

# The ISB Cancer Genomics Cloud

Sheila M Reynolds  
Institute for Systems Biology  
April 6<sup>th</sup> 2016

Bio IT World Conference & Expo



# ISB-CGC Team Members



Ilya Shmulevich  
Sheila Reynolds  
Phyliss Lee  
Michael Miller  
Kelly Iverson  
Abigail Hahn  
Zack Rodebaugh  
Kalle Leinonen  
Dave Gibbs  
Varsha Dhankani



Jonathan Bingham  
Nicole Deflaux  
Matt Bookman  
Jaclyn Kollar



David Pot  
Ross Casanova  
Sandeep Namburi  
Madelyn Reyes

This project has been funded in whole with Federal funds from the National Cancer Institute, National Institutes of Health, Department of Health and Human Services, under Contract No. HHSN261201400007C.



**CBIIT** Center for Biomedical Informatics  
and Information Technology



**NCI | NHGRI**

# THE CANCER GENOME ATLAS



## The Cancer Genome Atlas (TCGA)

“ ... a comprehensive and coordinated effort to accelerate our understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing”

# THE CANCER GENOME ATLAS

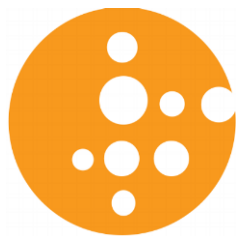


- Tissue Processing
  - High-quality sample prep, de-identified clinical data, imaging
- Research and Discovery
  - Genome sequencing & characterization centers
  - Genome data analysis centers
- Data Sharing
  - Public data repositories
  - open-access and controlled-access tiers

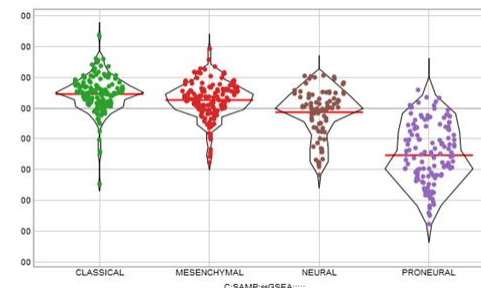
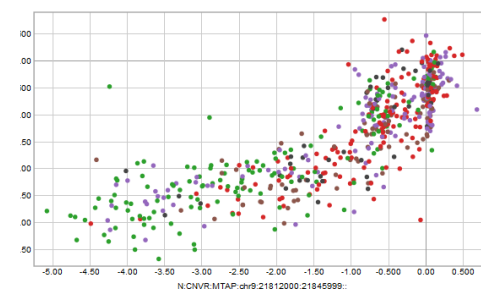
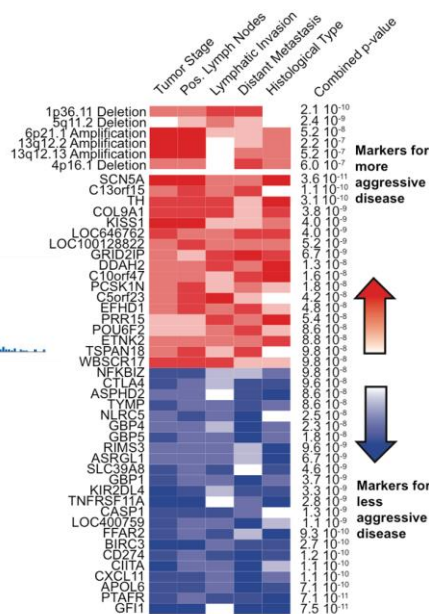
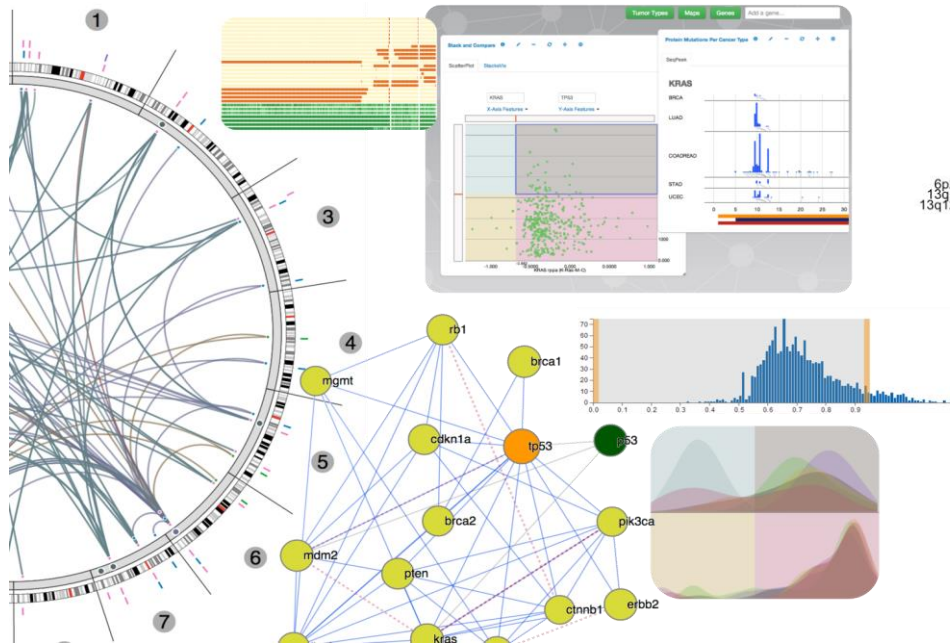
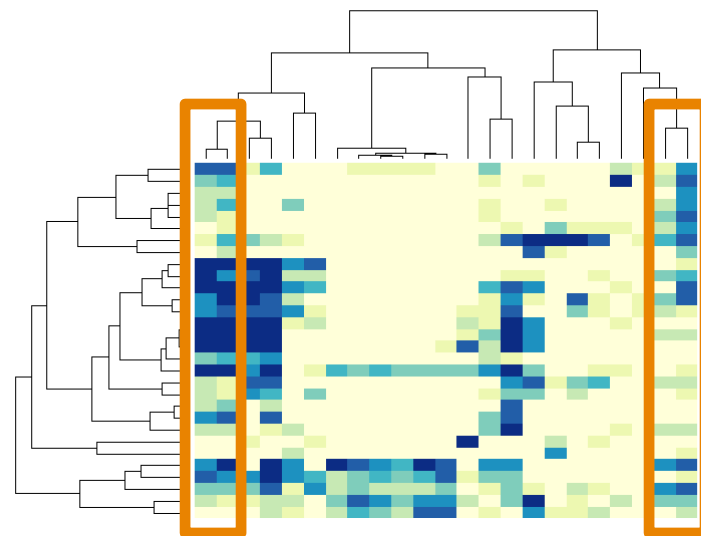
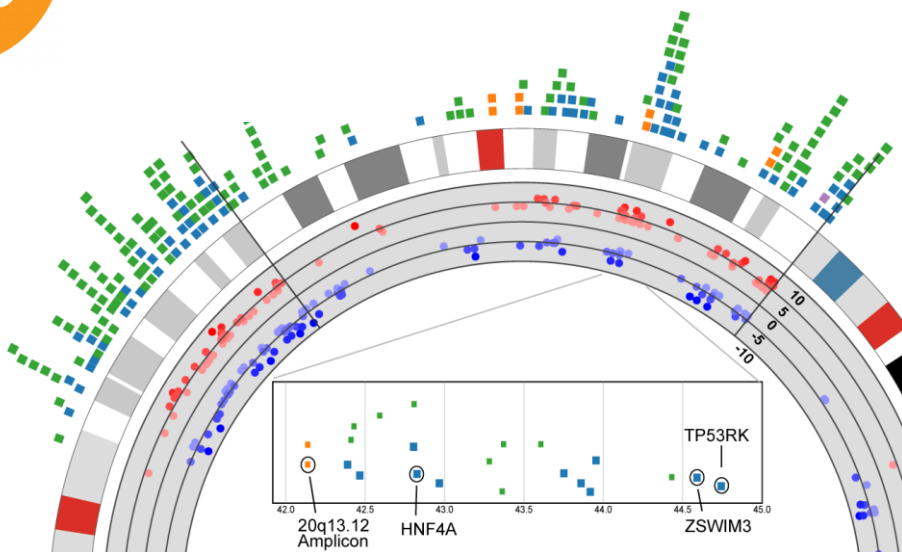
## The Cancer Genome Atlas (TCGA)

“ ... a comprehensive and coordinated effort to accelerate our understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing”

