THE UNIVERSITY OF CHICAGO | Center for Translational Data Science
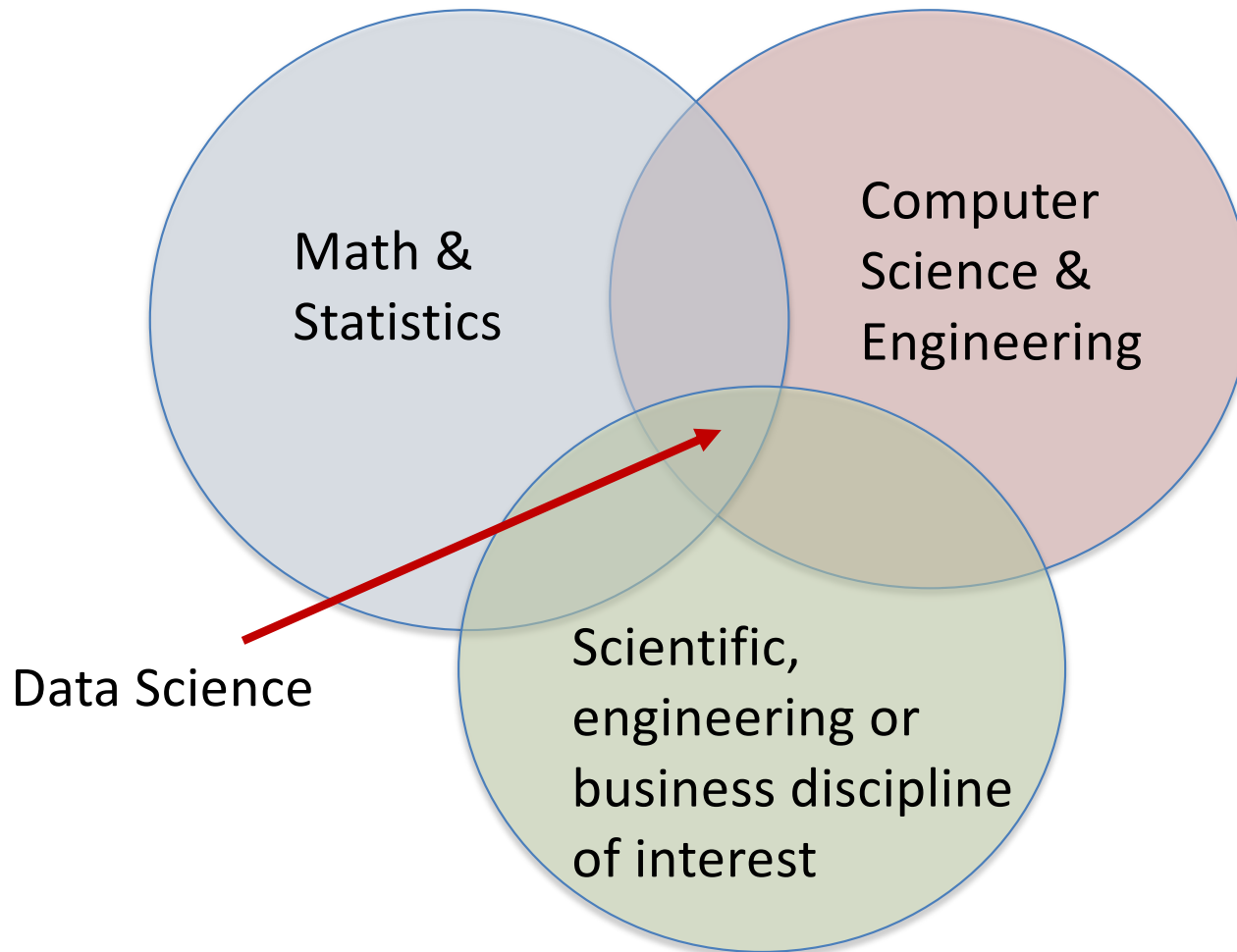
# Introduction to the Center