

Cancer Data Analytics on the ISB-CGC platform

April 9, 2021

www.isb-cgc.org



NATIONAL CANCER INSTITUTE
Cancer Research Data Commons



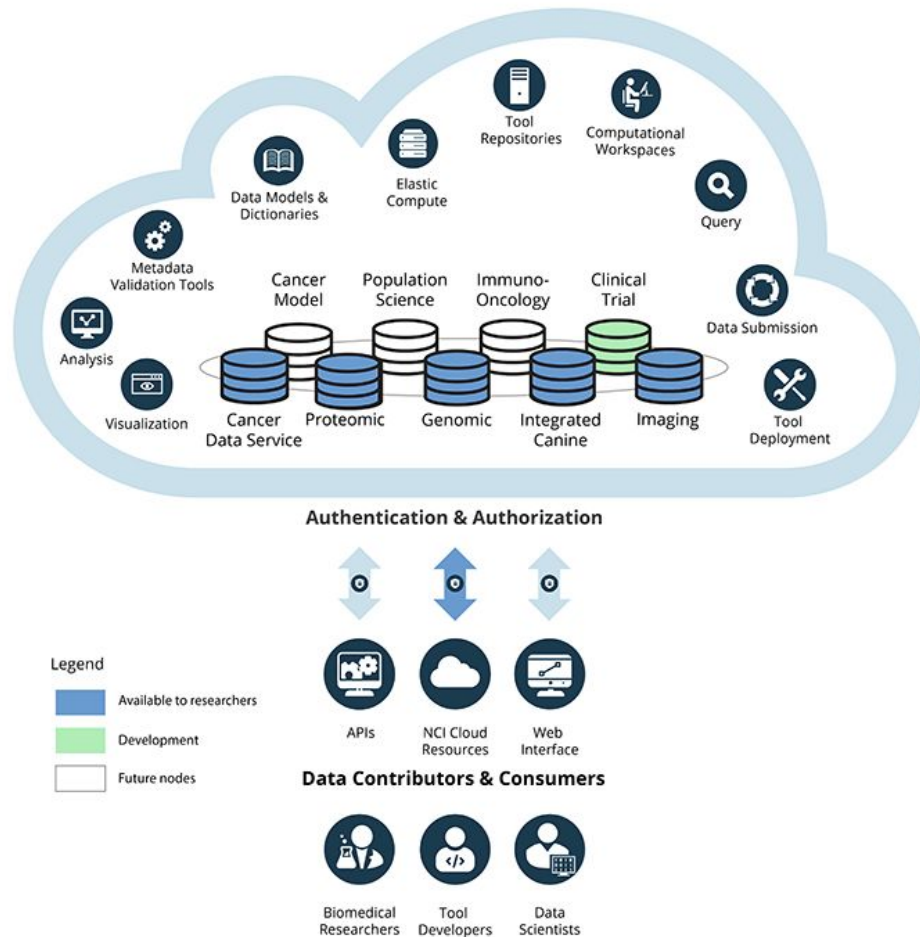
Agenda for Today's Presentation

- An introduction to the ISB-CGC role in the CRDC
- The data analysis avenues accessible via the ISB-CGC
 - Web tools
 - Google VMs
 - Google BigQuery
- Examples of analyses run with our tools

The CRDC Today

The CRDC is a **cloud-based data** science infrastructure that provides **secure access** to large, **comprehensive**, and expanding collections of **cancer research data**. Users can explore and use analytical and visualization tools for data analysis in the cloud.

The Cloud Resources have >3,000 on average active users per month from >500 institutions around the world



ISB-CGC provides Data as a Service (DaaS) solutions to the rapid growth of cancer data

Common problems of big data:

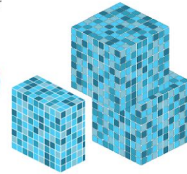
- Data discovery is onerous
- Transfer speeds become bottlenecks with scaling data size
- Availability of data can be tenuous

NATIONAL CANCER INSTITUTE THE CANCER GENOME ATLAS

TCGA BY THE NUMBERS

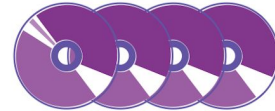
TCGA produced over

2.5
PETABYTES
of data



To put this into perspective, 1 petabyte of data is equal to

212,000
DVDs



TCGA data describes



33
DIFFERENT
TUMOR TYPES

...including

10
RARE
CANCERS

...based on paired tumor and normal tissue sets collected from



11,000
PATIENTS

...using

7

DIFFERENT
DATA TYPES



TCGA RESULTS & FINDINGS

The ISB-CGC landing page provides access to our resources

The screenshot shows the ISB-CGC landing page. At the top, there is a navigation bar with the ISB-CGC logo, links for 'Data Browsers', 'Resources', 'Documentation', 'About', and 'Help', and a 'Sign In' button. Below the navigation bar is a dark banner with the text 'Get started today! Contact us about setting up your own Google Cloud Platform Project with free cloud credits'. The main header area features the text 'A RESOURCE OF THE NCI CANCER RESEARCH DATA COMMONS', the 'ISB-CGC' logo, and the tagline 'Cancer Gateway in the Cloud'. Below this, it says 'Access, Explore and Analyze Large-Scale Cancer Data Through the Google Cloud'.

The page is divided into two main sections: 'Data Browsers' and 'Resources'.

Data Browsers

- BigQuery Table Search**: Browse BigQuery tables of metadata and molecular cancer data from the Genomic Data Commons and other sources. Jump directly to a table to perform discovery and computation via SQL. Includes 'Learn' and 'Launch' buttons.
- Data Browser**: Explore a comprehensive selection of cancer related data files in Google Cloud Storage Buckets, such as raw sequencing, cancer nucleotide variation, pathology or radiology images. Includes 'Learn' and 'Launch' buttons.
- Chromosomal Aberrations & Gene Fusions DB**: Browse the Mitelman Database of Chromosome Aberrations and Gene Fusions in Cancer which relates cytogenetic changes, in particular gene fusions, to tumor characteristics. Includes 'Learn' and 'Launch' buttons.

Resources

- Interactive Patient Cohort Exploration**: A web interface to build cohorts based on clinical demographics and molecular filters. Compare patient cohorts with various exploration tools including IGV viewer, image viewers, and analytical visualization. Includes 'Learn' and 'Launch' buttons.
- Programmatic Access**: Learn more about how to access and analyze cancer data through programmatic interfaces including Google Cloud virtual machines and APIs. Includes 'Learn' and 'Launch' buttons.
- Notebooks**: A collection of notebooks written in R and Python, to serve as both tutorials or analysis tools for a range of users; includes reproductions of Regulome Explorer functionality. Includes 'Statistical', 'Community', and 'GitHub' icons.
- Controlled Access Data**: Sign in to access controlled-access data on the Google Cloud. Authenticated users with proper dbGaP authorization only. Includes 'Learn' and 'Launch' buttons.

At the bottom of the page, there is a footer with the text: 'ISB-CGC is a component of the NCI Cancer Research Data Commons and has been funded in whole or in part with Federal funds from the National Cancer Institute, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN261201400008C and ID/IQ Agreement No. 17X146 under Contract No. HHSN261201500003I.' Below this is a copyright notice '© 2020, ISB', a 'PRIVACY POLICY' link, and social media icons for YouTube, Twitter, and GitHub.

www.isb-cgc.org

Our mission at ISB-CGC

To make NCI multi-omics cancer data as well as high-performance compute resources available via the Google Cloud Platform through multiple modes:

- Interactive web tools for cohort building and data discovery
- Easily accessible and query-able tables for multivariate data analysis
- Advanced pipeline and workflow execution on Google Cloud virtual machines

- Intentionally a “thin” layer on top of Google Cloud
 - Easy access to Google Tools
 - Direct access to NCI datasets
 - Raw GC bucket files
 - structured “BigQuery” tables for interactive data mining)

The data analysis avenues accessible via the ISB-CGC

Three entry points for exploring cancer data on ISB-CGC

ISB-CGC web tools

The screenshot shows the 'Create Cohort - Filters' page in the ISB-CGC web tool. The interface is divided into several sections: 'TCGA DATA', 'CCLE DATA', 'TARGET DATA', and 'USER DATA'. Under 'TCGA DATA', there are filters for 'PROGRAM', 'PROJECT SHORT NAME', 'DISEASE CODE', and 'VITAL STATUS'. The 'DISEASE CODE' section lists various cancer types with checkboxes, such as BRCA (3366), LUAD (1781), UCEC (1637), KIRC (1615), GBM (1573), and HNSC (1573). The 'VITAL STATUS' section has a 'Check All / Uncheck All' option. On the right side, there are 'Selected Filters' and 'Program Details' sections. The 'Program Details' section shows 'Total Number of Cases: 11353' and 'Total Number of Samples: 33460'. Below this, there are five 'Clinical Features' charts: 'Disease Code', 'Vital Status', 'Sample Type', 'Tumor Tissue Site', and 'Gender', each represented by a colorful stacked bar chart.

Google VMs

The screenshot shows a terminal window on a Google Cloud VM. The terminal output includes the following commands and their results:

```
https://ish.cloud.google.com/projects/ish-...
Linux x86_64
The programs included with the Debian GNU/Linux system have been configured to support the exact distribution terms for each individual file in /usr/share/doc/.
Debian GNU/Linux comes with ABSOLUTELY no warranty.
Last login: Thu Feb 6 22:31:46 2020
Copying gs://genomics-public-data/...
Operation completed over 1 objects/...
~/hg38_fasta$ ~/STAR/bin/Linux_x86_64/STAR --runThreadN 4 --runMode genomeGenerate --genomeDir ~/hg38_fasta/ --genomeFastaFiles GRCh38.
Feb 06 22:33:26 .... started STAR run
Feb 06 22:33:26 ... starting to generate Genome files
Feb 06 22:34:42 ... starting to sort Suffix Array. This may take a long time...
Feb 06 22:35:00 ... sorting Suffix Array chunks and saving them to disk...
```

BigQuery

The screenshot shows the BigQuery interface. On the left, there is a sidebar with 'Query history', 'Saved queries', 'Job history', 'Transfers', 'Scheduled queries', and 'BI Engine'. The main area shows a query that has been executed. The query results are displayed in a table with the following data:

Row	GTEX_tissueType	sample_barcode	TCGA_project	corr
1	Liver	TCGA-DD-A39V-11A	TCGA-LIHC	0.9213023777251851
2	Liver	TCGA-DD-A39Z-11A	TCGA-LIHC	0.9189148155140473
3	Liver	TCGA-DD-A3A1-11A	TCGA-LIHC	0.917682740669065

Some example use-cases of each major entry point

Interactive web-based exploration

- Select a subset of TCGA samples based on clinical or molecular characteristics
- Compare one cohort to another
- Upload a small private dataset to analyze in conjunction with TCGA data
- *etc...*

Interactive cancer data exploration and analysis

- Interactive data exploration in BigQuery
- Use R or Python to perform custom multivariate analyses
- Directly run statistical tests on data in BigQuery using custom functions
- *etc...*

Direct Command line Access to Google virtual machines

- Test new algorithm on hundreds or thousands of BAM or FASTQ files
- Run novel image segmentation method across whole-slide images
- *etc...*

Three entry points for exploring cancer data on ISB-CGC

ISB-CGC web tools

Dashboard WORKBOOKS PROGRAMS ANALYSES GENES & miRNAs VARIABLES COHORTS

Your Dashboard > Cohorts >

Create Cohort - Filters

Save As New Cohort

TCGA DATA CCLE DATA TARGET DATA USER DATA

CASE DATA MOLEC.

PROGRAM

PROJECT SHORT NAME

DISEASE CODE

- BRCA (3366)
- LUAD (1781)
- UCEC (1637)
- KIRC (1615)
- GBM (1573)
- HNSC (1573)
- 27 more

Check All / Uncheck All

VITAL STATUS

Selected Filters Clear All

Program Details

Total Number of Cases: 11353 Total Number of Samples: 33460

Clinical Features

Disease Code Vital Status Sample Type Tumor Tissue Site Gender

Google VMs

```
https://ssh.cloud.google.com/projects/...
Connected, host fingerprint: ssh-rsa
Linux 4.9.0-11-x86_64
The programs included with the Deb
Debian GNU/Linux comes with ABSOLU
permitted by applicable law.
last login: Thu Feb 6 22:31:46 20
Copying gs://genomics-public-data/
11 files|| 2.9 GiB / 2.9 GiB|
Operation completed over 1 objects
Primary assembly genome.fa --sdOverhang
Feb 06 22:33:26 .... started STAR run
Feb 06 22:33:26 ... starting to generate Genome files
Feb 06 22:34:42 ... starting to sort Suffix Array. This may take a long time...
Feb 06 22:35:00 ... sorting Suffix Array chunks and saving them to disk...
```

BigQuery

BigQuery

Query history

Saved queries

Job history

Transfers

Scheduled queries

BI Engine

Resources + ADD DATA

Search for your tables and data...

cgc-05-0050

isb-cgc

Run Save query Save view Schedule query More

Query complete (18.8 sec elapsed, 12.7 GB processed)

Job information Results JSON Execution details

Row	GTEx_tissueType	sample_barcode	TCGA_project	corr
1	Liver	TCGA-DD-A39V-11A	TCGA-LIHC	0.9213023777251851
2	Liver	TCGA-DD-A39Z-11A	TCGA-LIHC	0.9189148155140473
3	Liver	TCGA-DD-A3A1-11A	TCGA-LIHC	0.917682740669065

Interactive cohort-building using the ISB-CGC

ISB-CGC Data Browsers Resources Documentation About Help

Create Cohort - Filters

TCGA DATA CCLC DATA TARGET DATA

Hide attributes with 0 cases

CASE DATA MOLEC.

- PROGRAM
- PROJECT SHORT NAME
 - TCGA-BRCA (1,101 cases)
 - TCGA-GBM (617 cases)
 - TCGA-OV (608 cases)
 - TCGA-LUAD (586 cases)
 - TCGA-UCEC (560 cases)
 - TCGA-KIRC (537 cases)
- DISEASE CODE
- VITAL STATUS
- GENDER
- AGE AT DIAGNOSIS
- SAMPLE TYPE

28 more Check All / Uncheck All

ISB-CGC Data Browsers Resources Documentation About Help Sign In

Create Cohort - Filters

Log In To Save New Cohort

TCGA DATA CCLC DATA TARGET DATA BEATAML1.0 DATA USER DATA

CASE DATA MOLEC.

Selected Filters [Clear All](#)

Program Details

Total Number of Cases: 11,315 Total Number of Samples: 23,797

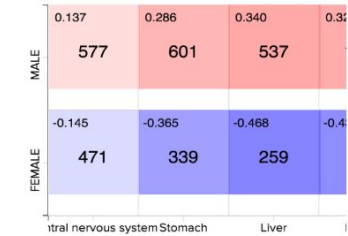
Clinical Features

Disease Code Vital Status Sample Type Tumor Tissue Site Gender

Show More

- PROGRAM
- PROJECT SHORT NAME
- DISEASE CODE
- VITAL STATUS
- GENDER
- AGE AT DIAGNOSIS
- SAMPLE TYPE
- TUMOR TISSUE SITE
- HISTOLOGICAL TYPE
- PATHOLOGIC STAGE
- NEOPLASM HISTOLOGIC GRADE
- BMI
- HPV STATUS
- RESIDUAL TUMOR
- RACE
- ETHNICITY

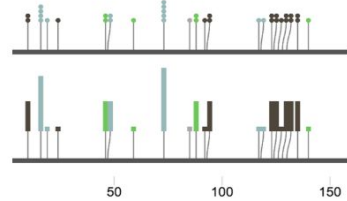
ISB-CGC has several interactive apps to generate quick plots



Cubby Hole Plot

Used to plot two categorical features. Boxes are colored by their related p-values.

