

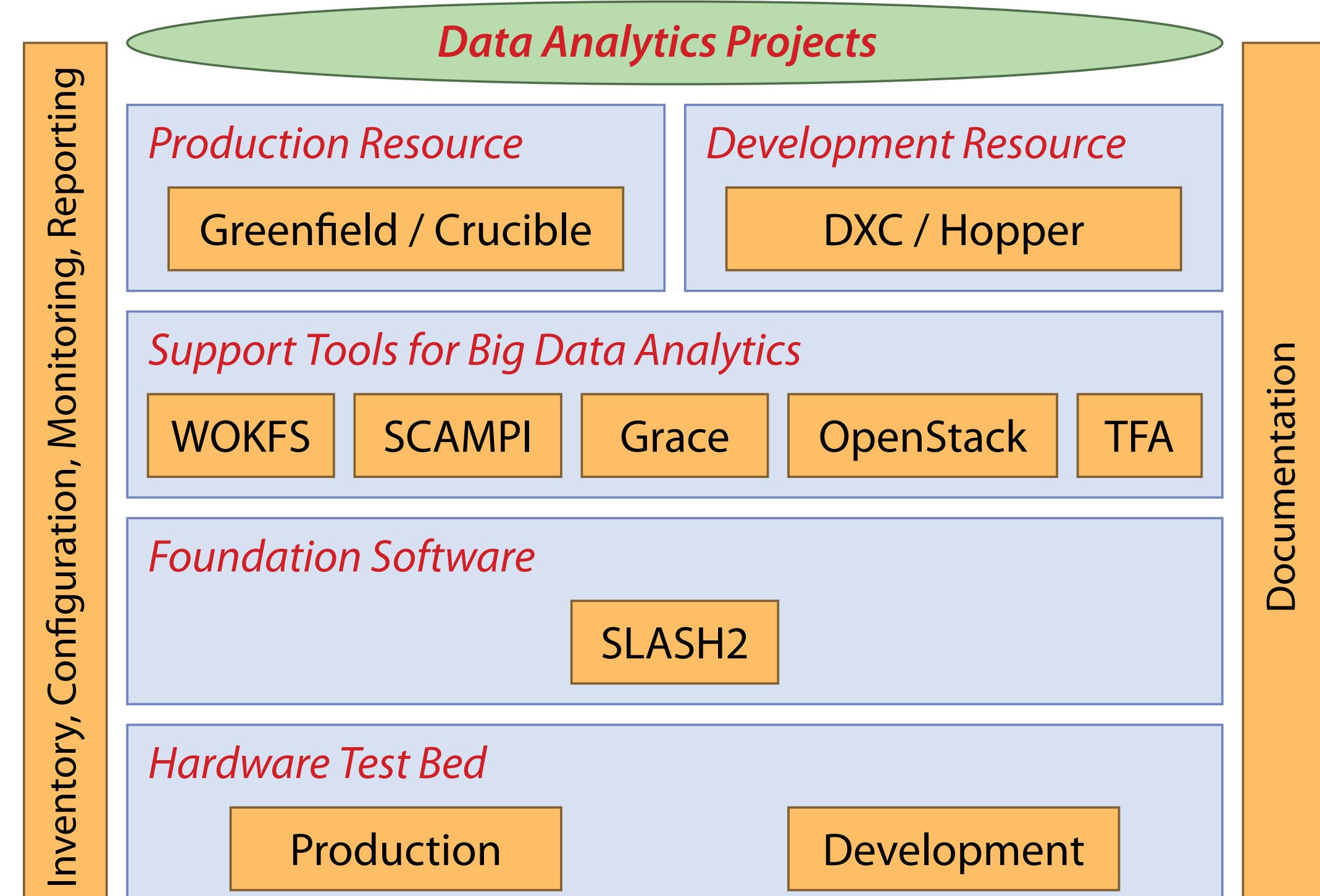
DATA EXACELL



Goals

- » Carry out an accelerated development pilot project to create, deploy, and test software building blocks and hardware implementing functionalities specifically designed to support data-analytic capabilities for data-intensive scientific research.
- » Implement and bring to production quality additional functionalities for distributed, disk-based data management.
- » Work with partners in diverse fields of science to provide scientific and technological drivers and system validation.

