

Standard (Alternative) data management tools

Daniel Traynor, QMUL
GridP 40, 10/04/2018

iRODS

- integrated Rule-Oriented Data System (iRODS). Open source project traces history back to 1995.
- Storage virtualisation of different disk and tape storage systems (sits on top of other storage systems e.g. posix, lustre, dcache). HSM capable. Load balancing. Concurrent access.
- A logical namespace across storage locations. Federation of sites.
- Controlled access through a number of different authentication mechanisms: Kerberos, secure username/password and Grid Security Infrastructure (GSI).
- A rule engine to automate data management and access according to defined policies.
- Client interfaces: web(idrop), cli, GUI (cyberduck!), gridFTP, webdav.

Policy examples

- Automatically replicate a file added to a collection into 3 geographically distributed sites.
- Periodically check integrity of files in a collection and repair/replace if needed/possible.
- Automatically pick a certain storage location based on user or collection or size or type.

iRODS - T2K

- Make data quality files available to world wide collaborators to be checked as soon as produced.
- This requires a data distribution system. iRODS rules used for access, replication and backup.
- Access the files is through a web or cli interface, with a read-only account

First iRODS Experience in a Neutrino Experiment

F. Di Lodovico^a, A. Hasan^b, Y. Iida^c, T. Sasaki^c

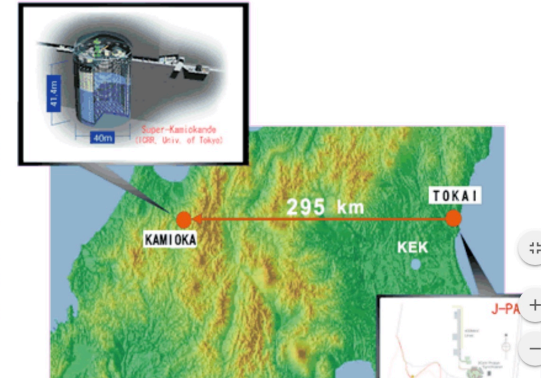
^aQueen Mary University of London, London, UK

^bUniversity of Liverpool, Liverpool, UK

^cHigh Energy Accelerator Research Organisation (KEK), Tsukuba, Japan

Abstract

The T2K experiment is a long-baseline neutrino experiment with the primary aim of observing muon into electron neutrino appearance. An important task is the determination of the quality of the collected data. In this paper we describe the first experience gained in managing the T2K data quality data for the near detector such that it is stored and accessible to collaborators world-wide in a timely fashion. We also describe a set of rules developed for the project that have much more general applicability. This is the first use of iRODS in a neutrino experiment.



<https://irods.org/uploads/2011/DiLodovico-Neutrino-T2K-paper.pdf>

iRODS - UCL

- The push towards making data openly, If data is going to be discoverable and effectively re-used it is essential that adequate metadata is generated that provides a full and accurate description of where the data came from and how it should be understood. As data repositories grow in size and complexity, a solution is needed which will connect everything together and enable researchers to search across multiple types of repository to identify and access relevant data.

The screenshot shows a webpage from the UCL Information Services Division. The header includes the UCL logo and navigation links: Home, Our services, How to guides, About ISD, Help & Support, and News. The breadcrumb trail is UCL Home / ISD / News / Managing research data with iRODS. The main article is titled 'Managing research data with iRODS' and is dated 11 January 2016. The text states: 'A new component, iRODS, has recently been added to the Research Data Services portfolio, which will make it easier for researchers to manage their data and prepare it for open access. We asked the Research Data Services team to tell us about iRODS and what it can do for UCL researchers.' Below the article is a section titled 'What is iRODS?' with the text: 'The Integrated Rule-Oriented Data System (iRODS) is the result of over 20 years' worth of effort to tackle many of the problems around research data management. Supported and maintained by the iRODS Consortium, of which UCL is a member. iRODS is open-source software that provides a comprehensive set of tools to support data'. On the right side, there is an 'RSS feeds' section with a list of feeds: Main ISD service news, Email & Calendar news, Wireless & Wired Networks news, Research IT Service news, Software & Hardware news, and Websites, Apps & Databases news. Below that is a 'Tweets by @uclisd' section showing a tweet from UCL ISD (@uclisd) stating: 'ISD will be closed from 5.30pm today (28/3/18) to 08:30am next Thursday (5/4/18) for the Easter break. This includes the ISD Service Desk email'.

