

¹Ontario Institute for Cancer Research, Toronto, ON, Canada; ²McGill University, Montreal, QC, Canada; ³University of Toronto, Toronto, ON, Canada; ⁴University of Ottawa, Ottawa, ON, Canada; ⁵Centre of Genomics and Policy, Montreal, QC, Canada; ⁶Simon Fraser University, Vancouver, BC, Canada; ⁷BC Cancer Agency Research Centre, Vancouver, BC, Canada

ABSTRACT

What is ICGC?

[illegible]

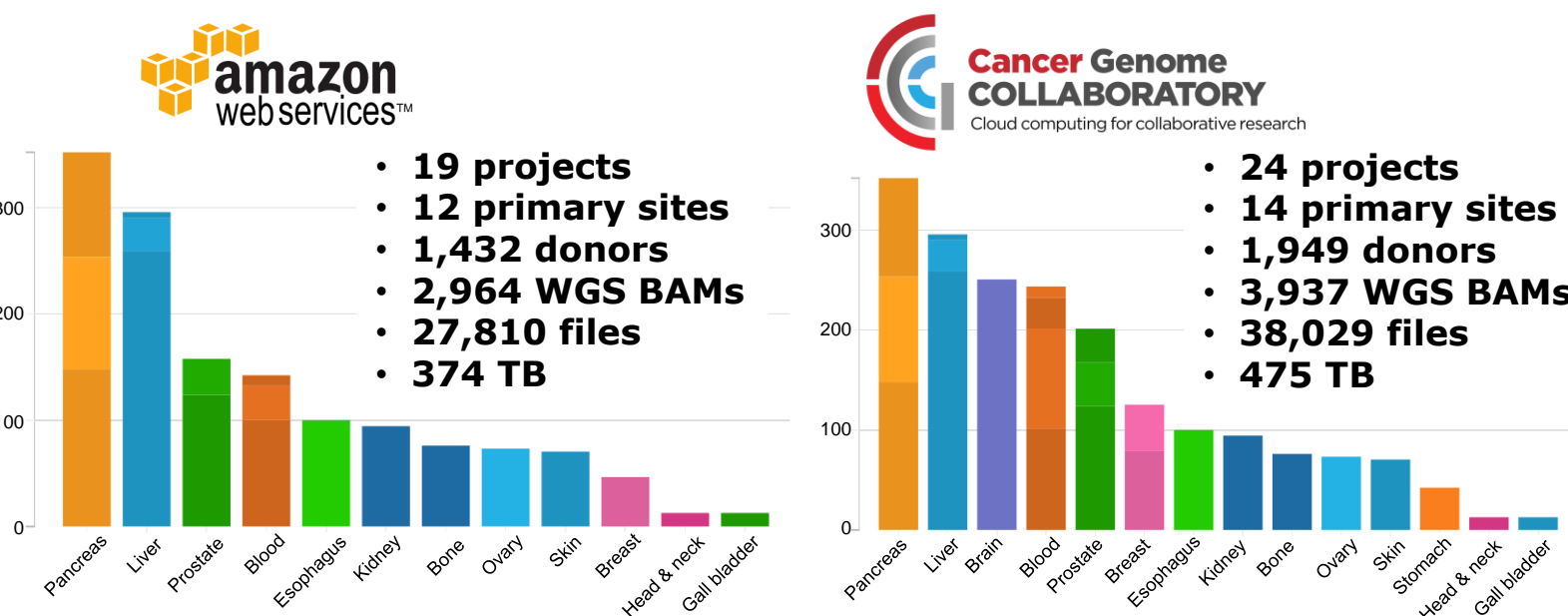
What is PCAWG?

Using a series of **14 academic and commercial compute clouds**, we were able to process this **800 terabyte data set** in just over a year's time. Given the improvements in the software that occurred over this period, the whole project would take less than 4 months on just a single commercial cloud if we were to start over. When the project is completed later in 2016, we will again use academic and compute clouds to publish the PCAWG data, its major results, and all the software used during the analysis, thereby allowing the research community to integrate PCAWG with their own data sets, and apply the same analytic procedures.

