# Poster 1067: The ISB Cancer Gateway in the Cloud (ISB-CGC): Access, explore and analyze large-scale cancer data through the Google Cloud



### Abstract

- Rapid growth of cancer data in recent decades has made data discovery and analysis difficult for cancer researchers. ISB-CGC's mission is to democratize access to large NCI cancer datasets.
- Data from NCI genomic and proteomic projects such as TCGA, TARGET, and CPTAC are populated into hundreds of tabular compute ready tables of mutations, gene expression, and protein abundance, which enable integrative data analysis in the cloud via SQL.
- The data can be accessed affordably from Google Cloud VMs where researchers can develop analysis pipelines in **Python, R, and Bioconductor**.
- BigQuery analyses are inexpensive and rapid even when scaled to petabyte sized inputs, for example we ran 6.6 billion correlations in 2.5 hours with a total cost of about one dollar.
- In this poster, we highlight some of the features of ISB-CGC and the hosted data, and present a collaborative example of using the cloud to **analyze cytogenetics data**.



sourced from the Genomic Data Commons. We populate the program metadata through the GDC API including subject clinical data, experimental strategies, etc. We then use these metadata to aggregate derived processed tables such as RNAseq, DNA methylation, etc.

	aliquot barcode	gene name 🔻	gene type		unstranded	-	fokm unstrar	nded 🔻	sample type	name 🔻		primary site	• •			
	TCGA-0R-A5J1-01A-11R-A29S	AKT2	protein coding			25902	35.	5365	Primary Tumo	or		Adrenal glar	nd			
	TCGA-0R-A5J1-01A-11R-A29S	IL 10RB	protein coding			37	0	1784	Primary Tumo	or		Adrenal glar	nd			
	TCGA-0R-A5J1-01A-11R-A29S-	CITA	protein coding		9305 80.7183		7183	Primary Tumor			Adrenal gland					
	TCGA-0R-A5.11-01A-11R-A29S-	CRBN	protein coding			1352 3.656		3 656	Primary Tumor			Adrenal gland				
	TCGA-OR-A5 11-01A-11R-A29S	C3orf33	protein_coding			189	2	0007	Primary Tum	or		Adrenal glar	nd nd			
	TCGA-OR-A5J1-01A-11R-A29S	PRPS2	protein_coding			1092	8	6382	Primary Tum	or		Adrenal glar	nd			
	TCGA-OR-A5J1-01A-11R-A29S	ZW10	protein_coding			573	4 2871 Primary Tu		Primary Tumo	iumor		Adrenal gland				
	TCGA-OR-A5J1-01A-11R-A29S	CYP2C8	protein_coding	protein coding		14	0 1188 Primary Tum		harv Tumor		Adrenal gland					
	TCGA-OR-A5J1-01A-11R-A29S	AL353751.1	IncRNA	-	1	Clinical/	1	Gene	Somatic	1	miRNA	DNA	Protein		Glycoproteom	Pho
	TCGA-OR-A5J1-01A-11R-A29S	PDK1	protein_coding		Projects GDC Metadata	Biospecime	File Metadata	Expression	Mutation	Copy Number	Expression	Methylation	Expression*	Acetylome	e	
	TCGA-OR-A5J1-01A-11R-A29S	ABHD14A	protein_coding		APOLLO	1	- V				-		1			
	TCGA-0R-A5J1-01A-11R-A29S	SUGT1-DT	IncRNA		CCLE	1	1	1		1	7					
	TCGA-0R-A5J1-01A-11R-A29S	OTUD5	protein_coding		CDDP EAGLE CGCI			1			0					
	TCGA-0R-A5J1-01A-11R-A29S	TNFSF10	protein_coding		CMI CPTAC	1		1			1					
					CTSP Exceptional Responders	1	1	1								
					FM	1	1	-								
				2	GENIE HCMI			1	1	1						
				U O	MATCH	1	1	1	1							
					MP2PRT	1	1			1						
C	SR CCC has	te a wide	vorioty		NCICCR	1		1							-	-
		is a wide	variety		ORGANOID	1	1	1								
					REBC		1				/					
	of derived a	and annot	ation		TCGA		1	1			1	1	1			
data our cataloque is always					TRIO	1	1									
					WCDT		1	1				-	-			-
		alogue le l	annayo		PDC metadata		1									
		andina			APOLLO	1	1									
	exha	anung.			CBTTC	1		2			e e e e e e e e e e e e e e e e e e e		1	-		-
					CPTAC	1	1							1	1	
Requests welcomed.				2	Georgetown Proteomics Research Program	1	1									
					ICPC	1	1						1			
					Tissue Biopsies	1	1									
					TCGA	1	1		_		v		-			-
					LITAN											
					HTAN Pancancer Atlas			1	1		1	1	1			-
					HTAN Pancancer Atlas Reactome	1			1	1	1	1	1			