UCL member of
iRODS consortium

iRODS - others

- eudat (european open science cloud): used in bsafe storage product. Provide federated access to data for collaborations.
- sanger: use iRODS to manage petabytes of genomics data with diverse scientific needs, rules for backup to slough,
- emedlab: possible use for allowing secure data access from VMs to local GPFS storage.
- Panasas(ActivStore array), HGST(Active Archive System),

CDMI

Cloud Data Management Interface

Industry standard implemented in several products.

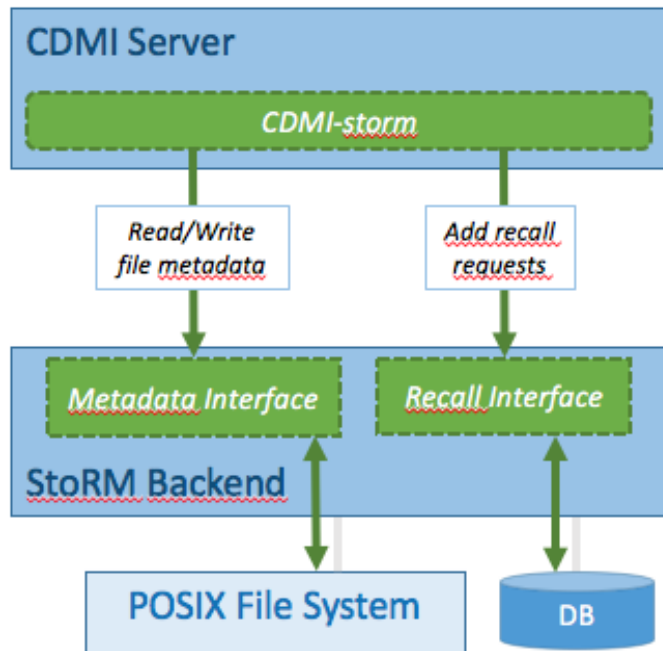
<https://www.snia.org/cdmi>

Cloud Data Management Interface

The **Cloud Data Management Interface (CDMI)** is a **SNIA** standard that specifies a protocol for self-provisioning, administering and accessing **cloud storage**.^[1]

CDMI defines **RESTful HTTP** operations for assessing the capabilities of the cloud storage system, allocating and accessing containers and objects, managing users and groups, implementing access control, attaching metadata, making arbitrary queries, using persistent queues, specifying retention intervals and holds for compliance purposes, using a logging facility, billing, moving data between cloud systems, and exporting data via other protocols such as **iSCSI** and **NFS**.

Transport security is obtained via **TLS**.



Functionality similar to Amazon's S3 cloud storage interface.

