# Acknowledgements

**The Global Seven Bridges Team**

**CGC users and collaborators**

Institute for Systems Biology

BROAD INSTITUTE

THE UNIVERSITY OF CHICAGO

Image courtesy: https://portal.gdc.cancer.gov/

# Agenda

- Background

- Access multi-omics datasets in the CGC

- Use Case 1: Identifying mutational burden in cancer - Analysis of TCGA datasets

- Use Case 2: microRNA biogenesis in cancer

- Questions/Discussion

# Background

# Explosion of genomics data with ease of sequencing



Big Data: Astronomical or Genomical? Stephens et al; PLoS Biol. 2015 Jul; 13(7): e1002195.

# Increasingly large datasets bring challenges to data analysis



www.cancer.gov/ccg

# Multi-omic data is critical for cancer research



Cancer is a complex disease!

Comprehensively understanding the full picture of a research question requires examining multiple modalities

Guillermo de Anda-Jáuregui and Enrique Hernández-Lemus, Computational Oncology in the Multi-Omics Era: State of the Art. Front. Oncol., 07 April 2020 | https://doi.org/10.3389/fonc.2020.00423

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# The Seven Bridges Cancer Genomics Cloud (CGC)

A Cloud Resource within the NCI Cancer Research Data Commons for secure storage, sharing & analysis of petabytes of public, multi-omic cancer datasets

NCI Cancer Research Data Commons (CRDC)

# Growth of the Cancer Genomics Cloud Ecosystem

TCGA Pilot Program announced

Awarded NCI Cancer Genomics Cloud (CGC) Pilot contract

SIMONS FOUNDATION

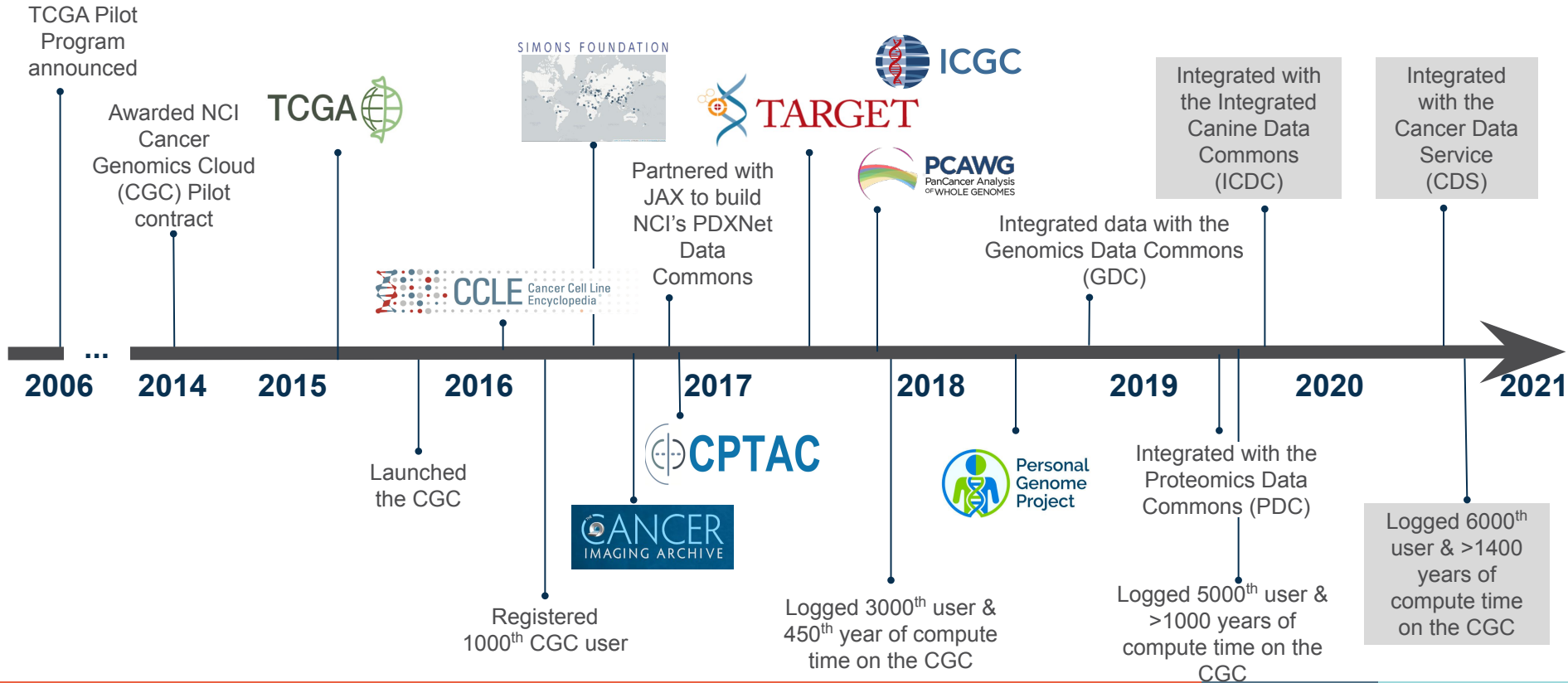TCGA

ICGC

TARGET

Integrated with the Integrated Canine Data Commons (ICDC)

