# Outline of today's agenda

- Introduction to ISB-CGC and the data commons
- ISB-CGC's approach to derived data
- Data exploration in BigQuery ('Excel-like data tables in the cloud')
- Hands on demonstration of the Google Cloud with focus on BigQuery

# The ISB-CGC homepage
## isb-cgc.org



Focus today is on analysis of derived data in BigQuery

Ways to use data and tools:
- Explore datasets
- Create cohorts
- Run pipelines
- Specialized DBs
  - Mitelman DB of Chromosomal Aberrations & Gene Fusions
  - The *TP53* DB
  - caNanoLab

The *TP53* Database compiles various types of data and information from the literature and generalist databases on human *TP53* gene variations related to cancer. The database is hosted by the National Cancer Institute (NCI) of the United States. The content reflects the R20, July 2019 version

Upcoming *TP53* meetings and conferences can be found in the Events menu.
Did you find any issues? Please submit a report.

## Functional / Structural Data ⊕

Explore functional and structural data and frequency statistics of all possible single nucleotide substitutions in *TP53* exonic sequences, other variants reported in human samples, and validated polymorphisms.

## Tumor Variants ⊕

Explore data for *TP53* tumor variants identified in human tumor samples. Includes data on the type and position of variants, detailed information on the tumor in which the variants have been found, and on various characteristics of the patients in which the tumor developed.

## Germline Variants ⊕

Explore data for individuals that are carriers of a *TP53* germline variant and families in which at least one family member has been identified as a carrier of a germline variant in the *TP53* gene.

## Cell Lines ⊕

Explore data for cell-lines that have been screened for *TP53* variant and have been published in the scientific literature, in the Sanger cell-line database, or the Broad Cancer cell-line Encyclopedia.

## Mouse Models ⊕

Explore data for mouse models with engineered *p53* that are compiled in the caMOD database or reported in the scientific literature.

## Experimentally Induced Variants ⊕

Explore data for variants in the human *TP53* gene obtained from mutagenicity assays in the Hupki mouse model (MEF cells treated with the indicated carcinogen agent) or in a yeast assay.

https://tp53.isb-cgc.org/

# caNanoLab

Enter keyword | Search

| RELATED LINKS | HOME | PROTOCOLS | SAMPLES | PUBLICATIONS | HELP | GLOSSARY | LOGIN |

## NCI

- caNanoLab Wiki
- ISA-TAB-Nano Wiki
- NCI CBIIT Home
- NCL Home
- NCI Nanodelivery Systems and Devices Home
- NCI Home
- Nanotechnology Working Group
- caNanoLab Curation

## EXTERNAL Disclaimer

- NBI
- nanoHUB
- SAFENANO
- OECD
- eNanoMapper

Logged in as guest

Associated Groups: Public

### Welcome to caNanoLab

Welcome to the cancer Nanotechnology Laboratory (caNanoLab) portal. caNanoLab is a data sharing portal designed to facilitate information sharing across the international biomedical nanotechnology research community to expedite and validate the use of nanotechnology in biomedicine. caNanoLab provides support for the annotation of nanomaterials with characterizations resulting from physico-chemical, *in vitro* and *in vivo* assays and the sharing of these characterizations and associated nanotechnology protocols in a secure fashion.

### Browse caNanoLab

| Data Type | Public Results |
|---|---|
| **Search Protocols** Search for nanotechnology protocols leveraged in performing nanomaterial characterization assays. | 211 |
| **Search Samples** Search for information on nanomaterials including the composition of the nanomaterial, results of physico-chemical, *in vitro*, *in vivo* and other characterizations, and associated publications. See also Advanced Sample Search | 1583 233 Sample Sources 7580 Characterizations 1738 Physico-chemical 2829 In Vitro 171 In Vivo 2842 Other |
| **Search Publications** Search for information on nanotechnology publications including peer reviewed articles, reviews, and other types of reports related to the use of nanotechnology in biomedicine. | 2404 |

## USER ACTIONS

No account is required to browse publicly available data.