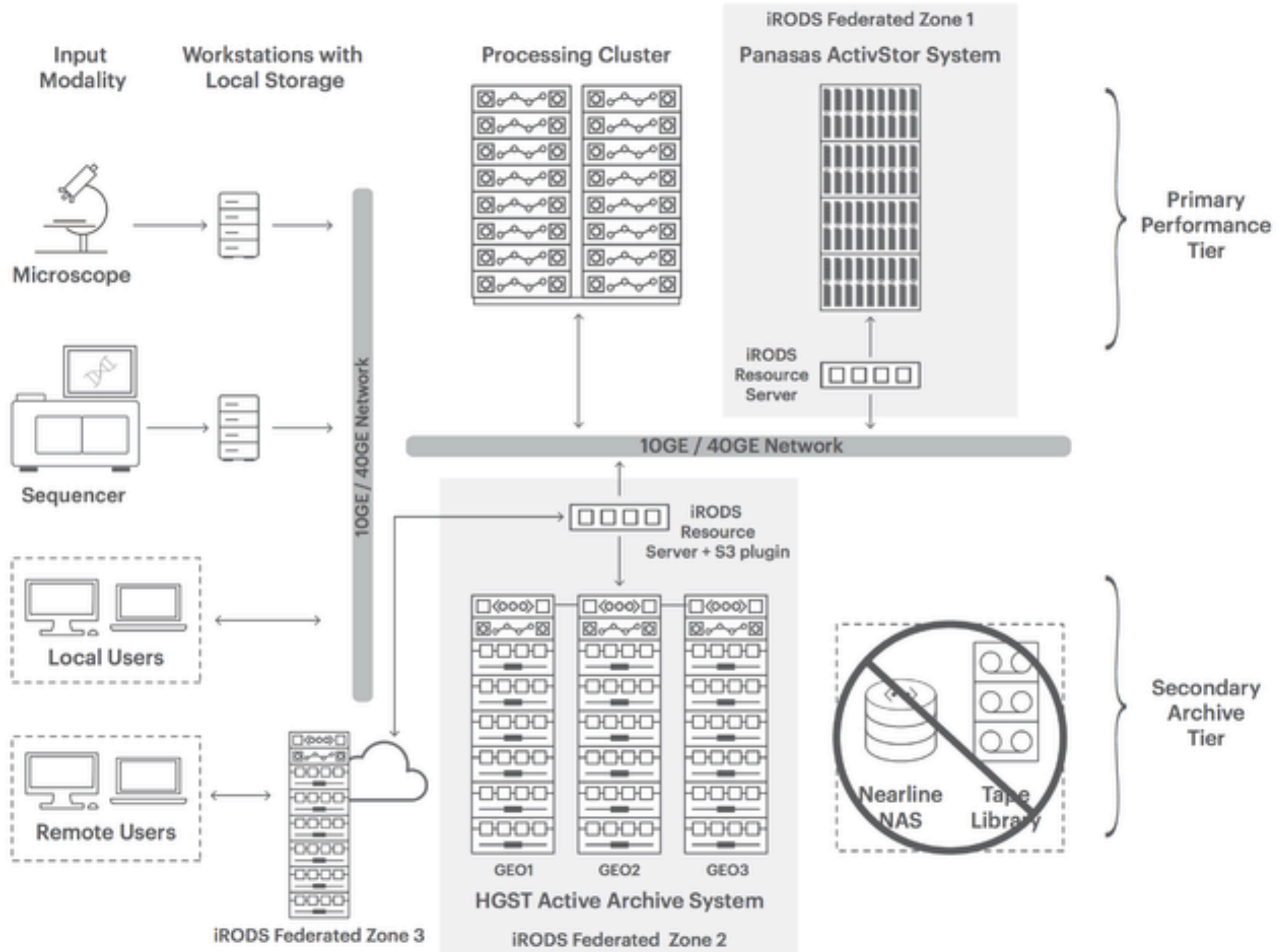
dcache and storm plugins

CDMI StoRM implements the Java Service Provider Interface **cdmi-spi** to provide a plugin for the **INDIGO DataCloud CDMI server** in order to support StoRM as storage back-end and allow users to negotiate the Quality of Service of stored data.

