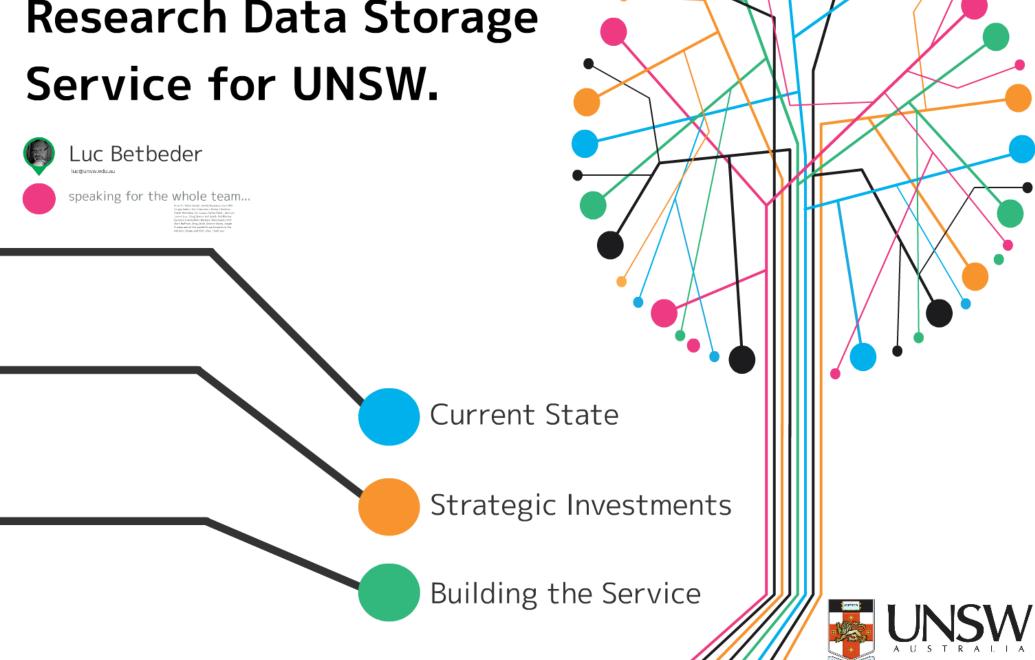


Building a long-term Research Data Storage Service for UNSW.



# Building a long-term Research Data Storage Service for UNSW.



Luc Betbeder

luc@upew.odu.au



speaking for the whole team...

From IT: Vishel Selsjal, Amery Nicoebeh, Chris Will, Sergey, Sochin, Seri Charcensri, Berhard Sentner, Dajan Martista, Jim Legapi, Denise Black., Ling vaccomes ton., Greg Sarayre and Inam), And the lary business state-falsers: Berhard Cinsielamiki, Prof. Mark Hoffman, Greg Leslie, Grainne Moron, Naude Frances and all the wonderful participants in the Advisory's Group and Platy sites. Transit you.



# Luc Betbeder

luc@unsw.edu.au



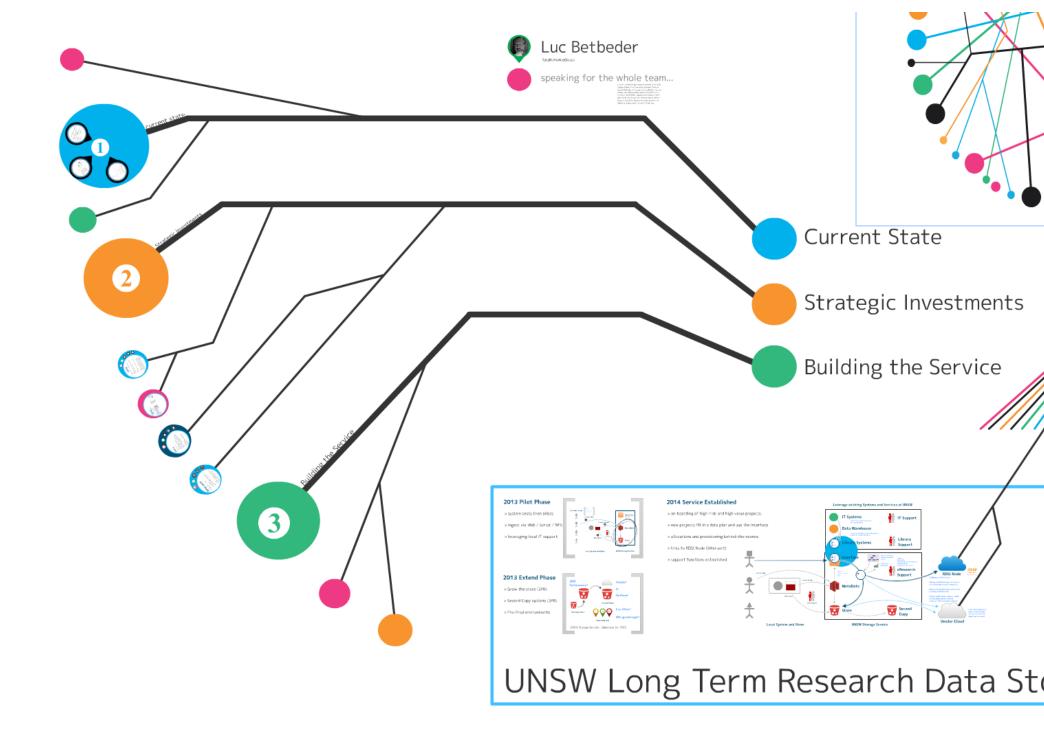
speaking for the whole team...

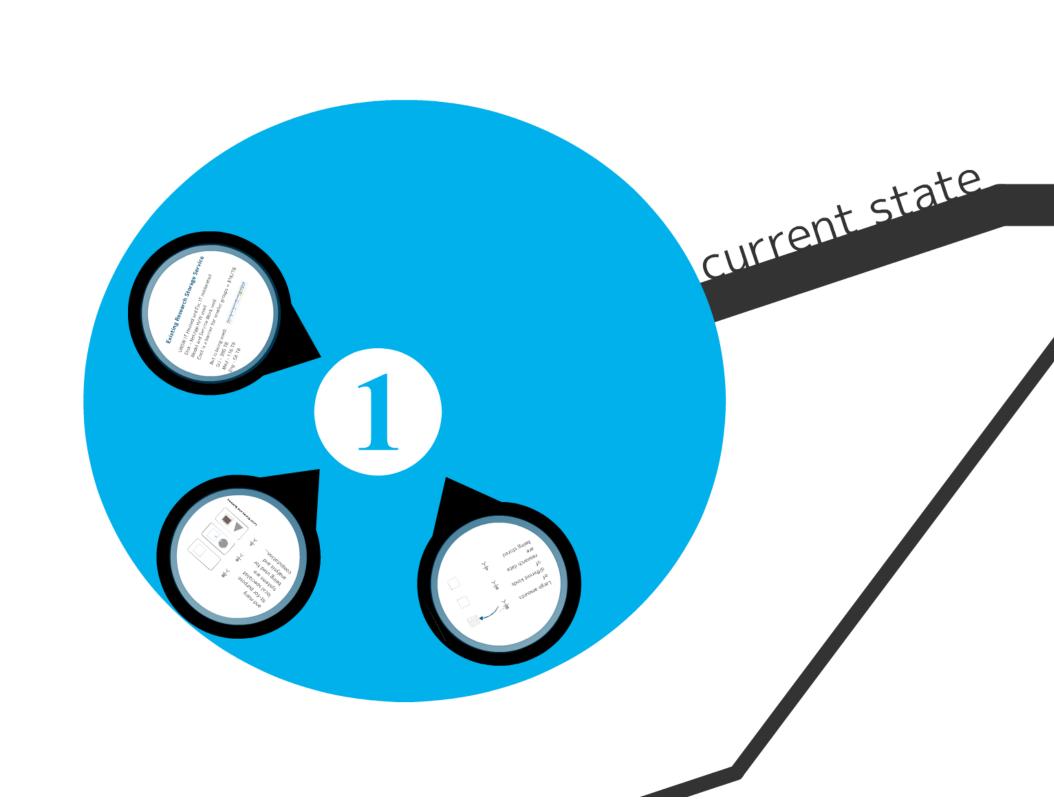
From IT: Vishal Sehgal, Amany Nuseibeh, Chris Will, Sergey Sashin, Seri Charoensri, Berhard Semtner, Dusan Munizaba, Jim Leeper, Denise Black... (and yes comms too... Greg Sawyer and team). And the key business stakeholders: Barbara Chmielewski, Prof Mark Hoffman, Greg Leslie, Grainne Moran, Maude Frances and all the wonderful participants in the Advisory Groups and Pilot sites. Thank you.