An example of the TCGA RNA expression table in our ecosystem. Data can be browsed in a columnar format, similar to Excel. You can also query the tables in SQL to perform operations such as filtering, set operations, or overview level statistics such as mean, minimum, and maximum.

 Synthetic Lethality Project
 ✓
 ✓

 Mitelman DB
 ✓
 ✓

 \* GDC protein expression is RPPA; PDC protein e ✓ pression is mass spec

ISB-CGC hosts more than one thousand tables from a large variety of large cancer initiatives all open access. We have a Search Tool on our Portal (isb-cgc.org) to assist in finding the right tables for your uses.

Fabian Seidl<sup>1</sup>; Jacob Wilson<sup>1</sup>; Boris Aguilar<sup>2</sup>; Lauren Hagen<sup>2</sup>; Deena Bleich<sup>1</sup>; Lauren Wolfe<sup>2</sup>; Poojitha Gundluru<sup>1</sup>; George White<sup>2</sup>; Suzanne Paquette<sup>2</sup>; Elaine Lee<sup>2</sup>; Danna Huffman<sup>1</sup>; David Pot<sup>1</sup>; William Longabaugh<sup>2</sup> <sup>1</sup>General Dynamics Information Technology; <sup>2</sup>Institute for Systems Biology

## **Co-analysis of multiple data types in BigQuery**

Within our BigQuery ecosystem it is possible to join between the many data types based on key fields such as case IDs, gene name, cancer type and more. These joins can generally be made in less than a minute for cents worth of cloud costs.



Derived data in BigQuery enables rapid data discovery and quick comparisons via "table joins". Below is an example of a set of queries slightly modified for simplicity joining gene expression data to clinical data and protein abundance.

#### SELECT <fields>

FROM `isb-cgc-bq.TCGA.RNAseq\_hg38\_gdc\_current` WHERE <conditionals>

The initial query selects specific fields from the table of interest, in this case upper quartile normalized fpkm.

### SELECT

<fields> FROM `isb-cgc-bq.TCGA.RNAseq\_hg38\_gdc\_current` JOIN `isb-cgc-bq.TCGA.clinical\_gdc\_current` WHERE <conditionals>

The second query performs a join between the expression and clinical tables, in this case selecting for cigarettes smoked per day.

### SELECT fiold

<116	105>
FROM	<pre>`isb-cgc-bq.TCGA.RNAseq_hg38_gdc_current`</pre>
JOIN	<pre>`isb-cgc-bq.TCGA.clinical_gdc_current`</pre>
JOIN	<pre>`isb-cgc-bq.TCGA.protein_expr_hg38_gdc_current`</pre>
WHERE	<conditionals></conditionals>

## **BigQuery is a powerful statistical tool**

One reason to host bioinformatic data from the Cancer Research Data Commons in BigQuery is to improve accessibility and ease of exploration. BigQuery also allows for custom functions written in SQL or JavaScript and features automated compute scaling that allow it to outperform traditional High Performance Compute clusters.



Tests (millions)

We have generated custom user defined functions for commonly used statistical tests and incorporated these tests into example notebooks.

