data Tsunami, we need to convert 'Big data' collections into HPD by
 - Aggregating data into seamless 'pre-processed' data products
 - Creating hyper-cubes and self describing data arrays



1964: 1KB = 2m of tape or ~ 20 cards

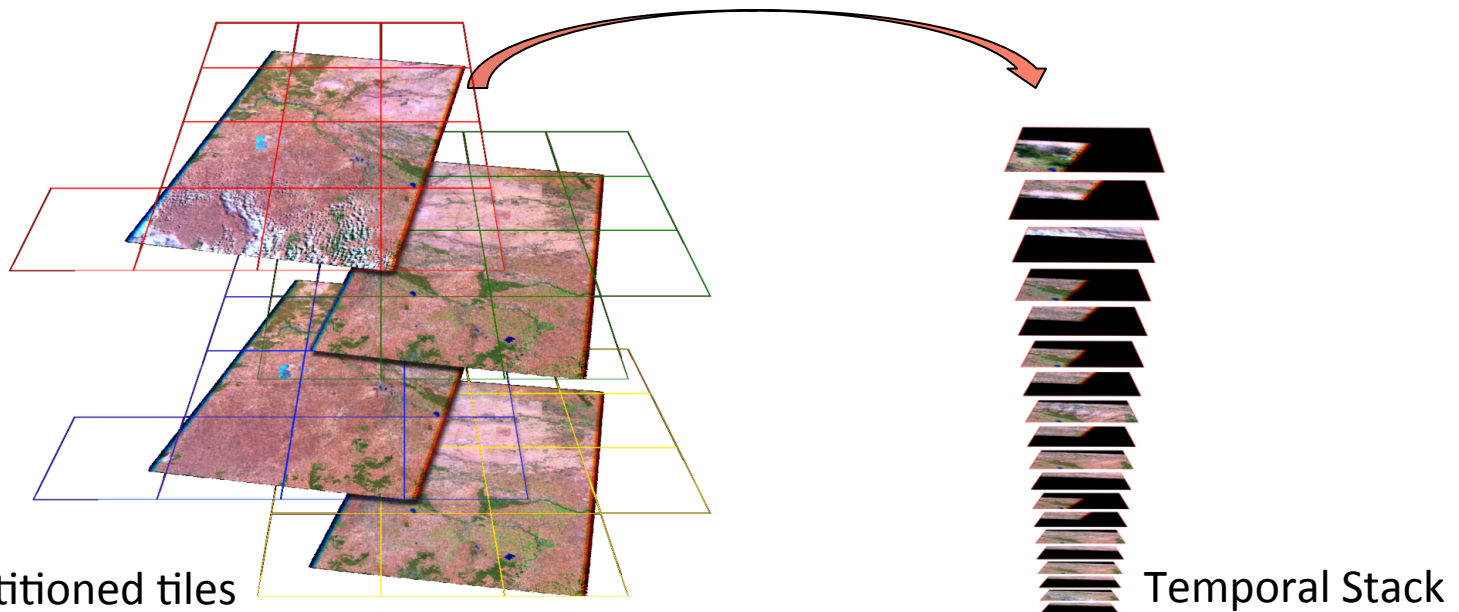


2014: a 4 GB Thumb drive =
~8000 Km of Tape
or ~83 million cards

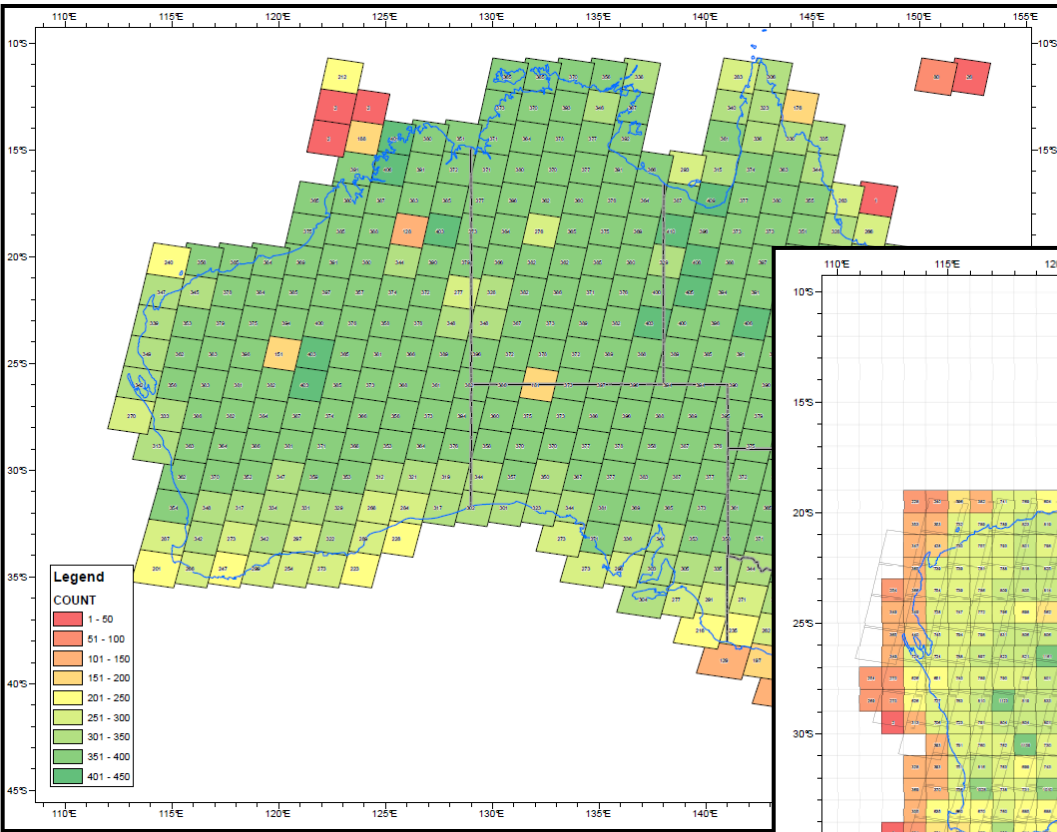


2014: 20 PB of modern storage =
~ 32 trillion metres of tape
~ 320 trillion cards

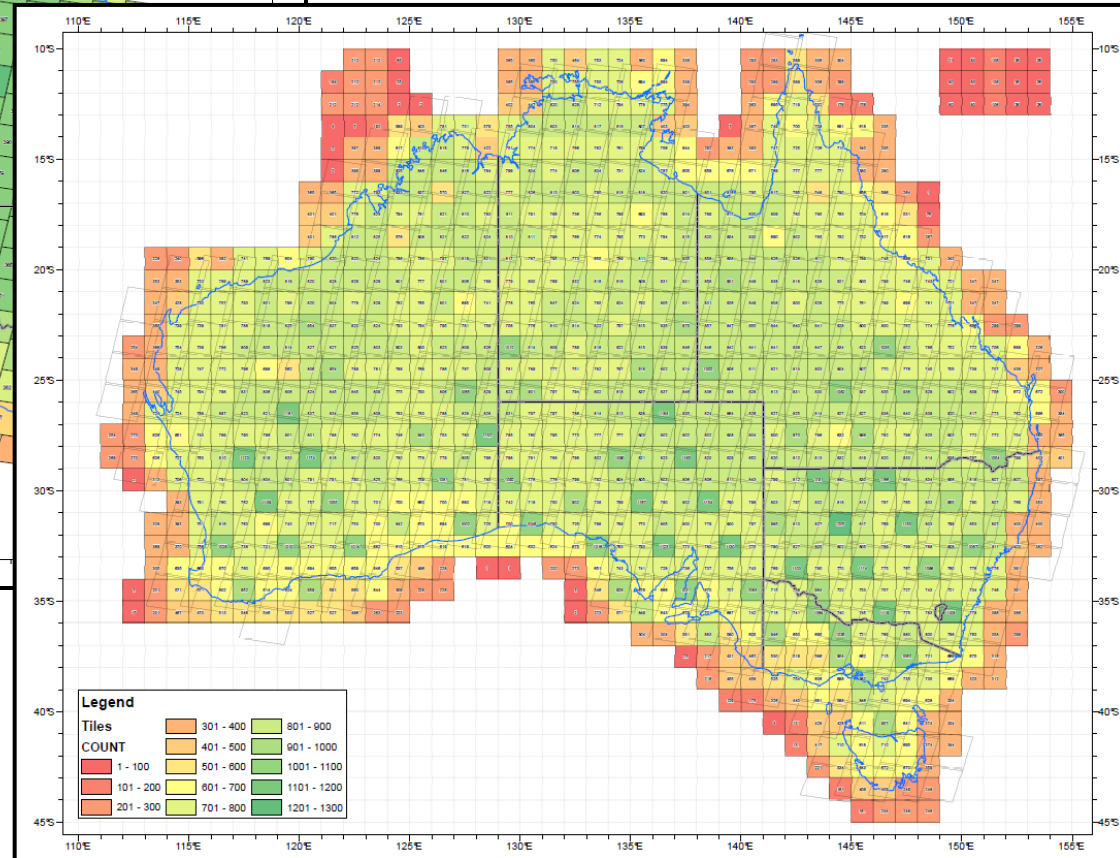
- The Landsat cube arranges 636,000 Landsat Source scenes spatially and temporally, to allow flexible but efficient large-scale in-situ analysis
- The data is partitioned into spatially-regular, time-stamped, band-aggregated tiles which are presented as temporal stacks.