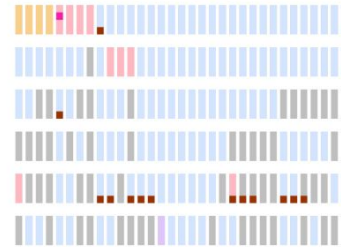
[Start A New Workbook Cubby Hole Plot](#)



SeqPeek

This visualization shows where somatic mutations have been observed on a linear representation of a specific protein. Each horizontal strip represents the protein, with data from different tumor types (aka cohorts or studies) shown stacked one on top of the other.

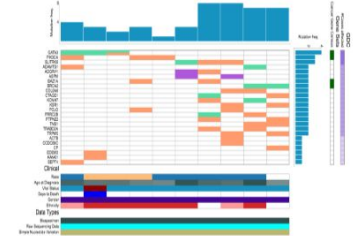
[Start A New Workbook SeqPeek](#)



OncoPrint

Used to plot multiple genomic alteration (somatic mutation) events across a set of samples by heatmap. OncoPrint is developed and provided by cBioPortal.

[Start A New Workbook OncoPrint](#)



OncoGrid

Used to view multiple genomic alteration (somatic mutation) events, clinical data, available files across a set of cases by interactive heatmap. OncoGrid library is developed at Ontario Institute for Cancer Research (OICR).

[Start A New Workbook OncoGrid](#)

File browser to browse the data hosted in cloud buckets

Integrated genome viewer
(view read pile-ups)

caMicroscope
(view histology)

OHIF
(view radiology)

Your Dash
File Browser

All Files IGV Pathology Images Pathology Reports Radiology Images

Build
HG19

- ▶ CASE
- ▶ DATA TYPE
- ▶ DATA CATEGORY
- ▶ EXPERIMENTAL STRATEGY
- ▶ DATA FORMAT
- ▶ PLATFORM
- ▶ DISEASE CODE

File Listing

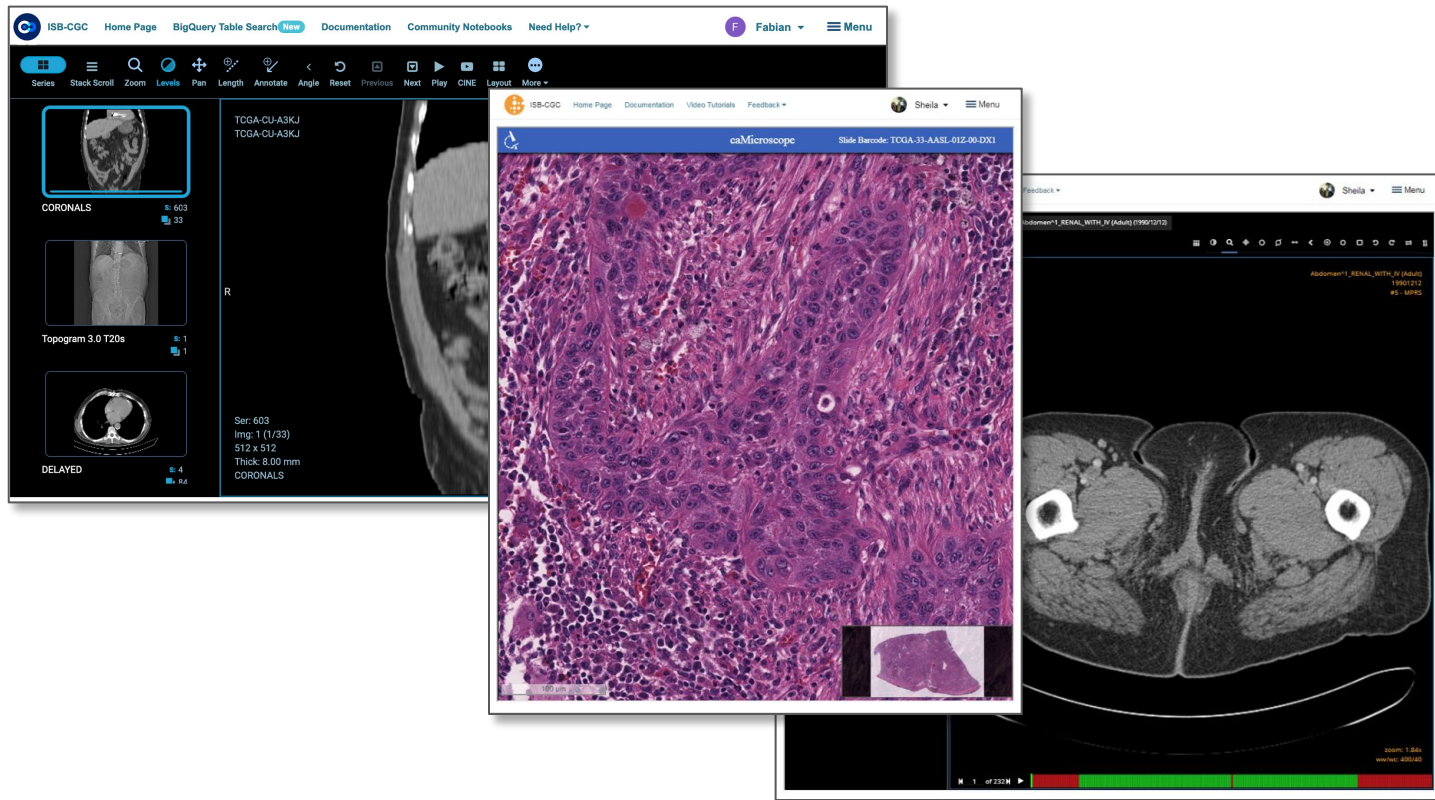
Showing 1 to 25 of 39692 entries

Show 25 entries Page Go Previous 1 2 3 ... 1588 Next

Choose Columns to Display ▾

Program	Case Barcode	File Name	Disease Code	Exp. Strategy	Platform	Data Category	Data Type	Data Format	File Size
TCGA	TCGA-OL-A660	SWEDE_p_TCGAb322_2... [GDC ID: 0686da7c-d103-...	BRCA	Genotyping array	Affymetrix SNP Array 6.0	Simple nucleotide variation	Genotypes	TXT	20.9 MB
TCGA	TCGA-OL-A660	SWEDE_p_TCGAb322_2... [GDC ID: 4bd19f77-9aa7-4...	BRCA	Genotyping array	Affymetrix SNP Array 6.0	Simple nucleotide variation	Genotypes	TXT	20.9 MB
TCGA	TCGA-OL-A660	UNCID_2171596.c7f5714... [GDC ID: b677ea35-d758-...	BRCA	RNA-Seq	Illumina HiSeq	Raw sequencing data	Aligned reads	BAM	7.8 GB
TCGA	TCGA-OL-A660	c61047b5e4ae38963735f... [GDC ID: 0a6db03e-748a-...	BRCA	WXS	Illumina HiSeq	Raw sequencing data	Aligned reads	BAM	4.9 GB
TCGA	TCGA-OL-A660	256cd674e76be0f163766b... [GDC ID: 72a31a7e-99df-4...	BRCA	WXS	Illumina HiSeq	Raw sequencing	Aligned reads	BAM	7.2 GB

Interactive image viewers allows for browsing of cancer imaging data



Mitelman database is also hosted by ISB-CGC

The screenshot displays the Mitelman Database website interface. On the left, a dark blue sidebar contains navigation links: Home, Search, Cases Cytogenetics, Gene Fusions, Clinical Associations, Recurrent Chromosome Aberrations, References, User Guide, About, and Contact. The main content area features a large, light-colored background with a pattern of white, interconnected, curved lines. Centered on this background is the text "Mitelman Database Chromosome Aberrations and Gene Fusions in Cancer" with a small icon of a chromosome. Below this, it states "This site has been funded by: National Cancer Institute, Swedish Cancer Society, Swedish Childhood Cancer Foundation". At the bottom, it mentions "This website is built and maintained by the ISB-CGC cloud project." To the right, a smaller inset shows a different view of the website with a search bar, navigation tabs (REFS CORNER, PEOPLE AND EVENTS, CONTACT US), a microscopic image of cells with green and red fluorescence, and a "p53 in the Clinics" book cover. A news section at the bottom of the inset contains a notice about the TP53 Database project transfer.

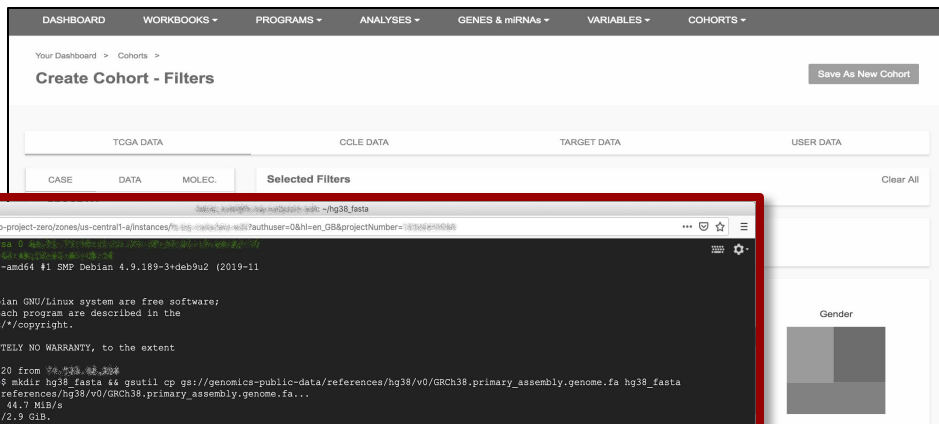
<https://mitelmandatabase.isb-cgc.org/>

NEWS
28/09/2020 - The IARC TP53 Database project will be transferred to another institution. The current version of the database (R20, July 2019) will remain accessible until February 2021 to allow users to download all available data. IARC thanks all contributors and data providers who made the success of this resource.

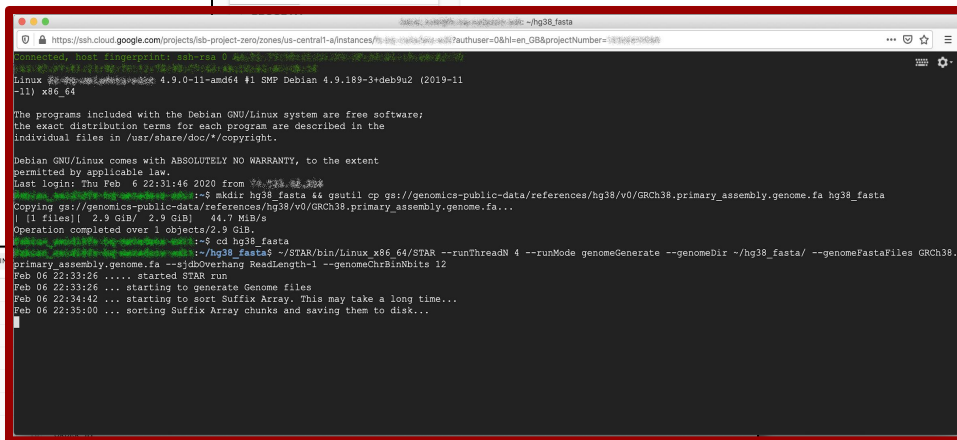
Also ongoing work towards hosting the TP53 database soon

Three entry points for exploring cancer data on ISB-CGC

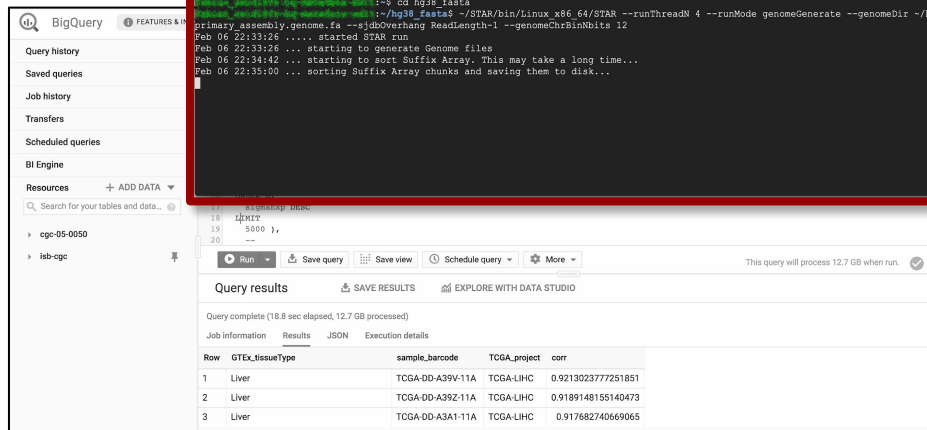
ISB-CGC web tools



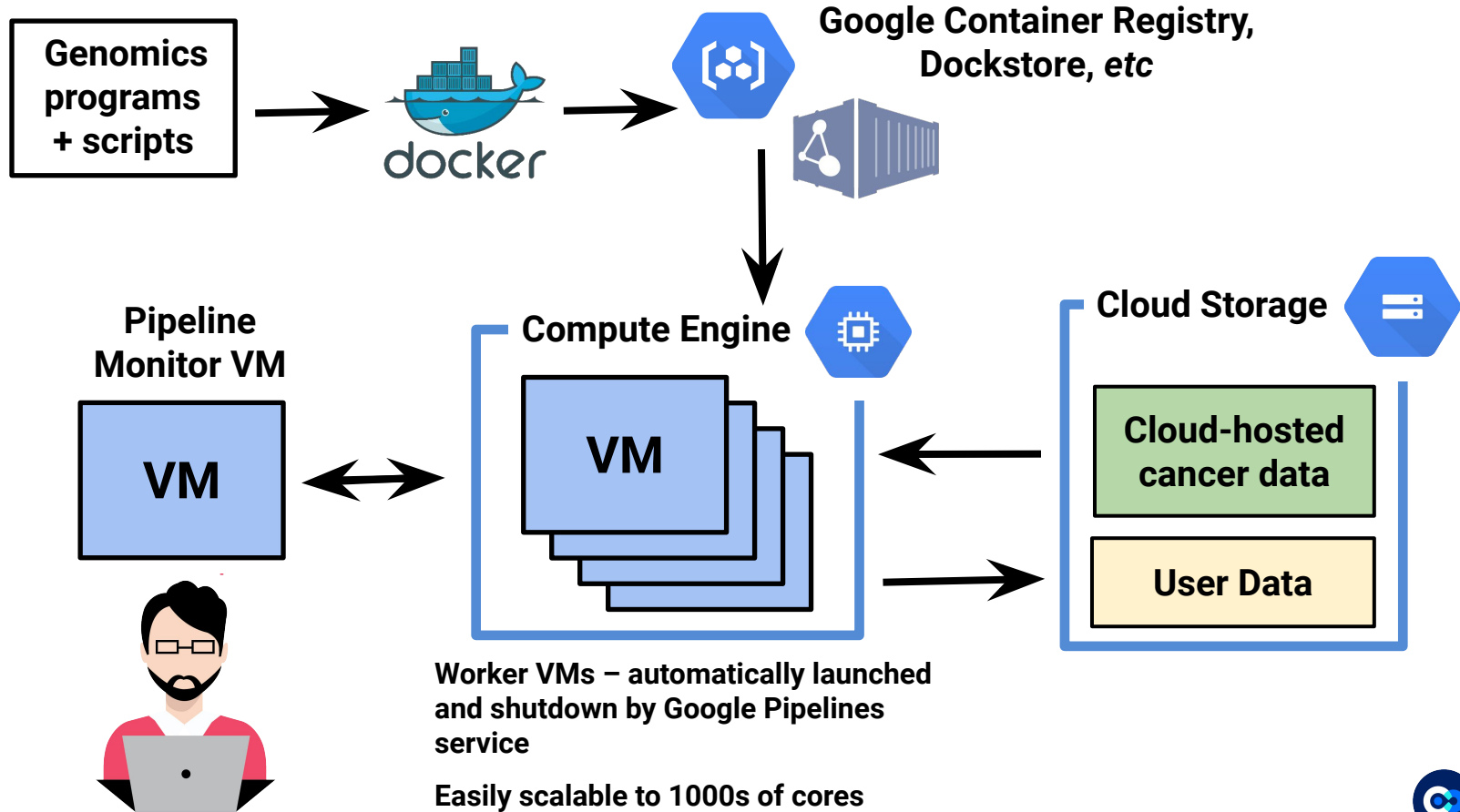
Google VMs



BigQuery



Advanced workflow execution on -omics data enabled by ISB-CGC



Google Cloud Platform Free Tier lets you perform small scale computes with no cost!