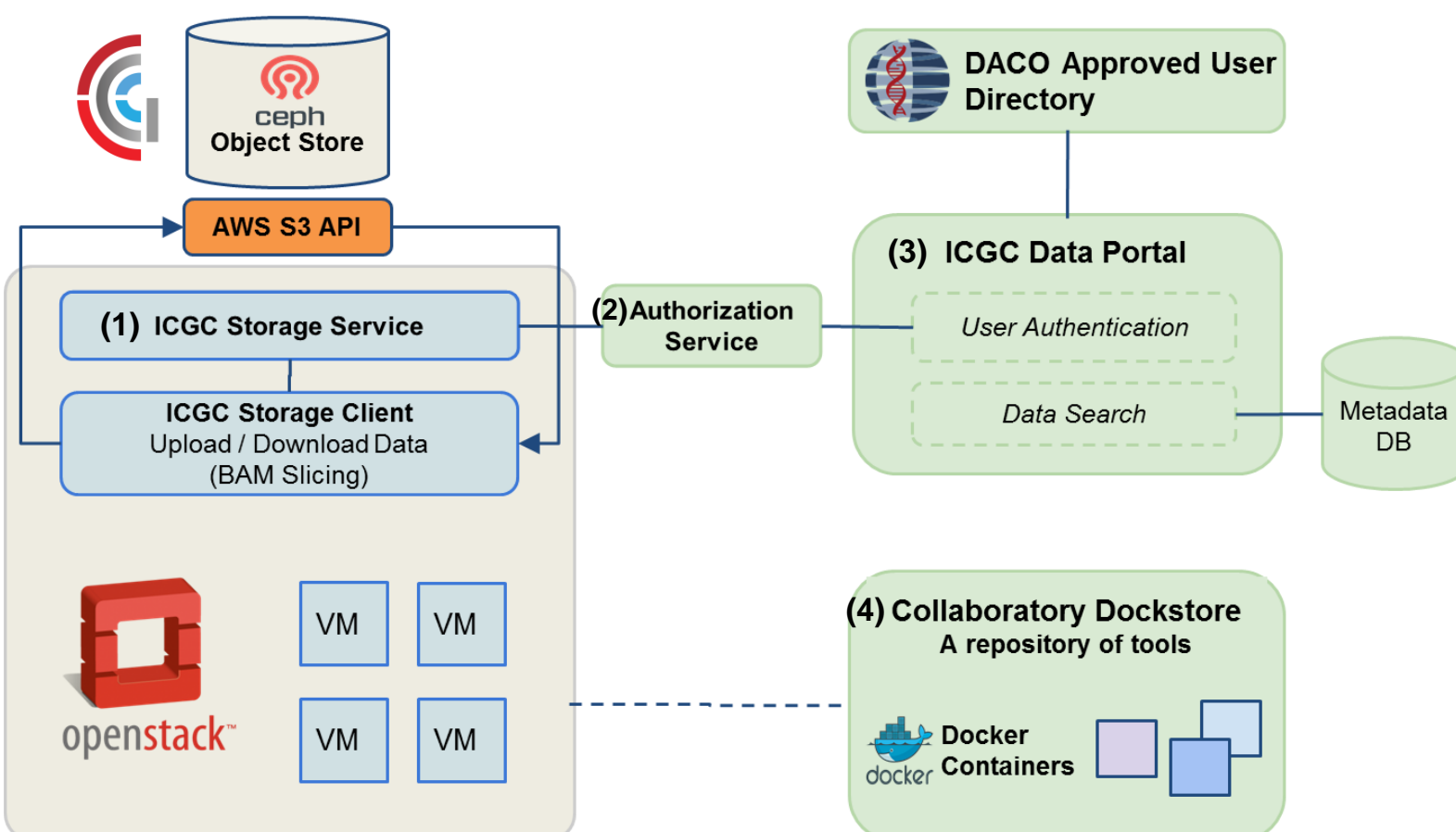
Why Cloud Computing?