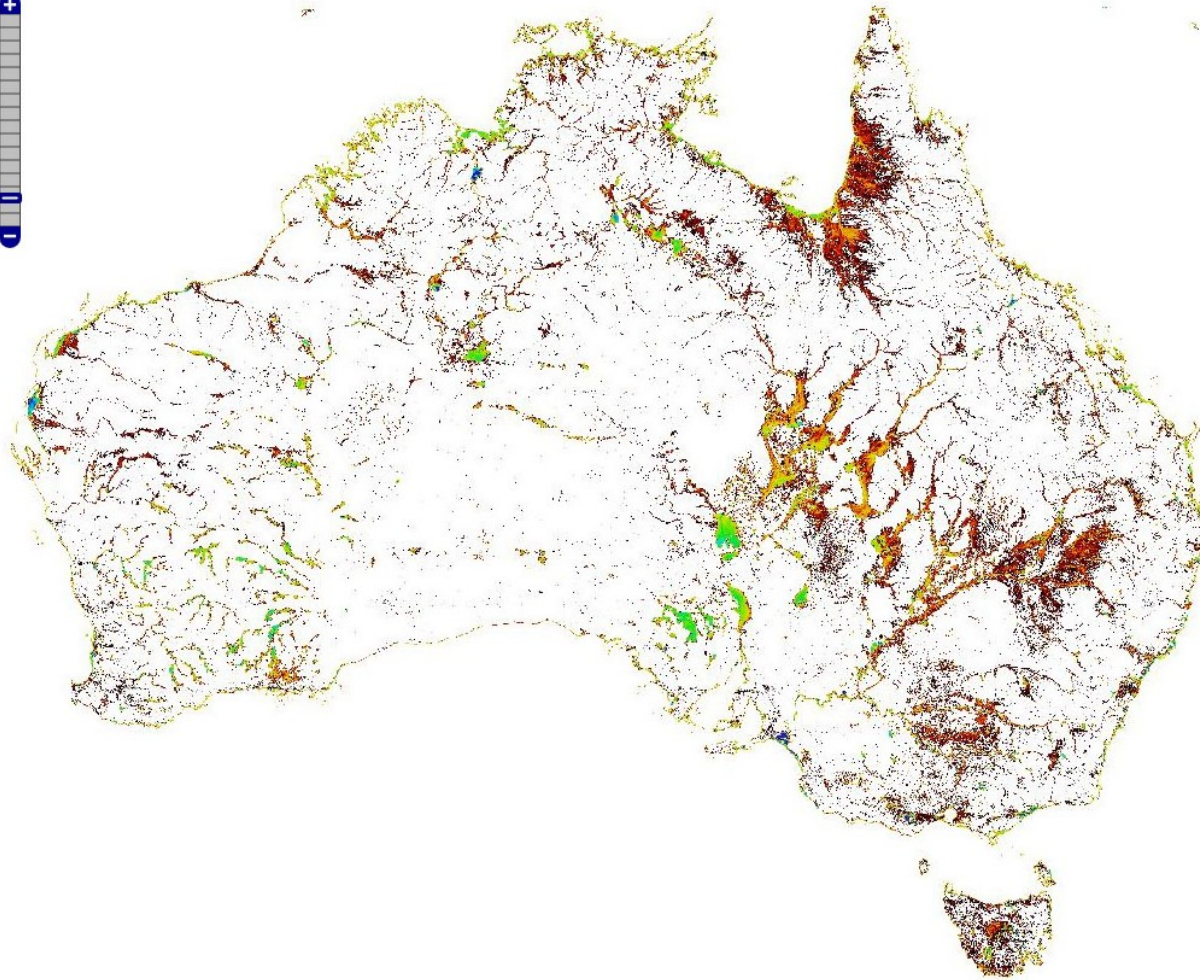
4M Spatially-Regular Time-Stamped Tiles (0.5 PB)



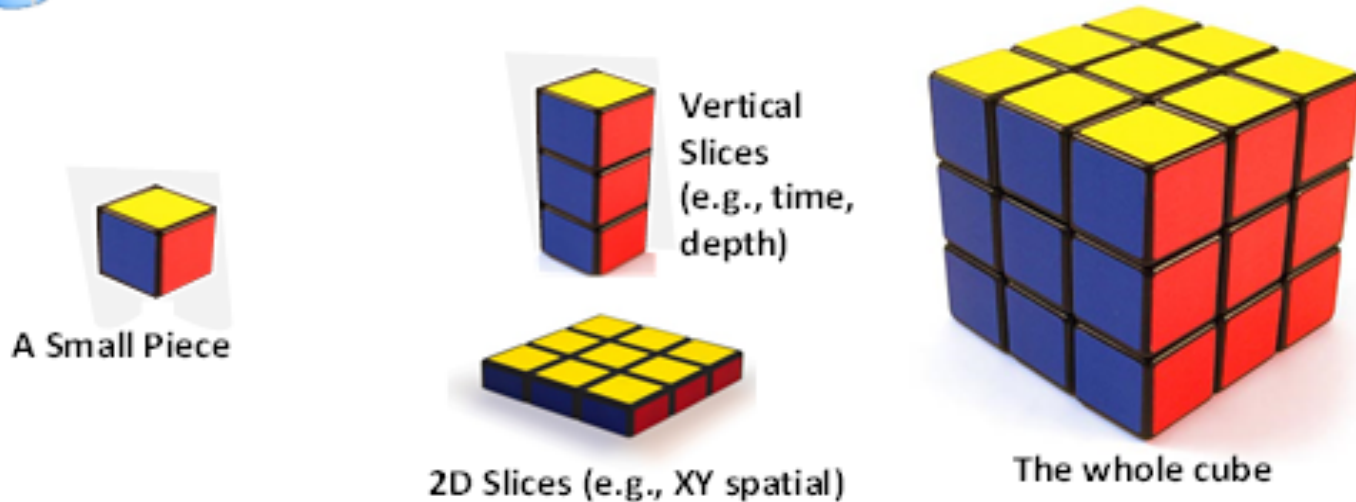
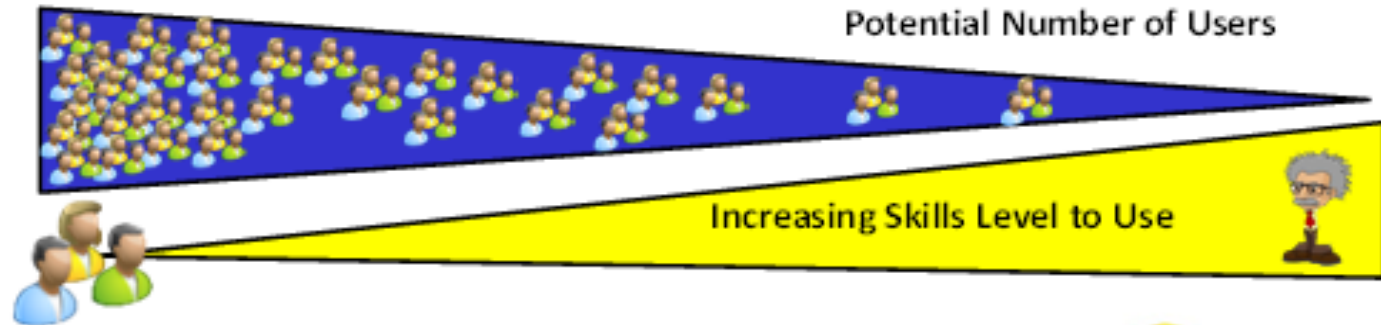
636,000 Landsat Source Scene Datasets
(~52 x 10¹² Pixels)