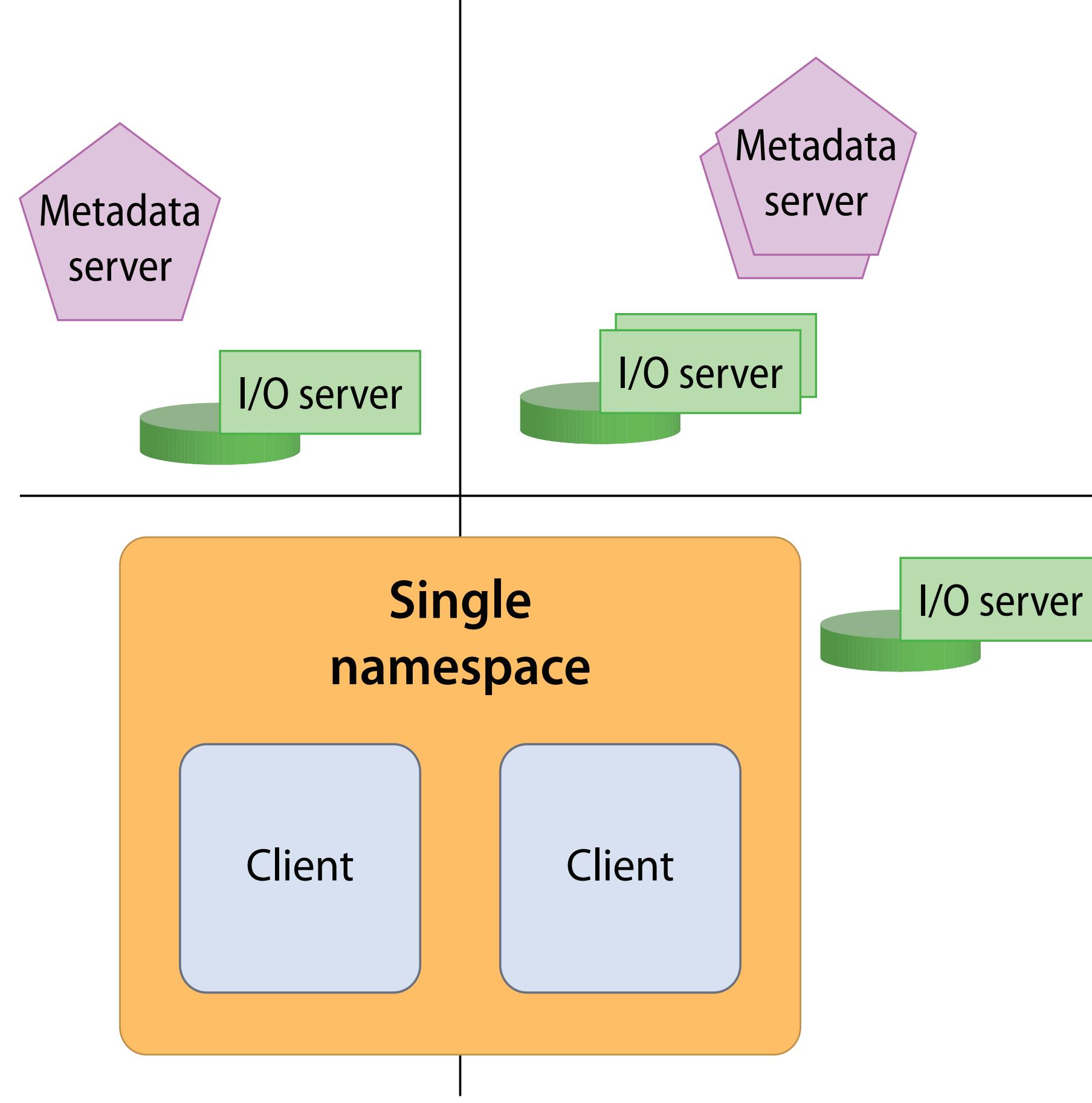
Data Infrastructure Building Blocks



- » **WOKFS** allows stackable FUSE layers without the penalty of having multiple FUSE mounts and repeated kernel entry & exit.
- » **SCAMPI** is a FUSE layer to restrict read/write/delete POSIX file system operations at a UID/GID level.
- » **GRACE** is a service platform allowing users with disabled accounts to perform file transfers and deletions.
- » **OpenStack** development includes procedures for deploying VMs to users and bare-metal booting, together with outreach.
- » **Two-Factor Authentication (TFA)** is implemented to enable enforcement for specific users with federated IDs.
- » **SLASH2** is a PSC-developed wide area file system.

SLASH2 Wide-Area File System

SLASH2 enhances distributed data access with optimal performance through replication and provides a single namespace across administrative domains.