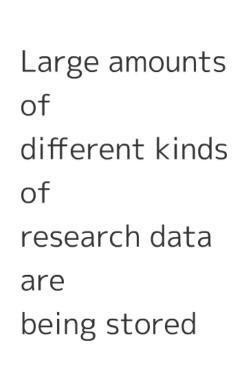
# whole team

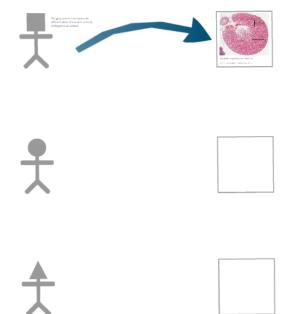
From IT: Vishal Sehgal, Amany Nuseibeh, Chris Will, Sergey Sashin, Seri Charoensri, Berhard Semtner, Dusan Munizaba, Jim Leeper, Denise Black... (and yes comms too... Greg Sawyer and team). And the key business stakeholders: Barbara Chmielewski, Prof Mark Hoffman, Greg Leslie, Grainne Moran, Maude Frances and all the wonderful participants in the Advisory Groups and Pilot sites. Thank you.

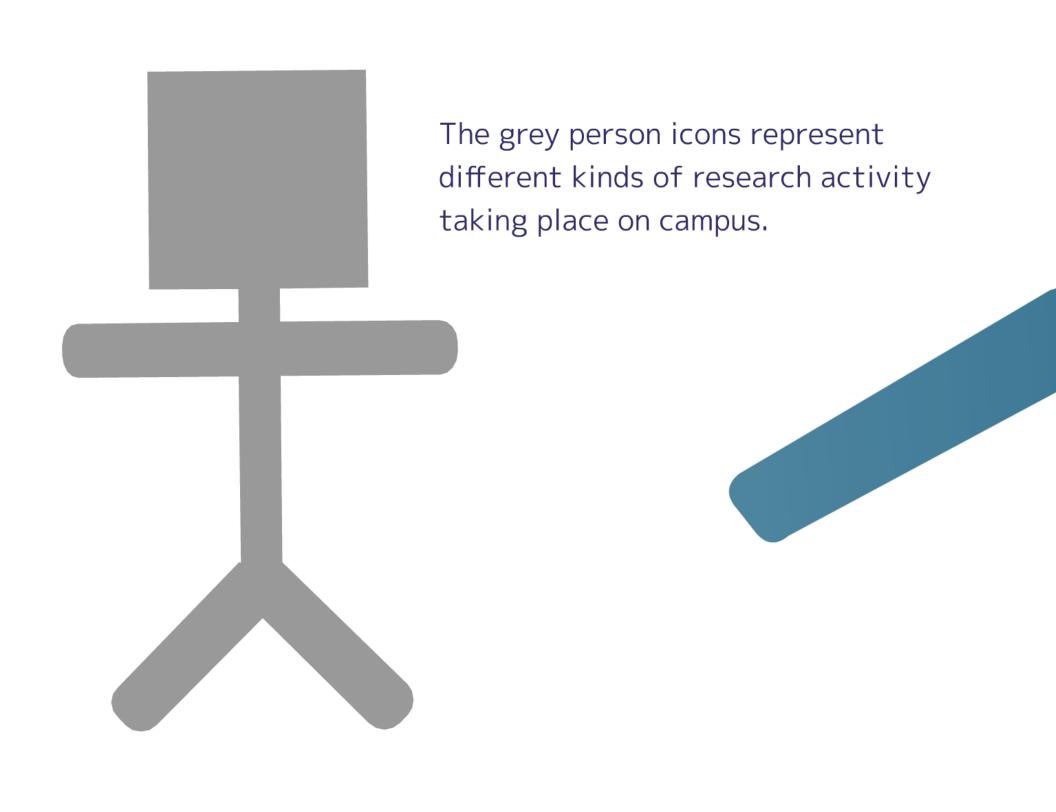




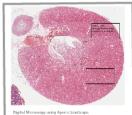












Creates super high resolution image files.

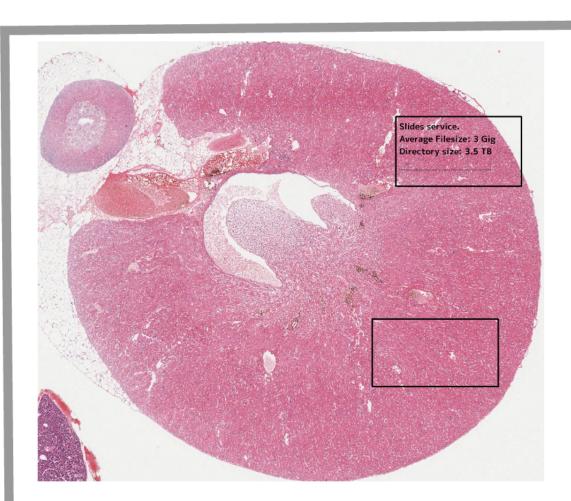












Digital Microscopy using Aperio ScanScope.