Integrated with the Cancer Data Service (CDS)

Partnered with JAX to build NCI's PDXNet Data Commons

PCAWG
PanCancer Analysis OF WHOLE GENOMES

Integrated data with the Genomics Data Commons (GDC)

CCLE Cancer Cell Line Encyclopedia

**2006** ... **2014** **2015** **2016** **2017** **2018** **2019** **2020** **2021**

Launched the CGC

CPTAC

Personal Genome Project

Integrated with the Proteomics Data Commons (PDC)

THE CANCER IMAGING ARCHIVE

Registered 1000th CGC user

Logged 3000th user & 450th year of compute time on the CGC

Logged 5000th user & >1000 years of compute time on the CGC

Logged 6000th user & >1400 years of compute time on the CGC

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# The CGC democratizes complex analyses in a FAIR data ecosystem

- A stable, secure, and highly customizable cloud storage and computing platform

- Promotes a **F**indable, **A**ccessible, **I**nteroperable and **R**eusable (FAIR) data ecosystem

- A user-friendly portal for collaborative analysis of petabytes of public data alongside private data

- An optimized venue for reproducible data analysis using validated tools and pipelines



| Easy data management | Secure collaboration & managed billing | Flexible & fully reproducible methods | Optimized bioinformatics algorithms | Scalable computation | Extensible & developer friendly tools |

Wilkinson, M., Dumontier, M., Aalbersberg, I. *et al.* The FAIR Guiding Principles for scientific data management and stewardship. *Sci Data* 3, 160018 (2016) doi:10.1038/sdata.2016.18
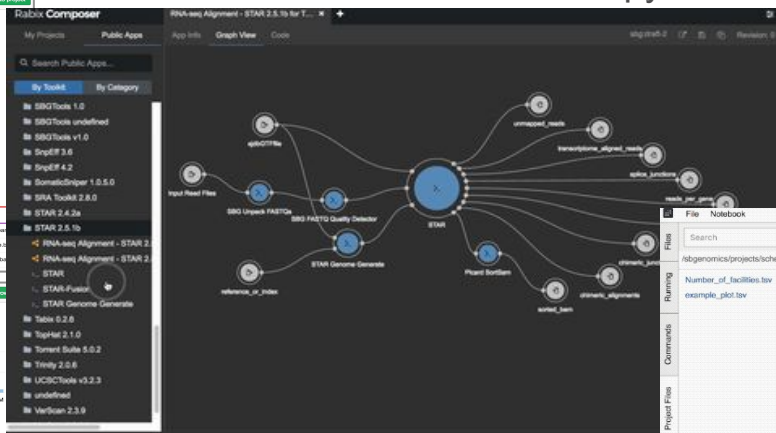
# Accelerating cancer research

- Detect aberrant splice junctions and splicing profiles across patient populations

- Identify neoantigens arising from novel gene fusion events

- Profile miRNA expression across patient populations

- Conduct  HLA typing to identify neoantigens

- Compare viral infection patterns across patient populations

- Detect novel gene fusions from RNA-Seq data

- Identify cis-regulatory region variants across patient populations

- ...and much more



CANCER GENOMICS CLOUD
SEVEN BRIDGES

# CGC provides an easy way to find and analyze data

Visually explore and access **3⁺ PB** of multi-omic public data through interactive query tools & APIs.

Use the **500⁺** cloud- and cost-optimized tools in our Public Apps library OR deploy custom tools using **Rabix Composer**, Jupyter notebooks or R packages

# Empowering a coordinating center on the CGC

**PDX Data Commons and Coordination Center JAX-Seven Bridges**

Collaborative and large-scale development and pre-clinical testing of targeted therapeutic agents in patient-derived models to advance the vision of cancer precision medicine.