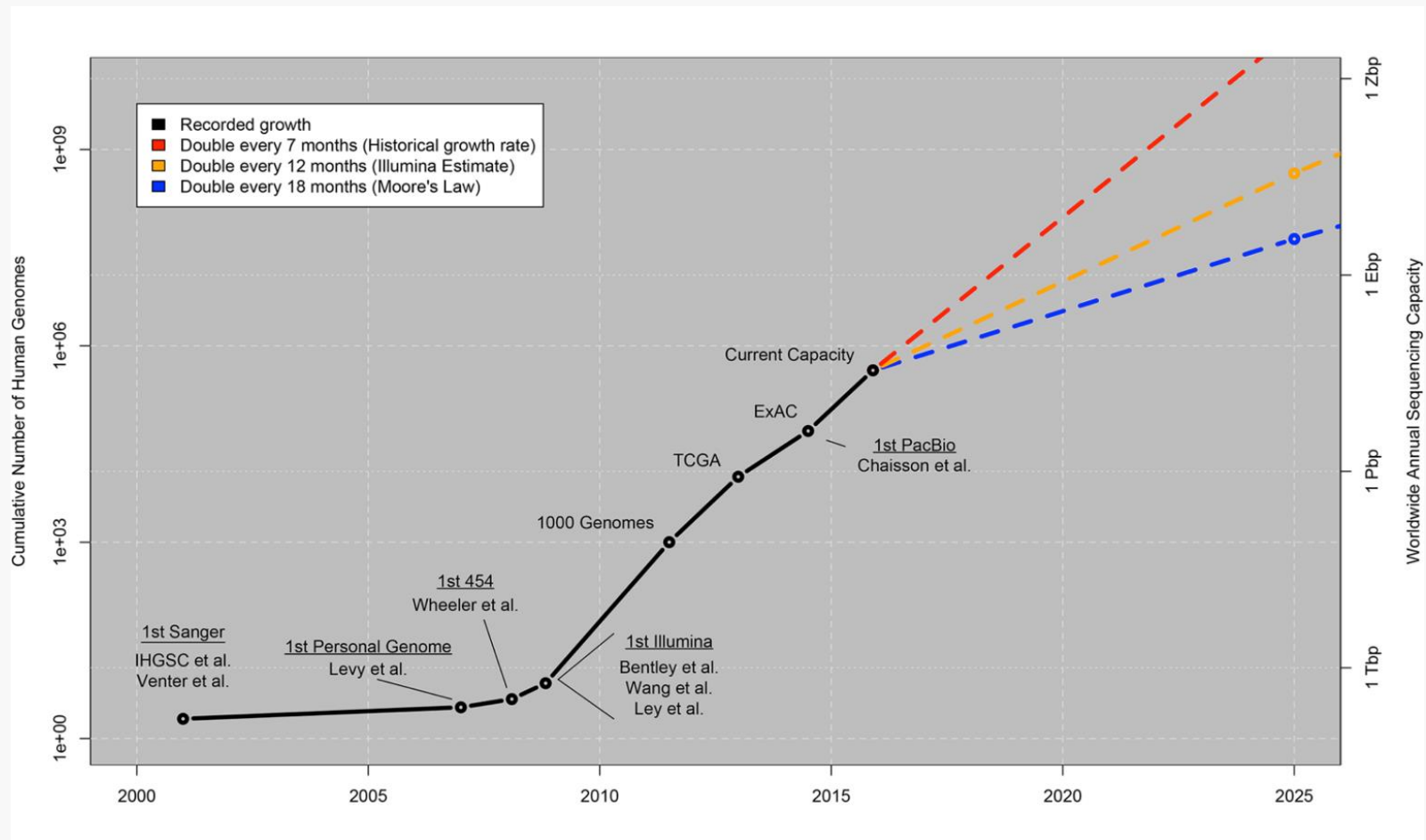




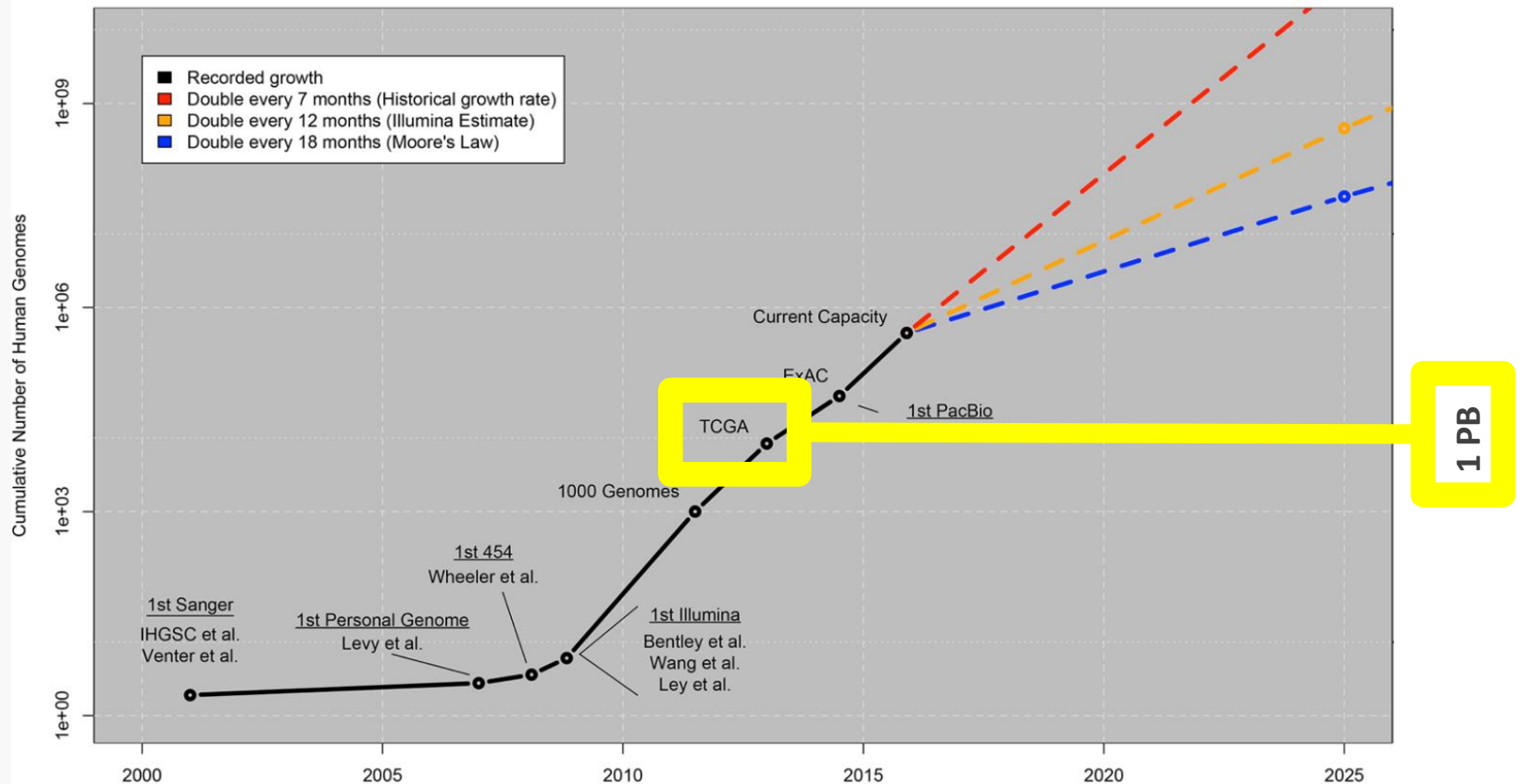
# ISB Genome Data Analysis Center



# The Challenge of Big Data



# The Challenge of Big Data





# NCI Cancer Genomics Cloud Pilots

Bringing data and computation together to create knowledge that accelerates cancer research and enables precision medicine

- **Goals:**

- Explore innovative methods for accessing and computing on genomic data
- Democratize access to NCI-generated data sets
- Cost-effective computational support to the cancer research community

- **Timeline:**

- 2015: Design and Build
- 2016: Community Evaluation



Source: <https://cbiit.nci.nih.gov/ncip/nci-cancer-genomics-cloud-pilots>



# Cancer Genomics Cloud

The ISB Cancer Genomics Cloud (ISB-CGC) is democratizing access to **TCGA** data and coupling it with unprecedented computational power to allow researchers to explore and analyze this vast data-space.

[Documentation](#)

[!\[\]\(a870788d6ed9b8fd294b7654a8c8526b\_img.jpg\) GitHub](#)

[Web App](#)

[Feedback](#)

# Cloud Paradigm Shift(s)

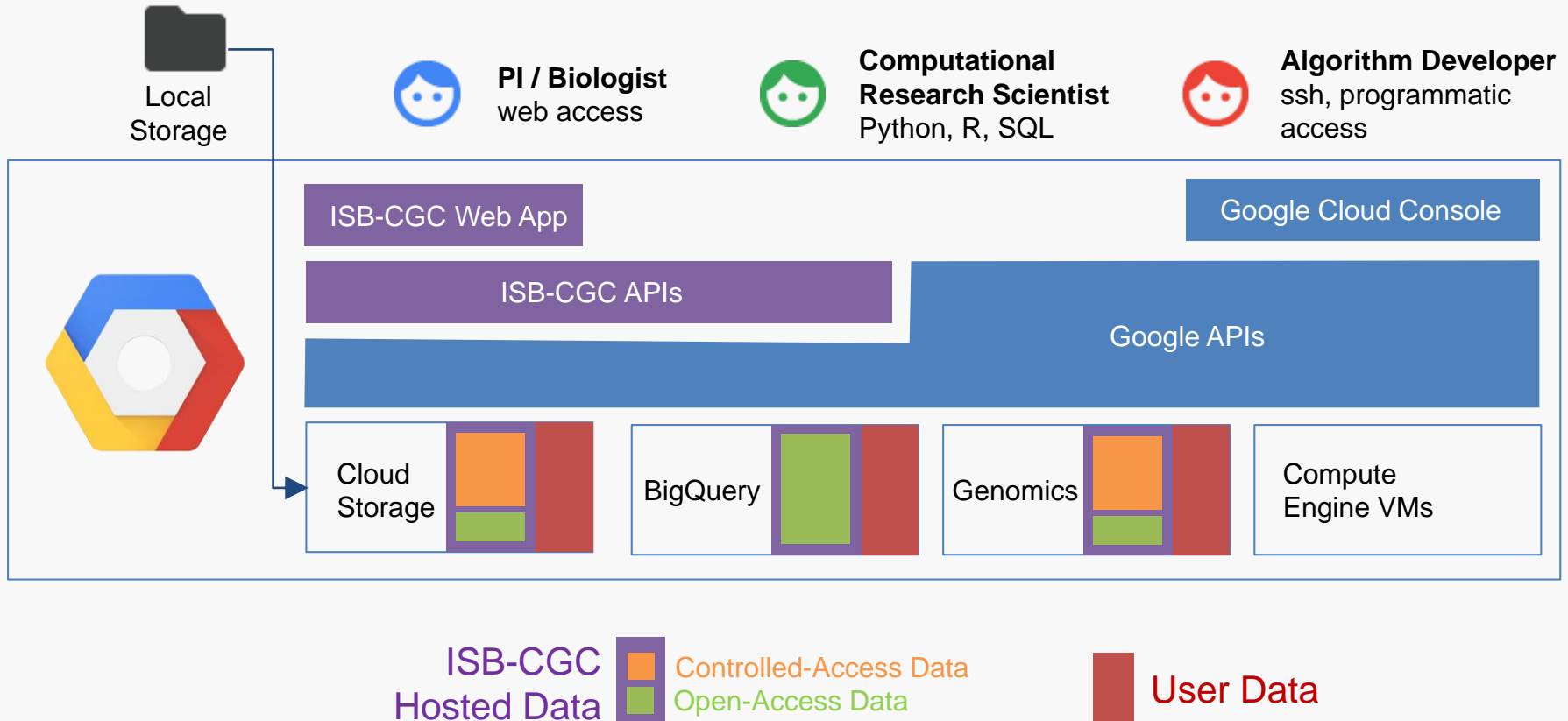
- Move data and existing pipelines to the cloud
  - all researchers access a single copy of the data
  - compute-power is “near” the data
  - pay only for minutes used
  - everyone saves time, money, and bandwidth
- Cloud-aware computing
  - rethink/redevelop approaches to fully leverage the power of the cloud
  - massively parallel, bursty, opportunistic computing