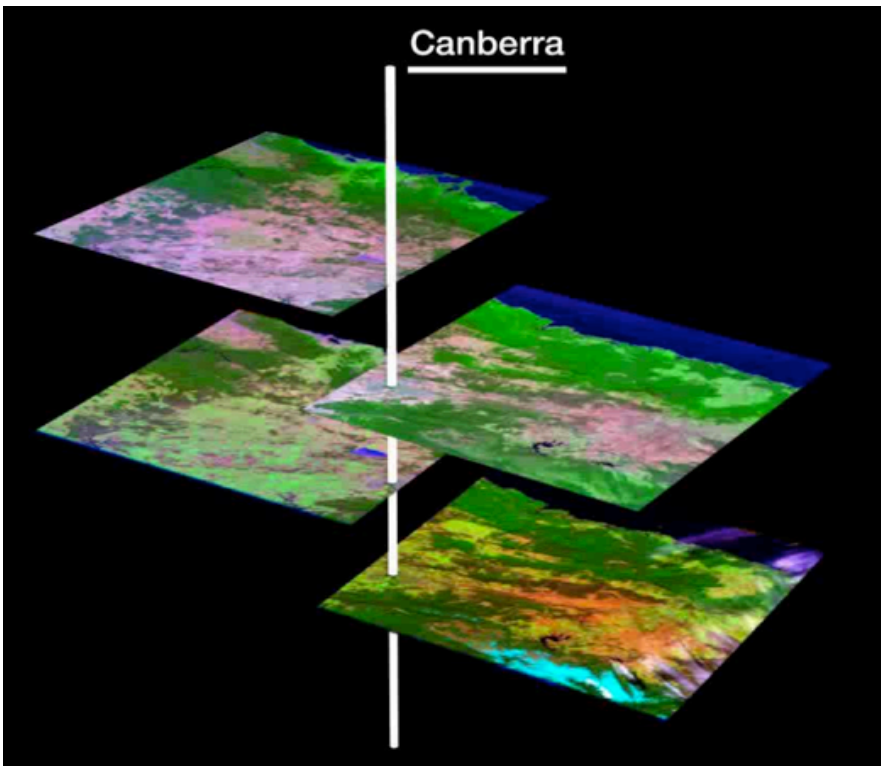


Water Detection from Space



- 15 Years of data from LS5 & LS7(1998-2012)
- 25m Nominal Pixel Resolution
- Approx. 133,000 individual source scenes in approx. 12,400 passes
- Entire archive of 1,312,087 ARG25 tiles => 21×10^{12} pixels can be processed in ~8 hours

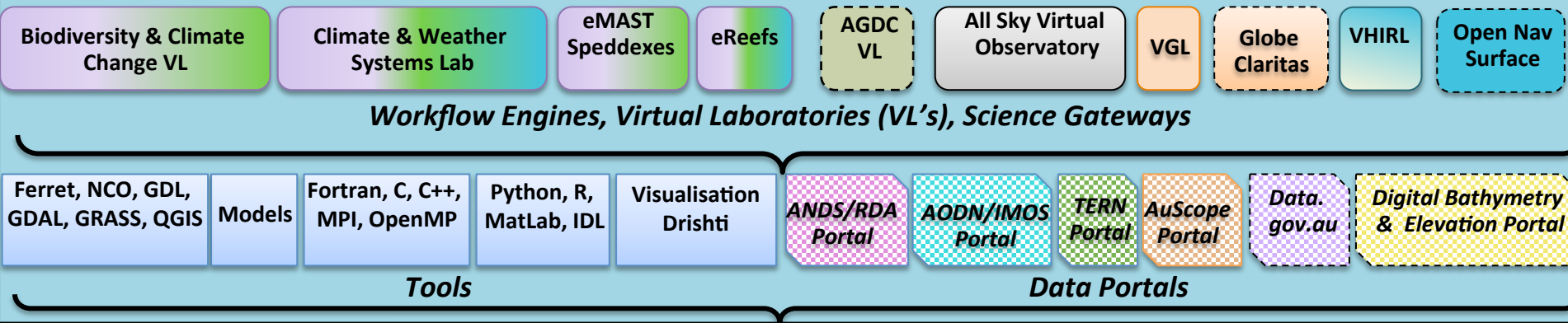




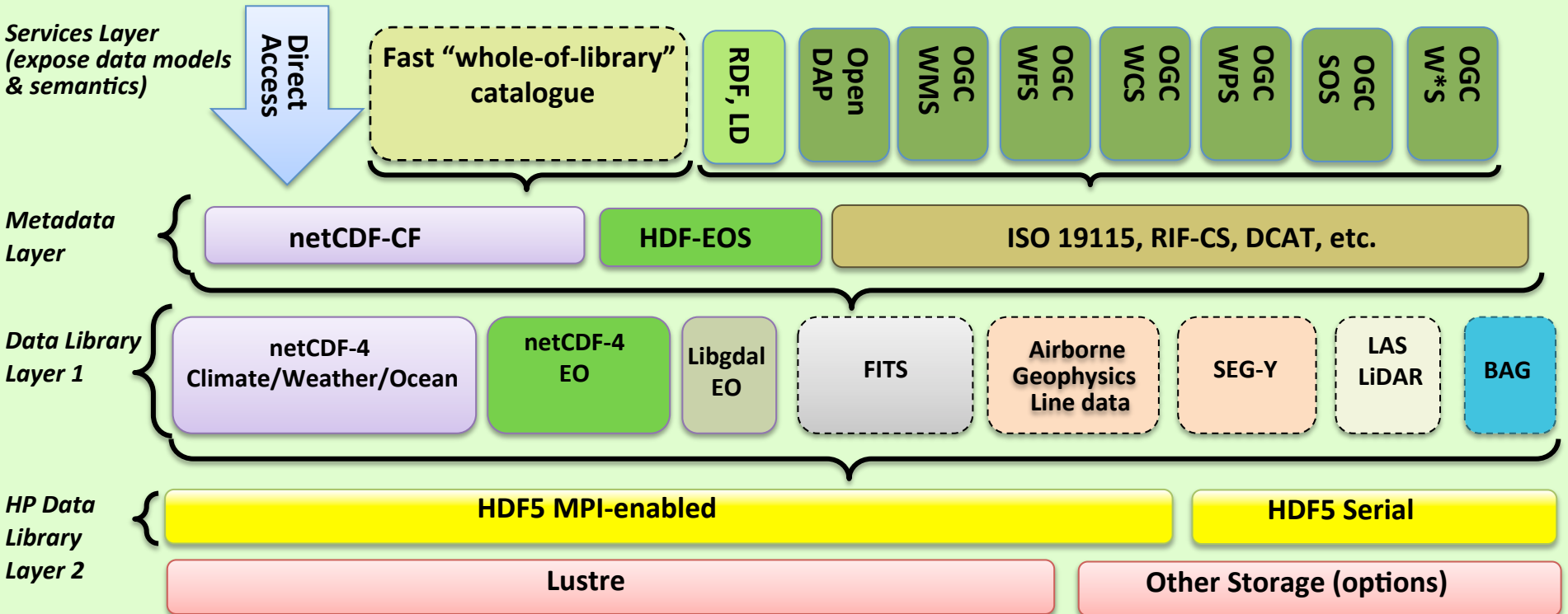
Do we enable individual scenes to be downloaded for locally hosted small scale analysis?
Or do we facilitate small scale analysis, in-situ on data sets that are dynamically updated?



Introducing the National Environmental Data Interoperability Research Platform (NERDIP)



National Environmental Research Data Interoperability Platform (NERDIP)



Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

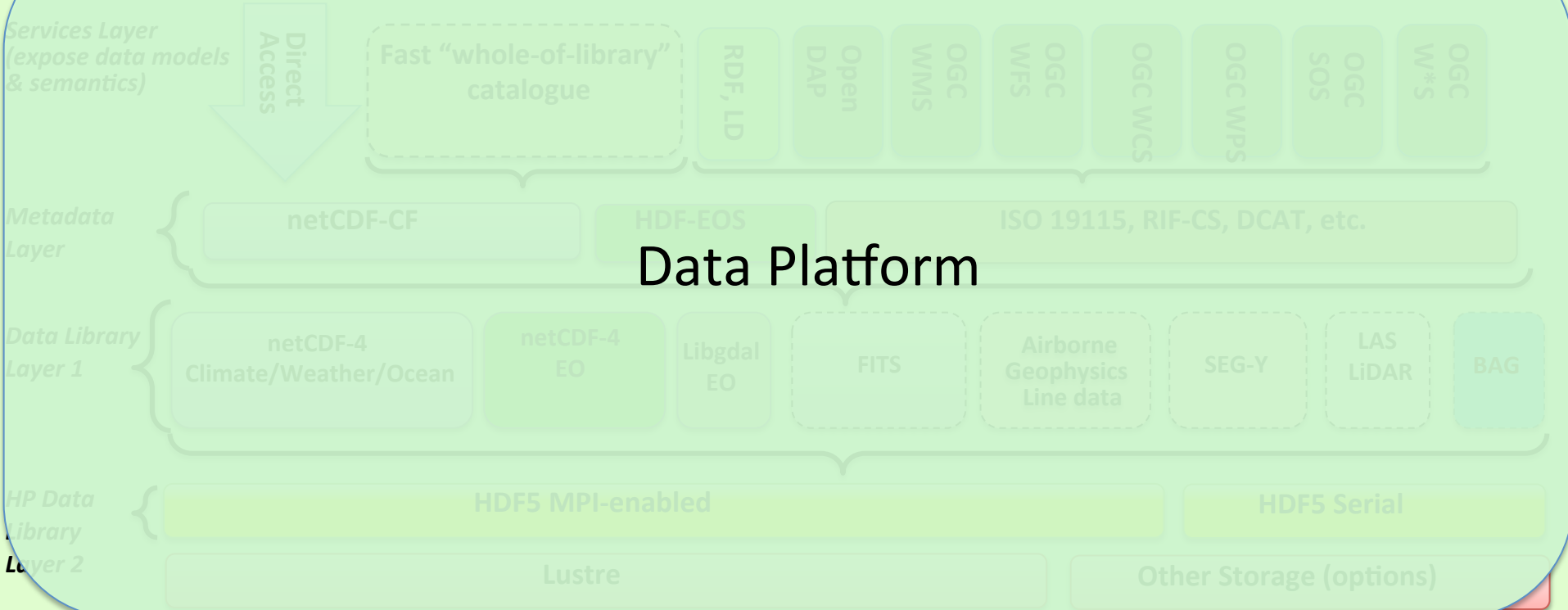
Ace Users

Tools

Data Discovery

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)



Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

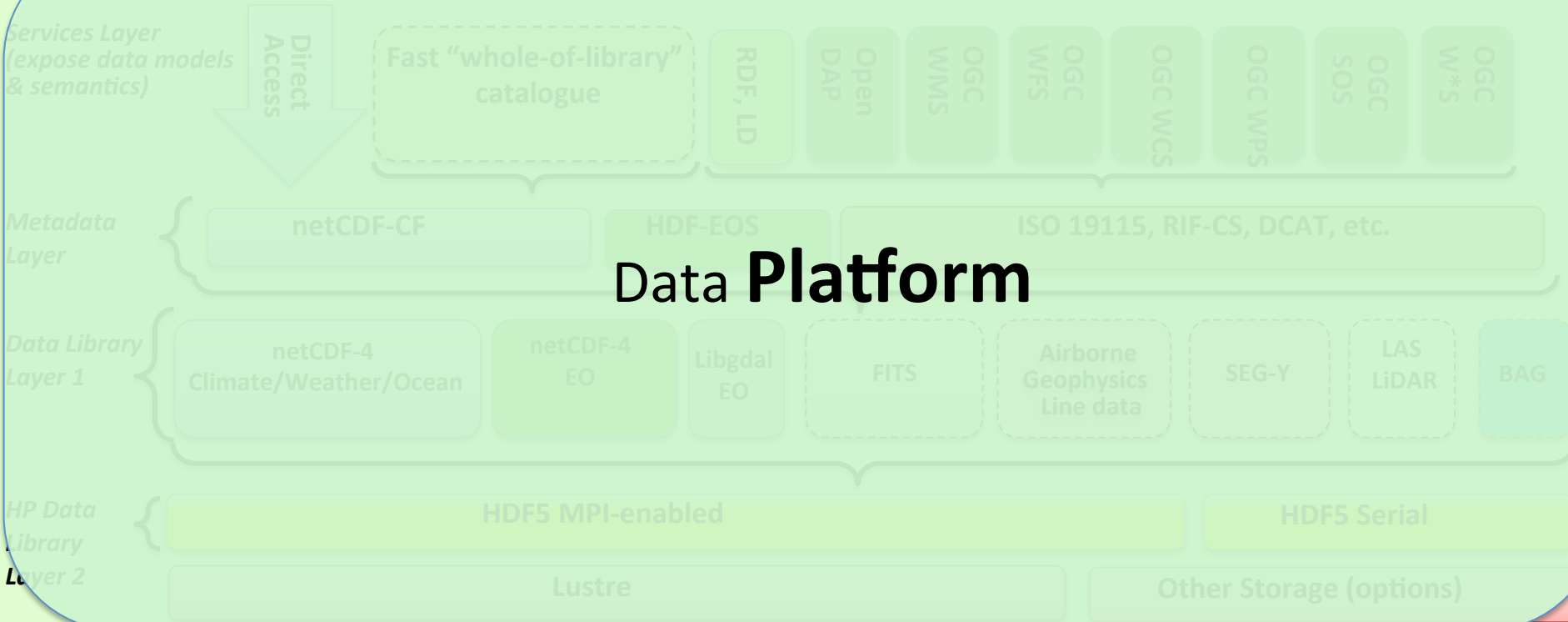
Ace Users

Tools

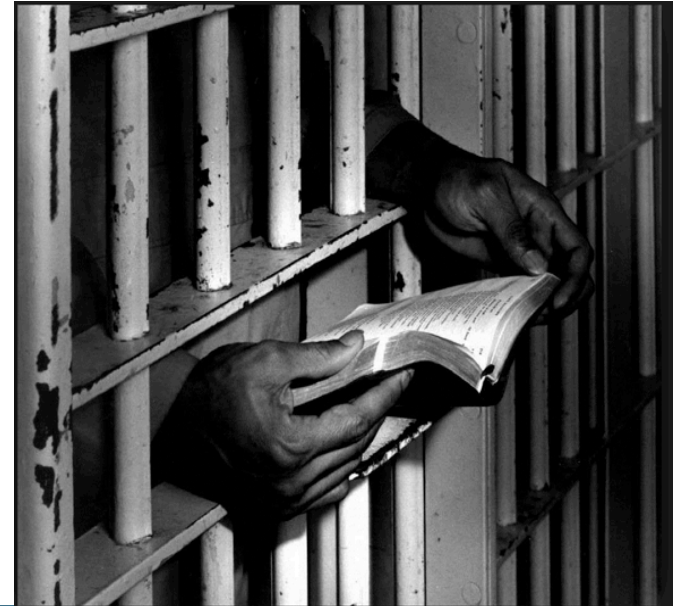
Data Portals

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)



- **Portals are for visiting, platforms are for building on**
- Portals present aggregated content in a way that invites exploration, but the experience is pre-determined by a set of decisions by the builder about what is necessary, relevant and useful.
- Platforms put design decisions into the hands of users: there are innumerable ways of interacting with the data
- Platforms offer many more opportunities for innovation: new interfaces can be built, new visualisations framed, ultimately new science rapidly emerges



Tim Sherratt <http://www.nla.gov.au/our-publications/staff-papers/from-portal-to-platform>

Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

Ace Users

Data Discovery

Tools

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)

*Services Layer
(expose data models & semantics)*

Fast "whole-of-library" catalogue

Data Platform

RDF, LD, Open DAP, W/M/S, OGC WFS, OGC WCS, OGC WPS, SOS, OGC W+S, ISO 19115, RIF-CS, DCAT, etc., FITS, Airborne Geophysics Line data, SEG-Y, LAS LiDAR, BAG, HDF5 MPI-enabled, HDF5 Serial, Other Storage (options)