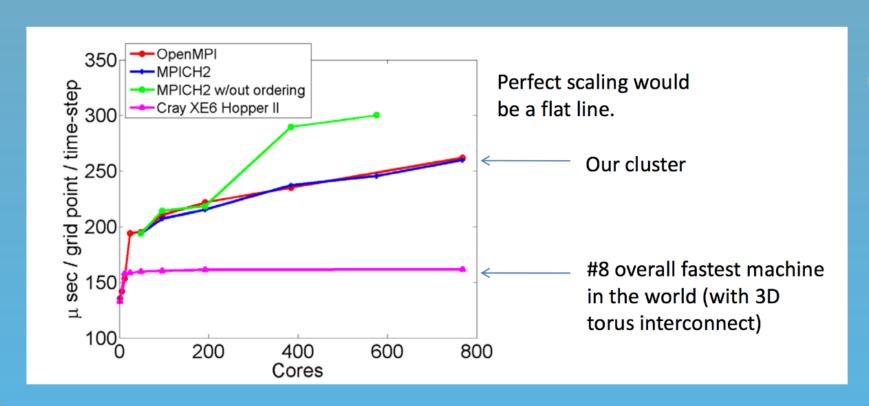
Creates super high resolution image files.

# Slides service. Average Filesize: 3 Gig Directory size: 3.5 TB

Forced to delete files and/or store them on external USB HDD.



# Computer cluster: Leonardi - 3000 nodes eg: Models of diesel spray and combustion in egines



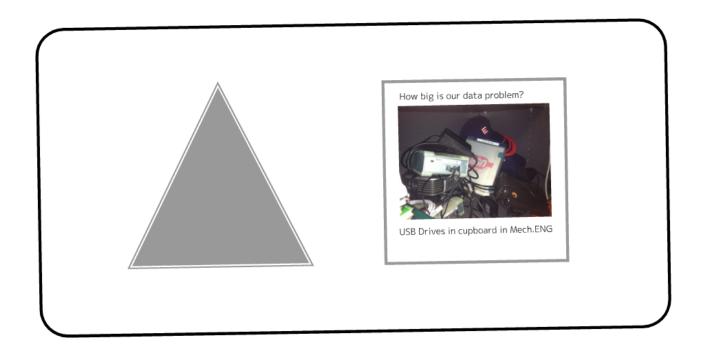
Computer cluster: Leonardi - 3000 nodes

Storage capacity: no-long term storage...

 $1 \times RAID5 (3 \times 2TB) + 1 HS (2TB) :$ /home (3TB) + /share/apps (750GB)

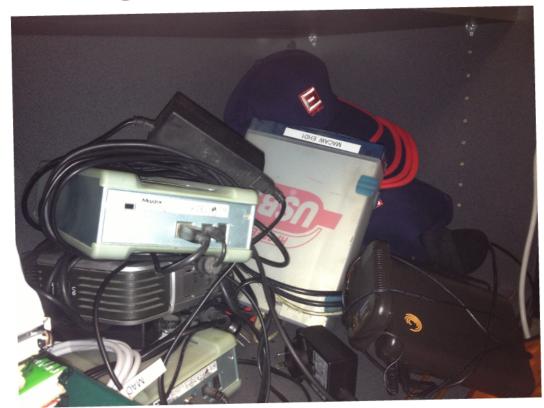
1 x RAID50 (5 x 5 x 2TB) + 1 HS (2TB) : /share/scratch (37TB)





### Local Stores and Systems

#### How big is our data problem?



USB Drives in cupboard in Mech.ENG



UNSW IT Hosted and Fac IT moderated
Disk - NetApp H/W used
Model and Service Work well
Cost is a barrier for smaller groups = \$1K/TB

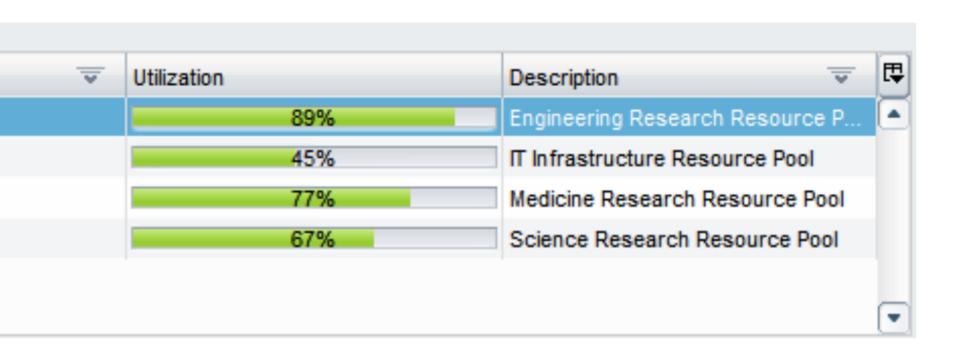
But is being used:

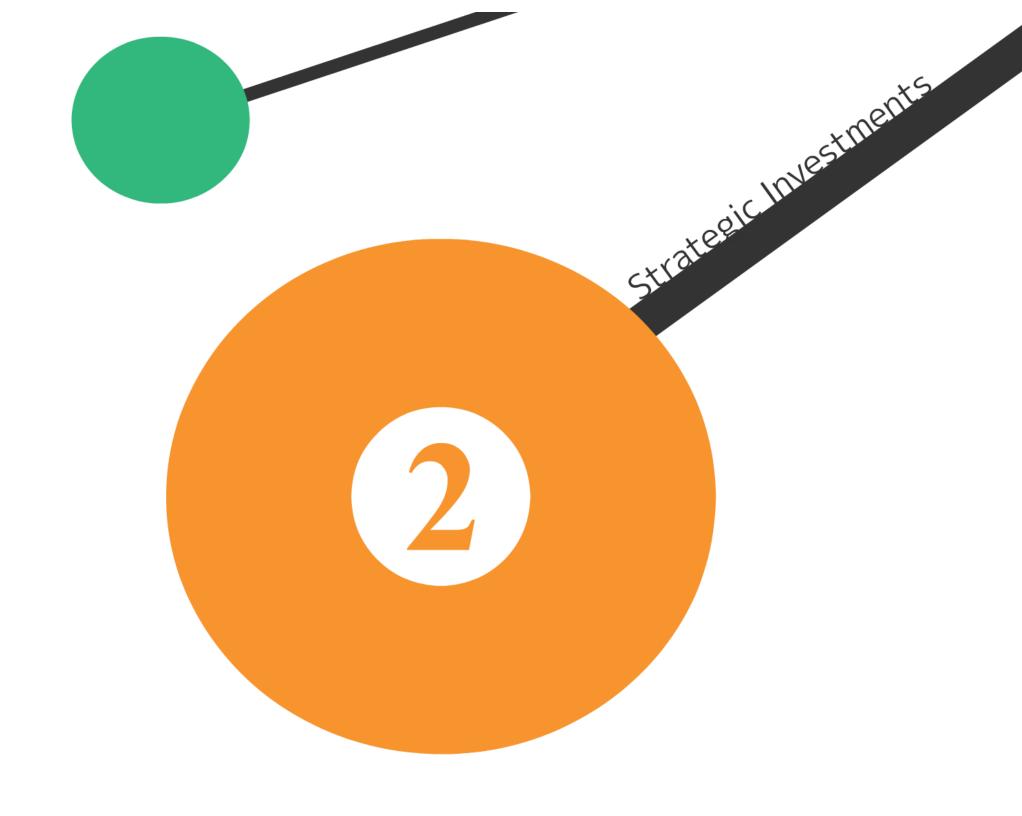


Sci - 380 TB

Med - 116 TB

Eng - 58 TB







#### 2012

Strategic Investment Planning (3 yr plan)

#### > process:

Align IT investment to UNSW "Business Domains" (ie: Research, Academic, Other)

#### > method:

In-domain prioritising > estimating > voting

#### > outcome:

Multi-year / multi-stream investment for research storage (To support "Research Practice", meet Policy obligations, reduce Risk and "Providing an excellent research environment, with cutting-edge facilities and equipment.")







