Hybrid Cloud and Local Data Lifecycle Management at a Medium-Scale Genomics Center

Sean M.L. Griffith, Benjamin Myers, Alice Sanchez, Michelle Mawhinney, Justin Paschall, Beth Marosy, David Mohr, Elvin Hsu, Kimberly F. Doheny

Department of Genetic Medicine // Center for Inherited Disease Research (CIDR) // Johns Hopkins Genomics (JHG) Johns Hopkins University School of Medicine, Baltimore, MD, USA



PromethION 24







Obstacles in Data Management

The latest crop of sequencing technologies produce even more data.







A Solution: Targeted Cloud Storage

To overcome the data deluge, we need more capacity and to ensure recovery, we need data replication. Cloud storage solves both problems, but, engaged naively, significantly raises the costs of storing data.



- Crying IT Person

NovaSeq X Plus

All organizations have an increasing need for recovery capabilities, such as following a cyber-attack or local disaster.



"I can't sleep anymore!" - Crying IT Person

An Ontology for Data Relationships

How is a set of genomics data constituted? We have more than just the raw instrument data to consider.

An Additional Assumption:

We should prepare for emergency situations, but they are unlikely. Replicated data will rarely need to be used for full recovery.

What is the solution?

- Store all cloud data in the 'deep freeze' or archival tier. This storage is much cheaper, but more expensive to retrieve, due to thawing costs.
- For active projects, replicate instrument data and perpetual project data, which are the minimal data set needed to reconstitute those projects in the case of data loss.
- For completed projects, replicate the final genomic data, along with perpetual project data, keeping essential and useful data available.
- Most importantly, develop, follow, and automate data lifecycles that apply to both instrument data, all types of projects, and any other regularly generated bulky data.

Technical Details - The BlobLobber Suite

All local data are stored on a dedicated Qumulo NAS, with cloud data written to Microsoft Azure Blob Storage. Data transfer from local storage to the cloud is automated via BlobLobber, a custom suite of tools built around the Azure CLI using a combination of shell scripts, Java, and Python, with an accompanying back-end housed in a MySQL database. BlobLobber is designed to transfer data from local storage to the cloud in a massively parallel fashion, using available HPC resources during periods of low occupancy.