LOGIN ID [                    ]

PASSWORD [                    ]

[ Login ]

## FEATURES

**caNanoLab provides access to information on:**
- Nanotechnology Protocols
- Nanomaterial Composition
- Nanomaterial Characterizations (physico-chemical, in vitro, in vivo)
- Nanomaterial Publications

**For additional information, see the caNanoLab** FAQ **or** User's Guide

## HOW TO

- How do I obtain a new caNanoLab Login ID and Password?
- How do I reset my current caNanoLab password?
- How do I save data to caNanoLab?
- How do I find nanotechnology protocols?
- How do I find nanotechnology publications?
- How can I search for nanomaterials, composition annotations, and characterizations?
- How can I add nanomaterial characterizations?
- Where can I get definitions for nanotechnology concepts?
- How do I incorporate caNanoLab into a data sharing plan?

## WHAT'S NEW

Release 3.1.9 of caNanoLab is not available for download yet, but you can download the prior version of caNanoLab (3.0) from the wiki home page.
**For information on caNanoLab releases, refer to the** caNanoLab Release Notes.

## KEEPING UP WITH caNanoLab

- Stay connected and provide feedback through the caNanoLab User Hub.
- What's New in caNanoLab?

| CONTACT US | PRIVACY NOTICE | POLICIES | DISCLAIMER |

Department of Health and Human Services | National Institutes of Health | National Cancer Institute | USA.gov | Vulnerability Disclosure

NIH...Turning Discovery Into Health®

caNanoLab Release 3.1.9 Build cananolab-3.1.9-0f58a02

https://cananolab.cancer.gov

# The three Data Commons host and control access to different types of cancer data

# ISB-CGC's approach to enabling data science in the cloud

- Moving Excel files into the cloud
- Derived molecular data available for query as you need, updated frequently
- Tooling examples provided to enable data mining and Machine Learning of your data
- Sharing of results with those you choose
- Maximum flexibility of scripting and compute for those who desire it

# ISB-CGC Focuses on Derived Data via BigQuery

| | Projects | Clinical/Biospecimen | File Metadata | Gene Expression | Somatic Mutation | Copy Number | miRNA Expression | DNA Methylation | Protein Expression* | Acetylome | Glycoproteome | Phosphoproteome | Ubiquitylome |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GDC | GDC Metadata | | X | | | | | | | | | | |
| | APOLLO | X | X | | | | | | X | | | | |
| | BEATAML1.0 | X | X | X | X | | | | | | | | |
| | CCLE | X | X | X | X | X | | | | | | | |
| | CDDP EAGLE | X | X | | X | X | | | | | | | |
| | CGCI | X | X | X | X | X | | | | | | | |
| | CMI | X | X | X | X | | | | | | | | |
| | CPTAC | X | X | X | X | X | X | | X | | | | |
| | CTSP | X | X | X | | | | | | | | | |
| | Exceptional Responders | X | X | X | X | | | | | | | | |
| | FM | X | X | | | | | | | | | | |
| | GENIE | X | X | | | | | | | | | | |
| | HCMI | X | X | X | X | X | | | | | | | |
| | MATCH | X | X | | | | | | | | | | |
| | MMRF | X | X | X | X | | | | | | | | |
| | MP2PRT | X | X | | | X | | | | | | | |
| | NCICCR | X | X | X | | | | | | | | | |
| | OHSU | X | X | X | | | | | | | | | |
| | ORGANOID | X | X | X | | | | | | | | | |
| | REBC | X | X | | | X | | | | | | | |
| | TARGET | X | X | X | X | X | X | | | | | | |
| | TCGA | X | X | X | X | X | X | X | X | | | | |
| | TRIO | X | X | | | | | | | | | | |
| | VAREPOP | X | X | | | | | | | | | | |
| | WCDT | X | X | X | | | | | | | | | |
| PDC | PDC metadata | | X | | | | | | | | | | |
| | APOLLO | X | X | | | | | | | | | X | |
| | BROAD | X | X | | | | | | | | | | |
| | CBTTC | X | X | | | | | | X | | | X | |
| | CPTAC | X | X | | | | | | | X | X | X | X |
| | Georgetown Proteomics Research Program | X | X | | | | | | | | | | |
| | ICPC | X | X | | | | | | X | | | X | |
| | Quantitative Digital Maps of Tissue | X | X | | | | | | | | | | |
| | TCGA | X | X | | | | | | | | | X | |