#### **Long-Term Storage**

A long-term "accessible archive" with metadata capability to support research practice at UNSW.

- > No direct charge to researchers or project.
- > Principles: large, functional, cheap, extensible, supportable, aligned, secure.

NOT a store for computation or analysis.

#### Interface / Portal

A portal for accessing and using the UNSW long-term store and other data storage services (e.g. RDSI node(s), vendor/cloud).

Linked to other UNSW systems (library, research projects, data warehouse, authentication etc.)

It is policy-driven. You create a data plan to use the store.

#### **Devices and Bookings**

We would like to link our research device booking system to the portal.

And eventually to provide tools to enable direct ingestion of research device data, device metadata and booking system metadata into store.

**Smart Devices** 

#### **Enhanced Metadata**

Over time we want to provide the Store with additional metadata, integration, search and collaboration functionality.

This may include federated search (searching across stores) and improving the ingest metadata capabilities of the store.

**Smarter Data** 

#### **Active Storage**

A self-service capability for researchers and research projects to access "active storage" (for compute) via the portal.

Active storage costs (unlike longterm storage) would be charged back to the projects.

Automation/Orchestration.

## 2012...

> procurement process.







#### Evaluation process based on these design principles.

	Design principles	
1	Resilient	Data must not be lost or corrupted. Data must be accessible for long time frames. Store outages should be minimal and planned.
2	Affordable	Long-term large-scale storage is expensive. The overall project must deliver large amounts of storage.
3	Supportable	Initial investment should be supportable within UNSW capabilities
4	Functional	Key functions must be available via initial investment. UNSW authentication, Drive mapping, Speed
5	Scalable	The capacity of the system is planned to grow from 1-5 Petabytes over 3 years. Initial system needs to be able to scale.
6	Extensible	The store will be extended with metadata, search, and collaboration features in 2014. Initial system needs to be able to be extended.
7	Aligned	The store should be aligned with national initiatives such as RDSI.

Object-Based Storage Devices (OSD)
Hierarchical Storage Management (HSM)
Expanding current Storage System
Cloud

> Very strong vendor responses.

#### HSM... Disk and Tape via SGI



DISK:

SGI IS5500 storage array.

0.5 PB

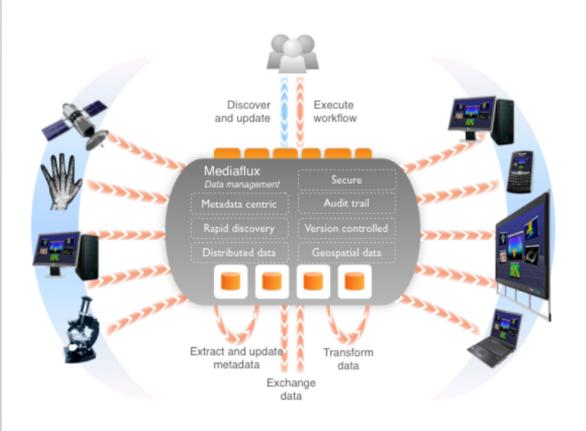


TAPE:

Oracle SL3000 with 700 Slots / 4 Drives

1.0 PB

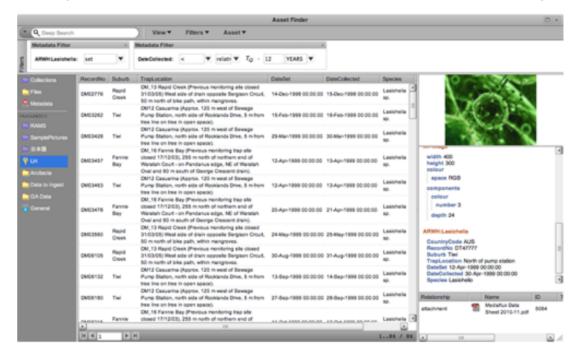
#### **SGI LiveArc** = Arcitecta MediaFlux



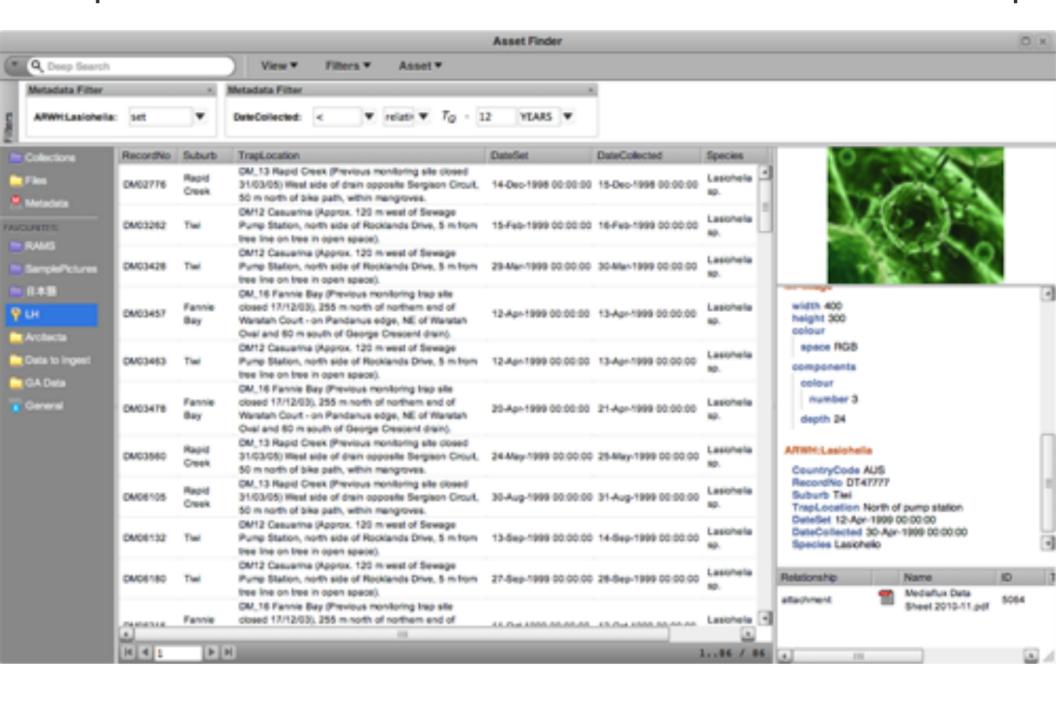
Ingest with MetaData
Search
Version
Script / API
Web Client (Java)



#### http://www.arcitecta.com/Products/Desktop



#### http://www.arcitecta.com/Products/Desktop





#### **Long-Term Storage**

- Installation and configuration
- Testing: July-Aug 💍
- Piloting: Sept-Dec
  - +Growing (to 3PB)
  - +Protecting (second copy)
  - +Staging (dev, UAT)

