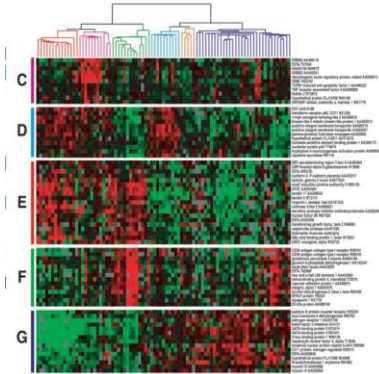# Introduction to the Gen3 Platform for Data Commons and Data Ecosystems

Phillis Tang
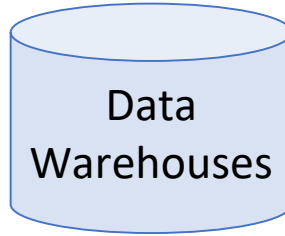
Center for Translational Data Science

University of Chicago
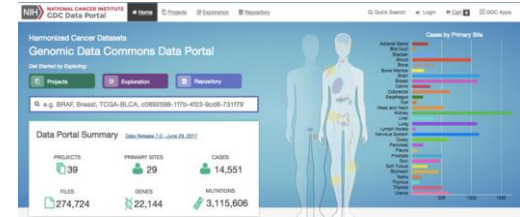
& Open Commons Consortium

**Databases** organize the data around a **project**.



Data Warehouses

**Data warehouses** organize the data for an **organization** (and are enabled by enterprise computing)



**Data Commons** organize the data for a scientific discipline, **community**, or field and are enabled by large scale cloud computing.

Multi-Discipline



Discipline



(Virtual) Organization
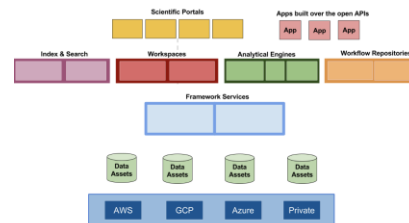


**Data Ecosystems**
2018 - 2028

Project



**Data Commons**
2014 - 2024

- Supports large data
- Workspaces
- **Common data models**
- **Core data services**
- **Data & Commons Governance**
- **Harmonized data**
- **Data sharing**
- **Reproducible research**

**Data Clouds**
2010 - 2020

- Supports large data & data intensive computing with **cloud computing**
- Researchers can analyze data with collaborative tools (**workspaces**) – so data does **not** have to be downloaded)

**Databases**
1982 - present

- Data repository
- Data catalogs
- Download data

- Interoperates **multiple data commons, databases, knowledge bases**, and other resources
- Supports **ecosystem of commons, portals, notebooks, applications & simulations** across multiple disciplines

# Genomic Data Commons - data exploration

# Data Access Control

**Cloud Bucket With Data**

- Bucket policy prevents access by unauthorized users

- Data access is logged for auditing and compliance
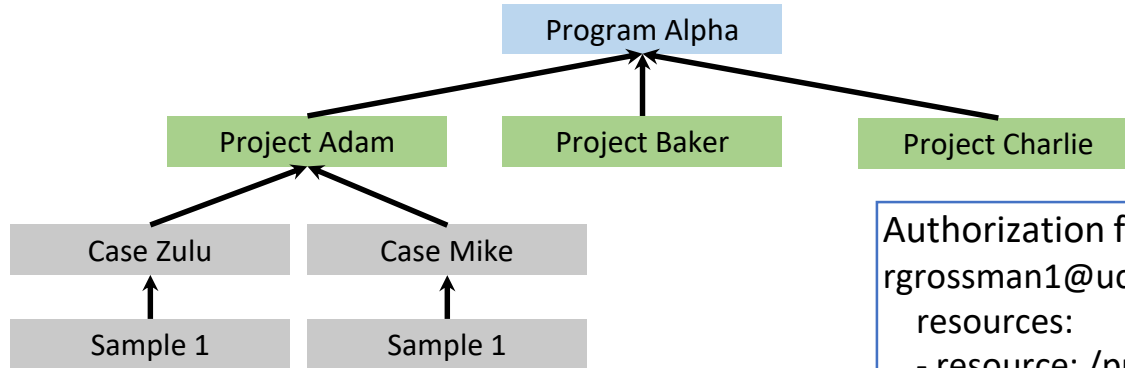
**Gen3 Auth**



- Gen3 Auth(Fence) provides Authentication and Authorization, and Data Access.