Metadata Layer

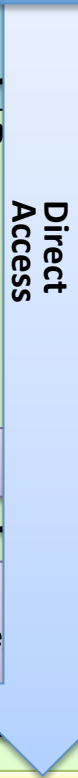
CDF-CF, HDF-EOS

Data Library Layer 1

Climate, F-4, netCDF-4 EO, Libgdal EO

HP Data Library Layer 2

HDF5 MPI-enabled, Lustre



APPLICATION

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

FOCUSSED DEVELOPERS

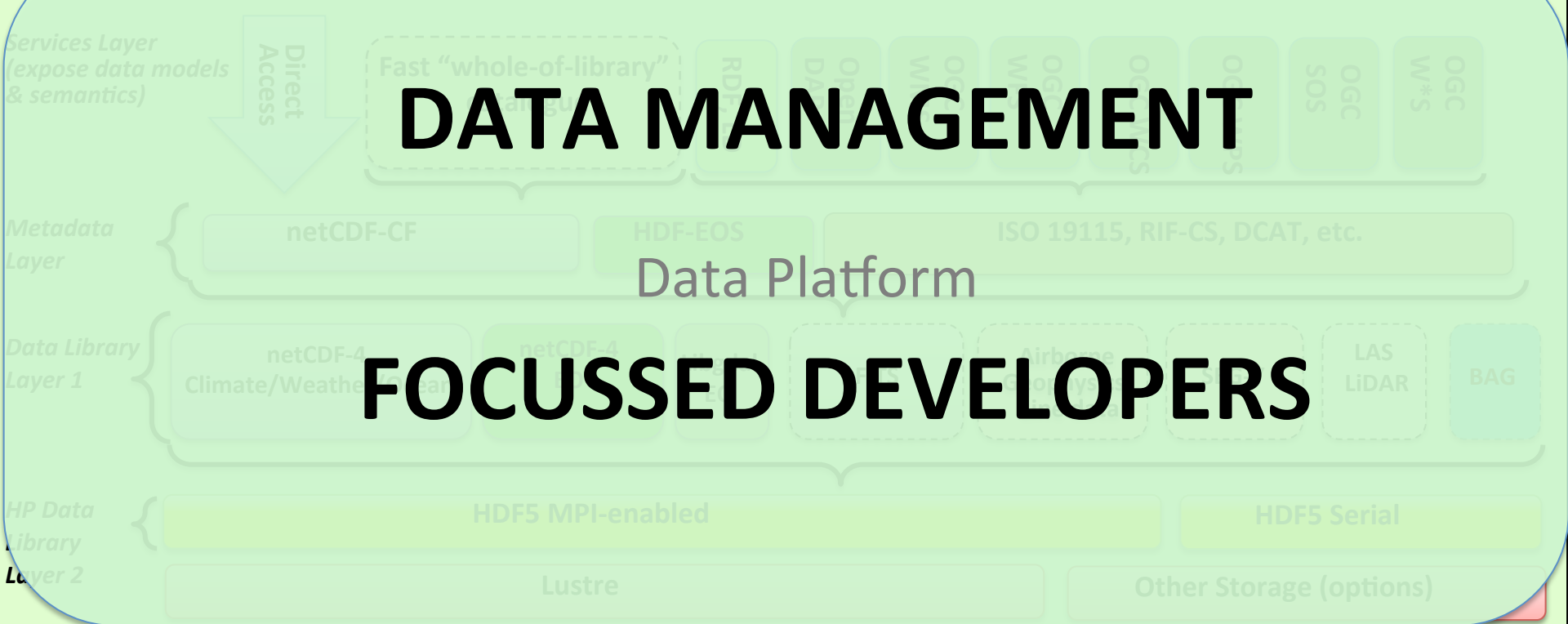
Tools

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)

DATA MANAGEMENT

FOCUSSED DEVELOPERS



Infrastructure to Lower Barriers to Entry

APPLICATION

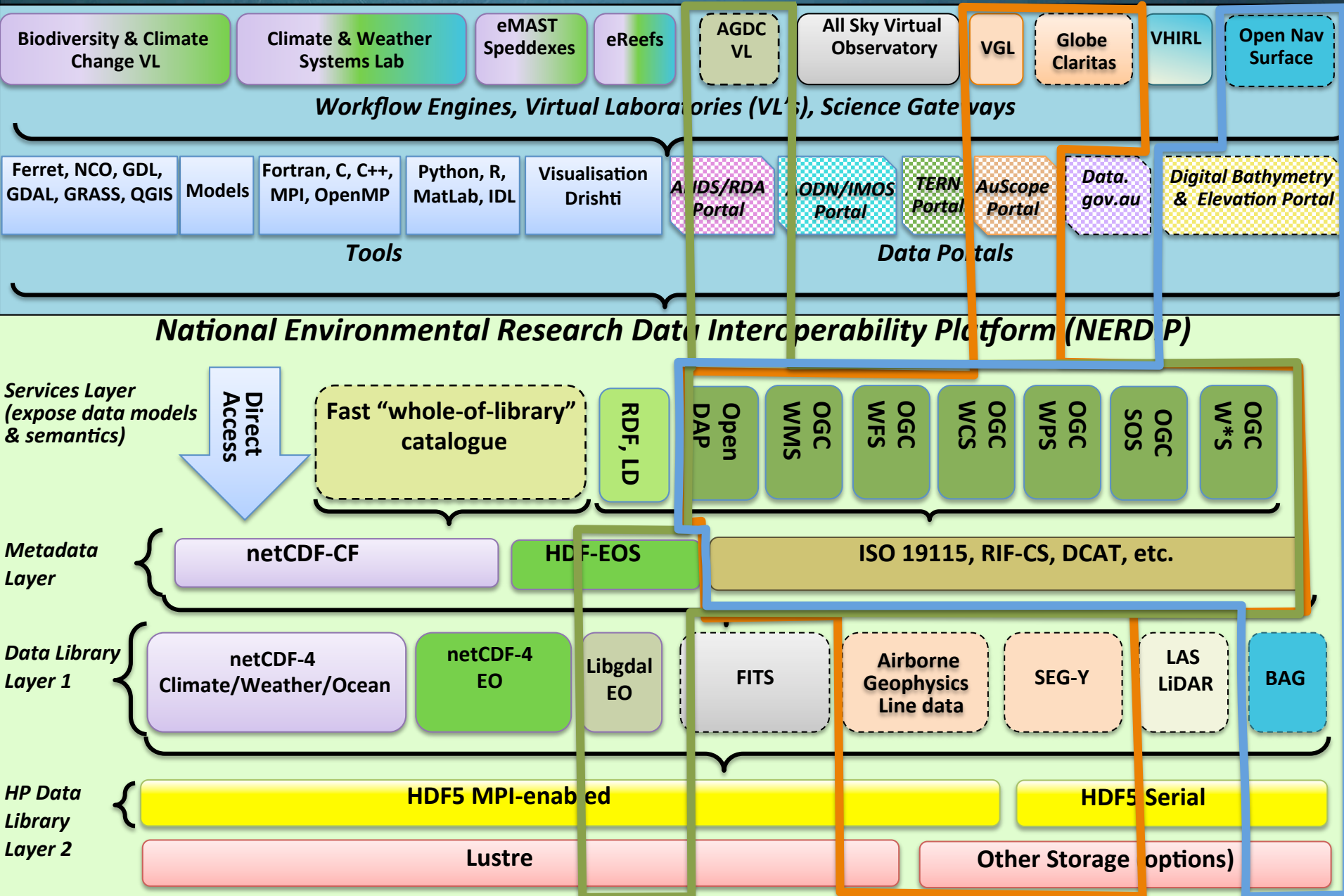
Ace Users

Data Discovery

FOCUSSED DEVELOPERS

DATA MANAGEMENT
FOCUSSED DEVELOPERS





APPLICATION

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

FOCUSSED DEVELOPERS

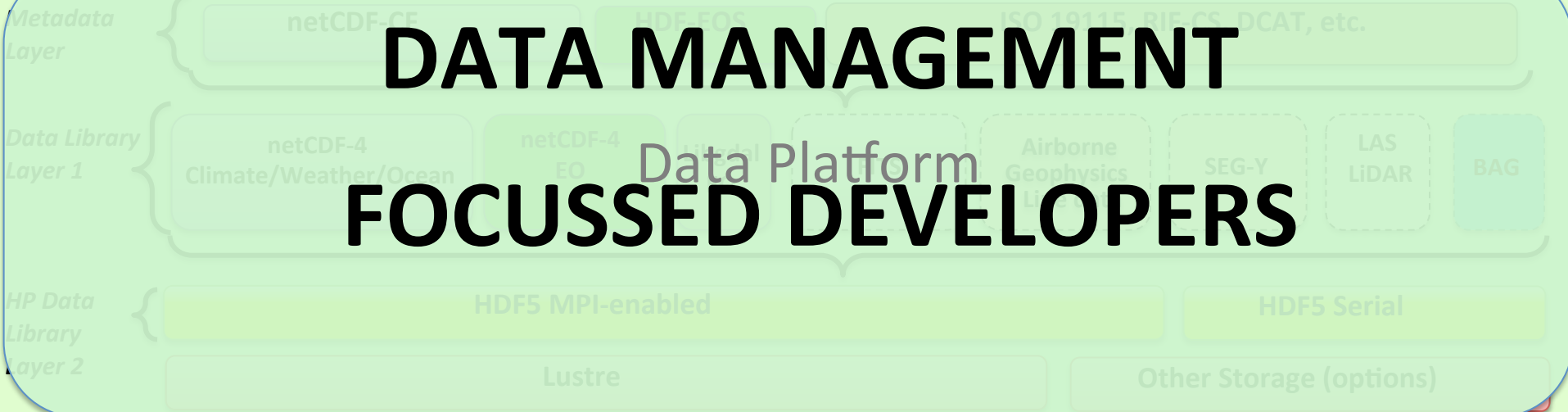
Tools

Data Portals

SERVICES INTERFACE

DATA MANAGEMENT

FOCUSSED DEVELOPERS



Infrastructure to Lower Barriers to Entry

APPLICATION

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

FOCUSSED DEVELOPERS

Tools

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)

Services Layer
(expose data models
& semantics)