# Our Mission

... is to make TCGA data, together with tools and compute-power, available and accessible to a broad range of users using multiple access modes:

- interactive web application
- scripting languages: R, Python, SQL
- direct programmatic access

Our Approach: build an open platform that can grow and evolve to satisfy a broad range of users and use-cases





# Data as a Service

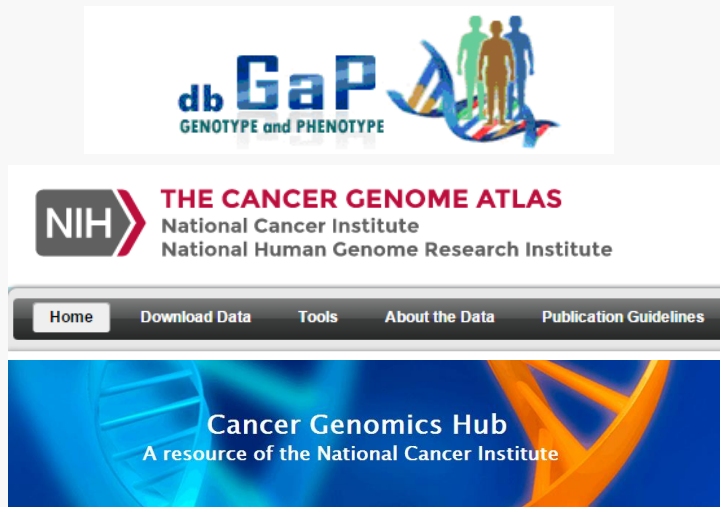
## Phase 1

- Low-level sequence and SNP data as *files* in **Cloud Storage**
- High-level data and annotations as *tables* in **BigQuery**

## Phase 2

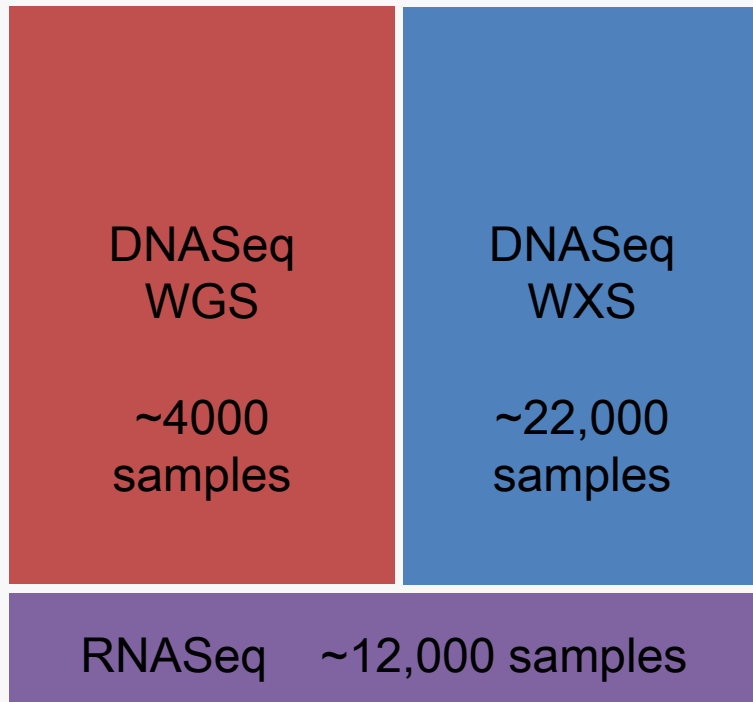
- Low-level sequence data in **Google Genomics**
- Variant calls in **Google Genomics** and **BigQuery**

# TCGA Data in the Cloud



# TCGA Size

>1 PB of sequence data  
(controlled access)

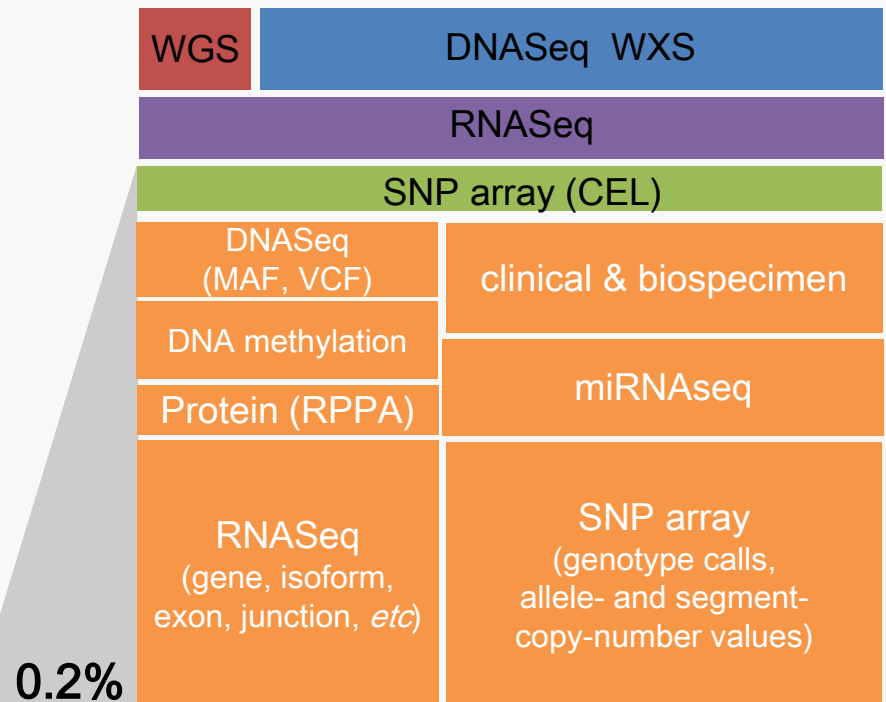


# TCGA Size & Complexity

>1 PB of sequence data  
(controlled access)



~400,000 files of  
heterogeneous data  
(mostly open-access)



Schema				
ParticipantBarcode		STRING	NULLABLE	Describe this field...
Study		STRING	NULLABLE	Describe this field...
Project		STRING	NULLABLE	Describe this field...
ParticipantUUID		STRING	NULLABLE	Describe this field...
TSSCode		STRING	NULLABLE	Describe this field...
age_at_initial_pathologic_diagnosis		INTEGER	NULLABLE	Describe this field...
anatomic_neoplas	Table Details: Somatic_Mutation_calls			
batch_number				
bcr	Schema			
clinical_M	ParticipantBarcode	STRING	NULLABLE	Describe this field...
clinical_N	Tumor_SampleBarcode	STRING	NULLABLE	Describe this field...
clinical_T	Tumor_AliquotBarcode	STRING	NULLABLE	Describe this field...
clinical_stage	Tumor_SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
colorectal_cancer	Normal_SampleBarcode	STRING	NULLABLE	Describe this field...
country	Normal_AliquotBarcode	STRING	NULLABLE	Describe this field...
vital_status	Normal_SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
days_to_birth	Study	STRING	NULLABLE	Describe this field...
days_to_death	Annotation_Transcript	STRING	NULLABLE	Describe this field...
days_to_last_know	CCLE_ONCOMAP_Total_Mutations_In_Gene	INTEGER	NULLABLE	Describe this field...
days_to_last_follo	COSMIC_Total_Alterations_In_Gene	INTEGER	NULLABLE	Describe this field...
days_to_initial_pat	Center	STRING	NULLABLE	Describe this field...
days_to_submitted	Chromosome	STRING	NULLABLE	Describe this field...
ethnicity	DNARepairGenes_Role	STRING	NULLABLE	Describe this field...
frozen_specimen_s	DbSNP_RS	STRING	NULLABLE	Describe this field...
gender	DbSNP_Val_Status	STRING	NULLABLE	Describe this field...
gleason_score_cor	DrugBank	STRING	NULLABLE	Describe this field...
histological_type	End_Position	INTEGER	NULLABLE	Describe this field...
history_of_colon_f	Entrez_Gene_Id	INTEGER	NULLABLE	Describe this field...
	GC_Content	FLOAT	NULLABLE	Describe this field...
	GENCODE_Transcript_Name	STRING	NULLABLE	Describe this field...
	GENCODE_Transcript_Status	STRING	NULLABLE	Describe this field...
	GENCODE_Transcript_Type	STRING	NULLABLE	Describe this field...
	GO_Biological_Process	STRING	NULLABLE	Describe this field...
	GO_Cellular_Componen			
	GO_Molecular_Function	Table Details: DNA_Methylation_beta		
	Gene_Type	Schema		
	Genome_Change	ParticipantBarcode	STRING	NULLABLE