\$1M

#### **Interface / Portal**

- Analysis
- Scope✓
- Architecture and Solution Design



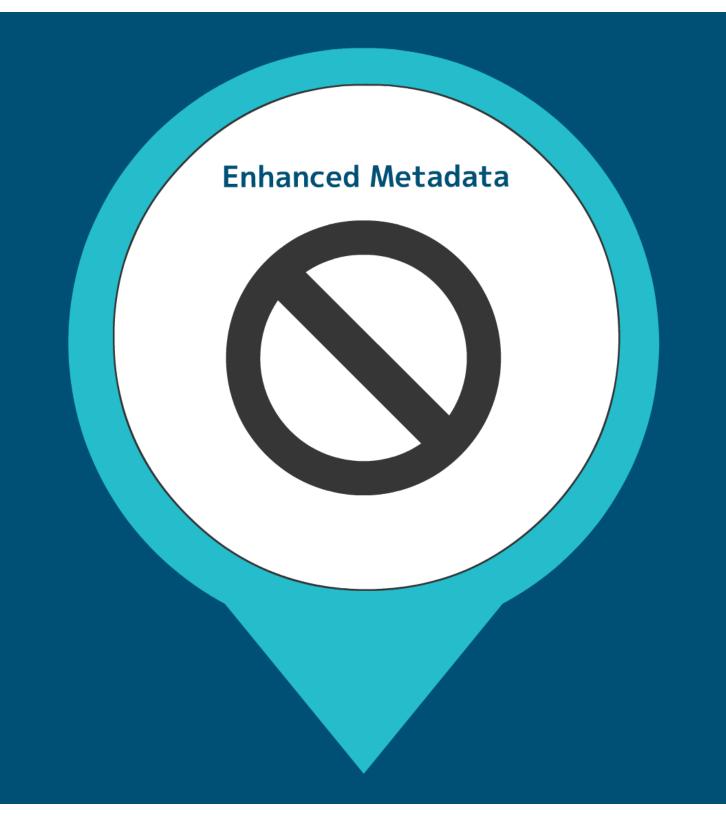
- Build: Aug-Dec
- Release: 2014

800K

#### **Devices (Bookings)**

- Analysis
- Scope
- Architecture and Solution Design

150K









#### 2014-2015



Run the service. Expand capacity / capability.

> space allocation Mix of existing targeted high risk and high value projects (100) and all new projects (1000).

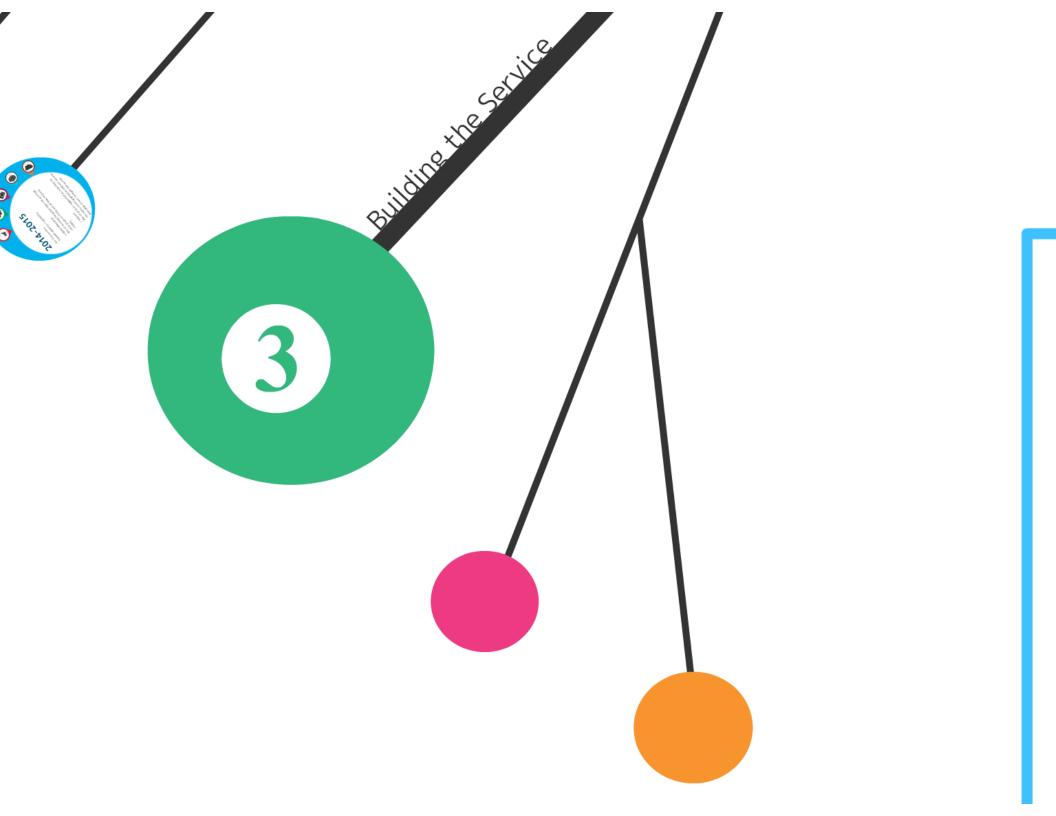
> support

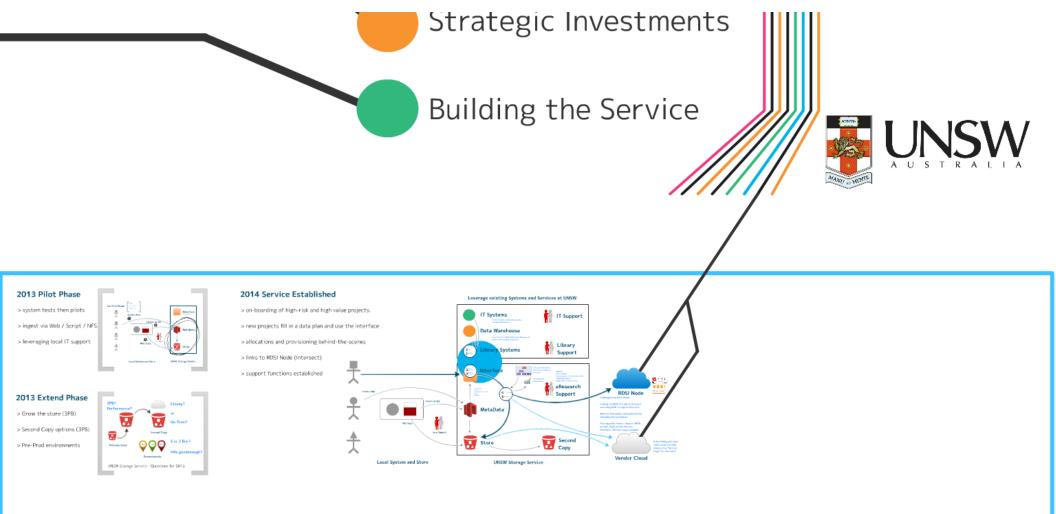
Mix of project-supported on-boarding for the messy existing projects and self-service data-plan-driven-through-the-portal.