- Data harmonized and securely shared
- Developed standardized PDX DNA-seq and RNA-seq workflows, available on the CGC
- Diverse models, metadata, and omics included



**https://portal.pdxnetwork.org/**

Enabled multiple high-impact publications
➔ Systematic Establishment of Robustness and Standards in Patient-Derived Xenograft Experiments and Analysis. *Cancer Research, March 2020*
➔ Conservation of copy number profiles during engraftment and passaging of patient-derived cancer xenografts. *Nature Genetics, January 2021*

# High impact publications on the CGC



**nature communications**

Explore our content | Journal information

nature > nature communications > articles > article

Article | Open Access | Published: 02 June 2020

## AGO-bound mature miRNAs are oligouridylated by TUTs and subsequently degraded by DIS3L2

Acong Yang, Tie-Juan Shao, Xavier Bofill-De Ros, Chuanjiang Lian, Patricia Villanueva, Lisheng Dai & Shuo Gu ✉

*Nature Communications* **11**, Article number: 2765 (2020) | Cite this article

2767 Accesses | 1 Citations | 11 Altmetric | Metrics

---

## CANCER RESEARCH
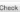
Home | About | Articles | For Authors | Alerts | News | COVID-19 | Search

Tumor Biology and Immunology

## Systematic Establishment of Robustness and Standards in Patient-Derived Xenograft Experiments and Analysis

Yvonne A. Evrard, Anuj Srivastava, Jelena Randjelovic; The NCI PDXNet Consortium, James H. Doroshow, Dennis A. Dean II, Jeffrey S. Morris, and Jeffrey H. Chuang

DOI: 10.1158/0008-5472.CAN-19-3101 Published June 2020    Check for updates

---

**Genome Medicine**

Home | About | Articles | Submission Guidelines

Research | Open Access | Published: 17 February 2020

## The pan-cancer landscape of prognostic germline variants in 10,582 patients

Ajay Chatrath, Roza Przanowska, Shashi Kiran, Zhangli Su, Shekhar Saha, Briana Wilson, Takaaki Tsunematsu, Ji-Hye Ahn, Kyung Yong Lee, Teressa Paulsen, Ewelina Sobierajska, Manjari Kiran, Xiwei Tang, Tianxi Li, Pankaj Kumar, Aakrosh Ratan & Anindya Dutta ✉

*Genome Medicine* **12**, Article number: 15 (2020) | Cite this article

2844 Accesses | 1 Citations | 78 Altmetric | Metrics

Oncogene
https://doi.org/10.1038/s41388-020-01507-5

**ARTICLE**

## Genetic alterations of *SUGP1* mimic mutant-*SF3B1* splice pattern in lung adenocarcinoma and other cancers

Samar Alsafadi [1,2] · Stephane Dayot [2] · Malcy Tarin [1] · Alexandre Houy [2] · Dorine Bellanger [2] · Michele Cornella [2] · Michel Wassef [3,4] · Joshua J. Waterfall [1,2] · Erik Lehnert [5] · Sergio Roman-Roman [1] · Marc-Henri Stern [2] · Tatiana Popova [2]

---

**nature genetics**

Explore our content | Journal information

nature > nature genetics > articles > article

Article | Published: 07 January 2021

## Conservation of copy number profiles during engraftment and passaging of patient-derived cancer xenografts

Xing Yi Woo, Jessica Giordano, Anuj Srivastava, Zi-Ming Zhao, Michael W. Lloyd, Roebi de Bruijn, Yun-Suhk Suh, Rajesh Patidar, Li Chen, Sandra Scherer, Matthew H. Bailey, Chieh-Hsiang Yang, Emilio Cortes-Sanchez, Yuanxin Xi, Jing Wang, Jayamanna Wickramasinghe, Andrew V. Kossenkov, Vito W. Rebecca, Hua Sun, R. Jay Mashl, Sherri R. Davies, Ryan Jeon, Christian Frech, Jelena Randjelovic, Jacqueline Rosains, Francesco Galimi, Andrea Bertotti, Adam Lafferty, Alice C. O'Farrell, Elodie Modave, Diether Lambrechts, Petra ter Brugge, Violeta Serra, Elisabetta Marangoni, Rania El Botty, Hyunsoon Kim, Jong-Il Kim, Han-Kwang Yang, Charles Lee, Dennis A. Dean II, Brandi Davis-Dusenbery, Yvonne A. Evrard, James H. Doroshow, Alana L. Welm, Bryan E. Welm, Michael T. Lewis, Bingliang Fang, Jack A. Roth, Funda Meric-Bernstam, Meenhard Herlyn, Michael A. Davies, Li Ding, Shunqiang Li, Ramaswamy Govindan, Claudio Isella, Jeffrey A. Moscow, Livio Trusolino, Annette T. Byrne, Jos Jonkers, Carol J. Bult, Enzo Medico ✉, Jeffrey H. Chuang ✉, PDXNET Consortium & EurOPDX Consortium  -Show fewer authors

*Nature Genetics* **53**, 86–99(2021) | Cite this article

618 Accesses | 42 Altmetric | Metrics

### Abstract

Patient-derived xenografts (PDXs) are resected human tumors engrafted into mice for preclinical studies and therapeutic testing. It has been proposed that the mouse host affects tumor evolution during PDX engraftment and propagation, affecting the accuracy of PDX