- » **Designed for wide-area operation:** no single point of failure | high recoverability: components start up in any order, can leave or return, and be added or removed | extensive end-to-end error checking | multiple metadata servers (MDS) | replication
- » **User space implementation:** file system in user space (FUSE) client | not tied to any specific UNIX kernel version or OS distribution | fast metadata import from local file system
- » **System administration:** configuration management; tools to communicate with services | networking | cross-administrative domain | scale | growth flux

Hardware Resources for Data Analytics

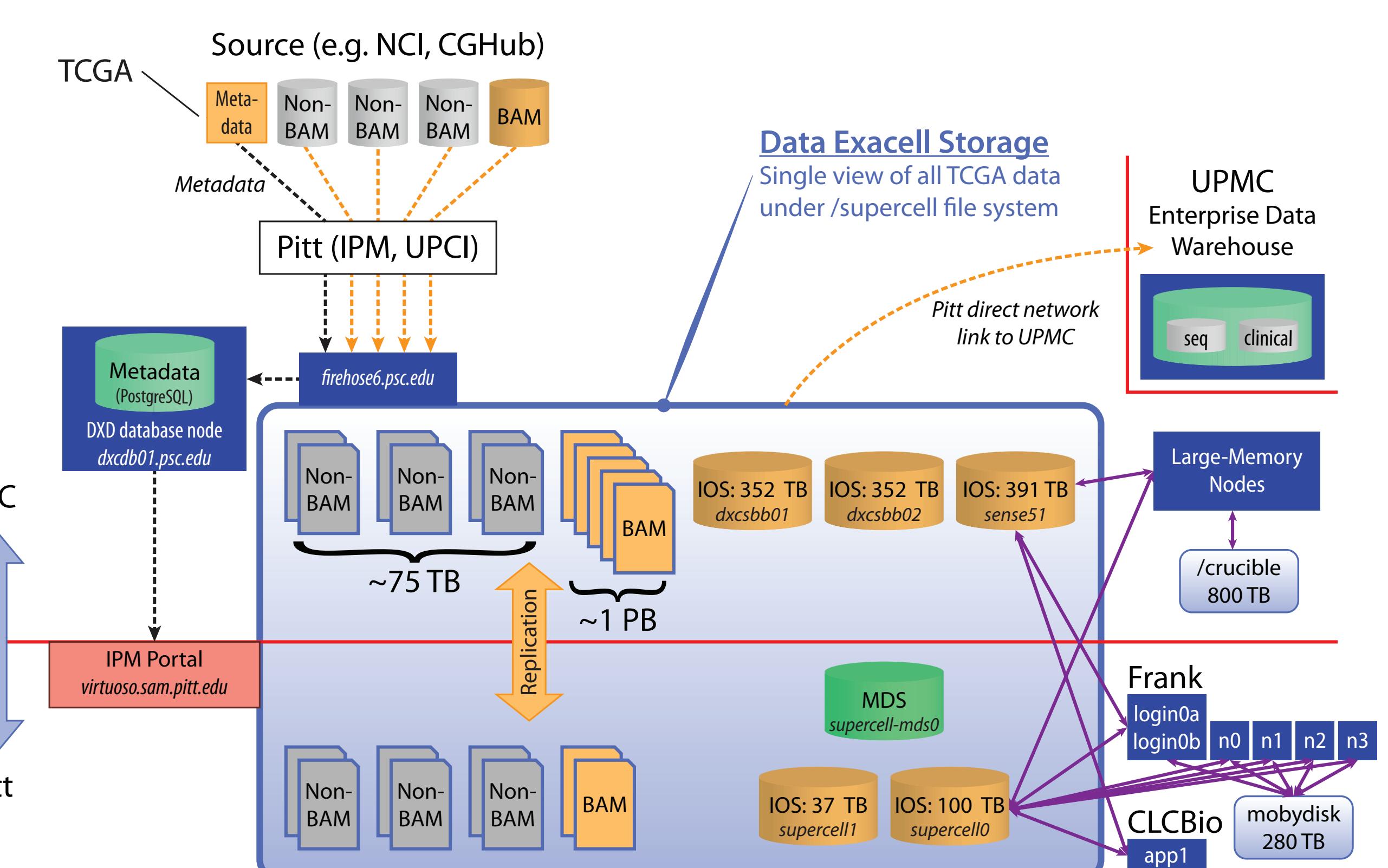
- » **Production:** **Greenfield** provides 12 TB and 3 TB large-memory nodes and a 1PB SLASH2 file system (**Crucible**) for researchers.
- » **Development:** Additional large-memory, accelerated, and standard nodes (**DXC**) provide development resources for facilities research and an additional file system (**Hopper**).

Pilot Applications

Partnering with research groups allows synergistic development of analytical capabilities with development of increasingly sophisticated mechanisms for data handling. Examples include the Pittsburgh Genome Resource Repository, Galaxy, the NRAO Green Bank Telescope Mapping Pipeline, the Causal Web, GenePattern, and text analytics.

Pittsburgh Genome Resource Repository (PGRR)

PGRR was developed using SLASH2 and other DXC infrastructure to effectively provide The Cancer Genome Atlas (TCGA) to researchers working in the life sciences.



Preliminary Results

SLASH2 and other data infrastructure building blocks have been extended, tested, and made production quality. They form a valuable foundation for data-analytic applications and higher layers of data infrastructure. Applications such as PGRR and the Causal Web have been built on them and then transitioned to production on Bridges.