for Translational Data Science (CTDS) and Some of the Data Commons It Develops
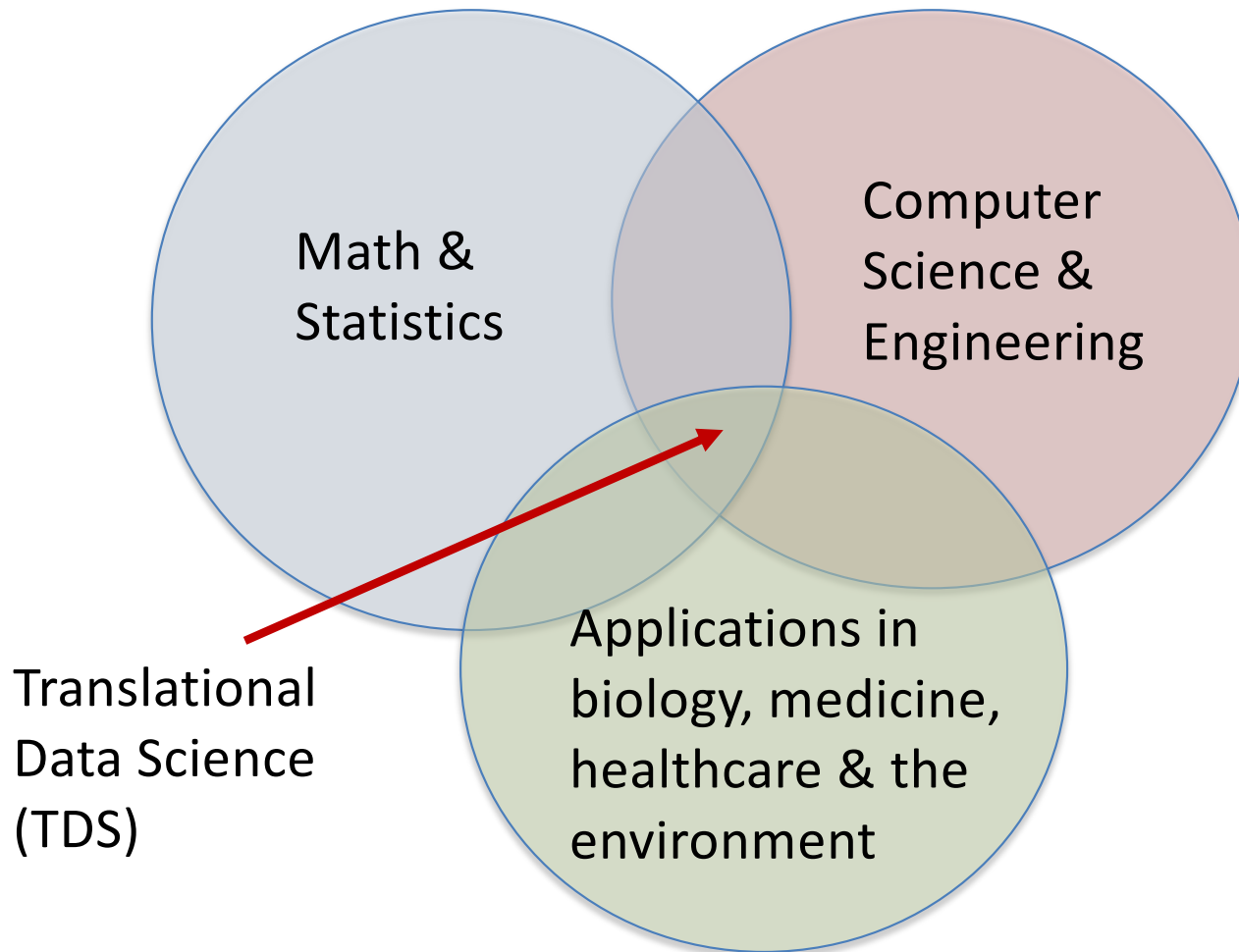
Robert Grossman