	Size	Transfer on dedicated 10G link	Transfer on shared University link
PCAWG data set today	0.8 PB	9 days	8 months
ICGC data set in 2018	5.0 PB	2 months	4 years

**Two Clouds Hosting PCAWG Data:
Amazon Web Services – commercial
Cancer Genome Collaboratory – academic**



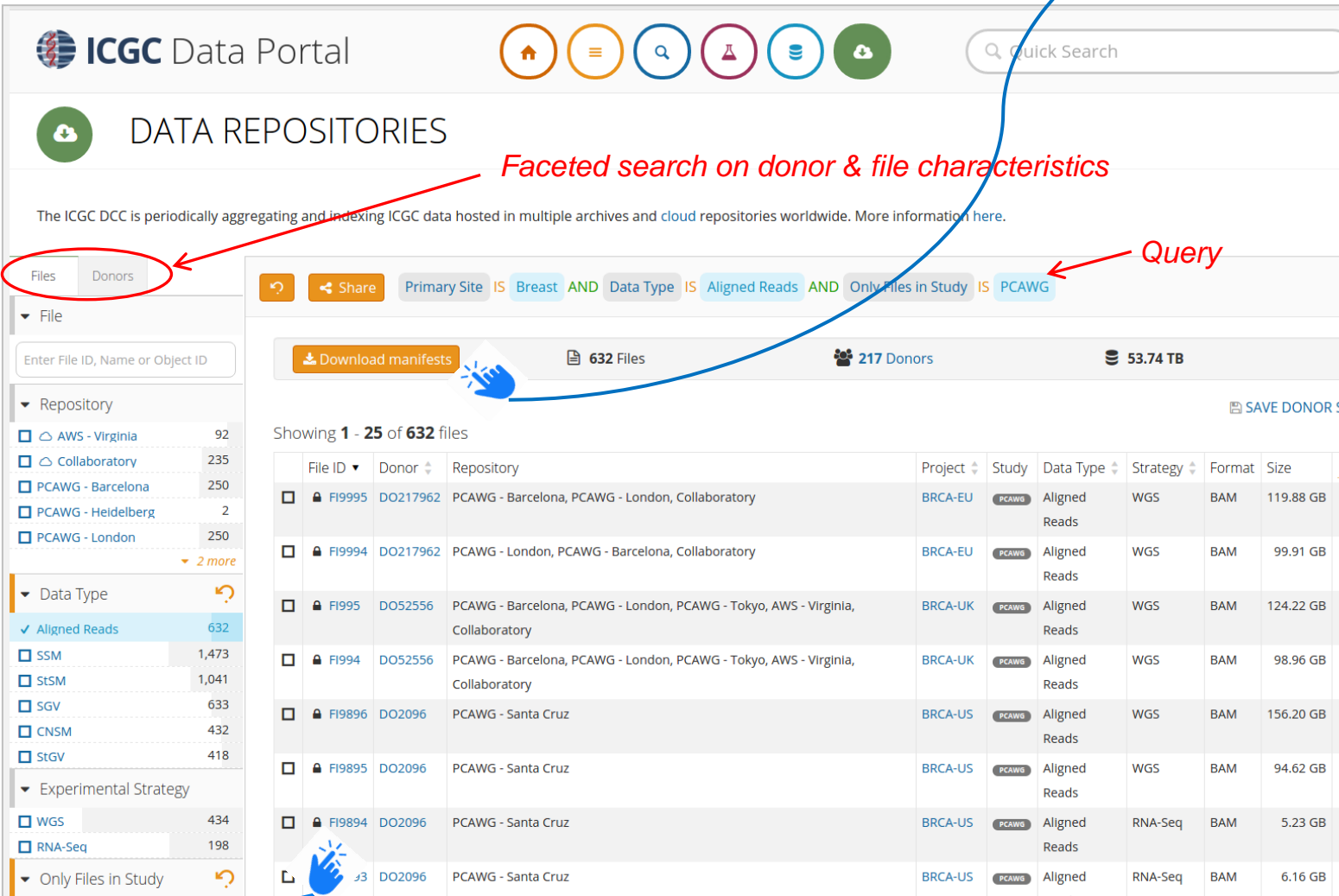
Infrastructure Overview



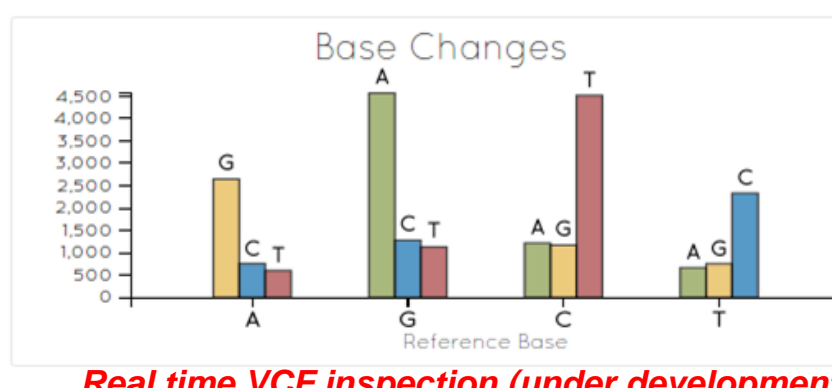
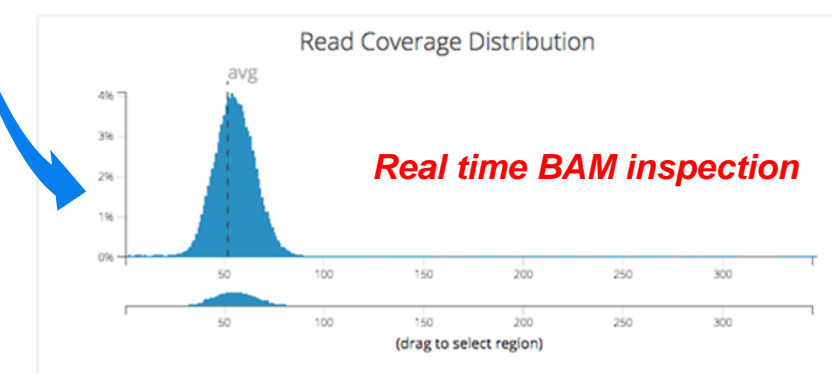
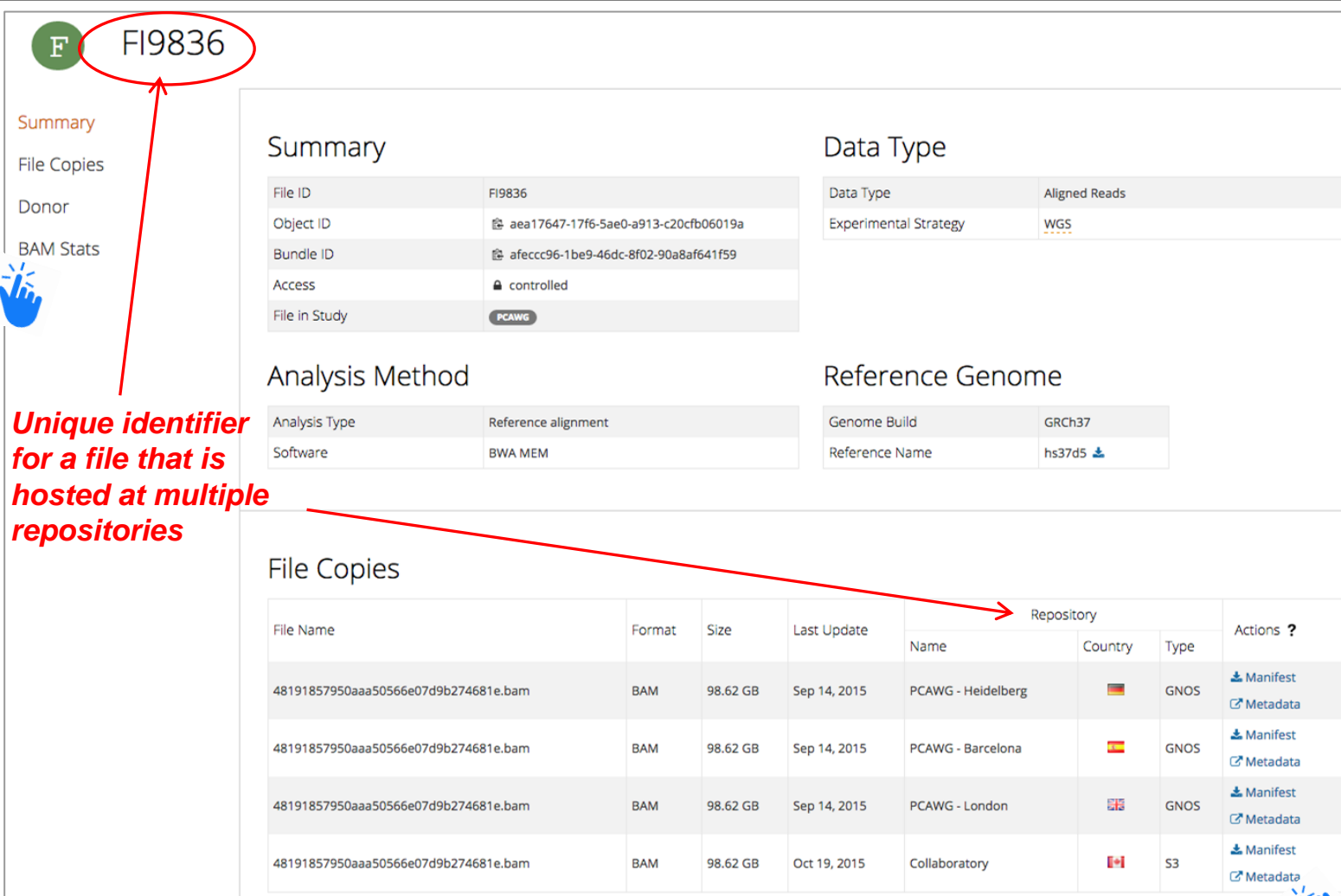
The "ICGC in the cloud" software infrastructure includes four main components: (1) a client-server storage module used to upload/download large data files to/from the Ceph object store; (2) an authentication and authorization module to ensure that only approved researchers from the ICGC Data Access Compliance Office (DACO; <https://icgc.org/daco>) can access ICGC control-tier data; (3) a user-friendly data search tool in the ICGC Data Portal and finally, (4) an open platform called Dockstore (dockstore.org) to share data analytic pipeline packaged in Docker containers. While the above figure depicts the infrastructure at Collaboratory, a similar system is used at AWS but with S3 in place of Ceph, and EC2 instances in place of OpenStack VMs.