<p>COMPUTE</p> <p>Cloud Run</p> <p>2 million</p> <p>Requests per month</p> <p>A fully managed environment to run stateless containers.</p> <p>▼</p>	<p>DATABASE</p> <p>Firestore</p> <p>1 GB</p> <p>Storage</p> <p>Scalable NoSQL, document database.</p> <p>▼</p>	<p>COMPUTE</p> <p>Compute Engine</p> <p>1</p> <p>F1-micro instance per month</p> <p>Scalable, high-performance virtual machines.</p> <p>▼</p>	<p>COMPUTE</p> <p>Compute Engine</p> <p>1</p> <p>F1-micro instance per month</p> <p>Scalable, high-performance virtual machines.</p> <hr/> <p>1 f1-micro instance per month (US regions only—excluding Northern Virginia [us-east4])</p> <hr/> <p>30 GB-months HDD</p> <hr/> <p>5 GB-months snapshot in select regions</p> <hr/> <p>1 GB network egress from North America to all region destinations per month (excluding China and Australia)</p> <p>▲</p>
<p>STORAGE</p> <p>Cloud Storage</p> <p>5 GB</p> <p>Months regional storage</p> <p>Best-in-class performance, reliability, and pricing for all your storage needs.</p> <p>▼</p>	<p>DATA ANALYTICS</p> <p>Pub/Sub</p> <p>10 GB</p> <p>Messages per month</p> <p>A global service for real-time and reliable messaging and streaming data.</p> <p>▼</p>	<p>COMPUTE</p> <p>Cloud Functions</p> <p>2 million</p> <p>Invocations per month</p> <p>A serverless environment to build and connect cloud services with code.</p> <p>▼</p>	
<p>COMPUTE</p> <p>Google Kubernetes Engine</p> <p>Clusters</p> <p>All size clusters</p> <p>One-click container orchestration via Kubernetes clusters, managed by Google.</p> <p>▼</p>	<p>COMPUTE</p> <p>App Engine</p> <p>28</p> <p>Instance hours per day</p> <p>Platform for building scalable web applications and mobile back ends.</p> <p>▼</p>	<p>MANAGEMENT TOOLS</p> <p>Stackdriver</p> <p>50 GB</p> <p>Logs with 30-day retention</p> <p>Monitoring, logging, and diagnostics applications on Google Cloud and AI</p> <p>▼</p>	
<p>DATA ANALYTICS</p> <p>BigQuery</p> <p>1 TB</p> <p>Queries per month</p> <p>Fully managed, petabyte scale, analytics data warehouse.</p> <p>▼</p>	<p>AI AND MACHINE LEARNING</p> <p>Vision AI</p> <p>1,000</p> <p>Units per month</p> <p>Label detection, OCR, facial detection and more.</p> <p>▼</p>	<p>AI AND MACHINE LEARNING</p> <p>Speech-to-Text</p> <p>60</p> <p>Minutes per month</p> <p>Speech-to-text transcription — the same powers Google's own products.</p> <p>▼</p>	

We have several workflow and pipeline tutorials at ISB-CGC

Tutorials and sample workflows designed to introduce and enable users on how to run workflows in CWL, Nextflow, Snakemake and WDL on GCP

https://isb-cgc.appspot.com/programmatic_access/

The screenshot displays the ISB-CGC website interface. At the top, there is a navigation bar with links for 'Data Browsers', 'Resources', 'Documentation', 'About', and 'Help'. The main content area is titled 'Pipelines and APIs' and includes a sub-header 'Getting Started' with the text 'Learn more about the different ways to access data programmatically'. Below this, there are three main sections: 'Tutorials for Workflow on Google Cloud' which features a workflow diagram with nodes like 'samtools_state_tool' and 'grep'; 'Comparison of Workflow Languages' which compares {wdl}, nextflow, and Snakemake; and 'ISB-CGC API' which lists various API endpoints such as 'GET /files/track/registration' and 'POST /access'. At the bottom, there is a 'Workflow Examples' table and a footer with contact information.

Name	Requirements	Language	Tools Used	Input	Output	Try this out
RNA-Seq	GCP, Miniconda	Snakemake	Hisat2, Samtools, Stringtie	Fasta, Fastq, GTF	Sam, Bam, GFF, TSV	Docs GitHub
RNA-Seq	GCP, CWLtool	CWL	Hisat2, Samtools, Stringtie	Fasta, Fastq, GTF	Sam, Bam, GFF, TSV	Docs GitHub
RNA-Seq	GCP, Nextflow	Nextflow	Hisat2, Samtools, Stringtie	Fasta, Fastq, GTF	Sam, Bam, GFF, TSV	Docs GitHub
GC-gather	GCP, Nextflow	Nextflow	Samtools	Bam	Text file with GC content	Docs GitHub
GC-gather	GCP, Miniconda	Snakemake	Samtools	Bam	Text file with GC content	Docs GitHub
GC-gather	GCP, CWLtool	CWL	Samtools	Bam	Text file with GC content	Docs GitHub
GC-gather	GCP, Cromwell	WDL	Samtools	Bam	Text file with GC content	Docs GitHub
Blast pipeline	GCP, Nextflow	Nextflow	Blast, Python	Fasta	Text file with contigs of interest	Docs GitHub
Blast pipeline	GCP, CWLtool	CWL	Blast, Python	Fasta	Text file with contigs of interest	Docs GitHub

Example RNASeq quantification workflow using a Docker image

```
STAR_alignment_run.sh — Desktop Add License
1 wget https://api.gdc.cancer.gov/data/25aa497c-e615-4cb7-8751-71f744f9691f \
2   $path_to_reference/gencode.v22.annotation.gtf.gz
3 gunzip gencode.v22.annotation.gtf.gz
4
5 path_to_reference=/home/fseidl/genome
6 mkdir $path_to_reference/star_index_oh75
7
8 sudo docker run --rm -v $path_to_reference:/data -t broadinstitute/gtex_rnaseq:v8 \
9   /bin/bash -c "STAR \
10  --runMode genomeGenerate \
11  --genomeDir /data/star_index_oh75 \
12  --genomeFastaFiles /data/GRCh38.d1.vd1.fa \
13  --sjdbGTFfile /data/gencode.v22.annotation.gtf \
14  --sjdbOverhang 75 \
15  --runThreadN 4"
16
17 sudo docker run --rm -v $path_to_reference:/data -t broadinstitute/gtex_rnaseq:v8 \
18   /bin/bash -c "rsem-prepare-reference \
19   /data/GRCh38.d1.vd1.fa \
20   /data/rsem_reference \
21   --gtf /data/gencode.v22.annotation.gtf \
22   --num-threads 4"
23
24 path_to_data=/home/fseidl/bams
25 input_bam=HG00182.mapped.ILLUMINA.bwa.FIN.low_coverage.20120522.bam
26 sample_id=HG00182
27 mkdir $path_to_data
28
29 gsutil cp gs://genomics-public-data/1000-genomes/bam/$input_bam $path_to_data
30
31 sudo docker run --rm -v $path_to_data:/data -t broadinstitute/gtex_rnaseq \
32   /bin/bash -c "src/run_SamToFastq.py /data/$input_bam -p $sample_id -o /data"
33
34 # STAR alignment
35 sudo docker run --rm -v $path_to_data:/data -v $path_to_reference:/genome -t broadinstitute/gtex_rnaseq:v8 \
36   /bin/bash -c "src/run_STAR.py \
37   /genome/star_index_oh75 \
38   /data/${sample_id}_1.fastq.gz \
39   /data/${sample_id}_2.fastq.gz \
40   ${sample_id} \
41   --threads 4 \
42   --output_dir /tmp/star_out && mv /tmp/star_out /data/star_out"
```



Row	HGNC_gene_symbol	gene_id	normalized_count
1	RELN	5649	5092.9882
2	DDX49	54555	1233.9143
3	OCSTAMP	128506	1.2346
4	RUFY4	285180	0.9481
5	SLC6A4	6532	1.0684
6	NOXA1	10811	14.8593
7	TAGLN2	8407	3785.6545
8	ZNF484	83744	40.8805
9	RNF217	154214	167.8584
10	RHOC	389	2670.3879
11	RNF219	79596	136.9329
12	MANEA	79694	797.0085
13	PALB2	79728	174.2287
14	MRPL35	51318	803.0303
15	IQSEC1	9922	2266.2474
16	FAM57B	83723	1.3774
17	CFLAR	8837	1108.9744
18	MAML2	84441	79.8898

Easily run batches of jobs in the cloud using dsub

Overview

`dsub` is a command-line tool that makes it easy to submit and run batch scripts in the cloud.

The `dsub` user experience is modeled after traditional high-performance computing job schedulers like Grid Engine and Slurm. You write a script and then submit it to a job scheduler from a shell prompt on your local machine.

Today `dsub` supports Google Cloud as the backend batch job runner, along with a local provider for development and testing. With help from the community, we'd like to add other backends, such as a Grid Engine, Slurm, Amazon Batch, and Azure Batch.

Getting started

You can install `dsub` from [PyPI](#), or you can clone and install from [github](#).

```
dsub \  
  --name kallisto_quant \  
  --project ${GS_PROJECT} \  
  --zones 'us-*' \  
  --image "nareshr/kallisto:v0.43" \  
  --input "KALIDX=${GS_BUCKET}/Homo_sapiens.GRCh37.cdna.all.kal.idx" \  
  --input "FASTQ=${GS_BUCKET}/All_CCLE_customDB.fastq" \  
  --output-recursive "KALOUT=${GS_BUCKET}/output" \  
  --logging ${GS_BUCKET}/log \  
  --min-cores 8 \  
  --command 'kallisto quant -i ${KALIDX} -o ${KALOUT} -b 100 --single -l 180 -s 20 -t 8 ${FASTQ}' \  
  --wait
```

Google Cloud Life Sciences

Features

Cost-optimized compute

Google Cloud's Healthcare and Life Sciences team has optimized the most popular methods—like [GATK](#), [DeepVariant](#), and [Sentieon](#)—to run on GCP.

Flexible machine sizes

Take advantage of [Compute Engine](#), our infrastructure as a service (IaaS), to run large-scale workloads on virtual machines and pay only for what you use.

Built for batch processing

[Preemptible VMs](#) for affordable batch processing on fault-tolerant workloads to save you time and money.

Fully integrated with GCP

Experience the power of Google Cloud's infrastructure with fast virtual machines, scalable storage, serverless data warehouses, and fully managed databases with GCP integration to tools like [Cloud Spanner](#) and [BigQuery](#).

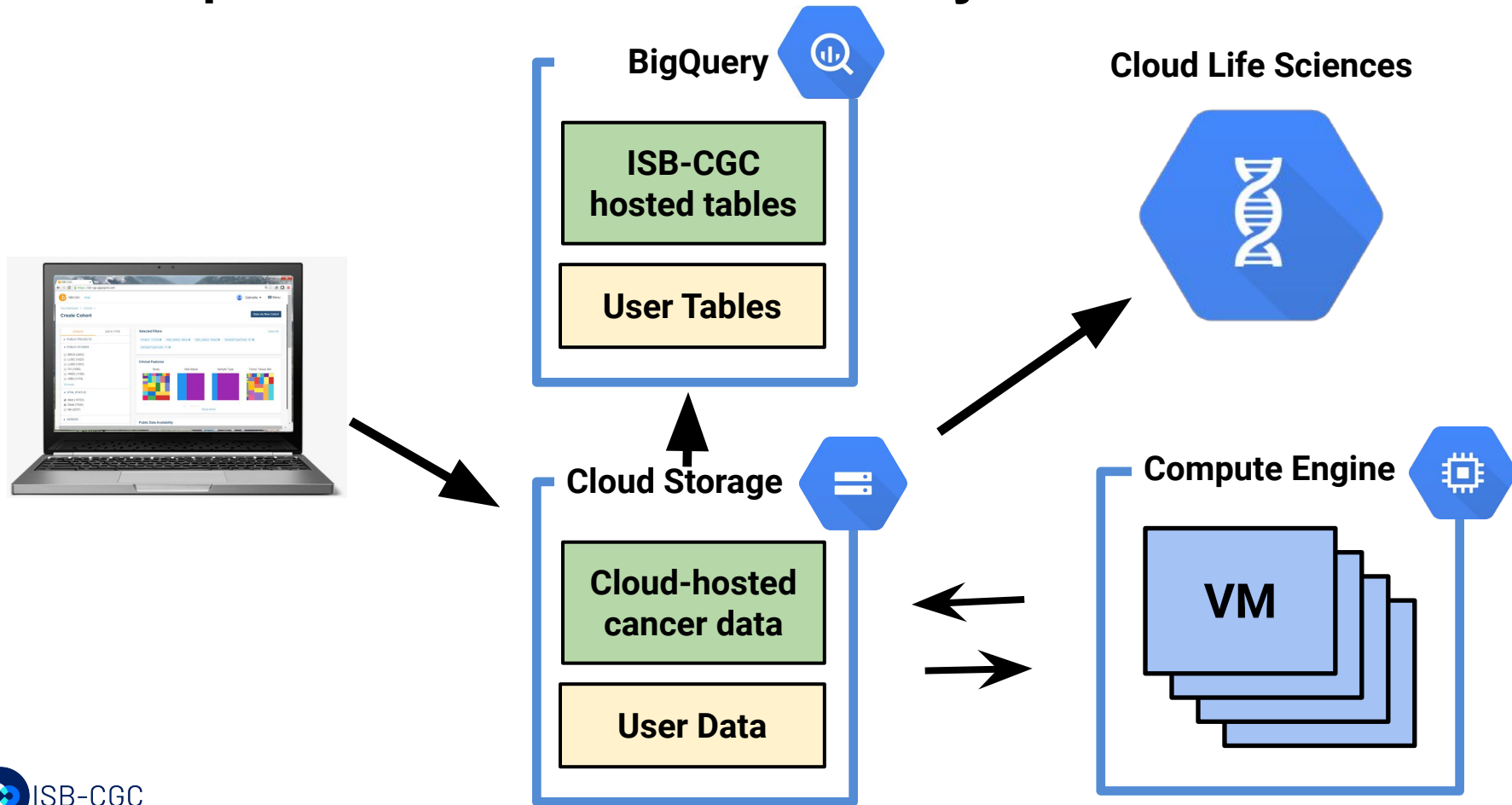
Open and interoperable

Use the tools and workflows you already know and enjoy support for open industry standards like [Global Alliance for Genomics and Health](#).