- Gen3 Auth works with multiple identify providers (IdP) including Google, and easily adaptable for any support OIDC provider
- This enables Single Sign On (SSO) compatibility with most systems

- Authorization for data access via internal Access Control List specified by the stakeholders

# Data Access Control

• Gen3 auth has a Role Based Access Control (RBAC) engine

**Gen3 Auth**

The RBAC engine understands the hierarchical nature of a users permissions, and can be used to determine if the user has access to a specific piece of data

```
                    Program Alpha

   Project Adam      Project Baker      Project Charlie

Case Zulu   Case Mike

Sample 1     Sample 1
```

Authorization for a user would then be stored as:
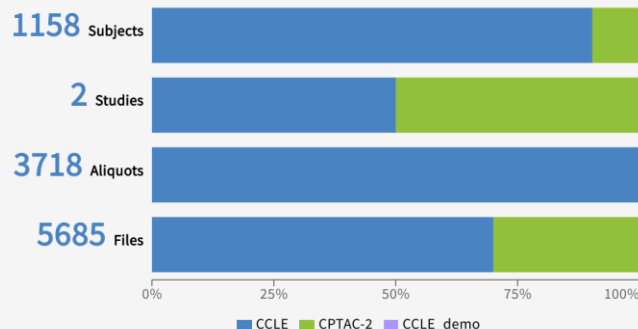rgrossman1@uchicago.edu:
   resources:
   - resource: /programs/alpha/projects/baker
     privilege: [create, read, read-storage, write-storage]
   - resource: /programs/alpha/projects/adam/cases/zulu
     privilege: [read, read-storage]

Giving write (submission) access to the Baker project and all nodes underneath it, while read access to only the Zulu case in the Adam project

# Data Access Control



**Query Gateway**

- Query gateway provides the potential to limit the queries that users can perform and control when results are returned.

Examples of queries:
Query1: StandardDeviation(variable) where STUDENTS_GENDER is MALE
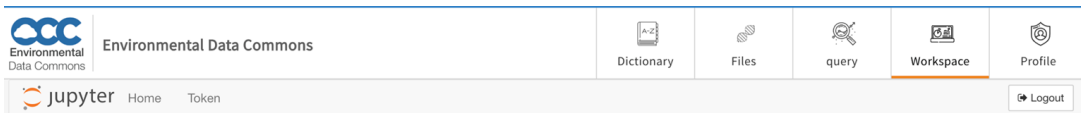
Blue = querying user can specify

Results returned only when # of students represented in the query > a threshold. *I.e. only return standard deviations when the query is computing it for at least 10 students.*
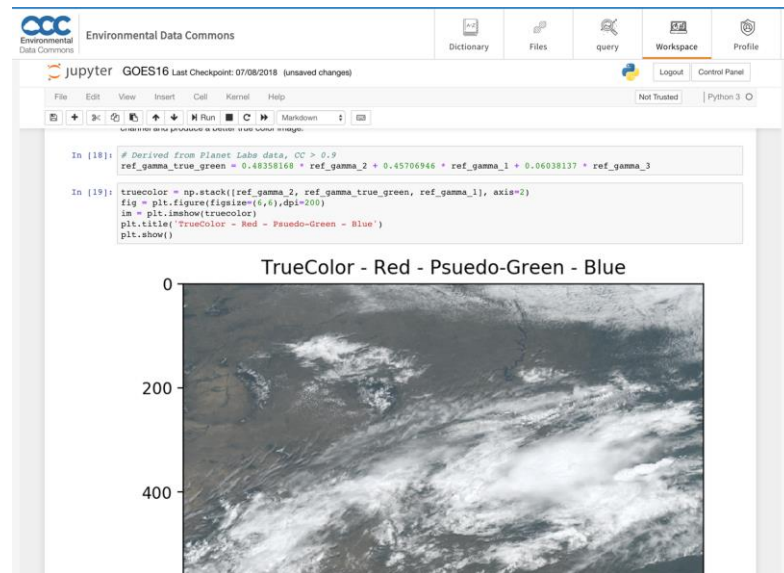
# Jupyter Notebooks

Jupyter

- Jupyter Notebooks are powerful tools for creating custom analysis over datasets
- Gen3 runs Jupyter Notebooks in a secure cloud environment helping to reduce the need to download data to laptops, etc.

# Data Ontologies

**Dictionary viewer**

- Gen3 dictionary viewer allows browsing data vocabularies within a particular data commons

# Data Ontologies

PFB

- Ontologies contain controlled vocabulary developed by a standards body.
- Data dictionaries contain references to the ontology terms allowing harmonization of differing data dictionaries

# Data Aggregation

Data & User Flow with Gen3