ISB-CGC hosts data from multiple well-known cancer datasets

ISB-CGC

# Data wrangling can be onerous, for example GDC has 24,944 individual transcriptome files for just TCGA

# The data commons host a wealth of data from 20 cancer initiative programs



Case count per Data Category

# ISB-CGC runs ETL pipelines to reduce the processing barrier of entry

| | |
|---|---|
| BEATAML | OHSU |
| CCLE | ORGANOID |
| CGCI | TARGET |
| CMI | TCGA |
| CPTAC | VAREPOP |
| CTSP | WCDT |
| FM | CBTN |
| GENIE | CPTAC |
| HCMI | ICPC |
| MMRF | Targetome |
| NCICCR | Reactome |
| Pan-Cancer Atlas | |
| Georgetown Proteomics Research Program | |
| Quantitative Digital Maps of Tissue Biopsies | |

## 500k+ Files of Heterogeneous Data

| | |
|---|---|
| WGS | DNASeq WXS |
| RNASeq (gene, isoform, exon, junction) | |
| SNP Array (CEL) | |
| DNASeq (MAF, VCF) | Clinical & Biospecimen |
| DNA Methylation | miRNASeq |
| Protein (RPPA) | SNP Array |

**Extract**

**Transform**

**Load**

**Download Data via Multiple Protocols**

**Convert & Standardize File Formats**

**Cloud Storage / Local VM Disk**

**QC & Normalize Data**

**Create Files for BigQuery Import**

**Final BigQuery Tables**

# The Google Cloud offers tools to simply host derived data by concatenating these files into a single BQ table



GDC case files

"one big csv"

case 1
case 2
case 3
case 4
case 5
⋮
case n

1
2
3
4
5
.
.
.
n

→ BigQuery

ISB-CGC

# BigQuery enables simple and efficient links between data types

# BigQuery enables simple and efficient linking between tables



```
SELECT
 <fields>
FROM `isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current` rna
WHERE <conditionals>
```

# BigQuery enables simple and efficient linking between tables



```
SELECT
 <fields>
FROM `isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current` rna
JOIN `isb-cgc-bq.TCGA.clinical_gdc_current` clin
WHERE <conditionals>
```

ISB-CGC

# BigQuery enables simple and efficient linking between tables

Query results

SAVE RESULTS ▾     EXPLORE DATA ▾     ↕

< | JOB INFORMATION | **RESULTS** | CHART | JSON | EXECUTION DETAILS | >

| Row | case_barcode ▾ | fpkm_uq_unstranded | exp__cigarettes_per_ | protein_expression |
|---|---|---|---|---|
| 1 | TCGA-86-A4P7 | 8.9616 | *null* | -0.6541634745 |
| 2 | TCGA-91-6829 | 14.4956 | 5.178082191780... | 0.399470062 |
| 3 | TCGA-91-6828 | 12.558 | *null* | -0.174617102 |
| 4 | TCGA-86-A4P8 | 1.947 | *null* | -0.4681542825 |
| 5 | TCGA-38-4629 | 40.1784 | 5.479452054794... | 0.739331331 |
| 6 | TCGA-38-6178 | 10.7342 | *null* | -0.013338784 |
| 7 | TCGA-78-7166 | 32.8658 | 2.082191780821... | -0.0290081565 |
| 8 | TCGA-78-7167 | 3.7353 | 3.506849315068... | -0.438169383 |

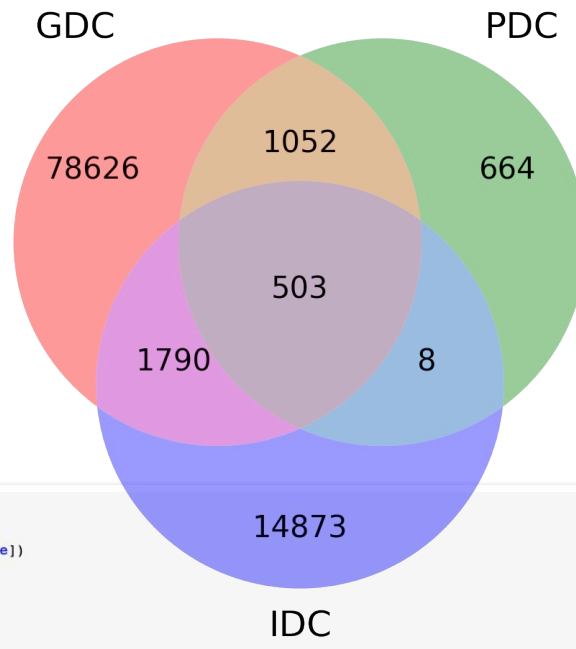Results per page: 50 ▾   1 – 50 of 378   |< < > >|

```
SELECT
 <fields>
FROM `isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current` rna
JOIN `isb-cgc-bq.TCGA.clinical_gdc_current` clin
JOIN `isb-cgc-bq.TCGA.protein_expression_hg38_gdc_current` prot
WHERE <conditionals>
```