Schema			
ParticipantBarcode	STRING	NULLABLE	Describe this field...
SampleBarcode	STRING	NULLABLE	Describe this field...
SampleTypeLetterCode	STRING	NULLABLE	Refer: <a href="https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm">https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm</a>
AliquotBarcode	STRING	NULLABLE	The Aliquot ID is an identifier/barcode of TCGA data. Refer: <a href="https://wiki.nci.nih.gov/display/TCGA/TCGA+bi">https://wiki.nci.nih.gov/display/TCGA/TCGA+bi</a>
Platform	STRING	NULLABLE	Refer: <a href="https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm">https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm</a>
Study	STRING	NULLABLE	TCGA disease type
Probe_Id	STRING	NULLABLE	illumina's CpG loci IDs. Refer: <a href="http://www.illumina.com/content/dam/illumina-marketing/documents/products/technote/technote_cpg_loci_identification.pdf">http://www.illumina.com/content/dam/illumina-marketing/documents/products/technote/technote_cpg_loci_identification.pdf</a>
Beta_Value	FLOAT	NULLABLE	The beta value ( $\beta$ ) is used to estimate the methylation level of the CpG locus using the ratio of intensities b

Schema			
ParticipantBarcode	STRING	NULLABLE	Describe this field...
SampleBarcode	STRING	NULLABLE	Describe this field...
SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
SampleType	STRING	NULLABLE	Describe this field...
Study	STRING	NULLABLE	Describe this field...
Project	STRING	NULLABLE	Describe this field...
SampleTypeCode	STRING	NULLABLE	Describe this field...
avg_percent_lymphocyte_infiltration	FLOAT	NULLABLE	Describe this field...
avg_percent_monocyte_infiltration	FLOAT	NULLABLE	Describe this field...
avg_percent_necrosis	FLOAT	NULLABLE	Describe this field...
avg_percent_neutrophil_infiltration	FLOAT	NULLABLE	Describe this field...
avg_percent_normal_cells	FLOAT	NULLABLE	Describe this field...
avg_percent_stromal_cells	FLOAT	NULLABLE	Describe this field...
avg_percent_tumor_cells	FLOAT	NULLABLE	Describe this field...
avg_percent_tumor_nuclei	FLOAT	NULLABLE	Describe this field...
batch_number	INTEGER	NULLABLE	Describe this field...
bcr	STRING	NULLABLE	Describe this field...
days_to_collection	FLOAT	NULLABLE	Describe this field...
days_to_sample_procurement	FLOAT	NULLABLE	Describe this field...
is_flow	STRING	NULLABLE	Describe this field...

max_perce	Schema			
max_perce	ParticipantBarcode	STRING	NULLABLE	Describe this field...
max_perce	SampleBarcode	STRING	NULLABLE	Describe this field...
max_perce	SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
max_perce	AliquotBarcode	STRING	NULLABLE	Describe this field...
max_perce	Study	STRING	NULLABLE	Describe this field...
max_perce	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...

Schema			
annotationId	INTEGER	NULLABLE	Describe this field...
annotationCategoryId	INTEGER	NULLABLE	Describe this field...
annotationCategoryName	STRING	NULLABLE	Describe this field...
annotationClassification	STRING	NULLABLE	Describe this field...
annotationNoteText	STRING	NULLABLE	Describe this field...
Study	STRING	NULLABLE	Schema
itemTypeName	STRING	NULLABLE	
ItemBarcode	STRING	NULLABLE	
AliquotBarcode	STRING	NULLABLE	
ParticipantBarcode	STRING	NULLABLE	
SampleBarcode	STRING	NULLABLE	
dateAdded	STRING	NULLABLE	
dateCreated	STRING	NULLABLE	Study
dateEdited	STRING	NULLABLE	Gene_Name
			Protein_Expression

Schema			
ParticipantBarcode	STRING	NULLABLE	Describe this field...
SampleBarcode	STRING	NULLABLE	Describe this field...
SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
AliquotBarcode	STRING	NULLABLE	Describe this field...
Study	STRING	NULLABLE	Describe this field...
Gene_Name	STRING	NULLABLE	Describe this field...
Protein_Expression	FLOAT	NULLABLE	Describe this field...
Protein_Name	STRING	NULLABLE	Describe this field...
Protein_Basename	STRING	NULLABLE	Describe this field...
_HiSeq_RSEM	TRING	NULLABLE	Describe this field...
	TRING	NULLABLE	Describe this field...
	TRING	NULLABLE	Describe this field...

Schema			
ParticipantBarcode	STRING	NULLABLE	Describe this field...
SampleBarcode	STRING	NULLABLE	Describe this field...
AliquotBarcode	STRING	NULLABLE	Describe this field...
Study	STRING	NULLABLE	Describe this field...
SampleTypeLetterCode	STRING	<b>Table Details: mRN</b>  <b>Schema</b>	
Platform	STRING		
original_gene_symbol	STRING		
HGNC_gene_symbol	STRING		
gene_id	INTEGER		
normalized_count	FLOAT	ParticipantBarcode	STRING
		SampleBarcode	STRING
		SampleTypeLetterCode	STRING

Schema				
ParticipantBarcode	STRING	NULLABLE	Describe this field...	
SampleBarcode	STRING	NULLABLE	Describe this field...	
SampleTypeLetterCode	STRING	NULLABLE	Describe this field...	
AliquotBarcode	STRING	NULLABLE	Describe this field...	
Study	STRING	NULLABLE	Describe this field...	
	STRING	NULLABLE	Describe this field...	
		NULLABLE	Describe this field...	
		NULLABLE	Describe this field...	
		NULLABLE	Describe this field...	
NG	NULLABLE	Describe this field...	NULLABLE	Describe this field...
NG	NULLABLE	Describe this field...	NULLABLE	Describe this field...
NG	NULLABLE	Describe this field...	NULLABLE	Describe this field...

Schema			
<b>ParticipantBarcode</b>	STRING	NULLABLE	Describe this field...
<b>SampleBarcode</b>	STRING	NULLABLE	Describe this field...
<b>AliquotBarcode</b>	STRING	NULLABLE	Describe this field...
<b>SampleTypeLetterCode</b>	STRING	NULLABLE	Describe this field...
<b>Study</b>	STRING	NULLABLE	Describe this field...
<b>Platform</b>	STRING	NULLABLE	Describe this field...
<b>mirna_id</b>	STRING	NULLABLE	Describe this field...
<b>mirna_accession</b>	STRING	NULLABLE	Describe this field...
<b>normalized_count</b>	FLOAT	NULLABLE	Describe this field...





### Schema

ParticipantBarcode	STRING	NULLABLE	Describe this field...
Study	STRING	NULLABLE	Describe this field...
Project	STRING	NULLABLE	Describe this field...
ParticipantUUID	STRING	NULLABLE	Describe this field...
TSSCode	STRING	NULLABLE	Describe this field...
age_at_initial_pathologic_diagnosis	INTEGER	NULLABLE	Describe this field...

### Schema

Field	Type	Nullable	Description
clinical_M	ParticipantBarcode	STRING	NULLABLE
clinical_N	Tumor_SampleBarcode	STRING	NULLABLE
clinical_T	Tumor_AlignotBarcode	STRING	NULLABLE
clinical_stage	Tumor_SampleTypeLetterCode	STRING	NULLABLE
colorectal_cancer	Normal_SampleBarcode	STRING	NULLABLE
country	Normal_AlignotBarcode	STRING	NULLABLE
vital_status	Normal_SampleTypeLetterCode	STRING	NULLABLE
days_to_birth	Study	STRING	NULLABLE
days_to_death	Annotation_Transcript	STRING	NULLABLE
days_to_last_know	CCLE_ONCOMAP_Total_Mutations_In_Gene	INTEGER	NULLABLE
days_to_last_follow	COSMIC_Total_Alterations_In_Gene	INTEGER	NULLABLE
days_to_initial_pat	Center	STRING	NULLABLE
days_to_submitted	Chromosome	STRING	NULLABLE
ethnicity	DNARepairGenes_Role	STRING	NULLABLE
frozen_specimen_s	Db SNP_RS	STRING	NULLABLE
gender	Db SNP_Val_Status	STRING	NULLABLE
gleason_score_cor	DrugBank	STRING	NULLABLE
histological_type	End_Position	INTEGER	NULLABLE
history_of_colon_f	Entrez_Gene_Id	INTEGER	NULLABLE

**Table Details:** DNA Methylation betas

Schema	ParticipantBarcode	STRING	NULLABLE	DATE
--------	--------------------	--------	----------	------

<b>ParticipantBarcode</b>	STRING	NULLABLE	Describe this field
<b>SampleBarcode</b>	STRING	NULLABLE	Describe this field
<b>SampleTypeLetterCode</b>	STRING	NULLABLE	Refer: <a href="https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm">https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm</a>
<b>AliquotBarcode</b>	STRING	NULLABLE	The Aliquot ID is an identifier/barcode of TCGA data. Refer: <a href="https://wiki.nci.nih.gov/display/TCGA/TCGA+bi">https://wiki.nci.nih.gov/display/TCGA/TCGA+bi</a>
<b>Platform</b>	STRING	NULLABLE	Refer: <a href="https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm">https://tcga-data.nci.nih.gov/datareports/codeTablesReport.htm</a>
<b>Study</b>	STRING	NULLABLE	TCGA disease type
<b>Probe_Id</b>	STRING	NULLABLE	illumina's CpG loci IDs. Refer: <a href="http://www.illumina.com/content/dam/illumina-marketing/documents/products/technotestechnote_cpg_loci_identification.pdf">http://www.illumina.com/content/dam/illumina-marketing/documents/products/technotestechnote_cpg_loci_identification.pdf</a>
<b>Beta_Value</b>	FLOAT	NULLABLE	The beta value ( $\beta$ ) is used to estimate the methylation level of the CpG locus using the ratio of intensities be

### Schema

<b>ParticipantBarcode</b>	STRING	NULLABLE	Describe this field...
<b>SampleBarcode</b>	STRING	NULLABLE	Describe this field...
<b>SampleTypeLetterCode</b>	STRING	NULLABLE	Describe this field...
SampleType	STRING	NULLABLE	Describe this field...