Direct
Access

Fast "whole-of-library"
catalogue

RDF, LD

Open
DAP

OGC
WMS

OGC
WFS

OGC
WCS

OGC
WPS

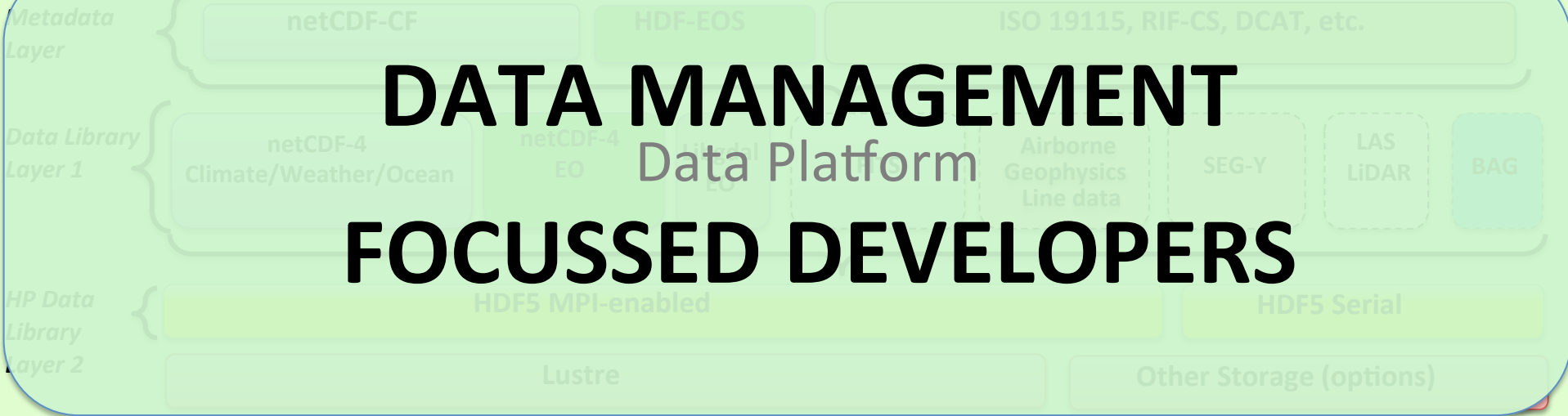
OGC
SOS

OGC
W*S

DATA MANAGEMENT

Data Platform

FOCUSSED DEVELOPERS



Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

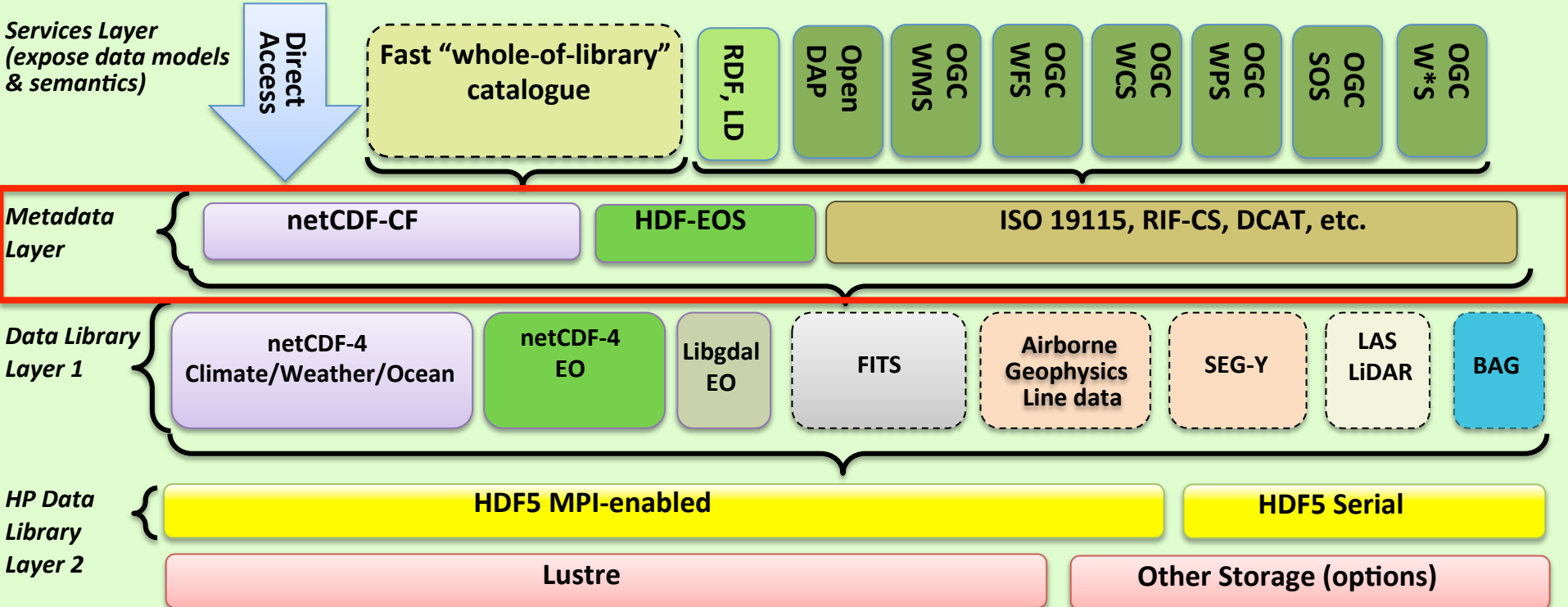
Ace Users

Data Discovery

Tools

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)



Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

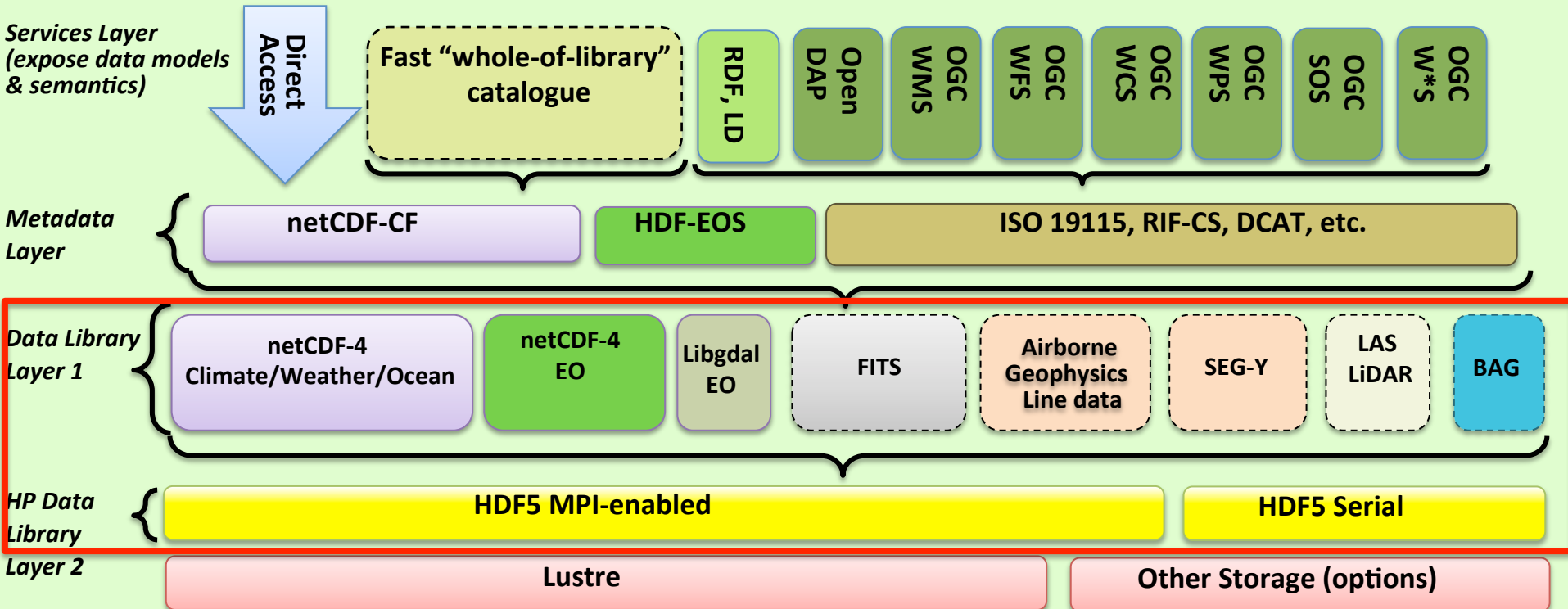
Ace Users

Tools

Data Discovery

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)