ISB-CGC

# There is a substantial overlap of data from the same cases across the Commons

```
1   sql_full = """
2   with gdc as (
3     with cases as (
4       SELECT
5         distinct
6           case_barcode,
7           case_gdc_id,
8           project_id
9       FROM `isb-cgc-bq.GDC_case_file_metadata.caseData_current`
10    ),
11    files as (
12      SELECT
13        case_gdc_id,
14        associated_entities__entity_submitter_id,
15        data_type
16      FROM `isb-cgc-bq.GDC_case_file_metadata.fileData_active_current`
17    )
18    SELECT
19      cases.case_barcode gdc_id,
20      cases.project_id,
21      array_agg(distinct files.data_type) gdc_data
22    FROM cases JOIN files ON cases.case_gdc_id = files.case_gdc_id
23    GROUP BY cases.case_barcode, cases.project_id
24  ),
25  pdc as (
26    SELECT
27      distinct case_submitter_id pdc_id
28    FROM `isb-cgc-bq.PDC_metadata.ali...
29  idc as (
30    SELECT
31      PatientID idc_id,
32      array_agg(distinct has_segmenta...
33      array_agg(distinct has_derived)
34      array_agg(distinct has_quantita...
35    FROM `canceridc-data.idc_current...
36    GROUP BY PatientID
37  )
38  SELECT * from gdc
39  FULL OUTER JOIN pdc ON gdc.gdc_id =
40  FULL OUTER JOIN idc ON gdc.gdc_id =
41
42  full_query = client.query(sql_full)
43  df = full_query.result().to_datafra...
44  df.head()
```

```
1   ids = []
2   for index, row in df.iterrows():
3       id = set([x for x in row[[0,3,4]] if x is not None])
4       if len(id) != 1: print(len(id))
5       id = id.pop()
6       ids.append(id)
7   df['id'] = ids
8   df.head()
```
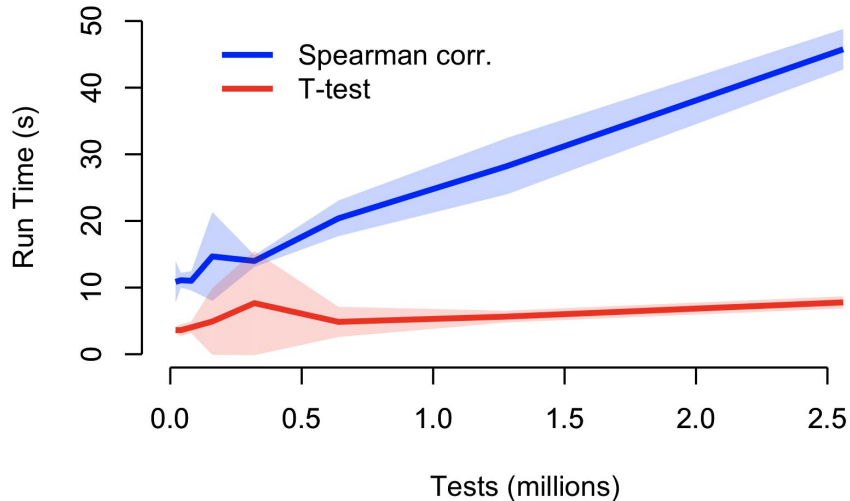
|   | gdc_id | project_id | gdc_data | pdc_id | idc_id | idc_segmentation | idc_derived | idc_quantitative | id |
|---|--------|-----------|----------|--------|--------|------------------|-------------|------------------|-----|
| 0 | TCGA-31-1951 | TCGA-OV | [Copy Number Segment, Masked Copy Number Segme... | None | None | [] | [] | [] | TCGA-31-1951 |
| 1 | TCGA-FS-A1Z7 | TCGA-SKCM | [Copy Number Segment, Masked Copy Number Segme... | None | None | [] | [] | [] | TCGA-FS-A1Z7 |
| 2 | TCGA-FS-A1ZQ | TCGA-SKCM | [Copy Number Segment, Masked Copy Number Segme... | None | None | [] | [] | [] | TCGA-FS-A1ZQ |
| 3 | TCGA-60-2716 | TCGA-LUSC | [Copy Number Segment, Masked Copy Number Segme... | None | TCGA-60-2716 | [False] | [False] | [False] | TCGA-60-2716 |
| 4 | TCGA-77-8153 | TCGA-LUSC | [Copy Number Segment, Masked Copy Number Segme... | None | None | [] | [] | [] | TCGA-77-8153 |