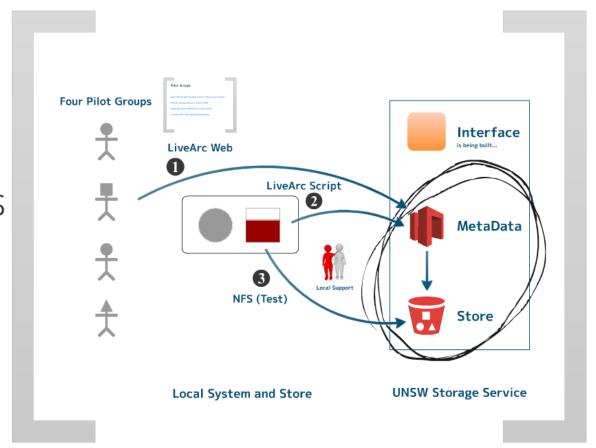




UNSW Long Term Research Data Storage Service

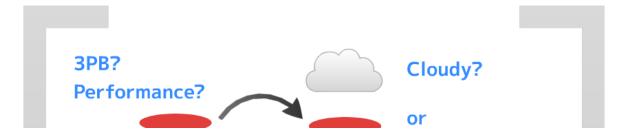
#### 2013 Pilot Phase

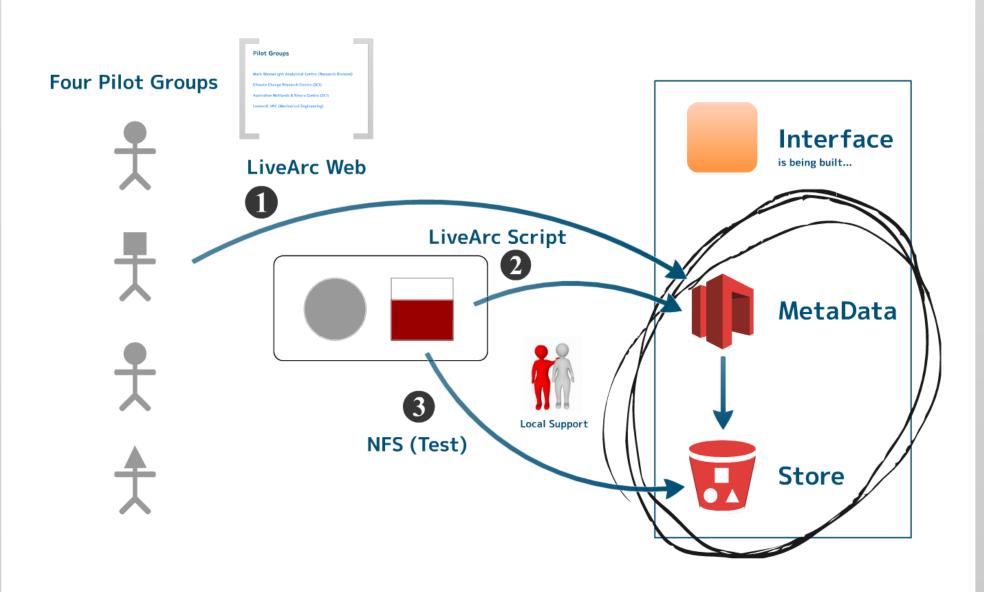
- > system tests then pilots
- > ingest via Web / Script / NFS
- > leveraging local IT support



#### 2013 Extend Phase

> Grow the store (3PB)





**Local System and Store** 

**UNSW Storage Service** 

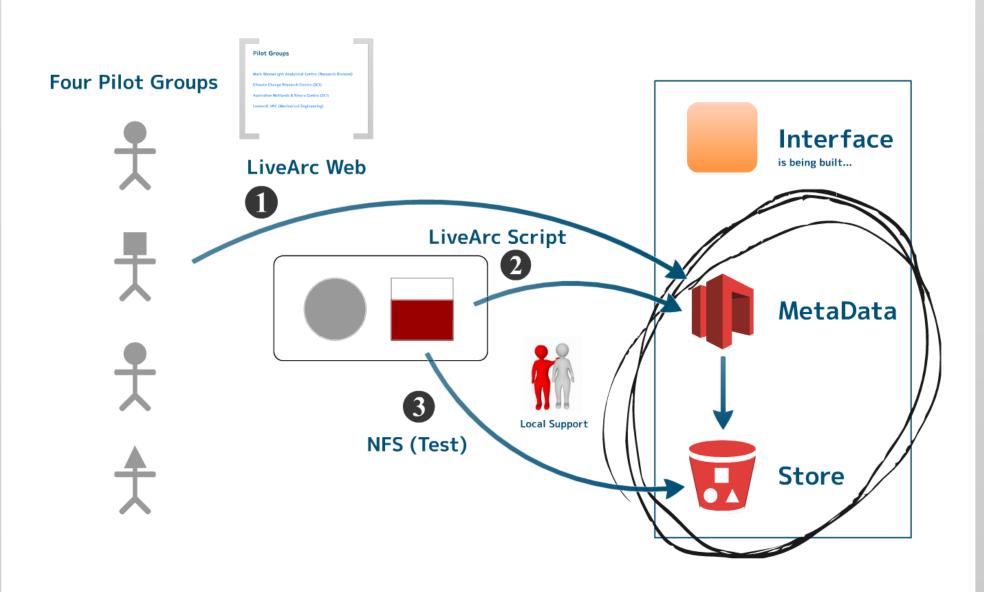
#### **Pilot Groups**

Mark Wainwright Analytical Centre (Research Division)

**Climate Change Research Centre (SCI)** 

**Australian Wetlands & Rivers Centre (SCI)** 

Leonardi HPC (Mechanical Engineering)