These notebooks as well as many others are a public teaching resource for researchers. isb-cancer-genomics-cloud.readthedocs.io/





Qu	ery results	ULTS 👻 📶 EXI	▼ m EXPLORE DATA ▼ ↓				
<	JOB INFORMATION	RESULTS	CHART	JSON EX	ECUTION DETAILS	>	
Row	case_barcode 💌	fpkr	n_uq_unstranded	expcigarettes_per_	protein_expression		
1	TCGA-86-A4P7		8.9616	nuli	-0.6541634745		
2	2 TCGA-91-6829		14.4956	5.178082191780	0.399470062		
3	3 TCGA-91-6828		12.558	nuli	-0.174617102		
4	TCGA-86-A4P8		1.947	nuli	-0.4681542825		
Ę	5 TCGA-38-4629		40.1784	5.479452054794	0.739331331		
e	5 TCGA-38-6178		10.7342	nuli	-0.013338784		
7	7 TCGA-78-7166		32.8658	2.082191780821	-0.0290081565		
8	3 TCGA-78-7167		3.7353	3.506849315068	-0.438169383		
		Results per pa	age: 50 🔻	1 – 50 of 378	I< < >	×	

Output of the final query joining between RNA expression, protein abundance, and a selected clinical field, ready for statistical tests and Machine Learning

Shown is the average run time required to calculate Spearman correlations and T-tests in BigQuery. Even **2.5 million** tests completed in less than 50 seconds in this trial.

BigQuery compute is cheap and we offer **\$300 cloud credits** for exploration and setup. Charges are based on data read rather than compute used. In a trial running 6.6 billion tests cost \$1.16.

Data type 1	Data type 2	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

The Mitelman Database of Chromosome Aberrations and Gene Fusions in Cancer is a free-access database devoted to chromosomes, genes, and cancer. The information in the Mitelman Database relates cytogenetic changes and their genomic consequences, in particular gene fusions, to tumor characteristics, based either on individual cases or associations.

### New Circos plots to visualize database contents.



# with gene fusions. opography Distribution



Example topography dis adenocarcinoma cases, and most common gen

Quarterly updates add thousands of new cases per year to the database. As of January 2025, the database consists of:

- 79,082 cytogenetic cases • 49,643 unique cytogenetic aberrations • 34,303 unique gene fusions 14,096 genes involved

- Rapid data exploration and quick statistics natively
- Derived data from well known reference NCI sets and annotations for co-analysis
- Fast links between unique data types
- Advanced statistical analyses using Python, R, SQL, and Bioconductor
- Rapidly able to expand to Machine Learning • Easy exploration of existing GDC and PDC data
- Access Virtual Machines for customized pipelines
- Multiple specialized databases such as Mitelman DB of Chromosomal Aberrations and Gene Fusions in Cancer

ISB-CGC is a component of the NCI Cancer Research Data Commons and has been funded with Federal funds from the National Cancer Institute, National Institutes of Health, Department of Health and Human Services. Contract No. HHSN261201400008C and ID/IQ Agreement No. 17X146 under Contract No. HHSN261201500003I.

If you have any questions about our resources or would like to collaborate please email us: Ask for help - we will get on a call. feedback@isb-cgc.org @isb\_cgc 😾

# **GENERAL DYNAMICS**

Information Technology

# isb-cgc.org

## **Recent upgrades**

### New charts and tables to view and compare morphology and topography distributions

	Fusion Gene	î↓ Tumor Morphology	î↓ Tumor Location	†↓ Count	↑↓
hology-Topography Distribution	TMPRSS2::ERG	Adenocarcinoma	Prostate	11	
Adverse invest	KIF5B::RET	Adenocarcinoma	Lung	8	
Adenocarcinoma - Dreast 9,905 9,905	TPM3::NTRK1	Adenocarcinoma	Lung	8	
1,607 enocarcinoma - Liver	EZR::ROS1	Adenocarcinoma	Lung	6	
Adenocarcinoma - Ovary 2,622 Adenocarcinoma - Prostate	IRF2BP2::NTRK1	Adenocarcinoma	Lung	6	
Adenocarcinoma - Lung	TPM3::NTRK1	Adenocarcinoma	Large intestine	6	
istributions for	CCDC6::RET	Adenocarcinoma	Lung	5	
he table displays the	CD74::ROS1	Adenocarcinoma	Lung	5	
e fusions.	ETV6::NTRK3	Adenocarcinoma	Lung	5	
	ETV6::NTRK3	Adenocarcinoma	Thyroid	5	

### **Summary and Conclusions**

- BigQuery is a tool with cloud-powered scaling of Microsoft Excel functionality
- Affordable storage and sharing of YOUR tabular data

### • Receive \$300 pilot funding and more available for your pilot project

### Funding