Ready for AI / ML

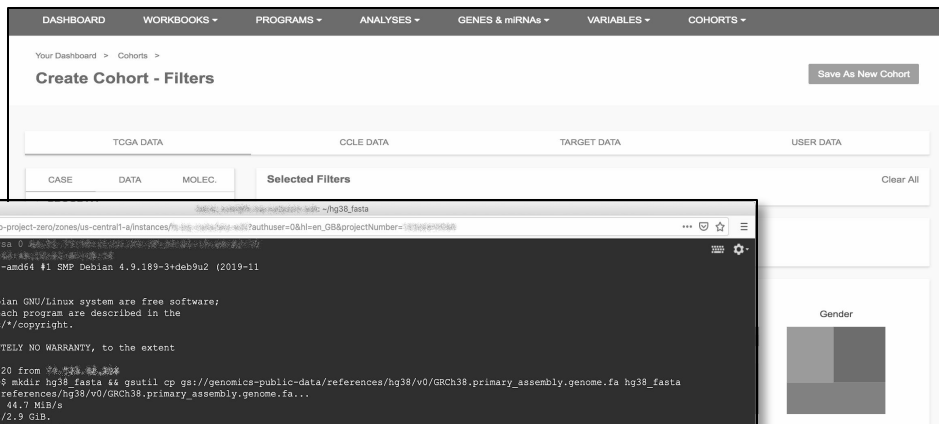
Bring your data closer to [public datasets](#) and advanced analytics that come along with GCP.

Example end-to-end workflow analysis on ISB-CGC

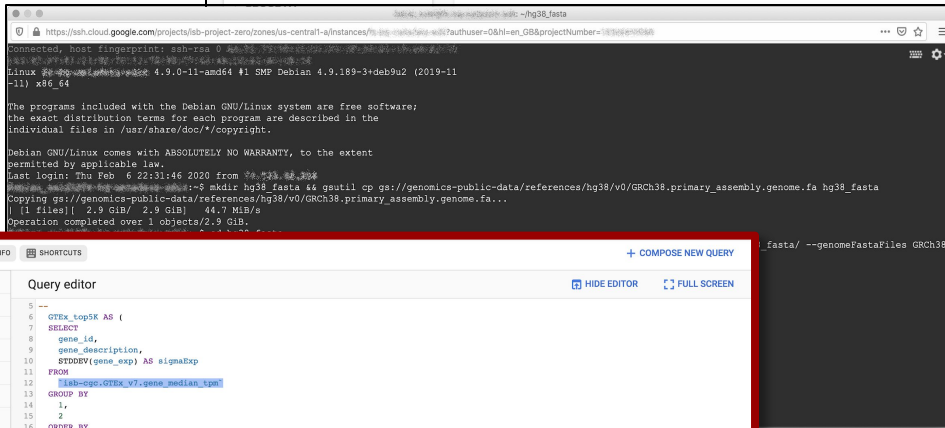


Three entry points for exploring cancer data on ISB-CGC

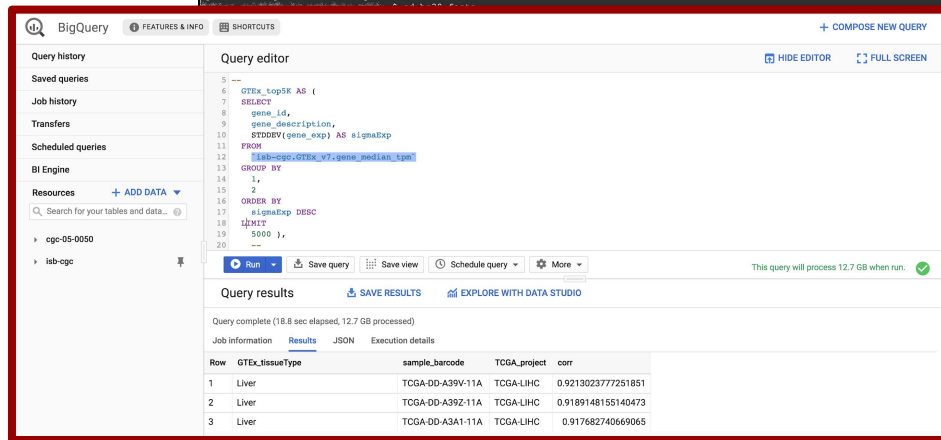
ISB-CGC web tools



Google VMs

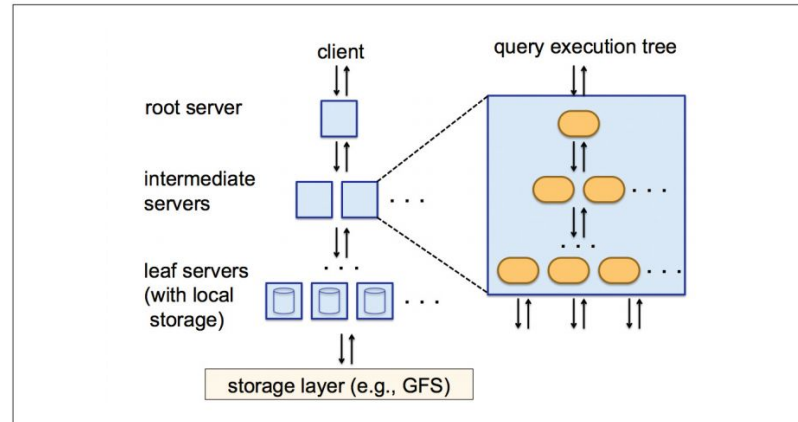


BigQuery



Attributes of Google BigQuery that make it ideal for cancer data analytics

- Columnar database ideal for storing tabular data (RNA sequencing etc.)
- Query speed is automatically scaled by multiprocessing
- Powerful SQL language interface, including user defined functions
- Can join tables based on shared variables
- ISB-CGC combines data of a similar type into single BigQuery tables
 - For example: ~150 individual MAF files were combined to generate a single table



Tree architecture of Dremel

You can think of individual BigQuery tables very simply as a tabular database

The screenshot displays the Google Cloud Platform BigQuery interface. The main content area shows a table named 'RNaseq_hg38_gdc_current'. The table has the following columns: primary_site, sample_barcode, aliquot_barcode, gene_name, gene_type, Ensembl_gene_id, Ensembl_gene_id_v, HTSeq_Counts, HTSeq_FPKM, and HTSeq_FPKM_UQ. The table contains 20 rows of data, each representing a different biological sample and its associated gene information.

primary_site	sample_barcode	aliquot_barcode	gene_name	gene_type	Ensembl_gene_id	Ensembl_gene_id_v	HTSeq_Counts	HTSeq_FPKM	HTSeq_FPKM_UQ
Bones, joints and articular cartilage of limbs	TARGET-40-0A418U-01A	TARGET-40-0A418U-01A-01R	ACD15489.13	sense_overlapping	ENSG00000270894	ENSG00000270894.1	39	0.735	19169.494
Bones, joints and articular cartilage of limbs	TARGET-40-PARBGW-01A	TARGET-40-PARBGW-01A-01R	GSTT2	polymorphic_pseudogene	ENSG00000099984	ENSG00000099984.9	18	0.3702	12565.4779
Bones, joints and articular cartilage of limbs	TARGET-40-0A410Q-01A	TARGET-40-0A410Q-01A-01R	TMED11P	unitary_pseudogene	ENSG00000215367	ENSG00000215367.9	16	0.2196	6770.8618
Bones, joints and articular cartilage of limbs	TARGET-40-PAUJML-01A	TARGET-40-PAUJML-01A-01R	KB-1836B5.1	sense_overlapping	ENSG00000260949	ENSG00000260949.1	19	0.9695	23251.0613
Bones, joints and articular cartilage of limbs	TARGET-40-PATEEM-01A	TARGET-40-PATEEM-01A-01R	IGHG4	IG_C_gene	ENSG00000211892	ENSG00000211892.3	454	5.601	121956.6821
Bones, joints and articular cartilage of limbs	TARGET-40-0A415B-01A	TARGET-40-0A415B-01A-01R	IGHV5-51	IG_V_gene	ENSG00000211966	ENSG00000211966.2	55	2.7753	59225.7578
Bones, joints and articular cartilage of limbs	TARGET-40-PASUJH-01A	TARGET-40-PASUJH-01A-01R	RP11-999E24.3	sense_overlapping	ENSG00000259969	ENSG00000259969.1	33	0.8218	21886.2353
Bones, joints and articular cartilage of limbs	TARGET-40-0A414E-01A	TARGET-40-0A414E-01A-01R	CTD-2201E18.3	sense_overlapping	ENSG00000177738	ENSG00000177738.3	74	0.9873	25266.5967
Bones, joints and articular cartilage of limbs	TARGET-40-PALFYN-01A	TARGET-40-PALFYN-01A-01R	RP11-517B11.7	sense_overlapping	ENSG00000261167	ENSG00000261167.1	24	0.4396	12773.4789
Bones, joints and articular cartilage of limbs	TARGET-40-PAPWVC-01A	TARGET-40-PAPWVC-01A-01R	RP11-166D19.1	sense_overlapping	ENSG00000255248	ENSG00000255248.5	1245	2.7999	70773.0387
Bones, joints and articular cartilage of limbs	TARGET-40-PAUTWB-01A	TARGET-40-PAUTWB-01A-01R	RP5-104218.7	sense_overlapping	ENSG00000261662	ENSG00000261662.1	76	1.6702	41929.5019
Bones, joints and articular cartilage of limbs	TARGET-40-PASUJH-01A	TARGET-40-PASUJH-01A-01R	RP11-389C8.2	sense_overlapping	ENSG00000261269	ENSG00000261269.1	68	0.6685	17803.6338
Bones, joints and articular cartilage of limbs	TARGET-40-PALWXX-01A	TARGET-40-PALWXX-01A-01R	MIR31HG	sense_overlapping	ENSG00000171889	ENSG00000171889.3	43	0.9736	25755.5202
Bones, joints and articular cartilage of other and unspecified sites	TARGET-40-0A410S-01A	TARGET-40-0A410S-01A-01R	TRBV19	TR_V_gene	ENSG00000211746	ENSG00000211746.3	17	0.3509	10221.4077
Bones, joints and articular cartilage of limbs	TARGET-40-PALZGU-01A	TARGET-40-PALZGU-01A-01R	CTB-109A12.1	sense_overlapping	ENSG00000251405	ENSG00000251405.2	16	0.1346	3071.1182
Bones, joints and articular cartilage of limbs	TARGET-40-PAMTCM-01A	TARGET-40-PAMTCM-01A-01R	CCNYL2	translated_unprocessed_pseudogene	ENSG00000182632	ENSG00000182632.13	31	0.0827	2387.338
Bones, joints and articular cartilage of limbs	TARGET-40-PASSLM-01A	TARGET-40-PASSLM-01A-01R	RP5-1126H10.2	sprime_overlapping_ncrna	ENSG00000272084	ENSG00000272084.1	24	0.1731	3943.9126
Bones, joints and articular cartilage of limbs	TARGET-40-PALZGU-01A	TARGET-40-PALZGU-01A-01R	TRAC	TR_C_gene	ENSG00000277734	ENSG00000277734.3	375	4.2533	97033.1155
Bones, joints and articular cartilage of limbs	TARGET-40-PAMYYJ-01A	TARGET-40-PAMYYJ-01A-01R	RP11-283I3.6	sense_overlapping	ENSG00000261799	ENSG00000261799.1	804	4.1921	115864.6689
Bones, joints and articular cartilage of limbs	TARGET-40-PALHRL-01A	TARGET-40-PALHRL-01A-01R	CDDC162P	unitary_pseudogene	ENSG00000203799	ENSG00000203799.9	15	0.1023	3259.2941
Bones, joints and articular cartilage of limbs	TARGET-40-PATMIF-01A	TARGET-40-PATMIF-01A-01R	RP13-122B23.8	sense_overlapping	ENSG00000260996	ENSG00000260996.1	20	0.1415	3684.5655
Bones, joints and articular cartilage of limbs	TARGET-40-PASEFS-01A	TARGET-40-PASEFS-01A-01R	RP11-13J10.1	sense_overlapping	ENSG00000269707	ENSG00000269707.1	28	1.1143	23450.6697
Bones, joints and articular cartilage of limbs	TARGET-40-PASYUK-01A	TARGET-40-PASYUK-01A-01R	CTD-2201E18.3	sense_overlapping	ENSG00000177738	ENSG00000177738.3	162	2.0041	48868.0707

The unique benefit is that you can mine these data quickly and cheaply

▶ RUN ⚙️ MORE ▾ 💾 SAVE ▾ ⋮ ✅ This query will process 33 GiB when run.

```
1 select avg(HTSeq__FPKM_UQ) as average_fpkm
2 from `isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current`
3 where aliquot_barcode = 'TCGA-3X-AAV9-01A-72R-A41I-07'
4 and gene_type = 'protein_coding'
5
```

Query results ⬇️ SAVE RESULTS 📊 EXPLORE DATA ▾

Query complete (1.9 sec elapsed, 33 GB processed)

Job information Results JSON Execution details

Row	average_fpkm
1	342956.0106327599

Analyze correlation between TCGA samples & GTEx tissue types quickly and cheaply

```
5 --
6 GTEx_top5K AS (
7 SELECT
8   gene_id,
9   gene_description,
10  STDDEV(gene_exp) AS sigmaExp
11 FROM
12   `isb-cgc.GTEx_v7.gene_median_tpm`
13 GROUP BY
14   1,
15   2
16 ORDER BY
17   sigmaExp DESC
18 LIMIT
19   5000 ),
20 --
```

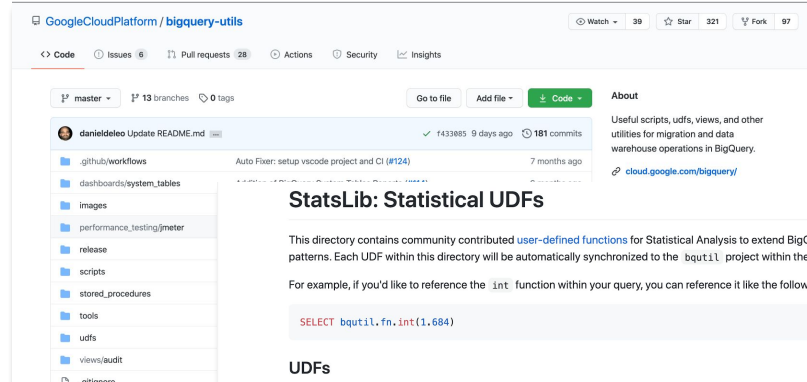
Query complete (18.8 sec elapsed, 12.7 GB processed)

Job information **Results** JSON Execution details

Row	GTEx_tissueType	sample_barcode	TCGA_project	corr
1	Liver	TCGA-DD-A39V-11A	TCGA-LIHC	0.9213023777251851
2	Liver	TCGA-DD-A39Z-11A	TCGA-LIHC	0.9189148155140473
3	Liver	TCGA-DD-A3A1-11A	TCGA-LIHC	0.917682740669065

Custom User Defined Functions allow users to expand beyond the default functionality

- We are contributing commonly used statistical methods as UDFs to the Google BQ StatsLib
- Users will be able to conduct the following statistical analyses directly within BigQuery



Data type	Data type	Statistical test/notebook
Gene expression	Clinical	Kruskal-Wallis score
Gene expression	Somatic mutation	T-test score
Gene expression	Gene expression	Spearman Correlation
Somatic mutation	Clinical	Chi Square test
Somatic mutation	Somatic Mutation	Fisher's exact test

StatsLib: Statistical UDFs

This directory contains community contributed [user-defined functions](#) for Statistical Analysis to extend BigQuery for more specialized usage patterns. Each UDF within this directory will be automatically synchronized to the `bqutil` project within the `fn` dataset for reference in queries.

For example, if you'd like to reference the `int` function within your query, you can reference it like the following:

```
SELECT bqutil.fn.int(1,684)
```

UDFs

- [kruskal_wallis](#)

Documentation

[kruskal_wallis\(arr\(struct\(factor STRING, val FLOAT64\)\)\)](#)

Takes an array of struct where each struct (point) represents a measurement, with a group label and a measurement value

The [Kruskal-Wallis test by ranks](#), Kruskal-Wallis H test (named after William Kruskal and W. Allen Wallis), or one-way ANOVA on ranks is a non-parametric method for testing whether samples originate from the same distribution. It is used for comparing two or more independent samples of equal or different sample sizes. It extends the Mann-Whitney U test, which is used for comparing only two groups. The parametric equivalent of the Kruskal-Wallis test is the one-way analysis of variance (ANOVA).