Center for Translational Data Science

University of Chicago

January 30, 2019
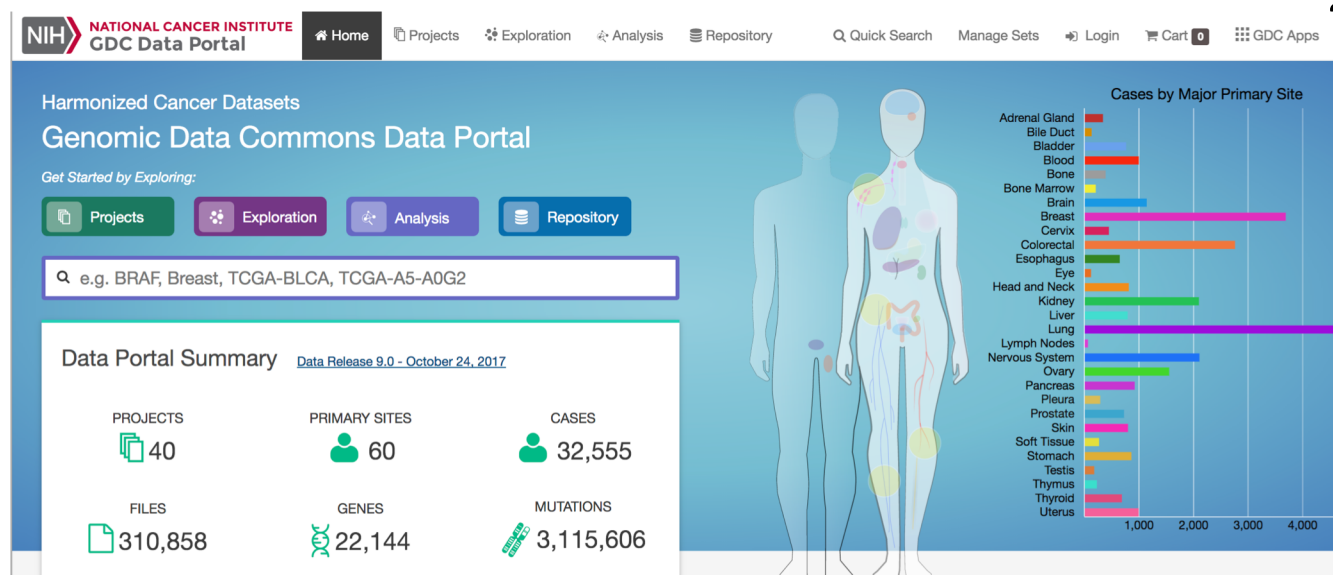
# Part 1: Introduction to the CTDS

Math & Statistics

Computer Science & Engineering

Scientific, engineering or business discipline of interest

Data Science

- As usual, we view data science as an approach that integrates math/stat, computer science & engineering, and its applications.

Math & Statistics

Computer Science & Engineering

Applications in biology, medicine, healthcare & the environment

Translational Data Science (TDS)

- Translation is about the human or societal impact of the data science.
- The challenge is **translating** a discovery in data science to have an impact.
- Translational data science is the discipline that supports this challenge.

# NCI Genomic Data Commons*



The GDC consists of a 1) data exploration & visualization portal (DAVE), 2) data submission portal, 3) data analysis and harmonization system system, 4) an API so third party can build applications.

*See: NCI Genomic Data Commons: Grossman, Robert L., et al. "Toward a shared vision for cancer genomic data." New England Journal of Medicine 375.12 (2016): 1109-1112.

The GDC makes available over 2.5 PB of data available for access via an API, analysis by cloud resources on public clouds, and downloading.
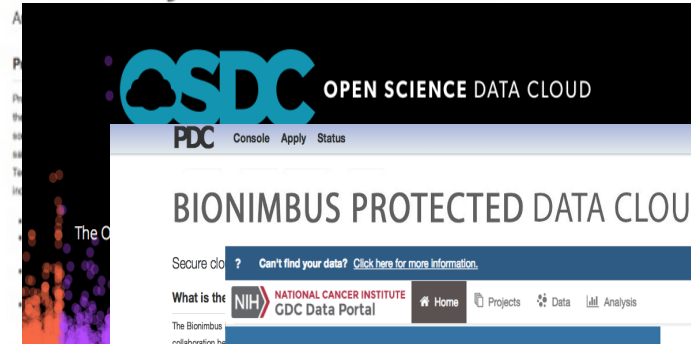
In an average month, the GDC is used by over 22,000 users and over 2 PB of data are downloaded.

The GDC is based upon an open source software stack that can be used to build other data commons.

OCC – NASA Project Matsu (2009)

OCC Open Science Data Cloud (2010)