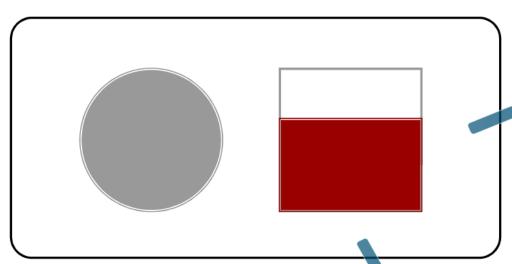


**Local System and Store** 

**UNSW Storage Service** 

1

#### **LiveArc Script**



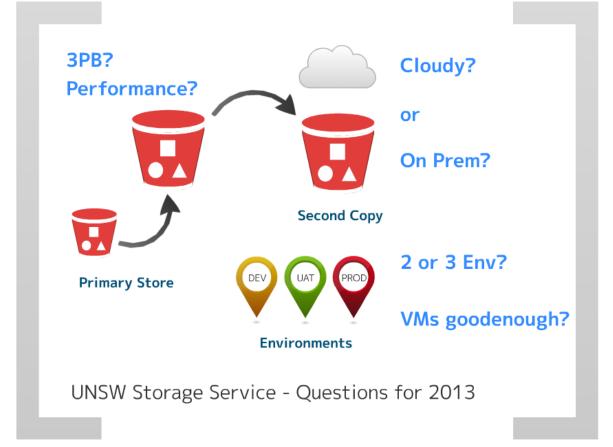
2

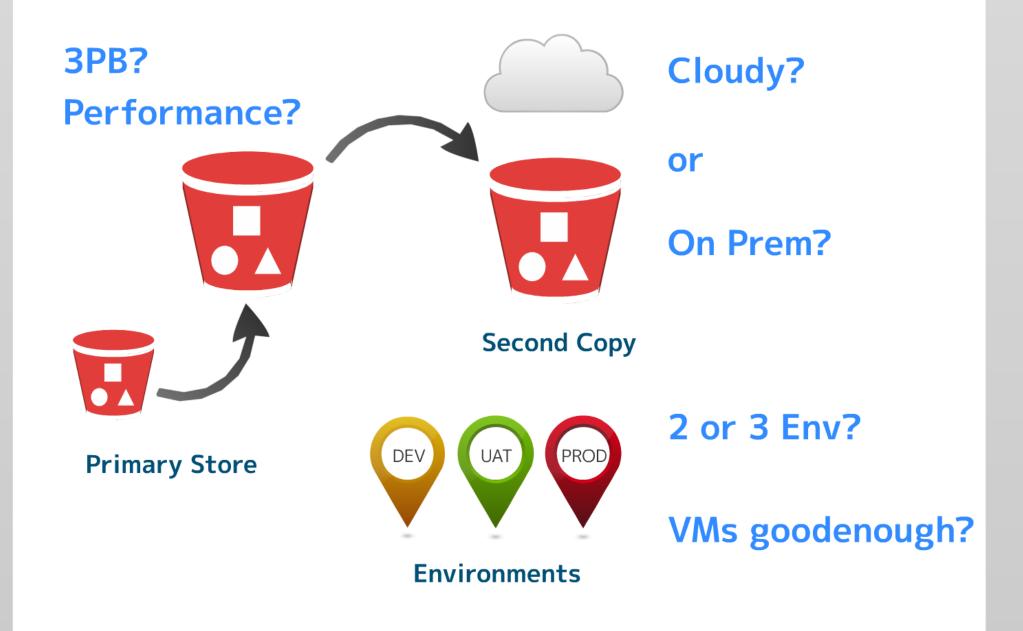
3 NFS (Test)



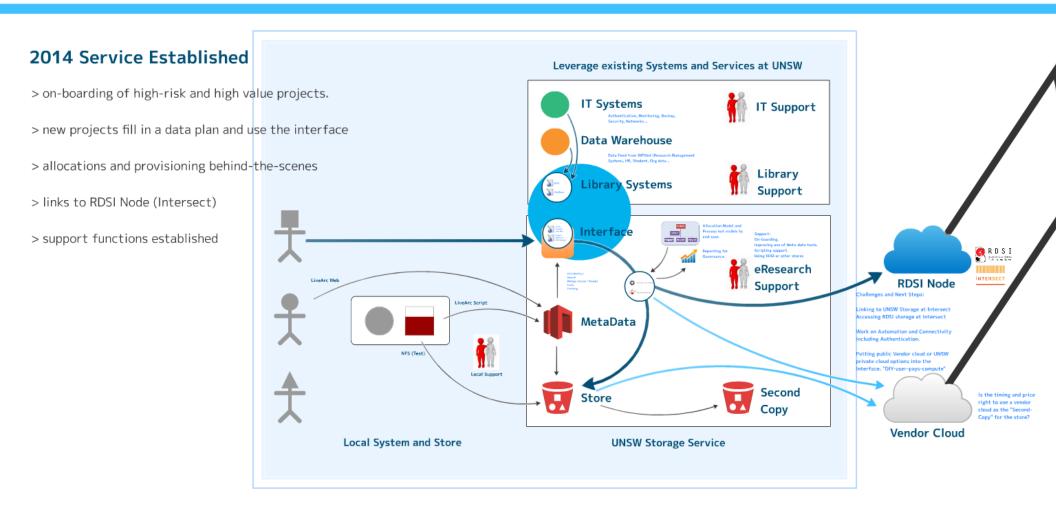
#### 2013 Extend Phase

- > Grow the store (3PB)
- > Second Copy options (3PB)
- > Pre-Prod environments





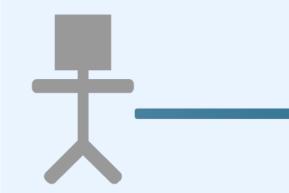
UNSW Storage Service - Questions for 2013



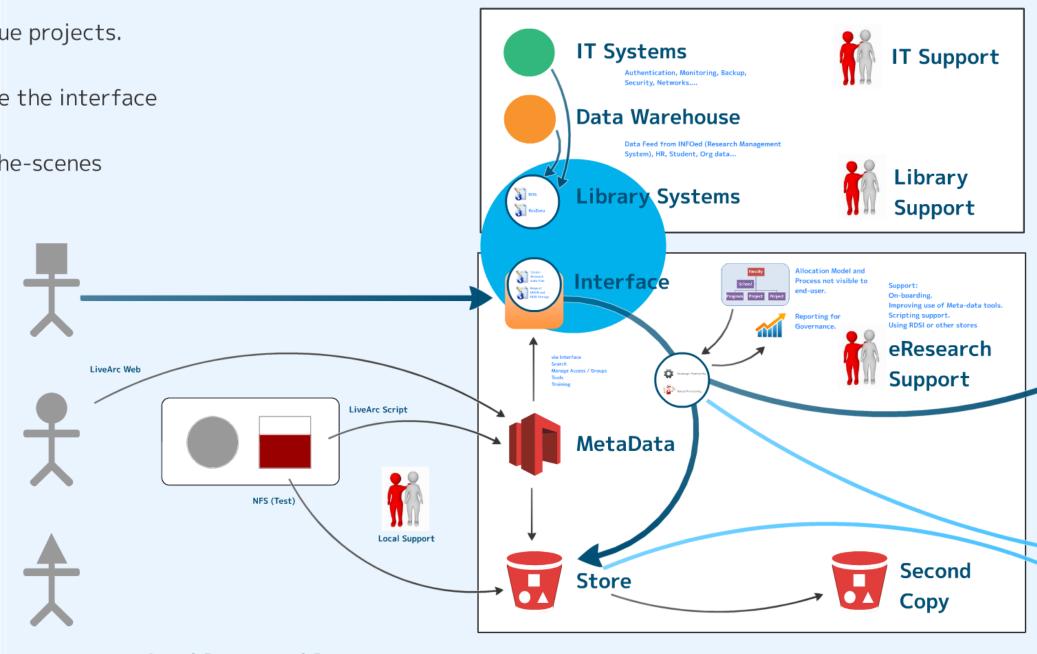
# ig Term Research Data

#### 2014 Service Established

- > on-boarding of high-risk and high value projects.
- > new projects fill in a data plan and use the interface
- > allocations and provisioning behind-the-scenes
- > links to RDSI Node (Intersect)
- > support functions established



#### Leverage existing Systems and Services at UNSW



**Local System and Store** 

**UNSW Storage Service** 



#### **Data Warehouse**

Data Feed from INFOed (Research Management System), HR, Student, Org data...



### **Library** Systems





#### **Interface**

Faculty
School
Program Project Project

Allocation Model and Process not visible to end-user.



Reporting for Governance.







#### Leverage existing Systems and Services at UNSW



Authentication, Monitoring, Backup, Security, Networks....