Infrastructure to Lower Barriers to Entry

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

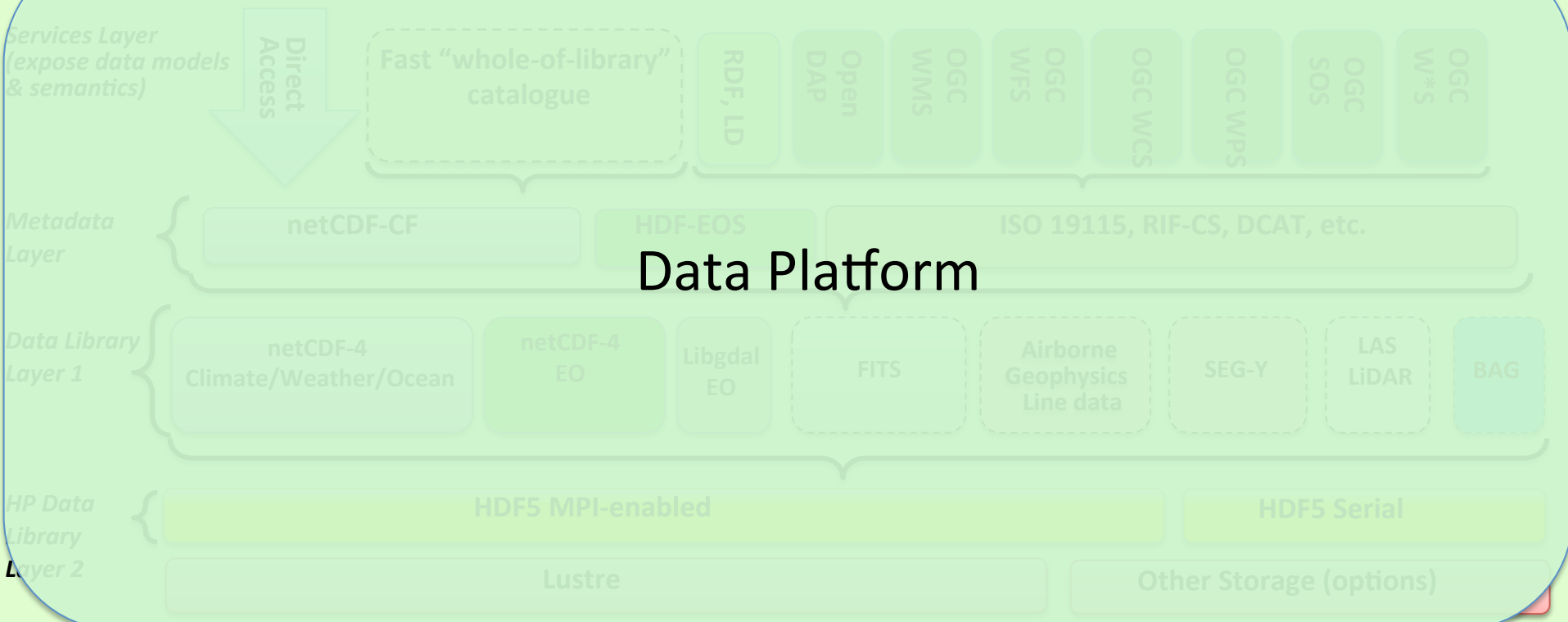
Ace Users

Data Discovery

Tools

Data Portals

National Environmental Research Data Interoperability Platform (NERDIP)



Biodiversity & Climate Change VL
Climate & Weather Systems Lab
eMAST Speddexes
eReefs
AGDC VL
All Sky Virtual Observatory
VGL
Globe Claritas
VHIRL
Open Nav Surface

Workflow Engines, Virtual Laboratories (VL's), Science Gateways

Ace Users

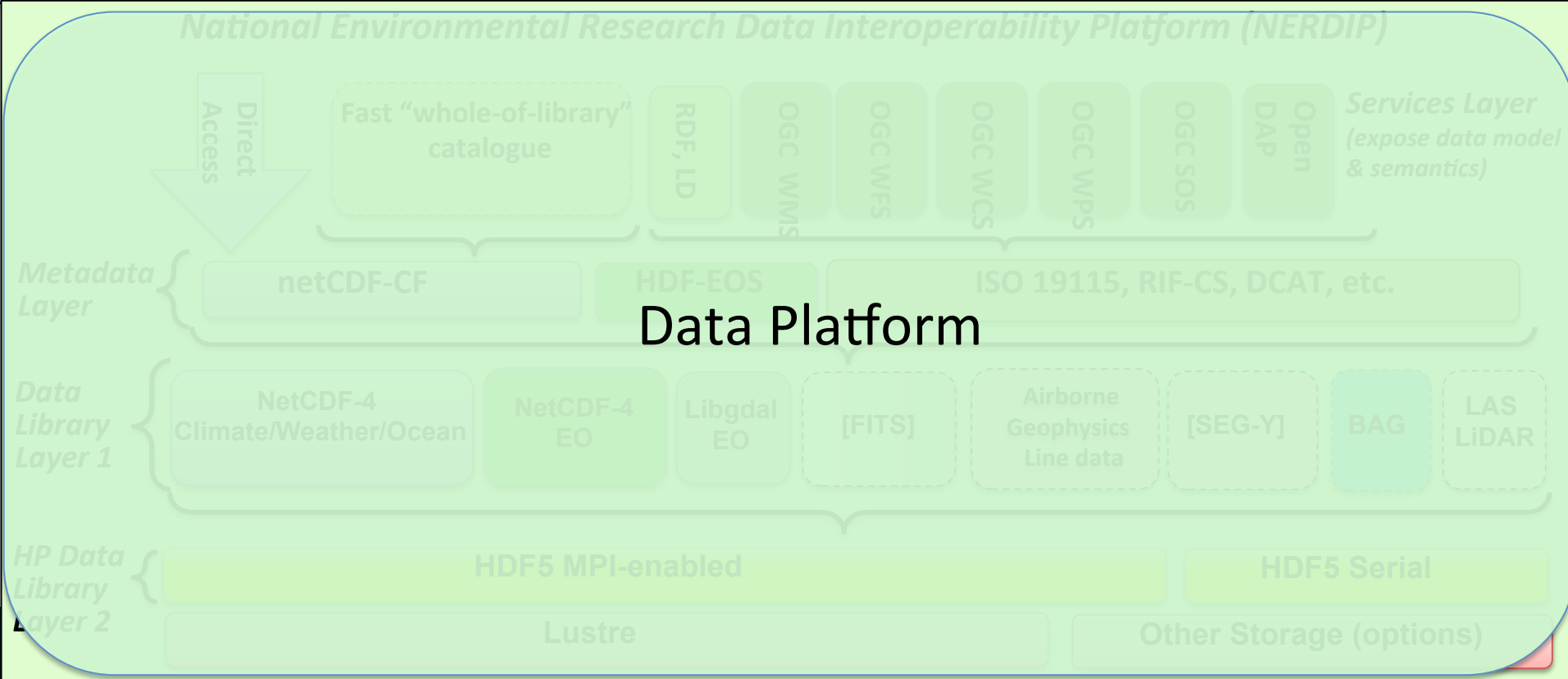
Ferret, NCO, GDL, GDAL, GRASS, QGIS
Models
Fortran, MPI, OpenMP
MatLab, IDL
Visualisation Drishti