GDC PDC IDC Venn diagram:
- GDC only: 78626
- GDC ∩ PDC: 1052
- PDC only: 664
- Center (GDC ∩ PDC ∩ IDC): 503
- GDC ∩ IDC: 1790
- PDC ∩ IDC: 8
- IDC only: 14873

# BigQuery is a powerful statistical tool that can run hundreds of millions of tests in seconds

Testing BigQuery compute time with statistical tests
- Millions of tests in 40 seconds
- 6.6 billion correlations for $1.16

# How to run statistics inside BigQuery

Big data is hard. Statistics is even harder. Doing statistics on big data is mind-blowingly hard. We are going to provide some tools to start you on a road to making statistics on big data, if not easy, at least possible.

Ross Thomson · *Follow*

8 min read · May 18, 2023

*Collaborators*:

**Ian Mathews,** Redivis; **Boris Aguilar,** Institute for Systems Biology

https://medium.com/@jrossthomson/how-to-run-statistics-inside-bigquery-95c0c6864f23

# I'll show hands on navigation of working in the Google Cloud Console

- BigQuery Search Tool
- Google Cloud Console
  - VMs and pipelines
  - Navigating BigQuery
- Create a project
- Upload a small table
- Introductory exploration
- Notebooks section

# isb-cgc.org

ISB-CGC

# The benefits of working with ISB-CGC in the cloud

- Multiple specialized databases such as Mittleman
- Easy exploration of existing GDC and PDC data
- Access Virtual Machines and controlled data for customized pipelines
- BigQuery as a tool with scaling Excel functionality
  - Affordable storage and sharing of tabular data
  - Data exploration and quick statistics
  - Derived data from well known reference NCI datasets and annotations
  - Fast links between diverse data types
  - Advanced statistical analyses using Python, R, and Bioconductor
  - Rapidly able to expand to Machine Learning

ISB-CGC

**✉ feedback@isb-cgc.org**

**@isb-cgc** 🐦

## ISB-CGC Office Hours

Do you need assistance with getting started? Questions on merging your research with cancer data in the cloud? Or possibly help with troubleshooting?

We have **virtual Office Hours on Tuesdays and Thursdays** for any questions on ISB-CGC functionality or data that you may have. We look forward to speaking with you.

| Day of the Week | Time | Host | Link |
|---|---|---|---|
| Tuesday | 2:00pm – 3:00pm Eastern | Poojitha Gundluru | http://meet.google.com/jkg-cxke-yzs |
| Thursday | 11:00am – 12:00pm Eastern | Poojitha Gundluru | http://meet.google.com/jai-kgkg-sii |

ISB-CGC

# The ISB-CGC team

ISB

**GENERAL DYNAMICS**
Information Technology

**Bill Longabaugh**
Suzanne Paquette
Bill Clifford
Elaine Lee
Mi Tian
Lauren Hagen
Boris Aguilar
Lauren Wolfe
Ilya Shmulevich

**David Pot**
Danna Huffman
Fabian Seidl
Jacob Wilson
Poojitha Gundluru
Prema Venkatesan
Deena Bleich

ISB-CGC