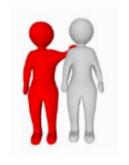


IT Suppo

#### **Data Warehouse**

Data Feed from INFOed (Research Management System), HR, Student, Org data...

**Library** Systems

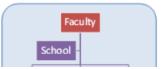


Library Support



ResData

**Interface** 



Allocation Model and Process not visible to

Support:



Create Research Data Plan



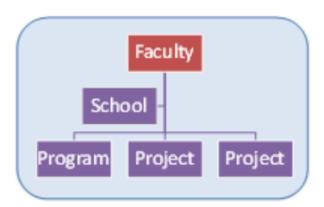
Request UNSW and RDSI Storage

## Interface

via Interface
Search
Manage Access / Groups
Tools
Training



## face



Allocation Model and Process not visible to end-user.



Reporting for Governance.





Allocation Model and Process not visible to end-user.

Reporting for Governance.



On-boarding.

Improving use of Meta-data tools.

Scripting support.

**Using RDSI** or other stores

# eResearch Support







**Challenges and Next Steps:** 

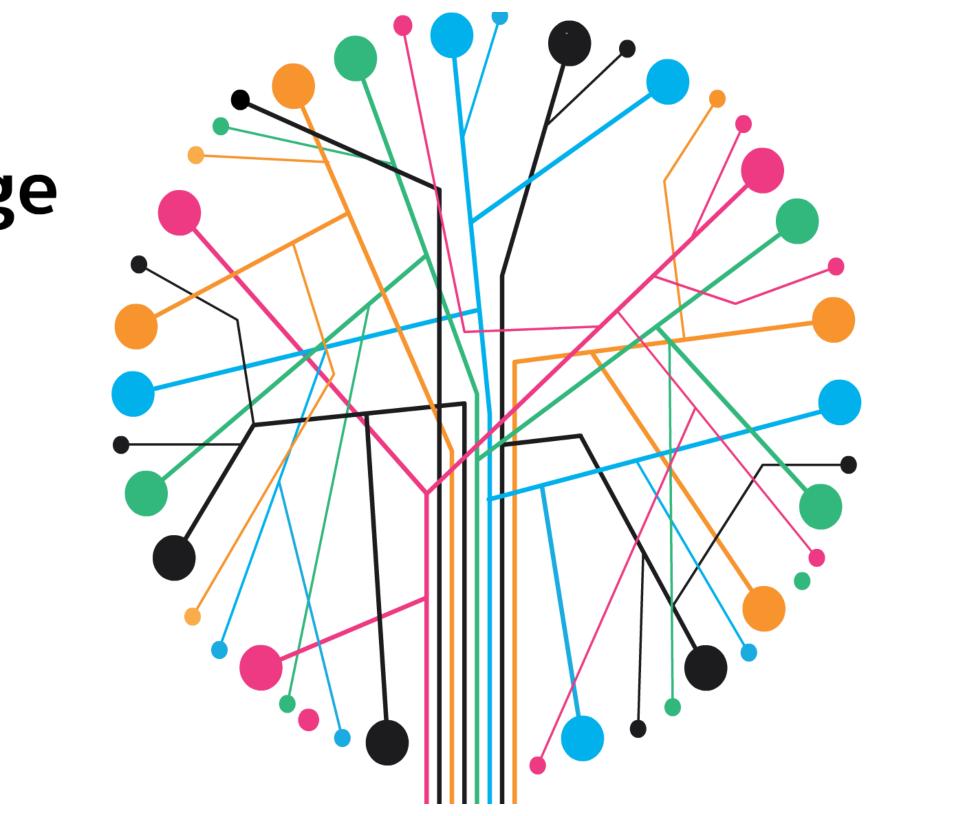
Linking to UNSW Storage at Intersect Accessing RDSI storage at Intersect

Work on Automation and Connectivity including Authentication.

Putting public Vendor cloud or UNSW private cloud options into the Interface. "DIY-user-pays-compute"

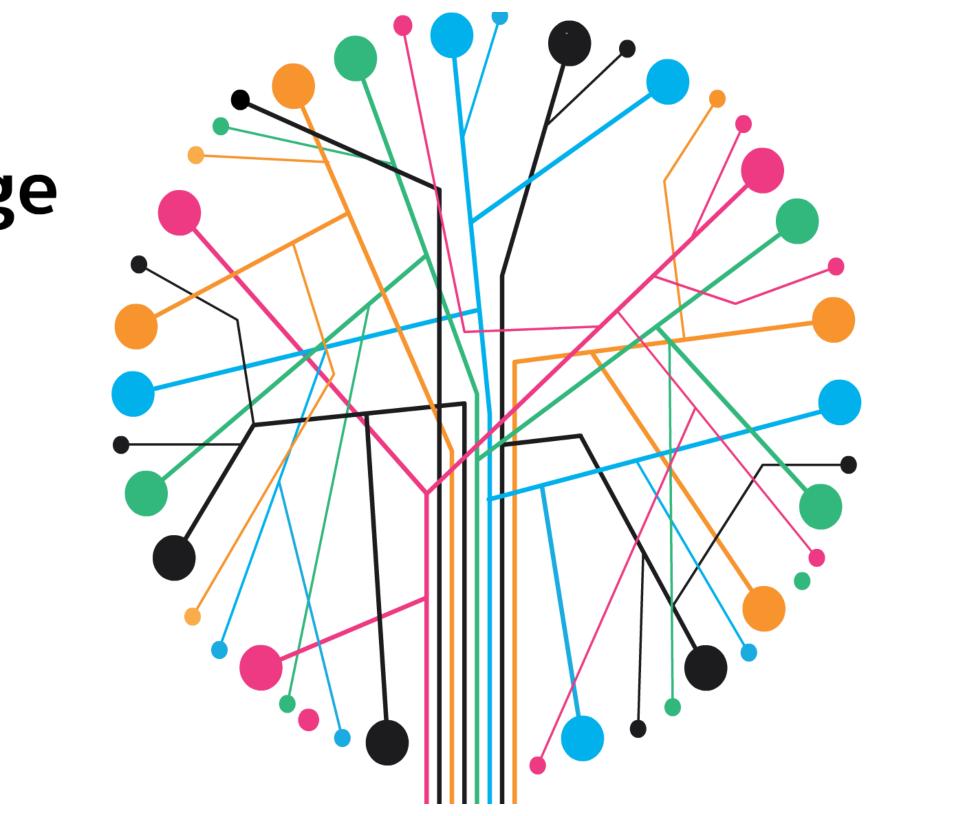
Is the timing and price right to use a vendor cloud as the "Second-Copy" for the store?

**Vendor Cloud** 





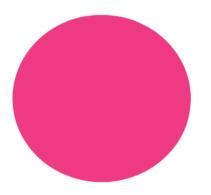






# Luc Betbeder

luc@unsw.edu.au



## speaking for the whole

From IT: Vishal Sehgal, Aman Sergey Sashin, Seri Charoens Dusan Munizaba, Jim Leeper, comms too... Greg Sawyer an business stakeholders: Barba Mark Hoffman, Greg Leslie, G Frances and all the wonderfu Advisory Groups and Pilot sin