- Input: array: struct <factor STRING, val FLOAT64>
- Output: struct<H FLOAT64, p-value FLOAT64, DOF FLOAT64>

Google Cloud Platform Free Tier lets you perform introductory queries at no cost

Category	Service	Free Tier Limit	Description
COMPUTE	Cloud Run	1 million requests per month	Serverless environment to run stateless apps.
DATABASE	Firestore	1 GB Storage	Scalable NoSQL document database.
COMPUTE	Compute Engine	1 F1-micro instance per month	Scalable, high-performance virtual machines.
STORAGE	Cloud Storage	5 GB	Global regional storage for all your storage needs.
DATA ANALYTICS	Pub/Sub	10 GB Messages per month	A global service for real-time and reliable messaging and streaming data.
COMPUTE	Cloud Functions	2 million Invocations per month	A serverless environment to build and connect cloud services with code.
COMPUTE	Google Kubernetes Engine	30 clusters	Managed Kubernetes clusters, managed by Google.
COMPUTE	App Engine	28 Instance hours per day	Platform for building scalable web applications and mobile back ends.
MANAGEMENT TOOLS	Stackdriver	50 GB Logs with 30-day retention	Monitoring, logging, and diagnostics for applications on Google Cloud and AWS.
DATA ANALYTICS	BigQuery	1 TB Queries per month	Fully managed, petabyte scale, analytics data warehouse.
AI AND MACHINE LEARNING	Vision AI	1,000 Units per month	Label detection, OCR, facial detection and more.
AI AND MACHINE LEARNING	Speech-to-Text	60 Minutes per month	Speech-to-text transcription -- the same that powers Google's own products.

DATA ANALYTICS

BigQuery

1 TB

Queries per month

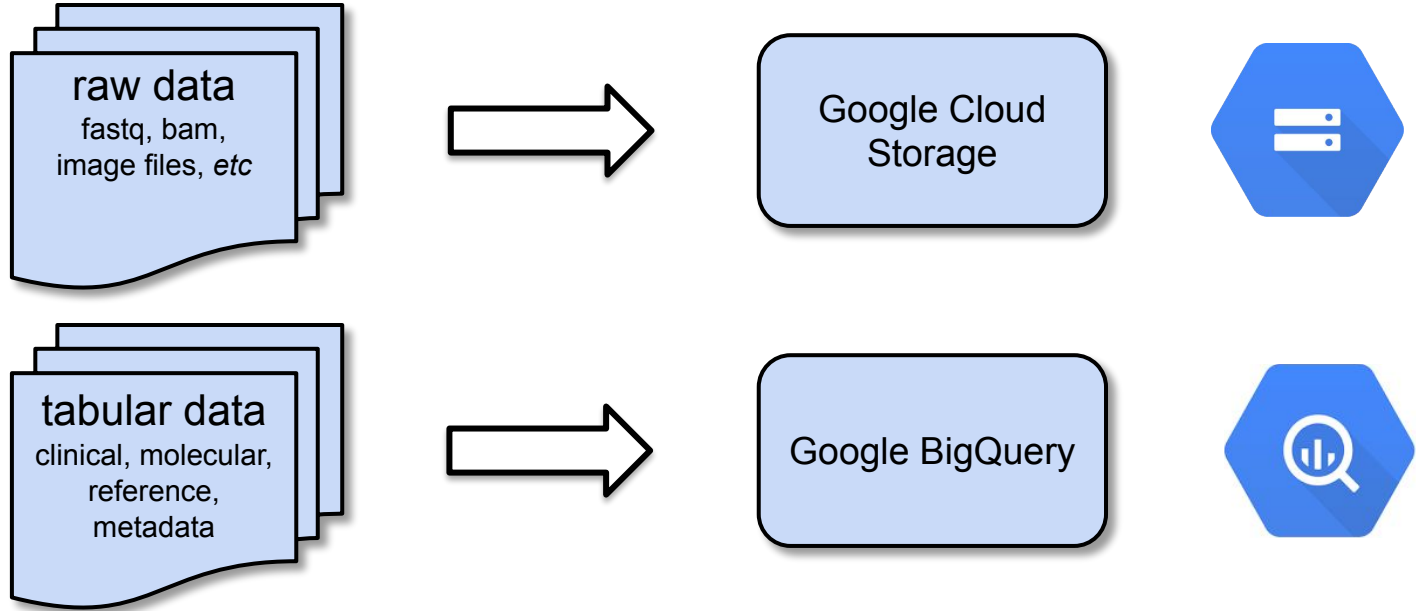
Fully managed, petabyte scale, analytics data warehouse.

1 TB of querying per month

10 GB of storage

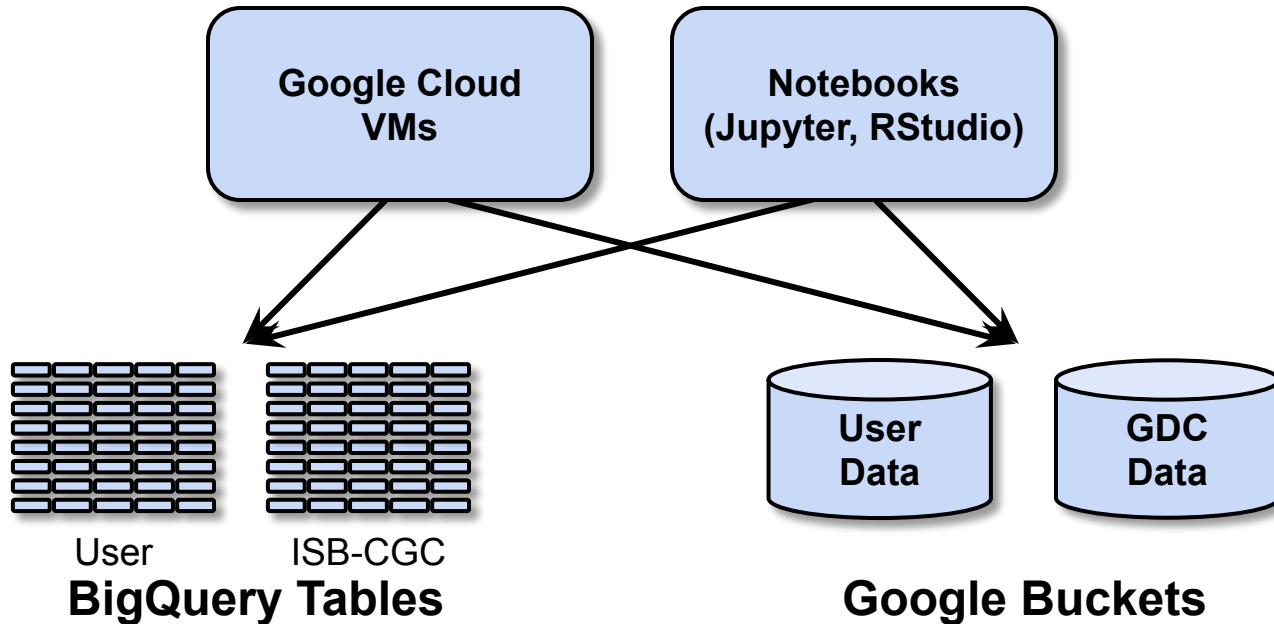
↑

How do users access data on ISB-CGC?



Multiple easy avenues for computing on data on ISB-CGC

ISB-CGC enables full command line access to analyze cloud hosted data via a collection of powerful tools and technologies along with the ability to install your own tools



In total, ISB-CGC hosts multiple TBs of cancer research data across 440 open-access BigQuery tables

- TCGA, TARGET, CCLE:
 - Clinical & Biospecimen information
 - Molecular data
 - mRNA, miRNA, protein expression
 - DNA copy-number
 - DNA methylation
 - Somatic mutations
- Reference information:
 - Gencode, Ensembl, miRBase, ...
 - hg19 to hg38 "liftover"
 - 450K DNA methylation annotations
 - Kaviar, COSMIC, ...
 - GDC metadata
- Other publicly available resources:
 - Google-hosted public datasets
 - Tute Genomics (ANNOVAR) annotations
- *Your* data!

Table Details: Biospecimen_data

Field	Type	Nullable	Description
code	STRING	NULLABLE	Describe this field
id	STRING	NULLABLE	Describe this field
siteCode	STRING	NULLABLE	Describe this field
	STRING	NULLABLE	Describe this field
	STRING	NULLABLE	Describe this field
	STRING	NULLABLE	Describe this field
id	STRING	NULLABLE	Describe this field
analyte_type	STRING	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field
analyte_type	FLOAT	NULLABLE	Describe this field

Table Details: Copy_Number_segments

Field	Type	Nullable	Description
ParticipantBarcode	STRING	NULLABLE	Describe this field
SampleBarcode	STRING	NULLABLE	Describe this field
SampleTypeLetterCode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
Study	STRING	NULLABLE	Describe this field
Platform	STRING	NULLABLE	Describe this field
Chromosome	STRING	NULLABLE	Describe this field
Start	INTEGER	NULLABLE	Describe this field
End	INTEGER	NULLABLE	Describe this field
Num_Probes	INTEGER	NULLABLE	Describe this field

Table Details: Annotations

Field	Type	Nullable	Description
annotation	INTEGER	NULLABLE	Describe this field
annotation_categorid	INTEGER	NULLABLE	Describe this field
annotation_categorname	STRING	NULLABLE	Describe this field
annotation_classification	STRING	NULLABLE	Describe this field
annotation_text	STRING	NULLABLE	Describe this field

Table Details: Protein_RPPA_data

Field	Type	Nullable	Description
Study	STRING	NULLABLE	Describe this field
AssayType	STRING	NULLABLE	Describe this field
ParticipantBarcode	STRING	NULLABLE	Describe this field
SampleBarcode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
SampleTypeLetterCode	STRING	NULLABLE	Describe this field
ParticipantBarcode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
Study	STRING	NULLABLE	Describe this field
SiteCreated	STRING	NULLABLE	Describe this field
SiteEdited	STRING	NULLABLE	Describe this field
Gene_Name	STRING	NULLABLE	Describe this field
Protein_Expression	FLOAT	NULLABLE	Describe this field
Position_Name	STRING	NULLABLE	Describe this field
Position_BaseName	STRING	NULLABLE	Describe this field

Table Details: mRNA_UNC_HiSeq_RSEM

Field	Type	Nullable	Description
ParticipantBarcode	STRING	NULLABLE	Describe this field
SampleBarcode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
Study	STRING	NULLABLE	Describe this field
SampleTypeLetterCode	STRING	NULLABLE	Describe this field
Platform	STRING	NULLABLE	Describe this field
original_gene_symbol	STRING	NULLABLE	Describe this field
HUGO_gene_symbol	STRING	NULLABLE	Describe this field
gene_id	INTEGER	NULLABLE	Describe this field
normalized_count	FLOAT	NULLABLE	Describe this field

Table Details: mRNA_BCSCGSC_HiSeq_RPKM

Field	Type	Nullable	Description
ParticipantBarcode	STRING	NULLABLE	Describe this field
SampleBarcode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
Study	STRING	NULLABLE	Describe this field
Platform	STRING	NULLABLE	Describe this field
original_gene_symbol	STRING	NULLABLE	Describe this field
HUGO_gene_symbol	STRING	NULLABLE	Describe this field
gene_id	INTEGER	NULLABLE	Describe this field
normalized_count	FLOAT	NULLABLE	Describe this field

Table Details: miRNA_expression

Field	Type	Nullable	Description
ParticipantBarcode	STRING	NULLABLE	Describe this field
SampleBarcode	STRING	NULLABLE	Describe this field
AlloportBarcode	STRING	NULLABLE	Describe this field
Study	STRING	NULLABLE	Describe this field
Platform	STRING	NULLABLE	Describe this field
miRNA_id	STRING	NULLABLE	Describe this field
miRNA_accession	STRING	NULLABLE	Describe this field
normalized_count	FLOAT	NULLABLE	Describe this field

SampleTypeLetterCode	STRING	NULLABLE	Refer: https://bga-data.ncbi.nih.gov/data-reports/codeTableReport.htm
	STRING	NULLABLE	The Alloport ID is an identifier for TCGA data. Refer: https://wiki.nci.nih.gov/display/TCGA/TCGA+
	STRING	NULLABLE	Refer: https://bga-data.ncbi.nih.gov/data-reports/codeTableReport.htm
	STRING	NULLABLE	TCGA disease type
	STRING	NULLABLE	Illumina's CpG loci IDs. Refer: http://www.illumina.com/content/illumina-marking-locations/illumina-marking-locations-technology/cpg_loci_identification.pdf
	FLOAT	NULLABLE	The beta value (B) is used to estimate the methylation level of the CpG locus using the ratio of intensities by



Google BigQuery



This is potentially overwhelming for data discovery, we have a table search UI to help

BigQuery Table Search

[ISB-CGC BigQuery Documentation](#) [ISB-CGC BigQuery Access Info](#) [Google BigQuery Console](#) [About BigQuery](#) [Release Notes](#)

Explore and learn more about available ISB-CGC BigQuery tables with this search feature.
Find tables of interest based on category, reference genome build, data type and free-form text search.

Status
CURRENT

Name

Program

Category
 CLINICAL BIOSPECIMEN DATA
 FILE METADATA
 GENOMIC REFERENCE DATABASE
 PROCESSED -OMICS DATA

Reference Genome
ALL

Source

Data Type

Experimental Strategy

[Reset All Filters](#)

[+ Show More Filters](#)

Show 10 entries

Columns CSV Download Search:

Name	Program	Category	Source	Data Type	Status	Rows	Created	Preview	Open
CCLE 2016 - AFFYU133 MICROARRAY	CCLE	PROCESSED -OMICS DATA	BROAD	GENE EXPRESSION	CURRENT	17,525,476	2/26/2016		
CCLE 2016 - COPY NUMBER SEGMENTS	CCLE	PROCESSED -OMICS DATA	BROAD	COPY NUMBER SEGMENT	CURRENT	760,192	2/27/2016		
CCLE 2016 - FASTQC METRICS	CCLE	PROCESSED -OMICS DATA	BROAD	FILE METADATA	CURRENT	1,249	3/28/2016		
CCLE 2016 - FILE METADATA	CCLE	PROCESSED -OMICS DATA	BROAD	FILE METADATA	CURRENT	1,915	3/29/2016		
CCLE 2016 - SAMPLE INFORMATION	CCLE	PROCESSED -OMICS DATA	BROAD	BIOSPECIMEN SUPPLEMENT	CURRENT	929	2/26/2016		
CCLE 2016 - SOMATIC MUTATION	CCLE	PROCESSED -OMICS DATA	BROAD	SOMATIC MUTATIONS	CURRENT	116,708	2/26/2016		
CCLE BIOSPECIMEN V0	CCLE	CLINICAL BIOSPECIMEN DATA	BROAD	BIOSPECIMEN SUPPLEMENT	CURRENT	954	4/4/2019		
CCLE CLINICAL V1	CCLE	CLINICAL BIOSPECIMEN DATA	BROAD	CLINICAL DATA	CURRENT	950	6/21/2019		
CCLE HG19 METADATA RELEASE 14	CCLE	FILE METADATA	BROAD	FILE METADATA	CURRENT	1,273	3/7/2019		
CLINVAR 20180401 GRCH37		GENOMIC REFERENCE DATABASE	CLINVAR	SOMATIC MUTATIONS	CURRENT	354,471	4/17/2018		

Showing 1 to 10 of 214 entries (filtered from 327 total entries)

Previous 1 2 3 4 5 ... 22 Next

Have feedback or corrections? Please email us at feedback@isb-cgc.org.

https://isb-cgc.appspot.com/bq_meta_search/

Table schemas are easily accessible through this UI

BigQuery Table Search

Explore and learn more about available ISB-CGC BigQuery tables with this search feature.
Find tables of interest based on category, reference genome build, data type and free-form text search.