[illegible]

Table Details: Copy Number segments

max_percent	ParticipantBarcode	STRING	NULLABLE	Describe this field
max_percent	SampleBarcode	STRING	NULLABLE	Describe this field
max_percent	SampleTypeLetterCode	STRING	NULLABLE	Describe this field
max_percent	Study	STRING	NULLABLE	Describe this field
max_percent	Platform	STRING	NULLABLE	Describe this field
max_percent	Chromosome	STRING	NULLABLE	Describe this field
max_percent	Start	INTEGER	NULLABLE	Describe this field
max_percent	End	INTEGER	NULLABLE	Describe this field

### Schema

<b>annotationId</b>	INTEGER	NULLABLE	Describe this field.
<b>annotationCategoryId</b>	INTEGER	NULLABLE	Describe this field.
<b>annotationCategoryName</b>	STRING	NULLABLE	Describe this field.
<b>annotationClassification</b>	STRING	NULLABLE	Describe this field.

Table Details: Protein RPPA data

Study	STRING	NUMERIC	Schema			
itemName	STRING	NUMERIC	ParticipantBarcode	STRING	NULLABLE	Describe this field...
itemBarcode	STRING	NUMERIC	SampleBarcode	STRING	NULLABLE	Describe this field...
AliquotBarcode	STRING	NUMERIC	SampleTypeLetterCode	STRING	NULLABLE	Describe this field...
dateAdded	STRING	NUMERIC	AliquotBarcode	STRING	NULLABLE	Describe this field...
dateCreated	STRING	NUMERIC	Study	STRING	NULLABLE	Describe this field...
dateEdited	STRING	NUMERIC	Gene_Name	STRING	NULLABLE	Describe this field...
			Protein_Expression	FLOAT	NULLABLE	Describe this field...
			Protein_Name	STRING	NULLABLE	Describe this field...
			Protein_Basename	STRING	NULLABLE	Describe this field...

Table Details: mRNA UNC HiSeq **RSEM**

Schema				TRIM	NULLABLE	Describe this field...
ParticipantBarcode	STRING	NULLABLE	Describe this field...	TRIM	NULLABLE	Describe this field...

Table Details: mRNA BCGSC HiSeq RPKM

Platform	STRING
HGNC_gene_symbol	STRING
gene_id	INTEGER
normalized_count	FLOAT

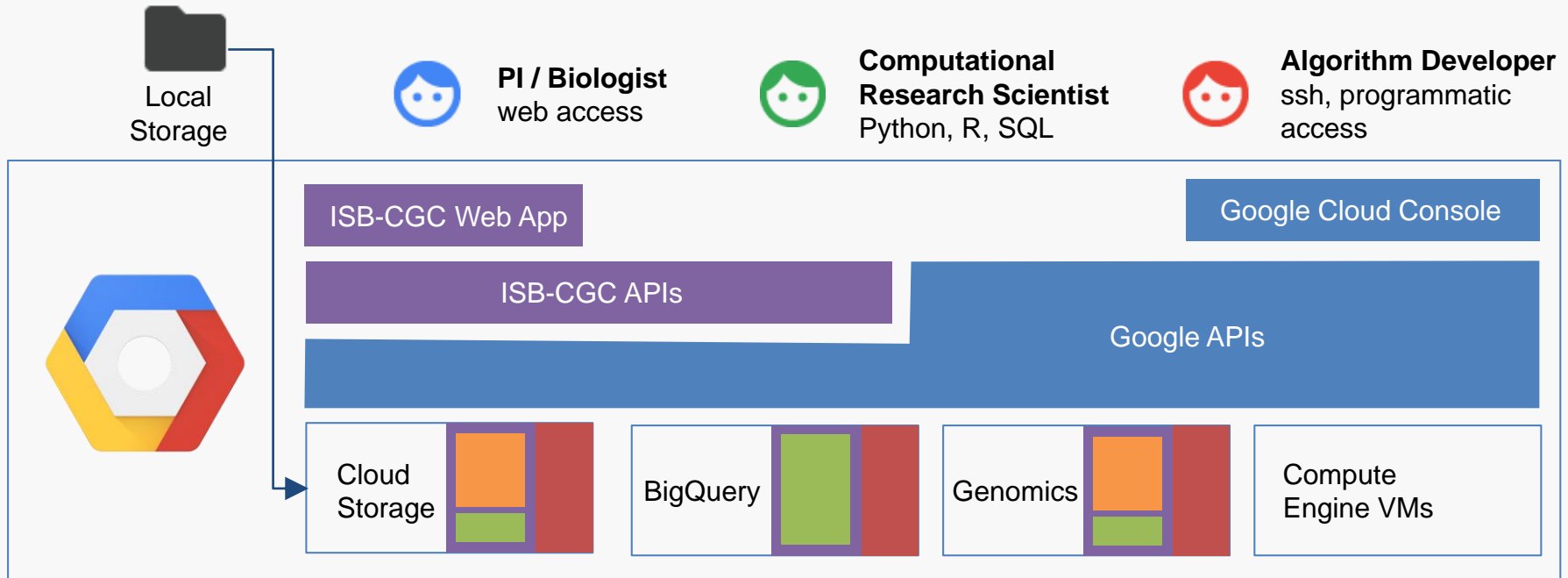
### Schema

ParticipantBarcode	STRING	NULLABLE	Describe this field.
SampleBarcode	STRING	NULLABLE	Describe this field.
SampleTypeLetterCode	STRING	NULLABLE	Describe this field.
AliquotBarcode	STRING	NULLABLE	Describe this field.
Study	STRING	NULLABLE	Describe this field.

Table Details: miRNA\_expression

Schema			
ParticipantBarcode	STRING	NULLABLE	Describe this field.
SampleBarcode	STRING	NULLABLE	Describe this field.
AliquotBarcode	STRING	NULLABLE	Describe this field.
SampleTypeLetterCode	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.