How to Search and Download Data?

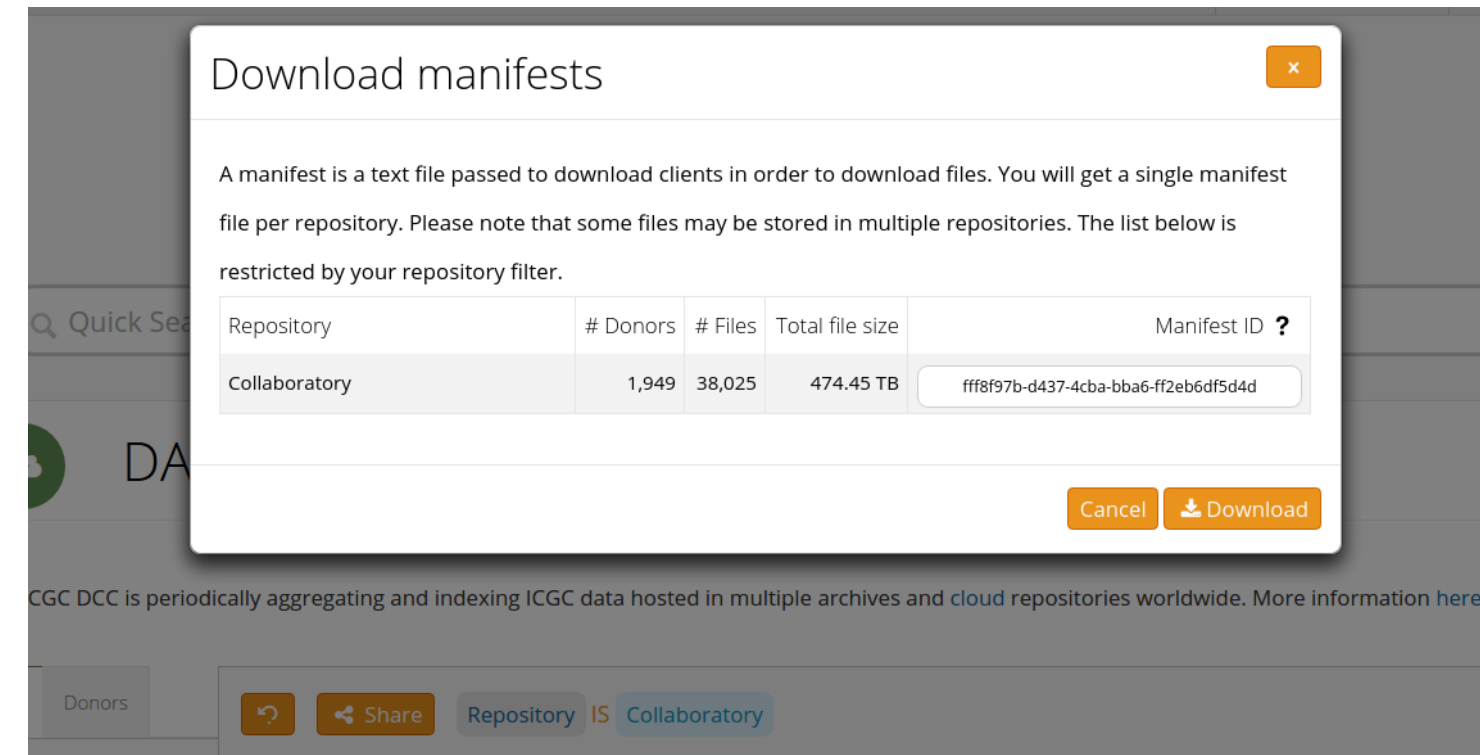
1. Perform faceted search at the ICGC Data Portal for data hosted at multiple repositories (GNOS, AWS, Collaboratory)