[ISB-CGC BigQuery Documentation](#) [ISB-CGC BigQuery Access Info](#) [Google BigQuery Console](#) [About BigQuery](#) [Release Notes](#)

Status
CURRENT

Name
[Text Input]

Program
Choose Programs...

Category
 CLINICAL BIOSPECIMEN DATA [?](#)
 FILE METADATA [?](#)
 GENOMIC REFERENCE DATABASE [?](#)
 PROCESSED -OMICS DATA [?](#)

Reference Genome
ALL

Source
Choose Sources...

Data Type
Choose Data Types...

Experimental Strategy
Choose Experimental Strategy...

[Reset All Filters](#)

[+ Show More Filters](#)

Show 10 entries

Columns CSV Download Search: [Text Input]

Name	Program	Category	Source	Data Type	Status	Rows	Created	Preview	Open
CCLE 2016 - AFFYU133 MICROARRAY	CCLE	PROCESSED -OMICS DATA	BROAD	GENE EXPRESSION	CURRENT	17,525,476	2/26/2016		
CCLE 2016 - COPY NUMBER SEGMENTS	CCLE	PROCESSED -OMICS DATA	BROAD	COPY NUMBER SEGMENT	CURRENT	760,192	2/27/2016		

Full ID isb-cgc.ccle_201602_alpha.Copy_Number_Segments [COPY](#) [OPEN](#)

Dataset ID ccle_201602_alpha

Table ID Copy_Number_Segments

Description Data was extracted from an older CCLE dataset from Google Genomics on February 2016. Copy number segment data are made available here.

Schema

Field Name	Type	Mode	Description
CCLE_name	STRING	NULLABLE	Cell line primary name, appended with a short name for the location of the cancer: e.g. TC71_BONE_HUPT4_PANGREAS, etc
Cell_line_primary_name	STRING	NULLABLE	The cell line primary name: e.g. TC71, NCI-60, etc
Platform	STRING	NULLABLE	Platform used to generate these data (Genome_Wide_SNP_6)
Chromosome	STRING	NULLABLE	Chromosome, possible values: chr1-22, and chrX
Start	INTEGER	NULLABLE	Start position
End	INTEGER	NULLABLE	End position
Num_Probes	INTEGER	NULLABLE	The num_probes field specifies the number of probes on the SNP chip that went into estimating the mean copy number for this segment
Segment_Mean	FLOAT	NULLABLE	Provides the log2(CN2) mean value estimate

Labels access: open data_type: copy_number_segment program: ccle reference_genome_0: hg18 source: broad category: processed_omics_data status: current

CCLE 2016 - FASTQC METRICS	CCLE	PROCESSED -OMICS DATA	BROAD	FILE METADATA	CURRENT	1,249	3/28/2016		
CCLE 2016 - FILE METADATA	CCLE	PROCESSED -OMICS DATA	BROAD	FILE METADATA	CURRENT	1,915	3/29/2016		
CCLE 2016 - SAMPLE INFORMATION	CCLE	PROCESSED -OMICS DATA	BROAD	BIOSPECIMEN SUPPLEMENT	CURRENT	929	2/26/2016		
CCLE 2016 - SOMATIC MUTATION	CCLE	PROCESSED -OMICS DATA	BROAD	SOMATIC MUTATIONS	CURRENT	116,708	2/26/2016		
CCLE BIOSPECIMEN V0	CCLE	CLINICAL BIOSPECIMEN DATA	BROAD	BIOSPECIMEN SUPPLEMENT	CURRENT	954	4/4/2019		
CCLE CLINICAL V1	CCLE	CLINICAL BIOSPECIMEN DATA	BROAD	CLINICAL DATA	CURRENT	950	6/21/2019		
CCLE HG19 METADATA RELEASE 14	CCLE	FILE METADATA	BROAD	FILE METADATA	CURRENT	1,273	3/7/2019		
CLINVAR 20180401 GRCH37		GENOMIC REFERENCE DATABASE	CLINVAR	SOMATIC MUTATIONS	CURRENT	354,471	4/17/2018		

Showing 1 to 10 of 214 entries (filtered from 327 total entries)

Previous 1 2 3 4 5 ... 22 Next




Have feedback or corrections? Please email us at feedback@isb-cgc.org.

https://isb-cgc.appspot.com/bq_meta_search/



Benefits of the ISB-CGC BigQuery Table Search

- Allows users to browse and learn more about available ISB-CGC BigQuery tables
- Each table has been curated to include detailed table and field descriptions as well as table labels
- Identify table(s) of interest by filtering (e.g. by reference genome build, data type, category) or via free-form text search
- Get a snapshot of table contents by previewing the first few (~10) lines
- Found a table of interested? An “open” button takes users directly to the GCP BigQuery Console.

	CCLL CLINICAL V1	CCLL	CLINICAL BIOSPECIMEN DATA	BROAD	CLINICAL DATA	CURRENT	950	6/21/2019		
---	------------------	------	---------------------------------	-------	---------------	---------	-----	-----------	---	---

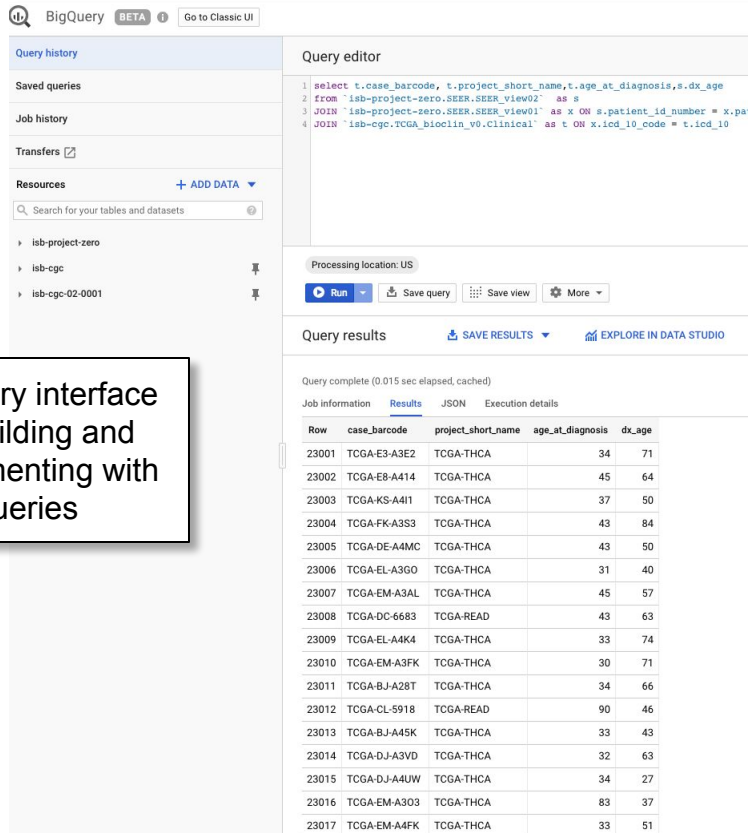
Use Google BigQuery to easily connect your research to public datasets

ISB-CGC and Other
Public Datasets



Private User Data
and Derived Results

Tables can be joined in BigQuery using SQL to draw connections amongst data



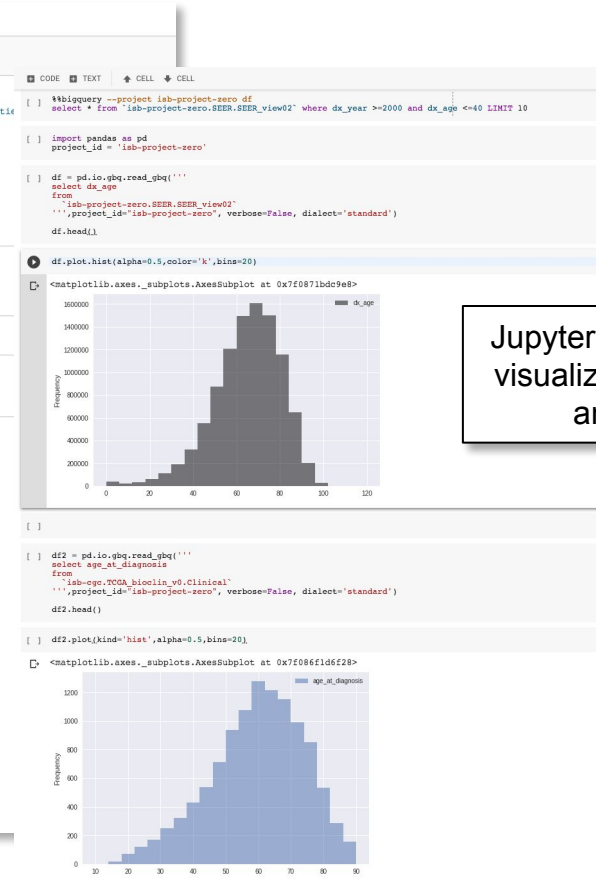
The screenshot shows the BigQuery interface. On the left, there is a sidebar with 'Query history', 'Saved queries', 'Job history', 'Transfers', and 'Resources'. The main area is the 'Query editor' with a SQL query:

```
1 select t.case_barcode, t.project_short_name, t.age_at_diagnosis, s.dx_age
2 from `isb-project-zero.SEER_SEER_View02` as s
3 JOIN `isb-project-zero.SEER_SEER_View01` as x ON x.patient_id_number = x.patien
4 JOIN `isb-cgc.TCGA_bioclin_v0.clinical` as t ON x.tcd_10_code = t.tcd_10
```

Below the editor, there are buttons for 'Run', 'Save query', 'Save view', and 'More'. The 'Query results' section shows a table with 17 rows and 5 columns: Row, case_barcode, project_short_name, age_at_diagnosis, and dx_age.

Row	case_barcode	project_short_name	age_at_diagnosis	dx_age
23001	TCGA-E3-A3E2	TCGA-THCA	34	71
23002	TCGA-E8-A414	TCGA-THCA	45	64
23003	TCGA-KS-A411	TCGA-THCA	37	50
23004	TCGA-FK-A3S3	TCGA-THCA	43	84
23005	TCGA-DE-A4MC	TCGA-THCA	43	50
23006	TCGA-EL-A3G0	TCGA-THCA	31	40
23007	TCGA-EM-A3AL	TCGA-THCA	45	57
23008	TCGA-DC-6683	TCGA-READ	43	63
23009	TCGA-EL-A4K4	TCGA-THCA	33	74
23010	TCGA-EM-A3FK	TCGA-THCA	30	71
23011	TCGA-BJ-A28T	TCGA-THCA	34	66
23012	TCGA-CL-5918	TCGA-READ	90	46
23013	TCGA-BJ-A45K	TCGA-THCA	33	43
23014	TCGA-DJ-A3VD	TCGA-THCA	32	63
23015	TCGA-DJ-A4UW	TCGA-THCA	34	27
23016	TCGA-EM-A3O3	TCGA-THCA	83	37
23017	TCGA-EM-A4FK	TCGA-THCA	33	51

BigQuery interface for building and experimenting with queries



The screenshot shows a Jupyter notebook with the following code cells:

```
%%bigquery --project isb-project-zero df
select * from `isb-project-zero.SEER_SEER_View02` where dx_year >= 2000 and dx_age <= 40 LIMIT 10
```

```
import pandas as pd
project_id = 'isb-project-zero'
```

```
df = pd.io.gbq.read_gbq('''
select dx_age
from
  isb-project-zero.SEER_SEER_View02
  ''', project_id='isb-project-zero', verbose=False, dialect='standard')
df.head()
```

```
df.plot.hist(alpha=0.5, color='k', bins=20)
```

The first histogram shows the frequency distribution of 'dx_age' for patients with dx_year >= 2000 and dx_age <= 40. The x-axis is 'dx_age' (0-120) and the y-axis is 'Frequency' (0-160000).

```
df2 = pd.io.gbq.read_gbq('''
select age_at_diagnosis
from
  isb-cgc.TCGA_bioclin_v0.clinical
  ''', project_id='isb-project-zero', verbose=False, dialect='standard')
df2.head()
```

```
df2.plot(kind='hist', alpha=0.5, bins=20)
```

The second histogram shows the frequency distribution of 'age_at_diagnosis' from the TCGA clinical data. The x-axis is 'age_at_diagnosis' (0-90) and the y-axis is 'Frequency' (0-1200).

Jupyter notebook to visualize and share analysis

BigQuery integrates with a variety of commonly used analysis tools



bigquery and
bigQueryR



googleAuthR



Pre-built VM images

IP[y]:

IPython



Cloud notebooks
and workspaces.

Cloud Datalab



Interactive step-by-step guide on BigQuery data analysis

× Cancer Data Discovery in the Cloud through ISB-CGC

1 Introduction

2 Getting Started

3 Exploring ISB-CGC BigQuery Tables

4 Query Cancer Data in the Cloud

5 Access data in BigQuery from R

6 Analyze data in BigQuery from R

7 Using Bioconductor packages on data in BigQuery

8 More Resources

1. Introduction



ISB-CGC, one of the National Cancer Institute's Cloud Resources, uniquely hosts cancer data including somatic mutations, copy number variations, gene and protein expressions, etc. from widely used cancer datasets including TCGA, TARGET and many more in Google BigQuery.

Google BigQuery is a massively-parallel analytics engine ideal for tabular data. ISB-CGC has combined data scattered over tens of thousands of files into easily accessible BigQuery tables. This novel approach allows our users to quickly analyze data from thousands of patients in ISB-CGC curated BigQuery tables.

- Users can explore and learn more about the ISB-CGC hosted BigQuery tables via an interactive [BigQuery Table Search User Interface](#). Users can find tables of interest based on category, reference genome build, data type and free-form text search.
- Users with Google Cloud Platform (GCP) projects can analyze patient, biospecimen, and molecular data from many NCI funded programs such as TCGA, TARGET, CCLL, GTEx all in ISB-CGC's BigQuery tables. SQL queries can be executed in many ways including through the Google Cloud BigQuery web console, Jupyter notebooks, and R scripts.