Comments

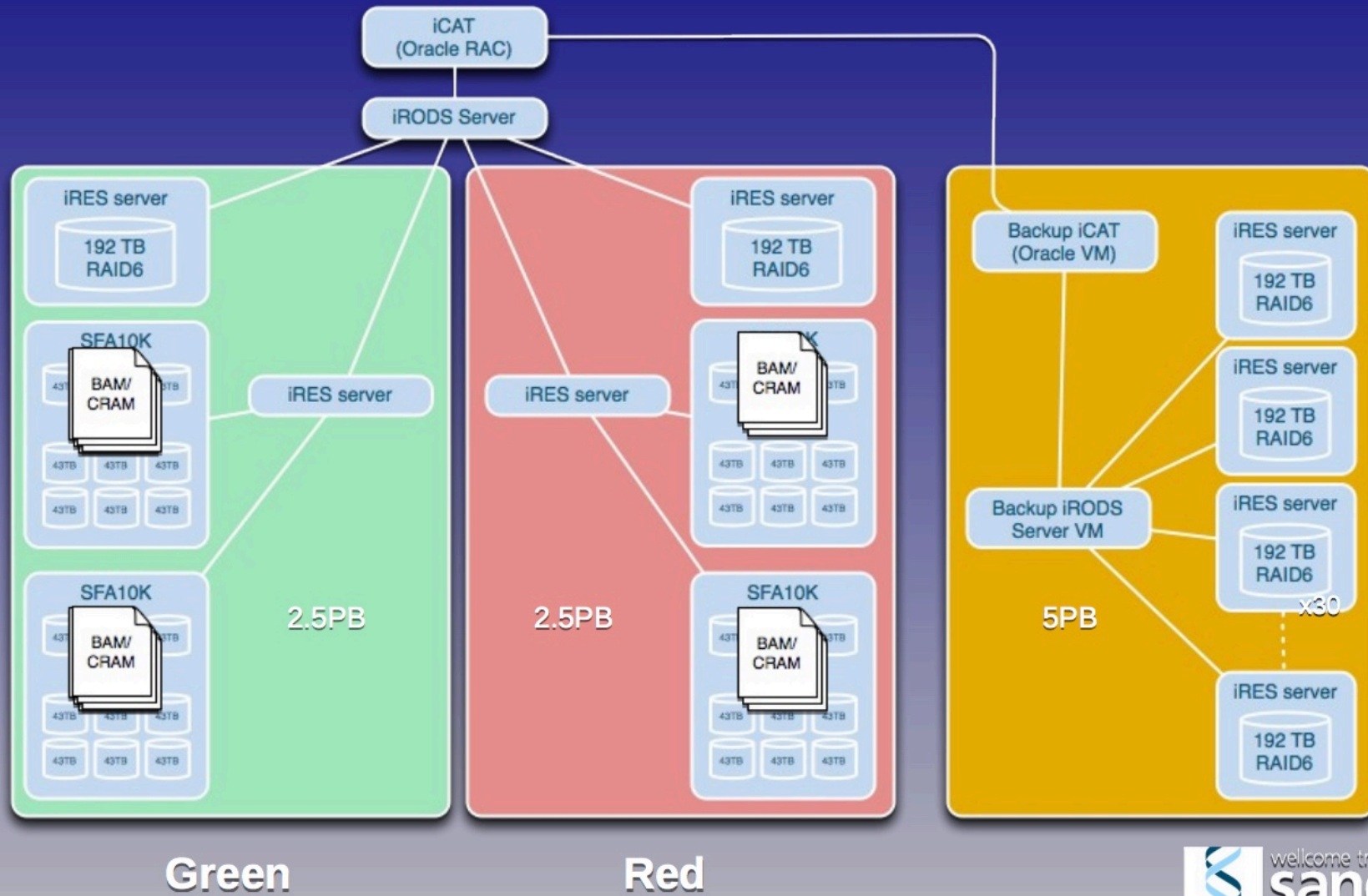
- iRODS mature product under active development and used widely for data management. Worth investigating further. Lustre and dcache plugins.
- iRODs may become a “must” to support a wider science base.
- CDMI pugins for dcache and StoRM make cloud and grid more interoperable. Support and usage not so clear.