2. View metadata and other details on the file entity page

[illegible]

3. Save the list of files as a manifest



4. Download the files using the ICGC Storage Client

- high-throughput multi-part download
- resumability after transfer interruption
- fault tolerant against data corruptions in transit
- MD5 checksum validation
- BAM slicing
- Filesystem in Userspace (FUSE)

Download manifest data

```
%: icgc-storage-client download --manifest 4jdyyqs099ew22
--output-dir data --output-layout bundle
```

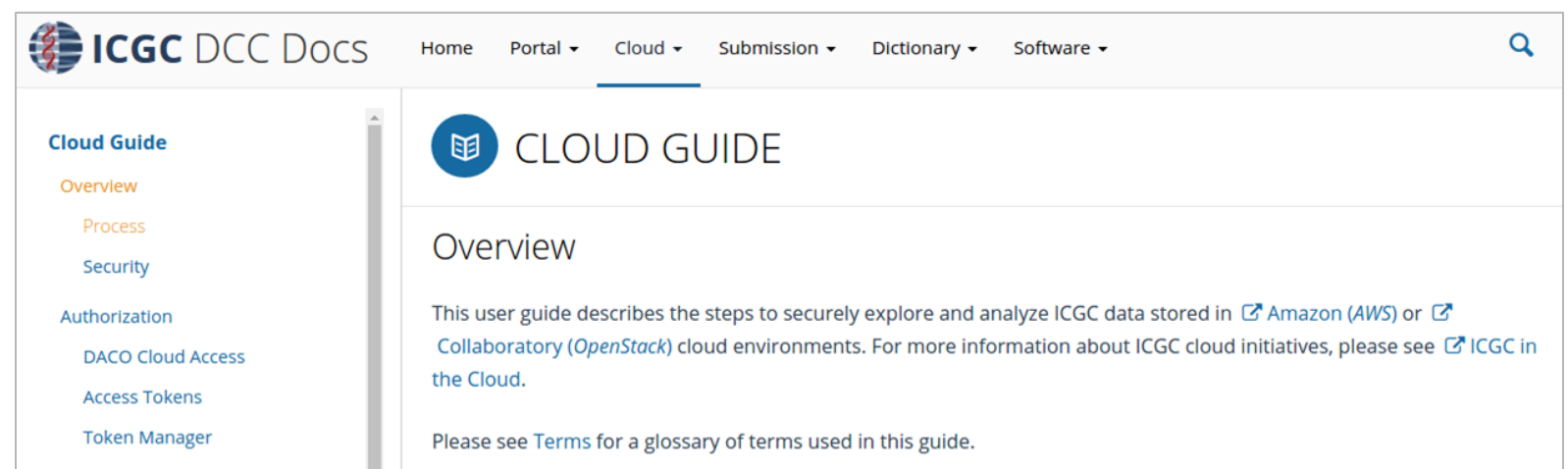
Download BAM slices

```
%: icgc-storage-client view --object-id ea17647-17f6-5ae0
--query 12:25357723-25403870
```

Mounting a manifest (FUSE)

```
%: icgc-storage-client mount --manifest 4jdyys099ew22 --
mount-point /tmp/
%: ls /tmp
```

Comprehensive User Guide on docs.icgc.org



Please contact help@cancercollaboratory.org for more information

Acknowledgements