**Gen1**

Bionimbus Protected Data Cloud* (2013)
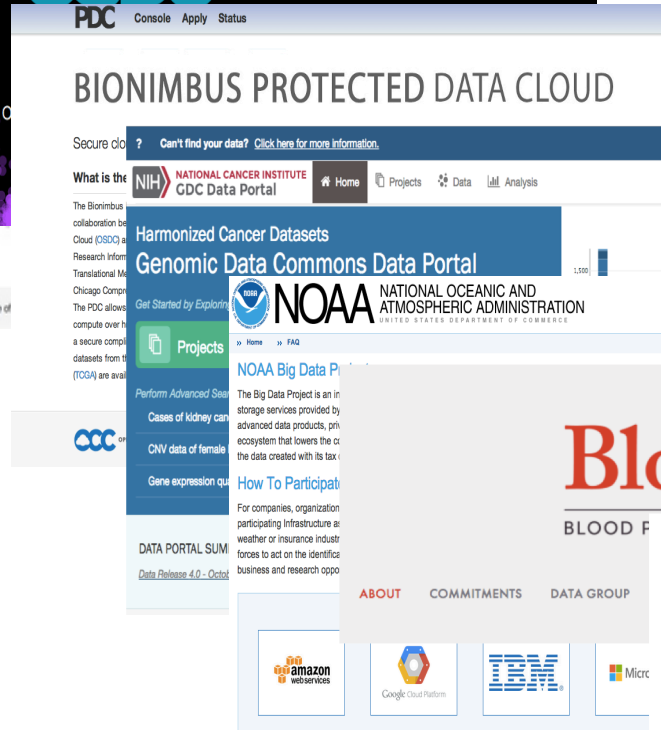
NCI Genomic Data Commons* (2016)

**Gen2**

OCC-NOAA Environmental Data Commons (2016)

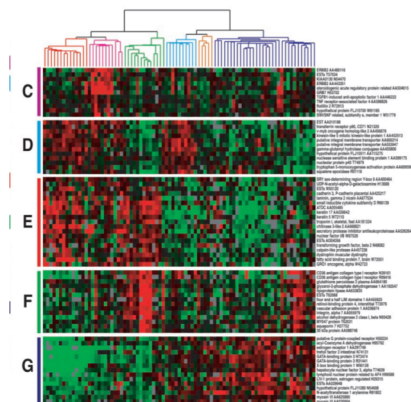OCC Blood Profiling Atlas in Cancer (2017)
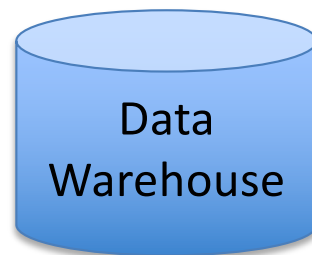
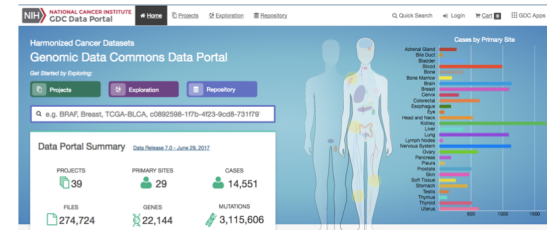Kids First Data Resource (2017)

Brain Commons (2017)

**Gen3**

*Operated under a subcontract from NCI / Leidos Biomedical to the University of Chicago with support from the OCC.

**Databases** organize the data for a project or department.



**Data warehouses** organize the data for an **organization**



**Data commons** organize the data for a research **discipline** or field

Multi-Discipline



Discipline



(Virtual) Organization



**Data Clouds**
2010 - 2020

Project



**Databases**
1982 - present

- Data repository
- Data catalogs
- Download data
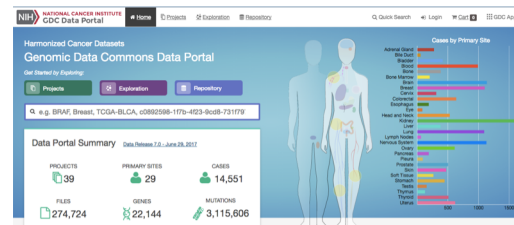
- Supports large data & data intensive computing with **cloud computing**
- Researchers can analyze data with collaborative tools (**workspaces**) – so data does **not** have to be downloaded)
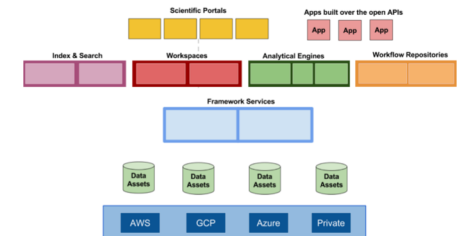
**Data Commons**
2014 - 2024

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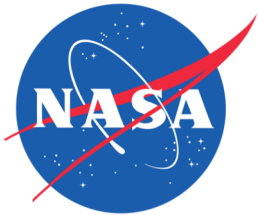- **Data sharing**
- **Reproducible research**