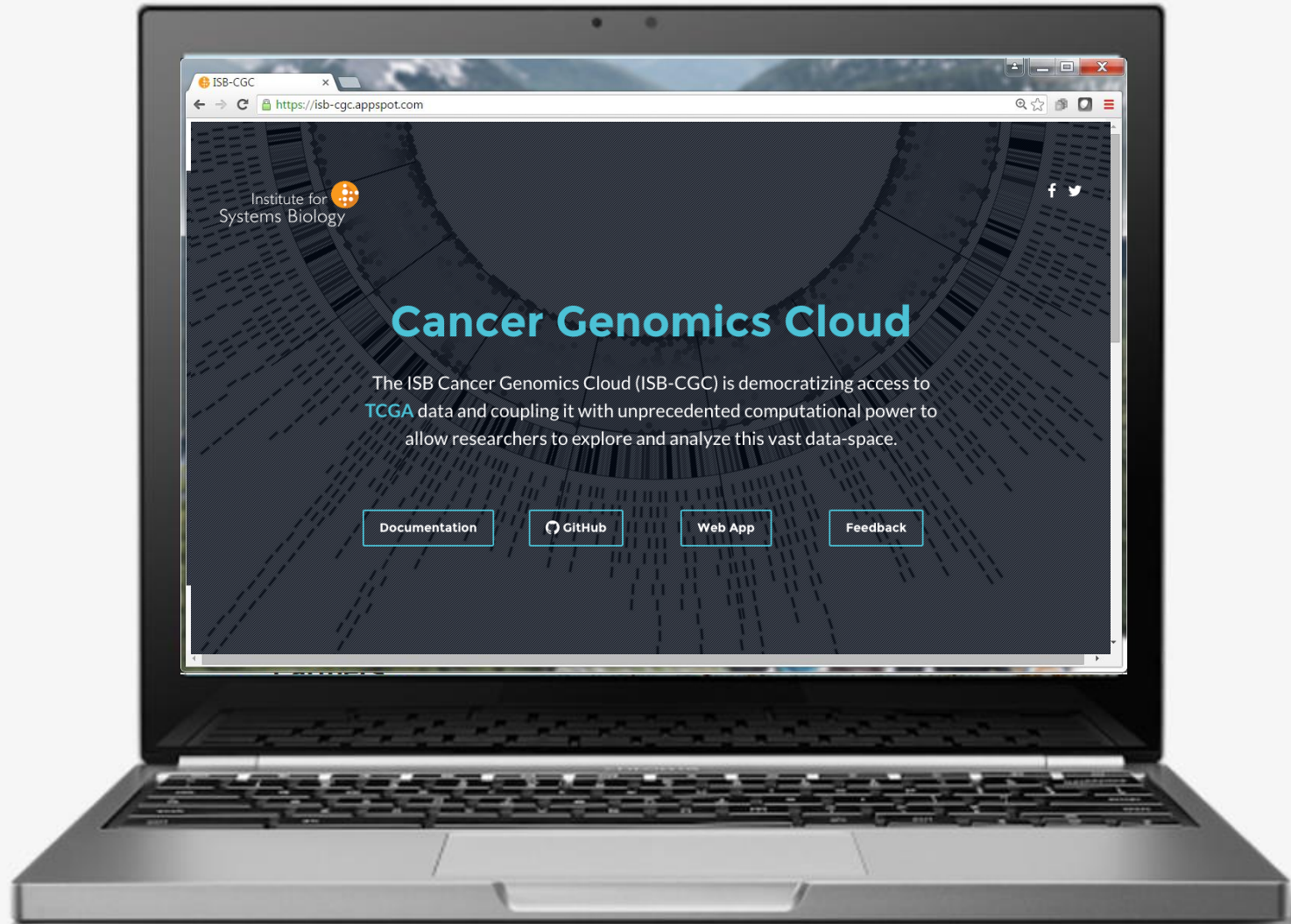
ISB-CGC  
Hosted Data



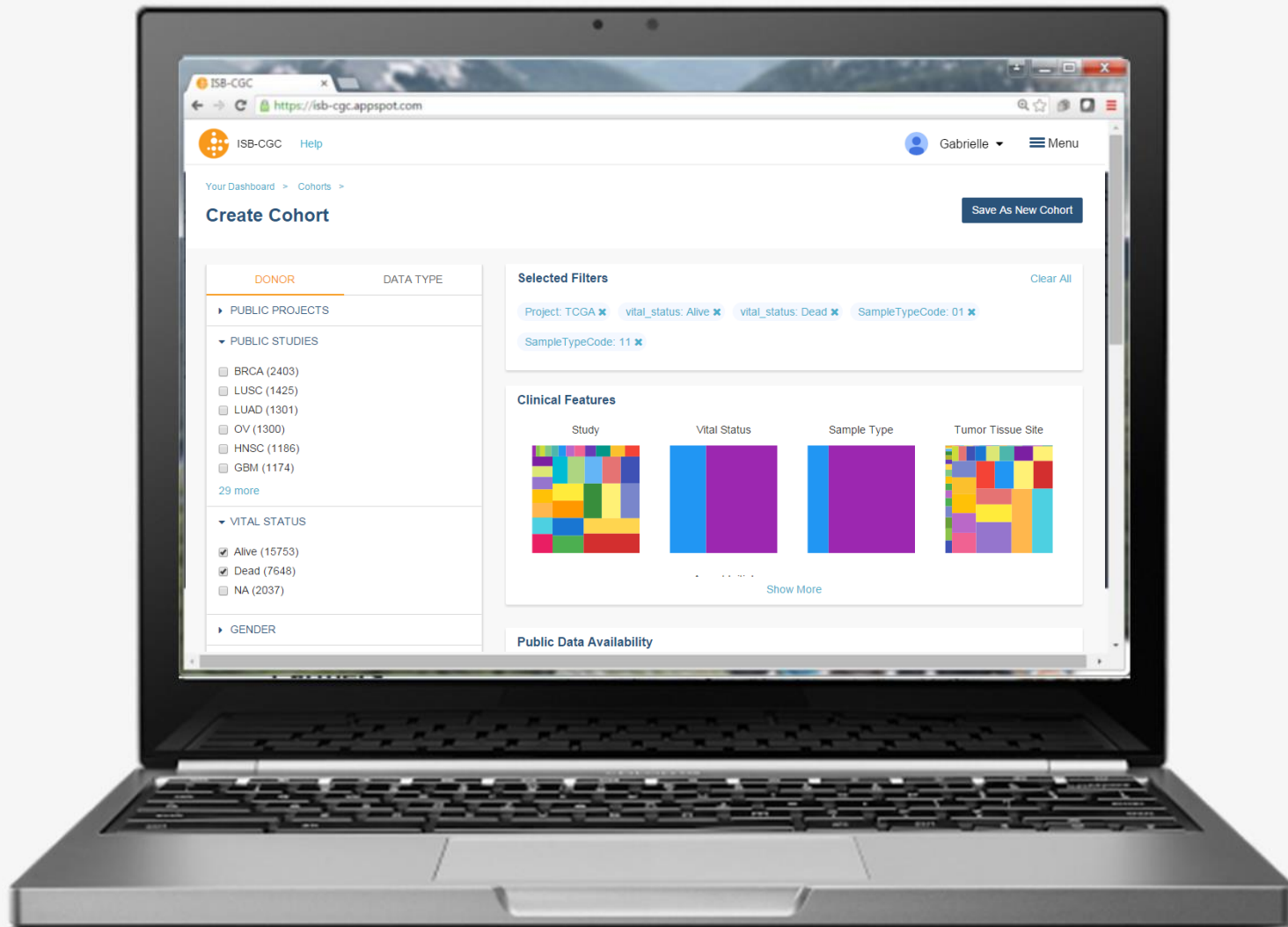
Controlled-Access Data  
Open-Access Data

User Data

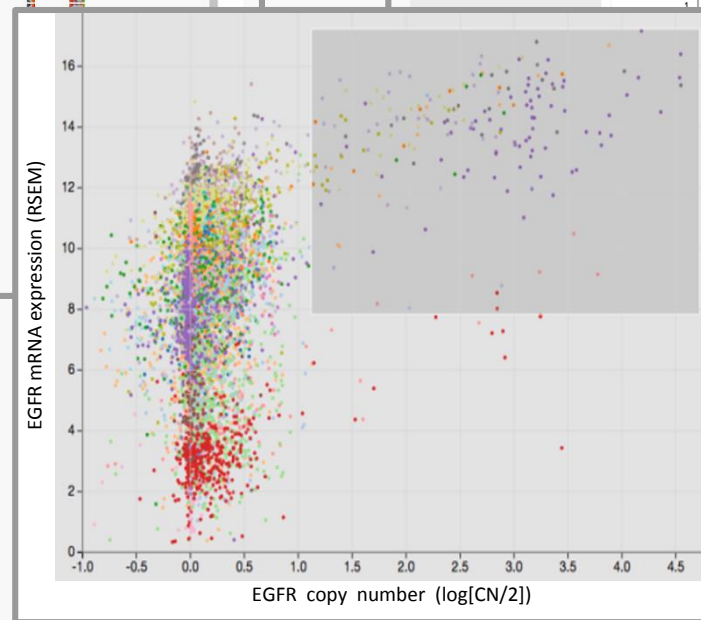
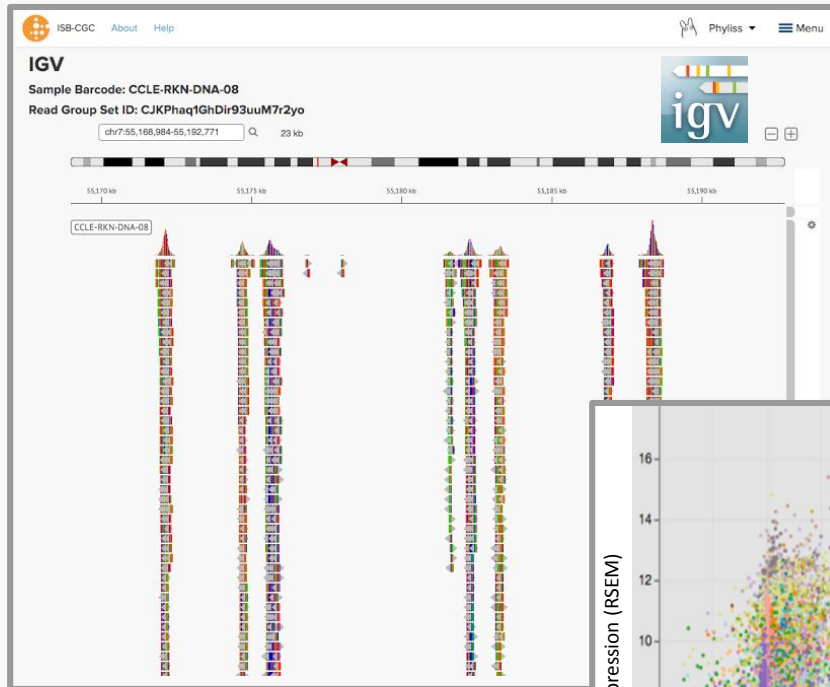
# Web access for the PI / Biologist:



# custom, interactive cohort definition



# Interactive data visualizations:





# Python, R, and SQL for the Computational Scientist:


IP[y]: IPython  
Interactive Computing



SQL







## ISB Cancer Genomics Cloud

The ISB-CGC is providing access to TCGA data and computation on the Google Cloud Platform.  
<http://www.isb-cgc.org>

[Repositories](#) [People](#) 32 [Teams](#) 5 [Settings](#)

### ISB-CGC-Webapp

JavaScript ★ 0 1

ISB CGC Webapp

Updated 22 hours ago



### ISB-CGC-data-proc

Python ★ 0 1

code for uploading cancer data into GCS and BigQuery

Updated 23 hours ago

### examples-R

HTML ★ 4 2

Analysis examples based on the ISB-CGC hosted TCGA data, using R and R Markdown.

Updated 23 hours ago


### examples-Python

★ 7 2

Analysis examples based on the ISB-CGC hosted TCGA data, using Python and IPython Notebooks.

Updated 3 days ago





## ISB Cancer Genomics Cloud

The ISB-CGC is providing access to TCGA data and computation on the Google Cloud Platform.  
<http://www.isb-cgc.org>

Repositories   People 32   Teams 5   Settings

### ISB-CGC-Webapp

ISB CGC Webapp  
Updated 22 hours ago

### ISB-CGC-data-proc

code for uploading cancer data into GCS and BigQuery  
Updated 23 hours ago

### examples-R

Analysis examples based on the ISB-CGC hosted TCGA data, using R and R Markdown.  
Updated 23 hours ago

### examples-Python

Analysis examples based on the ISB-CGC hosted TCGA data, using Python and IPython Notebooks.  
Updated 3 days ago

### examples-R

Analysis examples based on the ISB-CGC hosted TCGA data, using R and R Markdown.