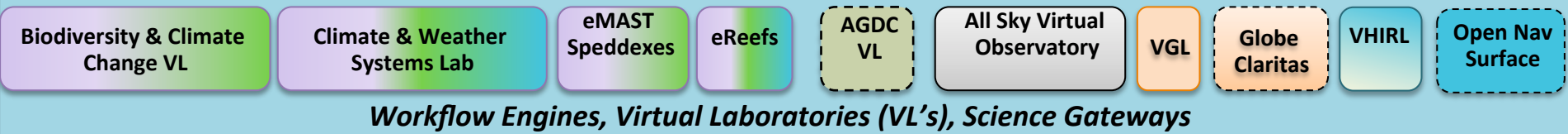
Tools

Data Discovery

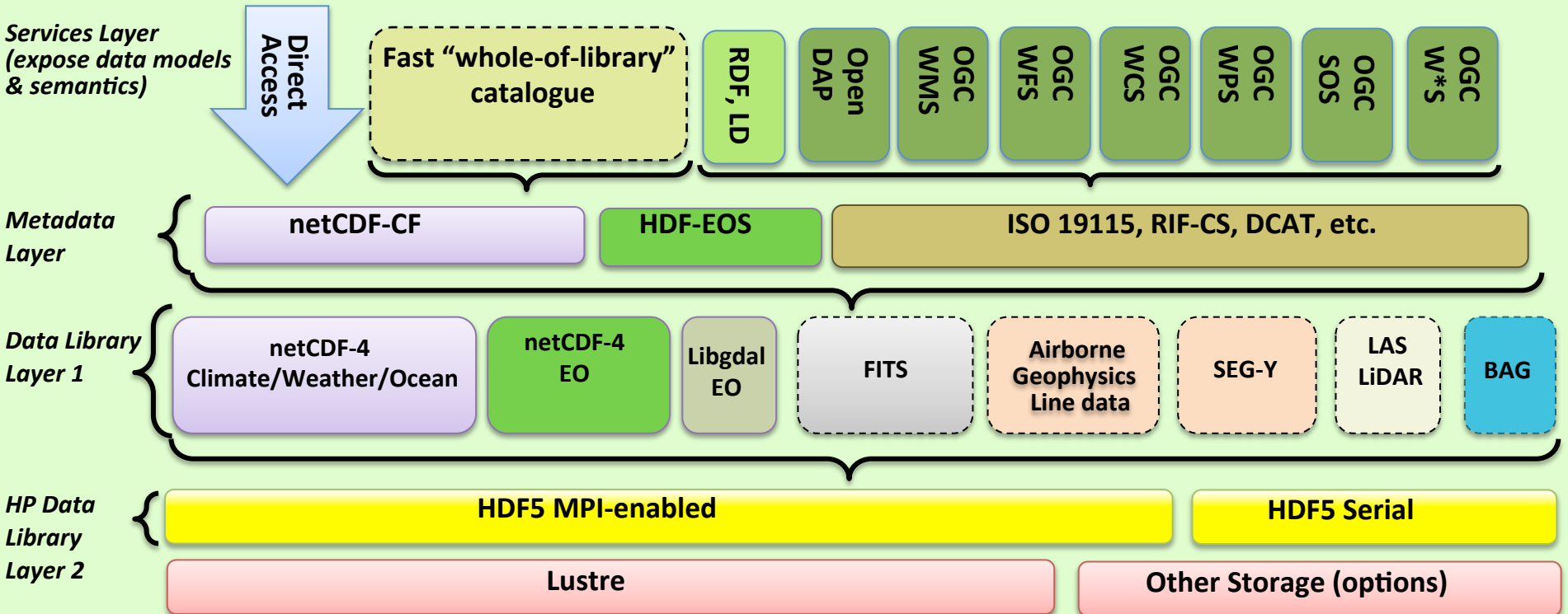
NDSD/RDA Portal
AODN/II Portal
Portal
Portal
gov.au
Digital Bathymetry & Elevation Portal

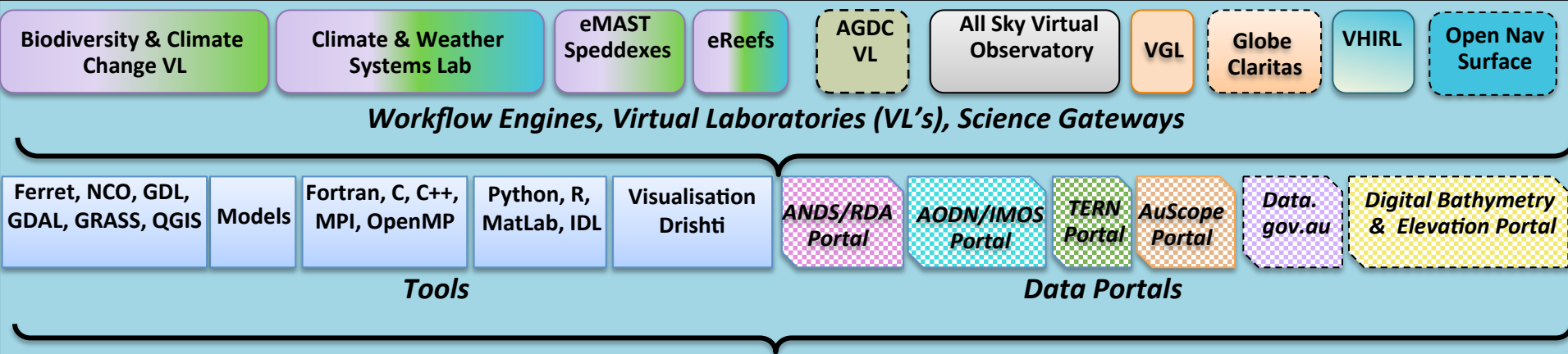
Data Portals



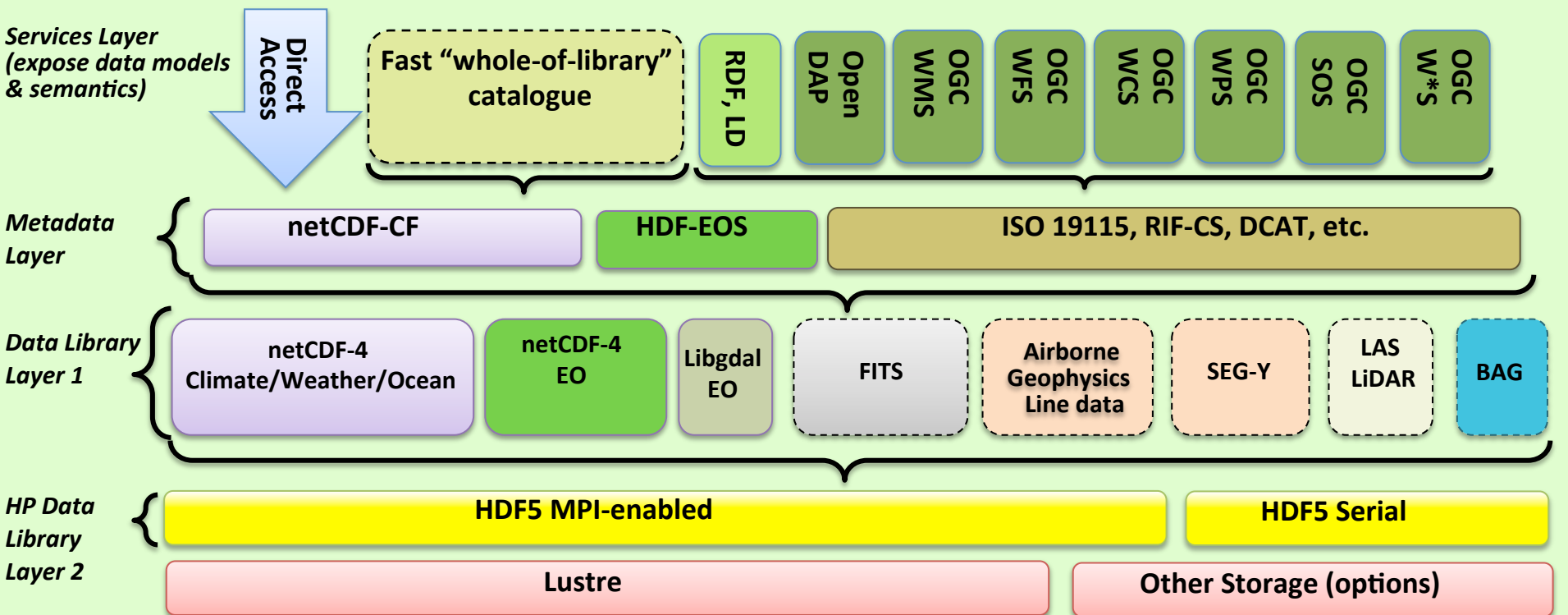


National Environmental Research Data Interoperability Platform (NERDIP)





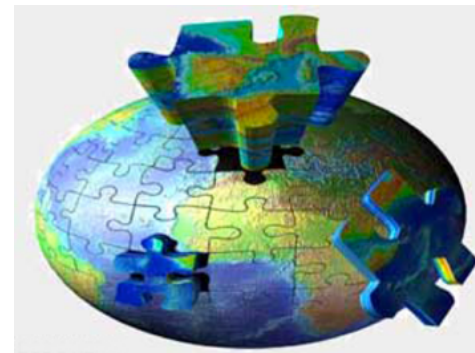
National Environmental Research Data Interoperability Platform (NERDIP)



- Data at scales of today have to be built as shared global facilities based around national institutions.
- Domain-neutral international standards for data collections and interoperability are critical for allowing complex interactions in HP environments both within and between HPD collections
- **No one can do it alone.** No one organisation, no one group, no one country has the required resources or the expertise.
- Shared collaborative efforts such as Research Data Alliance, the Earth Systems Grid Federation (ESGF), the Belmont Forum, EarthServer, the Oceans Data Interoperability Platform (ODIP), EarthCube, GEO and OneGeology are needed to realise the full potential of the new data intensive science infrastructures
- For it now takes a ‘village of partnerships’ to raise a ‘HPD data center’ in a Big Data World



https://www.sfwa.org/wp-content/uploads/2010/06/iStock_000012734413XSmall.jpg



<http://www.onegeology.org/>