**Data Ecosystems**
2018 - 2028

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

# Center for Translational Data Science – Selected Firsts



### OCC-NASA Project Matsu
**First cloud-based processing of satellite images**

### Bionimbus Protected Data Cloud
**First data cloud to earn NIH Trusted Partner status & to operate at FISMA Moderate**

### NCI Genomic Data Commons
**First genomic data commons**

### BRAIN Commons

### NCI DCF – First cancer data ecosystem

2009          2013          2016          2017

2010                                                          2018
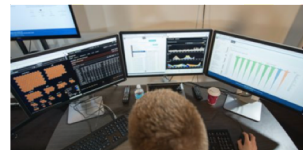
### Open Science Data Cloud
**First petabyte-scale data cloud (2011)**

### Open Commons Consortium
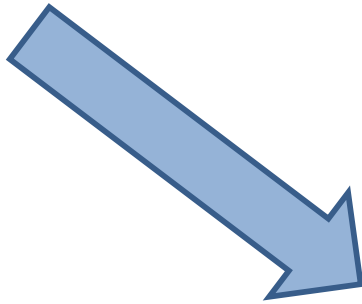**First set of open legal and government agreements**

### First Commons Services Operations Center (CSOC)

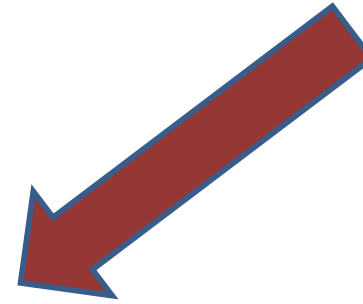### OCC-NOAA
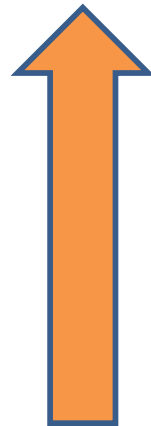**First environmental data commons**

# 2. Data Commons

IT infrastructure challenges
- Data size
- Security & compliance
- Policy restrictions

Growing importance of open data, open reproducible science & data ecosystems

Limited funding

IT infrastructure challenges

Data commons co-locate **data** with **cloud computing** infrastructure and commonly used **software services, tools & apps** for managing, analyzing and sharing data to create an **interoperable resource** for the research community.*

data commons
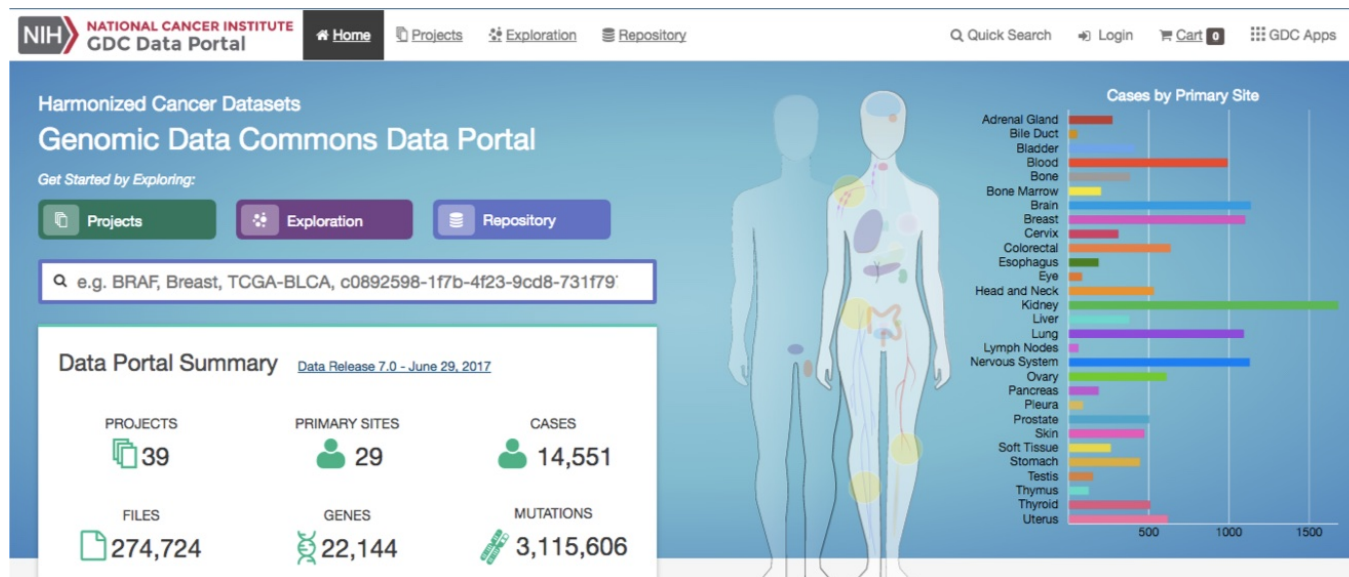
Growing importance of open data, open reproducible science & data ecosystems