To install

```
require(devtools) || install.packages("devtools")
install_github("isb-cgc", "examples-R", build_vignettes=TRUE)
```

To view and run the vignettes

```
help(package="ISBGCExamples")
```

There are vignettes for each TCGA data type, and more elaborate examples involving analyzing genomic data, correlating gene expression and methylation, and correlating protein and mRNA levels.

The vignettes as **R-markdown** can be found in the `examples/Rmarkdown` directory, which can serve as examples of using built-in BigQuery functions like Pearson correlation, or even how to implement more complex functions like Spearman's correlation. Queries can be simple character vectors, or standalone files. Results are returned as data frames using the `bigquery` package to interact with the servers.

The **SQL** files used in the vignettes can be found at `examples/Rmarkdown`. These are parsed and dispatched with arguments using the `DisplayAndDispatchQuery` function, found in the file of the same name in `examples-RR`.



# Bioconductor

OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

## The Comprehensive R Archive Network

### R Client for Google Genomics API

Bioconductor version: Release (3.2)

Provides an R package to interact with the Google Genomics API.

Author: Cassie Doll [aut], Nicole Deflaux [aut], Siddhartha Bagaria [aut, cre]


### bigquery: An Interface to Google's BigQuery API

Easily talk to Google's BigQuery database from R.

Version: 0.1.0  
Depends: R (≥ 3.1.0)  
Imports: [http](#), [jsonlite](#), [assertthat](#), [R6](#) (≥ 2.0.0), [dplyr](#) (≥ 0.3.0)  
Suggests: [testthat](#)  
Published: 2015-01-13  
Author: Hadley Wickham [aut, cre], RStudio [cph]







## ISB Cancer Genomics Cloud

The ISB-CGC is providing access to TCGA data and computation on the Google Cloud Platform.  
<http://www.isb-cgc.org>

Repositories   People 32   Teams 5   Settings

### ISB-CGC-Webapp

ISB CGC Webapp  
Updated 22 hours ago

### ISB-CGC-data-proc

code for uploading cancer data into GCS and BigQuery  
Updated 23 hours ago


### examples-R

Analysis examples based on the ISB-CGC hosted TCGA data, using R and R Markdown.  
Updated 23 hours ago

### examples-Python

Analysis examples based on the ISB-CGC hosted TCGA data, using Python and IPython Notebooks.  
Updated 3 days ago

IP[y]: IPython  
Interactive Computing



## CLOUD DATALAB<sup>BETA</sup>

An easy to use interactive tool for large-scale data exploration, analysis, and visualization.

README.md

### examples-Python

This repository contains analysis examples based on the ISB-CGC hosted TCGA data in BigQuery, using Python, IPython Notebooks, and Google Cloud Datalab.

#### Where to start?

You can find an overview of the BigQuery tables in this notebook and from there, we suggest that you look at the two "Creating TCGA cohorts" notebooks (part 1 and part 2) which describe and make use of the Clinical and Biospecimen tables. From there you can delve into the various molecular data tables as well as the Annotations table. For now these sample notebooks are intentionally relatively simple and do not do any analysis that integrates data from multiple tables but once you have a grasp of how to use the data, developing your own more complex analyses should not be difficult. You could even contribute an example back to our github repository! You are also welcome to submit bug reports, comments, and feature-requests as github issues.

#### How to run the notebooks

1. Launch your own Cloud Datalab instance in the cloud or run it locally.
2. Work through the introductory notebooks that are pre-installed on Cloud Datalab.
3. Run `git clone https://github.com/isb-cgc/examples-python.git` on your local file system to download the notebooks.
4. Import the ISB-CGC notebooks into your Cloud Datalab instance by navigating to the notebook list page and uploading them.

If you are running in the cloud, be sure to shut down Cloud Datalab when you are no longer using it. Shut down instructions and other tips are [here](#).



## Copy Number segments

The goal of this notebook is to in

This table contains all available T  
Genome Wide SNP6 array, as of  
recent archives (egbroad.mit.  
types was downloaded from the  
Each of these segmentation files  
During ETL the sample identifier  
the SDRF file in the associated n

In order to work with BigQuery,  
the name(s) of the table(s) you a

```
import gcp.bigquery as bq
cn_BQtable = bq.Table
```

From now on, we will refer to thi  
table name each time.

Let's start by taking a look at the

Bigquery schema --ta

name	type
ParticipantBarcode	STRING
SampleBarcode	STRING
SampleTypeLetterCode	STRING
AliquotBarcode	STRING
Study	STRING
Platform	STRING
Chromosome	STRING
Start	INTEGER
End	INTEGER
Num_Probes	INTEGER
Segment_Mean	FLOAT

Unlike most other molecular dat  
microRNAs, this data is produce  
sizes and positions of these segn

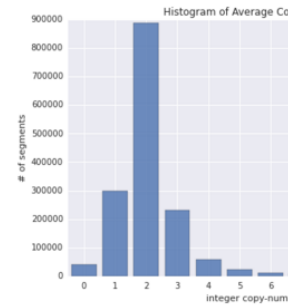
Now we'll use matplotlib to create some simple visual

```
import numpy as np
import matplotlib.pyplot as plt
```

For the segment means, let's invert the log-transform

```
%%sql --module getCNhist
SELECT
  lin_bin,
  COUNT(*) AS n
FROM (
  SELECT
    Segment_Mean,
    (2.*POW(2,Segment_Mean)) AS lin_
    INTEGER(((2.*POW(2,Segment_Mean)
  FROM
    st
  WHERE
    (End-Start+1)>1000 AND SampleLe
  GROUP BY
    lin_bin
  HAVING
    (n > 2000 )
  ORDER BY
    lin_bin ASC
```

```
CNhist = bq.Query(getCNhist,t=cn_BQ
bar_width=0.80
plt.bar(CNhist['lin_bin']+0.1,CNhist
plt.xticks(CNhist['lin_bin']+0.5,CN
plt.title('Histogram of Average Copy
plt.ylabel('# of segments');
plt.xlabel('integer copy-number');
```



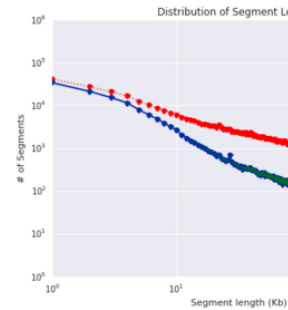
The histogram illustrates that the vast majority of the  
either side representing deletions (left) and amplificat

```
bin
ORDER BY
bin ASC

%%sql --module getSLhist_1k_amp
SELECT
  bin,
  COUNT(*) AS n
FROM (
  SELECT
    (END-Start+1) AS segLength,
    INTEGER((END-Start+1)/1000) AS b
  FROM
    st
  WHERE
    (END-Start+1)<1000000 AND SampleLe
  GROUP BY
    bin
  ORDER BY
    bin ASC
```

```
SLhistDel = bq.Query(getSLhist_1k_de
SLhistAmp = bq.Query(getSLhist_1k_am
```

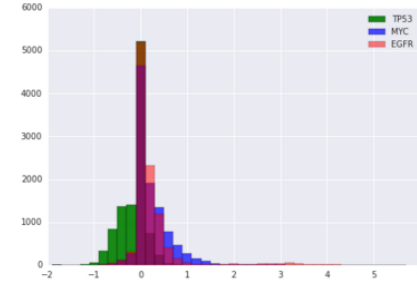
```
plt.plot(SLhist_1k['bin'],SLhist_1k[
plt.plot(SLhistDel['bin'],SLhistDel[
plt.plot(SLhistAmp['bin'],SLhistAmp[
plt.xscale('log');
plt.yscale('log');
plt.xlabel('Segment length (Kb)');
plt.ylabel('# of Segments');
plt.title('Distribution of Segment L
```



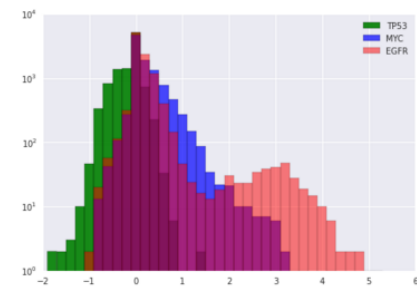
The amplification and deletion distributions are nearly  
from this graph that a majority of the segments less th  
lengths >100Kb are copy-number neutral.