In this tutorial, you will:

- Analyze gene expression and protein abundance differences between two types of TCGA kidney cancers, **Kidney Renal Papillary Carcinoma (TCGA-KIRP)** and **Kidney Chromophobe (TCGA-KICH)**.
- Build a cohort of patients with these cancer types and extract their respective gene expression and protein abundance data from ISB-CGC TCGA Google BigQuery tables.
- Connect to Google BigQuery tables from within R for data analysis and visualization

What you'll learn:

- How to explore the **ISB-CGC BigQuery Table Search User Interface**
- How to build and run queries in the **Google BigQuery Console**
- How to use **R notebooks** for data analysis and visualization
- How to use **Bioconductor packages** on data in **ISB-CGC BigQuery tables**

Next

Report a mistake

https://isb-cgc.appspot.com/how_to_discover/#0



Some examples of analyses run with our tools

BigQuery can be used to join tables of different data types

```
1  with gexp as
2  (
3    SELECT
4      project_short_name,
5      case_barcode,
6      gene_name,
7      avg(HTSeq__FPKM) as mean_gexp
8    FROM `isb-cgc.TCGA_hg38_data_v0.RNAseq_Gene_Expression`
9    WHERE
10     project_short_name in ('TCGA-KIRP', 'TCGA-KICH')
11     AND gene_type = 'protein_coding'
12   GROUP BY
13     project_short_name, case_barcode, gene_name
14 ),
15 pexp as
16 (
17   SELECT
18     project_short_name,
19     case_barcode, gene_name,
20     avg(protein_expression) as mean_pexp
21   FROM `isb-cgc.TCGA_hg38_data_v0.Protein_Expression`
22   WHERE
23     project_short_name in ('TCGA-KIRP', 'TCGA-KICH')
24   GROUP BY
25     project_short_name, case_barcode, gene_name
26 )
27 SELECT
28   gexp.project_short_name,
29   gexp.case_barcode,
30   gexp.gene_name,
31   gexp.mean_gexp,
32   pexp.mean_pexp
33 FROM gexp
34   inner join pexp
35   on
36   gexp.project_short_name = pexp.project_short_name
37   AND gexp.case_barcode = pexp.case_barcode
38   AND gexp.gene_name = pexp.gene_name
```

Query results [SAVE RESULTS](#) [EXPLORE DATA](#)

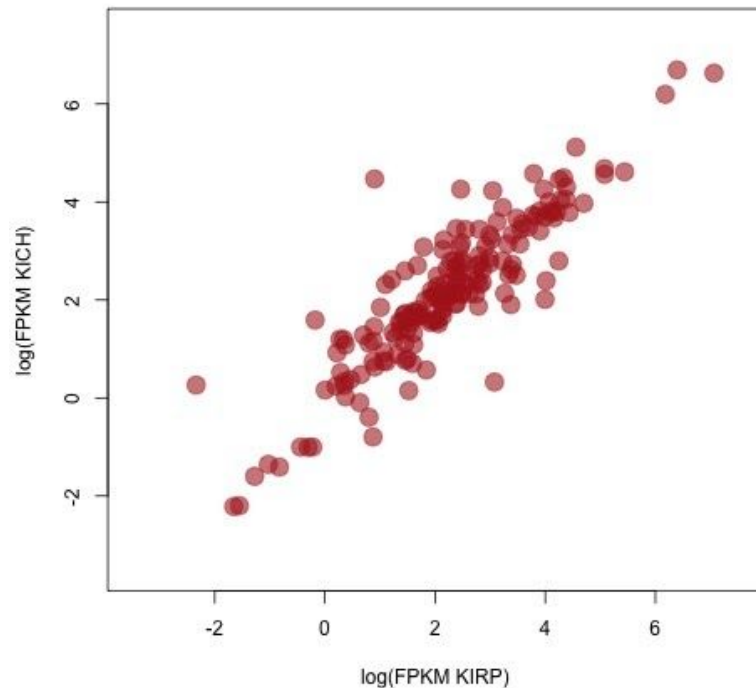
Query complete (3.6 sec elapsed, 36.5 GB processed)

[Job information](#) [Results](#) [JSON](#) [Execution details](#)

Row	project_short_name	case_barcode	gene_name	mean_gexp	mean_pexp
1	TCGA-KIRP	TCGA-BQ-5882	CASP9	4.36721774268	-0.77899630325
2	TCGA-KIRP	TCGA-BQ-5893	CASP9	3.51383227626	-0.00487536525000024
3	TCGA-KIRP	TCGA-UZ-A9PV	CASP9	5.18471681617	0.608095901
4	TCGA-KIRP	TCGA-HE-7130	CASP9	9.66588010427	1.26419392825
5	TCGA-KIRP	TCGA-BQ-7051	CASP9	4.960828359680001	0.16363261725
6	TCGA-KIRP	TCGA-B9-A5W7	CASP9	3.48913273215	0.46925972255
7	TCGA-KIRP	TCGA-B9-A69E	CASP9	2.31606330472	-0.79851830675
8	TCGA-KIRP	TCGA-IA-A83W	CASP9	7.0151441055	0.0390634209999998
9	TCGA-KIRP	TCGA-KV-A74V	CASP9	3.6511469304	0.58510363975
10	TCGA-KIRP	TCGA-DW-7963	CASP9	2.72821430732	-0.0676392402500001
11	TCGA-KIRP	TCGA-B1-A470	CASP9	4.38451034062	-0.30207994575
12	TCGA-KIRP	TCGA-SX-A7SM	CASP9	5.19366866041	-0.12840562825
13	TCGA-KIRP	TCGA-BQ-7061	CASP9	5.52320872773	-0.4549548795
14	TCGA-KIRP	TCGA-AL-A5DJ	CASP9	3.94040060142	-0.22713179975
15	TCGA-KIRP	TCGA-B9-A8YH	CASP9	4.40211877405	0.307956009
16	TCGA-KIRP	TCGA-J7-A8I2	CASP9	5.40615807819	0.9955787775
17	TCGA-KIRP	TCGA-IZ-8196	CASP9	4.76919558442	0.0767987447499998
18	TCGA-KIRP	TCGA-UZ-A9PR	CASP9	5.16376898893	0.41371152025
19	TCGA-KIRP	TCGA-BQ-5875	CASP9	3.445547761625	-0.35423371525
20	TCGA-KIRP	TCGA-SX-A7SS	CASP9	3.32930470316	-0.12148286825
21	TCGA-KIRP	TCGA-DW-7840	CASP9	4.47098482015	-0.1197760605

This query can be executed and the results downloaded and plotted from R

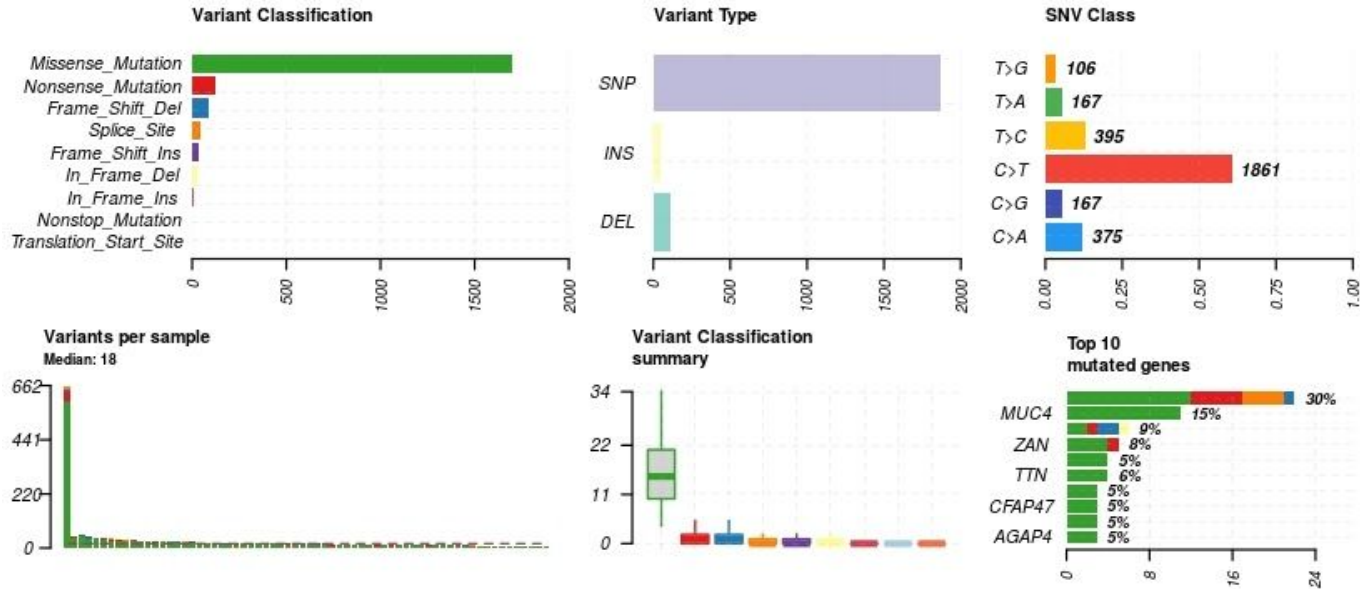
```
1
2 # Load data into a data frame
3 expression_data <- bq_table_download(bq_project_query (project, query = sql_expression))
4
5 expression_data$id <- paste(expression_data$project_short_name, expression_data$case_barcode,
6 sep='.')
7 cases <- unique(expression_data$id)
8
9 # Transform the expression_data data frame, so that columns are samples, rows are genes.
10 list_exp <- lapply(cases, function(case){
11   temp <- expression_data[expression_data$id == case, c('gene_name', 'mean_gexp')]
12   names(temp) <- c('gene_name', case)
13   return(temp)
14 })
15 gene_exps <- Reduce(function(x, y) merge(x, y, all=T, by="gene_name"), list_exp)
16
17 exp_p <- gene_exps[,grep('KIRP', names(gene_exps))]
18 exp_c <- gene_exps[,grep('KICH', names(gene_exps))]
19 plot(log(rowMeans(exp_p)), log(rowMeans(exp_c)),
20      xlab='log(FPKM KIRP)', ylab='log(FPKM KICH)',
21      xlim=c(-3.5,7.5), ylim=c(-3.5,7.5), pch=19, cex=2,
22      col=rgb(178,34,34,max=255,alpha=150))
```



Lots of data types can directly be plugged into Bioconductor pipelines

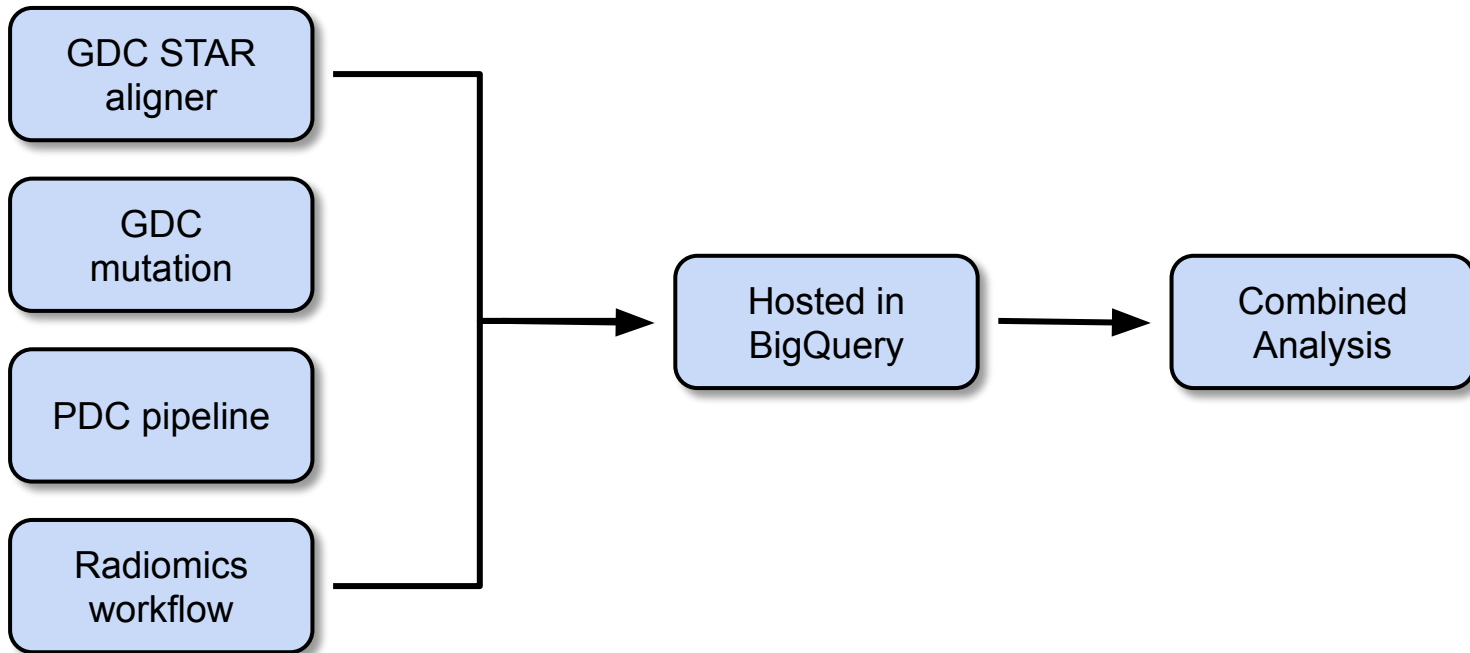
```
1 # Use BigQuery to load TCGA somatic mutation data for our cancers of interest.
2 sql_kich <- "
3 SELECT
4   Hugo_Symbol,
5   Chromosome,
6   Start_Position,
7   End_Position,
8   Reference_Allele,
9   Tumor_Seq_Allele2,
10  Variant_Classification,
11  Variant_Type,
12  sample_barcode_tumor
13 FROM
14   `isb-cgc.TCGA_hg38_data_v0.Somatic_Mutation`
15 WHERE
16   project_short_name = 'TCGA-KICH'"
17
18 # Put data into a dataframe
19 maf_kich <- bq_table_download(bq_project_query (project, query = sql_kich))
20
21 # Rename column 9 to the field name required by maftools.
22 colnames(maf_kich)[9] <- "Tumor_Sample_Barcode"
23
24 # Convert data frames to maftools objects.
25 kich <- read.maf(maf_kich)
26
27 # Maftools plots
28 ▼ plotmafSummary(maf = kich, rmOutlier = TRUE,
29                addStat = 'median', dashboard = TRUE,
30                titvRaw = FALSE)
```


A plot of MAF characteristics using the Bioconductor maftools library



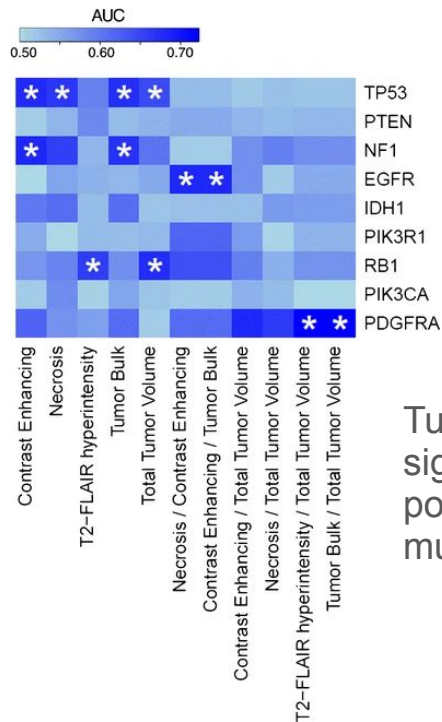
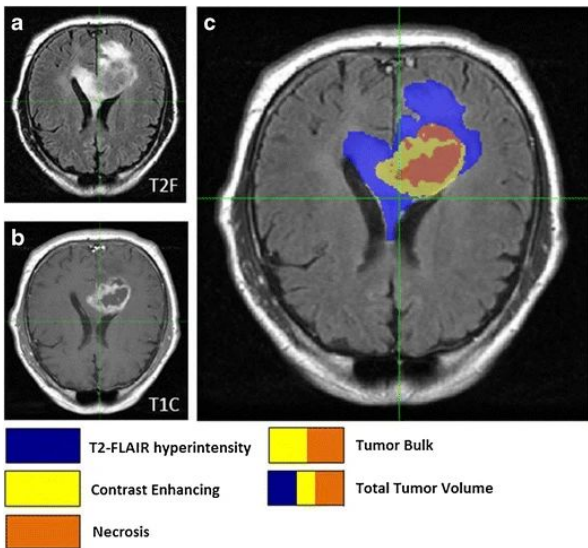
Connecting existing genomics, proteomics, and radiomics data from the respective data commons using BigQuery

Existing data from



Collaboration between ISB-CGC Cloud Resource and IDC

Tumor identification and characterization by machine learning



Tumor features have significant predictive power about the mutations present

Gutman, et al.
Diagnostic Neuroradiology 2015

We have created workflows to join data and run statistical tests

Generated significance values from the correlation between gene expression and tumor features

Subset of somatic mutations for each case id to be used for association analyses

Using an initial feature set from Bakas, et al. *Nature Scientific Data*, 2017 to pilot analyses

The screenshot shows a data explorer interface with a search bar and a list of pinned projects on the left. The main area displays a table titled 'tcga_gbm_radiomicFeatures'. The table has columns for Row ID, Date, VOLUME_ET, VOLUME_NET, VOLUME_ED, and VOLUME_TC. The data rows show various TCGA case IDs and their corresponding volume measurements.