Limited funding

*Robert L. Grossman, Allison Heath, Mark Murphy, Maria Patterson and Walt Wells, A Case for Data Commons Towards Data Science as a Service, IEEE Computing in Science and Engineer, 2016.   Source of image: The CDIS, GDC, & OCC data commons infrastructure at a University of Chicago data center.

# NCI Genomic Data Commons*



- The GDC was launched in 2016 with over 4 PB of data.

- Used by 1500 - 3000+ users per day and over 36,000 researchers each month.

- Based upon an open source software stack that can be used to build other data commons.

*Source: NCI Genomic Data Commons: Grossman, Robert L., et al. "Toward a shared vision for cancer genomic data." New England Journal of Medicine 375.12 (2016): 1109-1112.

# Apps 1 and 2: Data Portals to Explore and Submit Data

# App 3: Analysis & Harmonization of all Submitted Data with a Common Set of Bioinformatics Pipelines



- MuSE
  (MD Anderson)
- VarScan2 (Washington Univ.)
- SomaticSniper
  (Washington Univ.)
- MuTect2
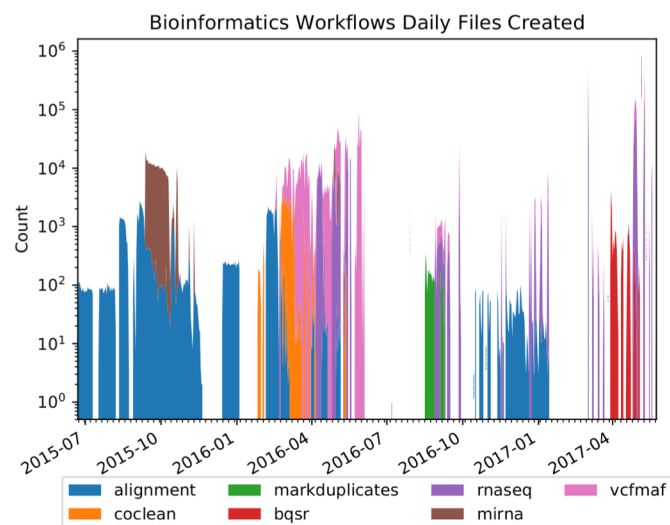  (Broad Institute)

Source: Zhenyu Zhang, et. al. and the GDC Project Team, Uniform Genomic Data Analysis in the NCI Genomic Data Commons, to appear.

# System 4: An API to Support User Defined Applications and Notebooks to Create a Data Ecosyst



GDC developed apps

Third party apps

API URL          Endpoint          Optional Entity ID          Query parameters

https://gdc-api.nci.nih.gov/files/5003adf1-1cfd-467d-8234-0d396422a4ee?fields=state

# Commons



Commons are **resources** that are held in common (and not owned privately) that a group or **community** manage for individual and collective benefit.

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



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 support an
information platform.

# 3. Gen3 Data Commons

# A Gen3 Data Commons Platform in Six Steps



Gen3.org
(an open source platform for developing & operating data commons)

1. Define a data model.
2. Use the Gen3 software to auto-generate the data commons and associated API.
3. Import data into the commons using Gen3 import application.
4. Use Gen3 to explore your data and create synthetic cohorts.
5. Use cloud based platforms, tools and workspaces to analyze the synthetic cohosts.
6. Develop your own container-based workflows, applications and Jupyter Notebooks.