And now we'll take a look at histograms of the average copy-number for these three genes. TP53 (in green) shows a significant number of partial deletions (CN<0), while MYC (in blue) shows some partial amplifications -- more frequently than EGFR, while EGFR (pale red) shows a few extreme amplifications (log2(CN/2) > 2). The final figure shows the same histograms on a semi-log plot to bring up the rarer events.

```
binwidth = 0.2
binvals = np.arange(-2+(binwidth/2.), 6-(binwidth/2.), binwidth)
plt.hist(tp53CN['avgCN'],bins=binvals,normed=False,color='green',alpha=0.9,label='TP53');
plt.hist(mycCN ['avgCN'],bins=binvals,normed=False,color='blue',alpha=0.7,label='MYC');
plt.hist(egfrCN['avgCN'],bins=binvals,normed=False,color='red',alpha=0.5,label='EGFR');
plt.legend(loc='upper right');
```



```
plt.hist(tp53CN['avgCN'],bins=binvals,normed=False,color='green',alpha=0.9,label='TP53');
plt.hist(mycCN ['avgCN'],bins=binvals,normed=False,color='blue',alpha=0.7,label='MYC');
plt.hist(egfrCN['avgCN'],bins=binvals,normed=False,color='red',alpha=0.5,label='EGFR');
plt.yscale('log');
plt.legend(loc='upper right');
```

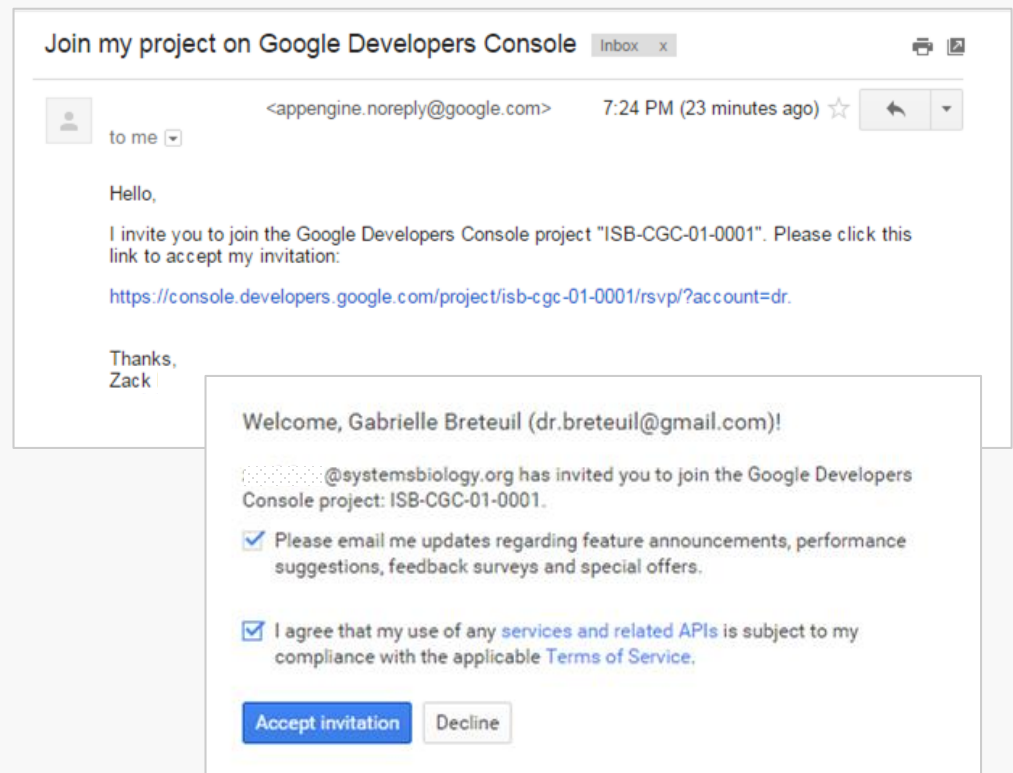




# Programmatic access for the Algorithm Developer:

Your own Google Cloud Project, with automatic access to:

- Cloud Storage
- BigQuery
- Google Genomics
- Compute Engine
- Container Engine
- Cloud Dataflow



# The ISB-CGC API provides programmatic access to the same functionality as the web-app and more:

## Cloud Endpoints API (backed by App Engine)

- authenticate from the command-line
- make requests to Endpoints API, *eg*:
  - get list of my cohorts
  - get cohort details
  - save a new cohort
  - get list of data files

```
plantain:~ kiverson$ python isb_curl.py https://isb-cgc.appspot.com/v1/cohorts_list?cohort_id=12
{
  "count": "2",
  "items": [
    {
      "name": "OV 80+",
      "filter_value": "OV",
      "comments": "None",
      "id": "12",
      "perm": "OWNER",
      "filter_name": "Study",
      "last_date_saved": "2015-11-19 00:05:12",
      "email": "kiverson@systemsbiology.org",
      "kind": "cohort_api#cohortsItem"
    },
    {
      "name": "OV 80+",
      "filter_value": "Over_80",
      "comments": "None",
      "id": "12",
      "perm": "OWNER",
      "filter_name": "Study",
      "last_date_saved": "2015-11-19 00:05:12",
      "email": "kiverson@systemsbiology.org",
      "kind": "cohort_api#cohortsItem"
    }
  ]
}
```

The screenshot shows the Google APIs Explorer interface. At the top is the Google logo and a search bar. Below it, the 'APIs Explorer' section is active, showing a list of services. The 'cohort\_api' service is selected, and its endpoints are listed in a table. The table has two columns: the endpoint name and its description. The endpoints listed are:

Endpoint	Description
cohort_api.cohort_endpoints.cohorts.cohort_patients_samples_list	Takes a cohort id as a required parameter and returns information about the participants and samples in a particular cohort. Authentication is required. User must have either READER or OWNER permissions on the cohort.
cohort_api.cohort_endpoints.cohorts.datafilenamekey_list_from_cohort	Takes a cohort id as a required parameter and returns cloud storage paths to files associated with all the samples in that cohort, up to a default limit of 10,000 files. Authentication is required. User must have READER or OWNER permissions on the cohort.
cohort_api.cohort_endpoints.cohorts.datafilenamekey_list_from_sample	Takes a sample barcode as a required parameter and returns cloud storage paths to files associated with that sample.

# Workflows & Pipelines on VM Clusters

- Grid Engine / Elasticcluster
  - use ISB-CGC API combined with Grid Engine running on a Compute Engine (GCE) cluster
- Common Workflow Language (CWL)
  - provision and configure GCE VM for use with CWL
  - create and run CWL workflows
- Kubernetes + GKE (Container Engine)
  - command line interface for launching containerized workflows
- Google Genomics “Pipelines API”
  - released in February
  - Docker-based
  - encapsulates VM-provisioning and data-staging

## Table Of Contents

- The ISB Cancer Genomics Cloud
  - Contents

## Next topic

About the ISB Cancer Genomics Cloud

## This Page

Show Source

## Quick search

Enter search terms or a module, class or function name.

# The ISB Cancer Genomics Cloud

Welcome to the ISB-CGC Documentation on Read the Docs.

Here you will find information describing the features of the ISB-CGC platform, tips on how to use it, and details about the data that we are hosting on the Google Cloud Platform.

The diagram illustrates the ISB-CGC platform architecture. It shows a central layer with 'ISB-CGC Web App' and 'ISB-CGC APIs' (purple) and 'Google APIs' (blue). Below this is a layer of services: 'Cloud Storage', 'BigQuery', 'Genomics', and 'Compute Engine VMs'. A 'Local Storage' icon is connected to 'Cloud Storage'. 'Google Cloud Console' is shown as a management interface. Above the diagram, three user roles are listed with their access methods: 'PI / Biologist' (web access), 'Computational Research Scientist' (Python, R, SQL), and 'Algorithm Developer' (ssh, programmatic access). A legend at the bottom identifies data types: 'ISB-CGC Hosted Data' (Controlled-Access Data in orange, Open-Access Data in green), and 'User Data' in red.

**Local Storage**

**PI / Biologist**  
web access

**Computational Research Scientist**  
Python, R, SQL

**Algorithm Developer**  
ssh, programmatic access

**ISB-CGC Web App**

**ISB-CGC APIs**

**Google APIs**

**Google Cloud Console**

**Cloud Storage**

**BigQuery**

**Genomics**

**Compute Engine VMs**

**ISB-CGC Hosted Data**

- Controlled-Access Data
- Open-Access Data

**User Data**

The ISB-CGC aims to serve the needs of a broad range of cancer researchers ranging from scientists or clinicians who prefer to use an interactive web-based application to access and explore the rich TCGA dataset, to computational scientists who want to write their own custom scripts using languages such as R or Python, accessing the data through APIs, and to algorithm developers who wish to spin up thousands of virtual machines to analyze hundreds of terabytes of sequence data.

This documentation is a work-in-progress, please let us know how we can improve it! [feedback@isb-cgc.org](mailto:feedback@isb-cgc.org)

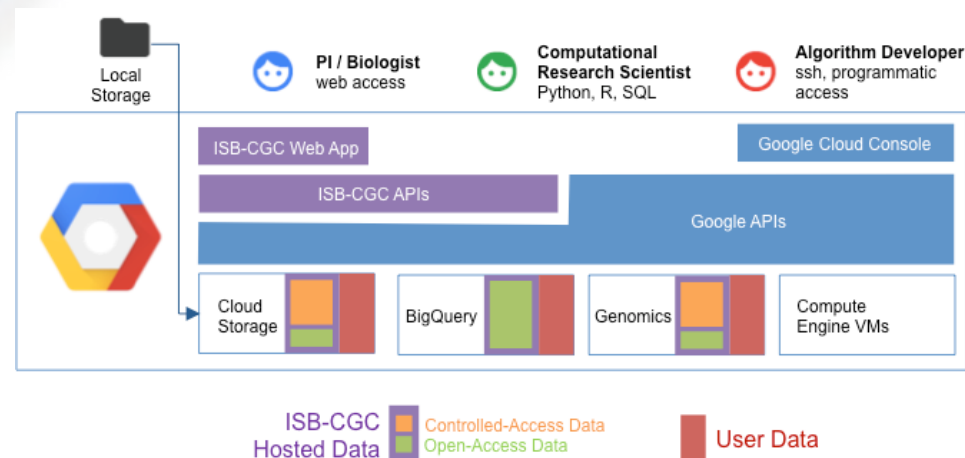
– the ISB-CGC team

v: latest

# ISB Cancer Genomics Cloud

# Thank You

[www.isb-cgc.org](http://www.isb-cgc.org)  
[info@isb-cgc.org](mailto:info@isb-cgc.org)



# ISB Cancer Genomics Cloud

# Questions?

[www.isb-cgc.org](http://www.isb-cgc.org)  
[info@isb-cgc.org](mailto:info@isb-cgc.org)