backup

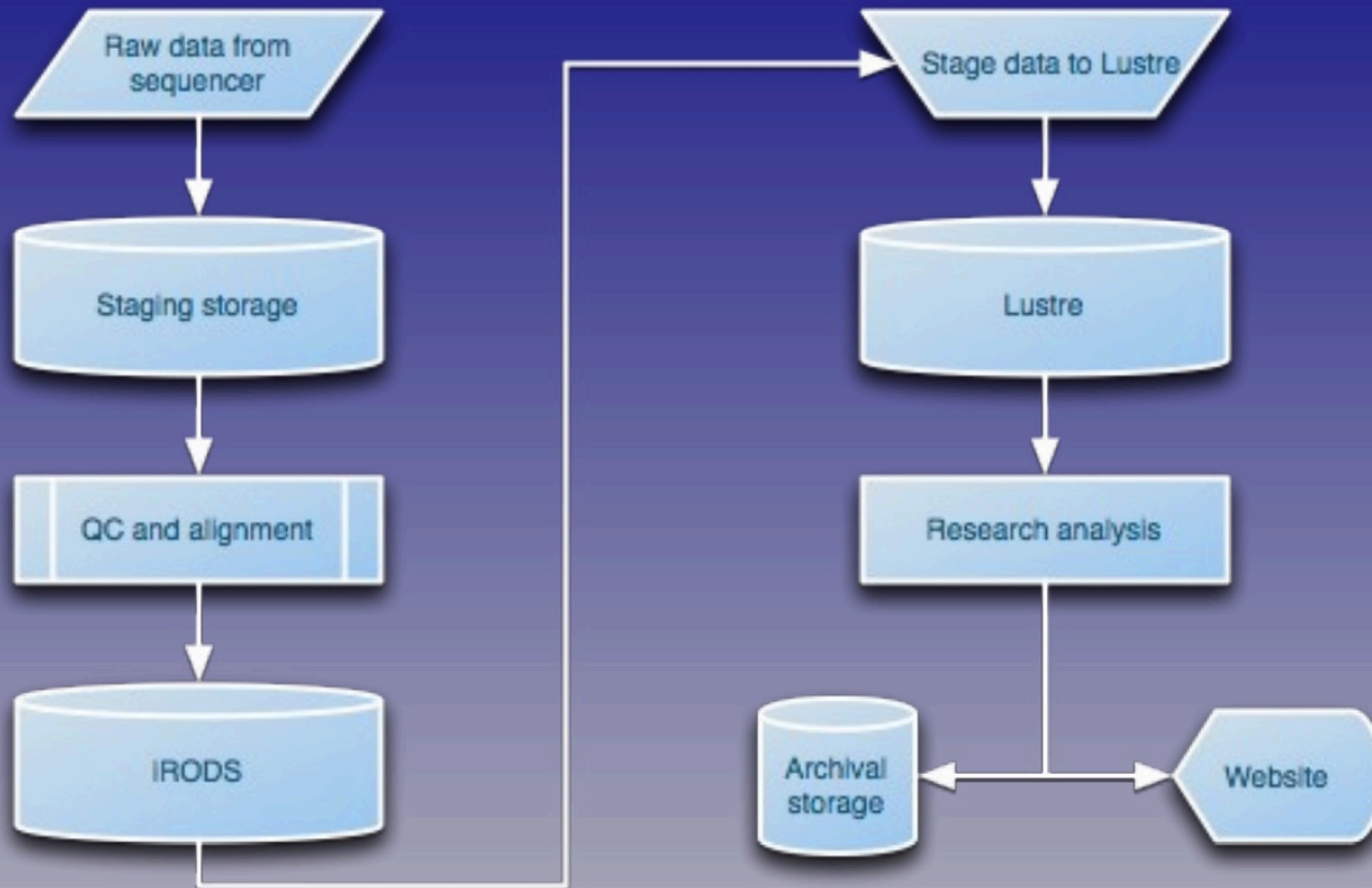


https://www.theregister.co.uk/2016/10/07/building_irods_for_saving_life_sciences_back/

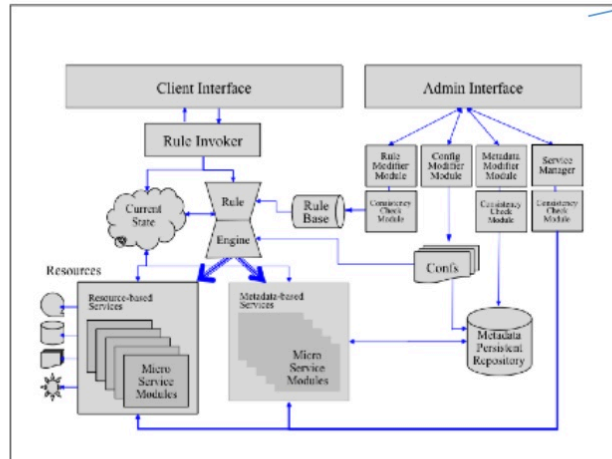
Capacity Expansion and Offsite Migration



Sanger NGS Data Flow



iRods – Running inside OpenStack VMs



Provides a proper research data management layer with an API-based user interface

Make data on GPFS available to VMs