Row ID	Date	VOLUME_ET	VOLUME_NET	VOLUME_ED	VOLUME_TC
1	TCGA-02-0006	1996.08.23	1662	384	36268
2	TCGA-02-0009	1997.06.14	4362	4349	15723
3	TCGA-02-0046	1998.11.28	35721	4212	103826
4	TCGA-02-0047	1998.12.15	50410	11908	160681
5	TCGA-02-0064	1999.08.08	26241	9905	19635
6	TCGA-02-0068	2000.05.16	13224	11978	74657
7	TCGA-02-0070	2000.07.10	8994	3169	25665
8	TCGA-02-0075	1999.09.24	13892	26109	94553
9	TCGA-02-0085	1999.01.30	25236	3277	66878
10	TCGA-02-0087	1999.12.13	3852	42877	59124
11	TCGA-02-0102	1997.12.15	1037	204	11711
12	TCGA-02-0106	1998.10.30	46904	12627	33130
13	TCGA-02-0116	1997.03.22	35019	1665	30118
14	TCGA-06-0119	2003.12.26	64595	15826	83133
15	TCGA-06-0122	2004.09.14	17941	736	48147
16	TCGA-06-0130	2001.09.11	9138	2149	30952
17	TCGA-06-0137	2001.12.24	49748	7195	129318
18	TCGA-06-0138	2002.11.25	26938	15565	80426

```
2 #get gene expression data
3 rnaseq AS (
4   select case_barcode, sample_barcode, gene_name, MAX(HTSeq__Counts) gexp
5   from
6     isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current as rna
7   WHERE
8     gene_type = 'protein_coding'
9     AND project_short_name = 'TCGA-GBM'
10  GROUP BY gene_name, case_barcode, sample_barcode
11 ),
12 # join gene exp with image features
13 combined AS (
14   SELECT gene_name, feature, gexp, value ,
15         (RANK() OVER( PARTITION BY gene_name,feature ORDER BY gexp )) + (COUNT(*) OVER
16         (RANK() OVER( PARTITION BY gene_name,feature ORDER BY value )) AS ranky
17 FROM rnaseq
18 JOIN isb-cgc-ids-collaboration.Analysis.unpivoted_tcga_gbm_radiomicFeatures img
19 ON img.ID = rnaseq.case_barcode
20 AND STARTS_WITH(img.feature, "VOLUME")
21 ),
22 # spearman correlation
23 correlation AS (
24   SELECT gene_name, feature, COUNT(gexp) as n, CORR(rankx,ranky) as corr
25 FROM combined
26 GROUP BY gene_name, feature
27 )
28 #p-value computation
29 SELECT feature, gene_name, n, corr,
30        'cgc-05-0042.functions.jstat_ttest'(t, n-2, 2) as p
31 FROM (
32   SELECT *,
33         ABS(corr)*SQRT( (n-2)/((1-corr*corr))) AS t
34 FROM correlation
35 WHERE ABS(corr) < 1.0
36 )
37 ORDER BY p
```

The screenshot shows a table titled 'GBM_gene_matrix'. The table has columns for Row, gene, TCGA_DB_S275, TCGA_DH_A66D, TCGA_DH_A66F, TCGA_DU_6392, and TCGA_DU_6402. The data rows show various ENSG gene IDs and their corresponding mutation counts across different TCGA datasets.

Row	gene	TCGA_DB_S275	TCGA_DH_A66D	TCGA_DH_A66F	TCGA_DU_6392	TCGA_DU_6402
4401	ENSG00000113621	0	0	0	0	0
4402	ENSG00000172508	0	0	0	0	0
4403	ENSG00000165028	0	0	0	0	0
4404	ENSG00000166159	0	0	0	0	0
4405	ENSG00000159212	0	0	0	0	0
4406	ENSG00000061794	0	0	0	0	0
4407	ENSG00000173812	0	1	0	0	0
4408	ENSG00000183011	0	0	0	0	0
4409	ENSG00000165175	0	0	0	0	0
4410	ENSG00000162595	0	0	0	0	0
4411	ENSG00000188846	0	0	0	0	0
4412	ENSG00000143353	0	0	0	0	0
4413	ENSG00000100350	0	0	0	0	0
4414	ENSG00000138382	0	0	0	0	0
4415	ENSG00000141295	0	0	0	0	0
4416	ENSG00000101152	0	0	0	0	0
4417	ENSG00000180483	0	0	0	0	0
4418	ENSG00000143545	0	0	0	0	0
4419	ENSG00000166974	0	0	0	0	0

With more effort these cross-omics analyses could lead to novel insights

Using a simple t-test run via SQL in BigQuery we were able to replicate a significant association between a mutated MUC16 gene and an increased tumor volume as measured by radiology.

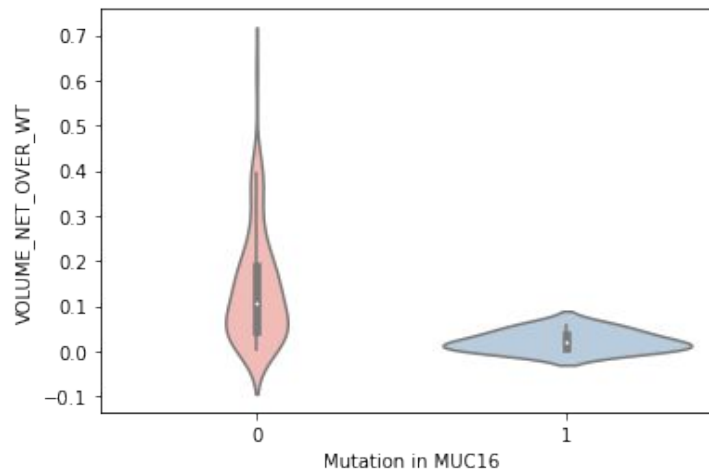
Query results [SAVE RESULTS](#) [EXPLORE DATA](#)

Query complete (2.5 sec elapsed, 68.2 MB processed)

Job information **Results** JSON Execution details

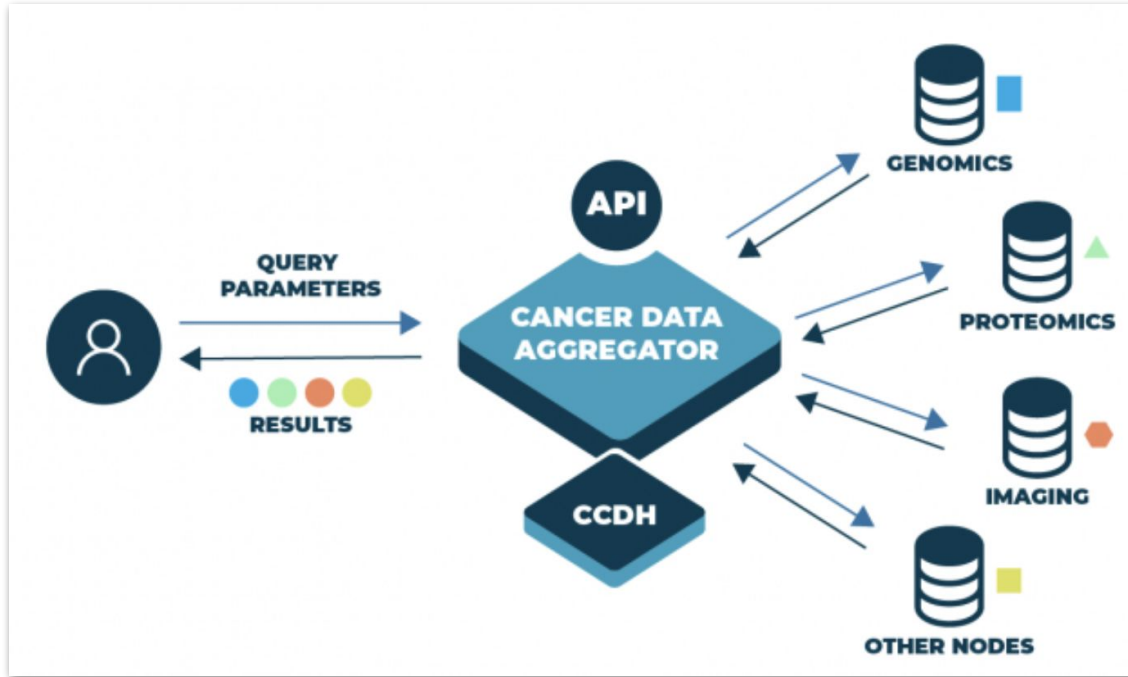
Row	Study	radiomic_feature	gene	n1	n0	avg1	avg0	pvalue
1	GBM	VOLUME_NET_OVER_WT	MUC16	6	96	0.02313391	0.139958265	1.5318989464578824E-11
2	GBM	VOLUME_NET_OVER_ED	MUC16	6	96	0.036119709	0.333068168	6.005824166041314E-10
3	GBM	VOLUME_NET_OVER_BRAIN	MUC16	6	96	0.00180942	0.011286297	6.260805815986468E-8
4	GBM	VOLUME_NET	MUC16	6	96	2660	16385.6875	1.2425312760676515E-7
5	GBM	VOLUME_ET_over_TC	MUC16	6	96	0.90486	0.691931177	2.036526115878079E-5
6	GBM	VOLUME_NET_over_TC	MUC16	6	96	0.095138	0.308068906	2.0366175438960117E-5

Rows per page: 100 1 - 100 of 210 First page <



The Cancer Data Aggregator (CDA) is aimed to bridge data discovery between the Data Commons

Federated realtime search across nodes



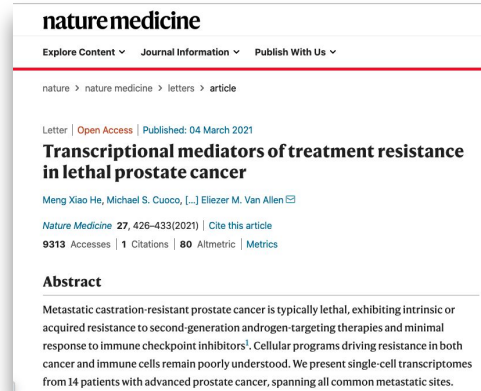
Some more workflows done by ISB-CGC end-users

Multiple PanCancer Atlas projects, including:

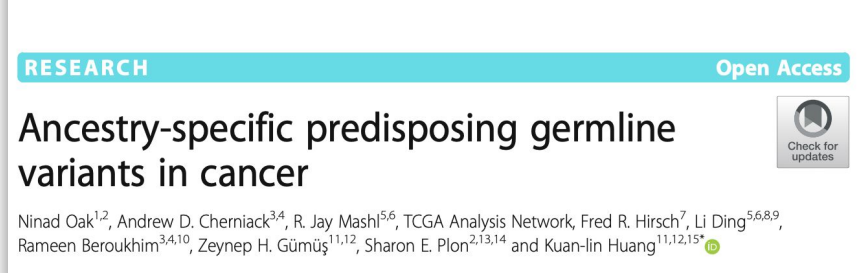
- Germline-variant calling
- Fusion gene analysis
- T-cell and B-cell receptor analysis
- viral DNA screening
- MYC pathway analysis (BQ)
- 8-oxoG filtering (MC3 project)

Other end-user projects include:

- SMC-RNA Dream challenge (supporting both the organizers and many participants)
- tumor-specific alternative polyadenylation
- ML algorithm evaluation & benchmarking
- RNA seq alignment to novel transcriptome(s)
- mRNA expression quantitation
- targeted de-novo assembly
- structural variations (WGS + SNP6 data)
- metagenomics / cancer analysis
- statistical meta-analysis of miRNAs in cancer
- code/tutorial development
- GDC hg38 TCGA miRNA QC (w/ BCGSC)



Oak et al. *Genome Medicine* (2020) 12:51
<https://doi.org/10.1186/s13073-020-00744-3>



*with many other manuscripts
and grants currently in
progress or submitted*

nature genetics

Explore Content | Journal Information | Publish With Us | Subscribe

nature > nature genetics > letters > article

Letter | Published: 17 August 2020

Extrachromosomal DNA is associated with oncogene amplification and poor outcome across multiple cancers

Hoon Kim, Nam-Phuong Nguyen, Kristen Turner, Sihan Wu, Amit D. Gujar, Jens Luebeck, Jihe Liu, Viraj Deshpande, Utkrish Rajkumar, Sandeep Namburi, Samir Kumar B. Amin, Eunhee Yi, Francesca Menghi, Johannes H. Schulte, Anton G. Henssen, Howard Y. Chang, Christine R. Beck, Paul S. Mischel, Vineet Bafna & Roel G. W. Verhaak

Genome Medicine

Resources to help you get started on ISB-CGC

- Our documentation can be found at:
<https://isb-cancer-genomics-cloud.readthedocs.io/>
- Our workflow tutorials can be found at:
https://isb-cgc.appspot.com/programmatic_access/
- Finally, a quick guide to BigQuery analysis is at:
https://isb-cgc.appspot.com/how_to_discover/#0

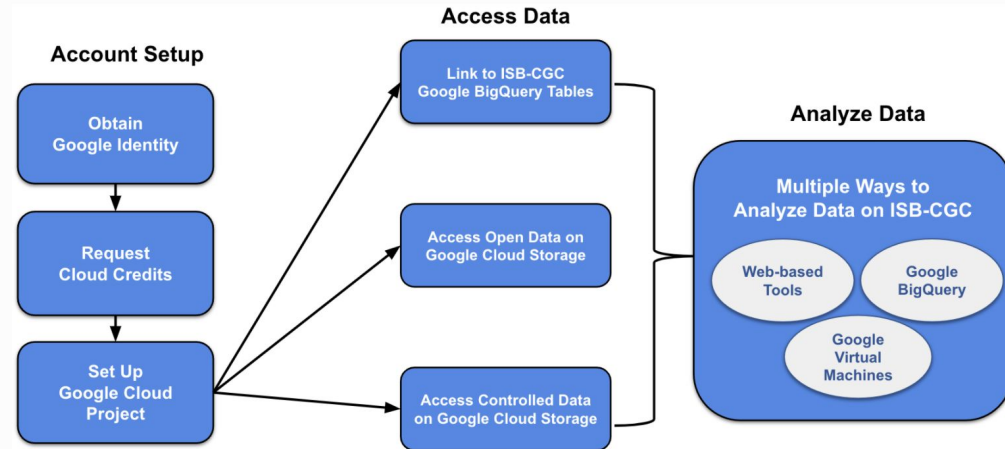
For feedback/suggestions or need help getting starting, contact us at feedback@isb-cgc.org

Contact us about setting up your own Google Cloud Platform Project!

feedback@isb-cgc.org

Quick-Start Guide

ISB-CGC provides both interactive (through a [web application](#)) and programmatic access to data hosted by institutes such as the Genomic Data Commons (GDC) of the National Cancer Institute (NCI), and the Wellcome Trust Sanger Institute, leveraging many aspects of the Google Cloud Platform. To get started, you'll need a Google Cloud Project. Additionally, to access controlled data, you'll also need [dbGaP authorization](#).



ISB-CGC Team



Bill Longabaugh
Suzanne Paquette
David Gibbs
Jennifer Dougherty
Bill Clifford
Elaine Lee
Lauren Hagen
Boris Aguilar
Mi Tian
Lauren Wolfe
Ilya Shmulevich

GENERAL DYNAMICS
Information Technology

David Pot
Madelyn Reyes
Kawther Abdilleh
Fabian Seidl
Deena Bleich
Owais Shahzada

Questions?