# (Selected) GEN3 Data Commons

**BRAIN Commons**
Total Files: 16,733
Total Size: 270.14 GB

**GenoMEL** the Melanoma Genetics Consortium
Total Files: 4,008
Total Size: 20.77 TB

**BloodPAC** BLOOD PROFILING ATLAS IN CANCER
Total Files: 1,855
Total Size: 295.73 GB

**ACCOUNT**
Total Files: 1,952
Total Size: 3.77 TB

**NHLBI DATA STAGE**
Total Files: 71,368
Total Size: 344.03 TB

Gabriella Miller **Kids First** PEDIATRIC RESEARCH PROGRAM Data Resource Center
Total Files: 134,531
Total Size: 2.33 PB

**NIAID DATA HUB**
Total Files: 156,368
Total Size: 304.59 GB

NIH NATIONAL CANCER INSTITUTE Cancer Research Data Commons
Total Files: 1,688,568
Total Size: 2.2 PB

**Environmental Data Commons**
Total Files: 12,317,334
Total Size: 28.93 TB

# From Data Commons to Data Ecosystems of Interoperating Data Commons



1. Build data commons over hosted Data Commons Framework Services
2. Interoperate your data commons with other DCFS compliant data commons.

# 4. Suggested Guidelines for Foundations for Data Sharing

# Key Issues

- Data sharing is not as simple as copying data or data & metadata from a repository
- Data Governance and Commons Governance is required
- Data standards
- It is important to fund both the commons and the data scientists that curate the data and build applications for the end users
- Data commons don't care about the specific data, as long as the support the required data types and applications.
- Each research group doesn't need to build their own commons (this is called multi-tenancy).

# Sharing Data with Data Commons – the Main Steps

1. **Require data sharing.** Put data sharing requirements into your grant agreements. We can work out some common language.
2. **Build a commons.** Lead, co-lead or join a data commons, fund it, and develop an operating plan, governance structure, and a sustainability plan.
3. **Populate the commons.** Provide resources to your researchers to get the data into data commons.
4. **Interoperate with other commons.** Fund your commons developers and operators to interoperate with other commons that can accelerate research discoveries.
5. **Support commons use.** Support applications that ask for support to build apps over commons.

# The Components of a Data Commons

**Sponsor (One or More Foundations)**

**Third party open source apps**

**Third party vendor apps**

**Commons Operations Center**

**Sponsor funded apps**

**Data Commons Governance & Standards**

**Open Source Software for Data Commons**

**Data managed by the data commons**

**On Premise Clouds**

**Public Clouds**

Object-based storage with access control lists

Scalable light weight workflow

Database services

Community data products

Portals for accessing & submitting data

**Data Commons 1**

**Data Commons 2**

APIs

Apps

Notebooks

Apps

**Apps & Notebooks**

Workspaces

Workspaces

**Workspaces**

Commons Framework Services (Digital ID, Metadata, Authentication, Auth., etc.) that support multiple data commons.

**Data Commons Framework Services**

# Building the BloodPAC Data Commons
# (Examplar of Principles 1 & 2)



Adapt the data model to your project

Set up & configure the data commons (CSOC)

Put in place the OCC data governance model

Research groups submit data

Clean and process the data following the standards

New research discoveries

Researchers use the commons for data analysis

# Three Principles For Foundations Funding Research

1. Require that researchers share the data generated by research that you fund.
2. Foundations should provide the computing infrastructure and bioinformatics resources that is required to support data sharing.
3. The data commons supported by Foundations should themselves share data and interoperate with other data commons.

# Questions?

rgrossman.com
@bobgrossman

# References

# For more information:

- To learn more about data commons: Robert L. Grossman, et. al. A Case for Data Commons: Toward Data Science as a Service, Computing in Science & Engineering 18.5 (2016): 10-20. Also https://arxiv.org/abs/1604.02608
- To large more about large scale, secure compliant cloud based computing environments for biomedical data, see: Heath, Allison P., et al. "Bionimbus: a cloud for managing, analyzing and sharing large genomics datasets." Journal of the American Medical Informatics Association 21.6 (2014): 969-975. This article describes Bionimbus Gen1.
- To learn more about the NCI Genomic Data Commons: Grossman, Robert L., et al. "Toward a shared vision for cancer genomic data." New England Journal of Medicine 375.12 (2016): 1109-1112. The GDC was developed using Bionimbus Gen2.
- To learn more about BloodPAC, Grossman, R. L., et al. "Collaborating to compete: Blood Profiling Atlas in Cancer (BloodPAC) Consortium." Clinical Pharmacology & Therapeutics (2017). BloodPAC was developed using the GDC Community Edition (CE) aka Bionimbus Gen3

# Contact Information

Robert L. Grossman
rgrossman.com

@BobGrossman
robert.grossman@uchicago.edu

THE UNIVERSITY OF CHICAGO | Center for Translational Data Science

ctds.uchicago.edu