---

**CANCER GENOMICS CLOUD**
SEVEN BRIDGES

# Participating in open standards groups helps make us more FAIR

# How do I get an account on the CGC?

- Sign up with your email
  - **https://www.cancergenomicscloud.org/**
- Option to connect with eRA Commons to access controlled data
- **$300 of pilot funding** to get your project started
- Comprehensive online documentation and training resources
- Technical support from a team of scientists, bioinformaticians, and engineers

# Access multi-omics datasets in the CGC

# Access and search large public datasets on the CGC

| Dataset | Description | Experimental setup | File types |
|---|---|---|---|
| **TCGA** | Rich dataset of tumor and normal tissues from 11,000 patients, covering 33 cancer types | WES, RNAseq, miRNAseq, methylation, genotyping, ATACseq, imaging, WGS, .. | BAM, VCF, MAF, TXT, TSV, SVS, XML |
| TARGET | Dataset of genomic changes in childhood cancers | RNASeq, WGS, WES, miRNAseq | BAM, MAF, TSV, VCF, XLSX, TXT |
| CANCER IMAGING ARCHIVE | Imaging data from many 21 tumor types | Imaging | DCM |
| CPTAC | Proteomics of 10 tumor types and associated genomic data | Proteomics, WGS, WES, RNAseq | BAM, TSV, VCF, mzML.gz, mzid.gz, raw, tar.gz |
| International Cancer Genome Consortium | Consortium of many datasets, 20 studies on CGC | WGS, RNASeq | BAM, VCF |
| CCLE Cancer Cell Line Encyclopedia | Dataset of 1457 cancer cell lines | WGS, WES, RNAseq | BAM |
| SIMONS FOUNDATION | Genome sequencing of 130 populations | WGS | BAM, VCF |
| Personal Genome Project | Crowdsourced genomics, datasets from 10 individuals | WGS, WGBS, RNAseq, methylation | BAM, FASTQ, IDAT, TBI, VCF |
| HUMAN CELL ATLAS | Single-cell genomics of healthy tissues | scRNASeq | FASTQ |

# CGC connects with several CRDC data repositories

Coming soon!

# Use Case 1: Identifying mutational burden in cancer - Analysis of TCGA datasets

Image courtesy: https://portal.gdc.cancer.gov/

# Typical User Flow

**Create a Project**

**Find datasets of interest**

**Bring/Build tools or workflows**

**Analyze**

Organizational unit within the CGC

Many ways to find and bring in data:
- Data Browser
- Desktop uploader
- Command line uploader
- Volumes

Tools, workflows, and software packages
- Public Apps Gallery
- Tools or workflows wrapped in CWL
- R packages
- Python libraries

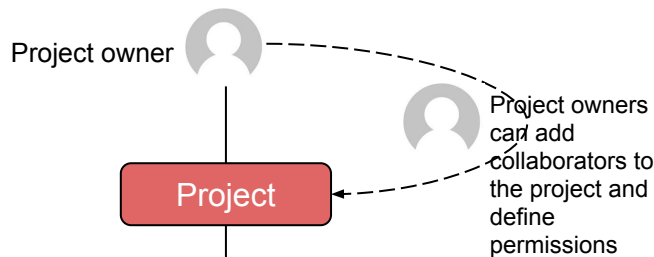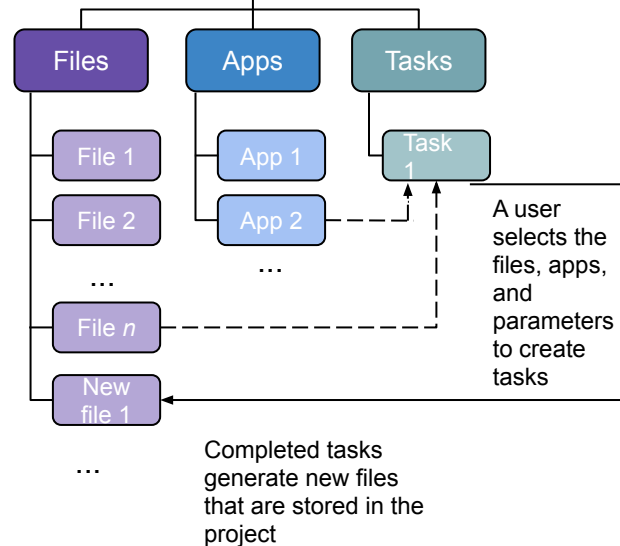Specify how an analysis will be run
- Task page
- Notebooks in RStudio or JupyterLab

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Projects organize files, methods, and results

Project owner



Project owners can add collaborators to the project and define permissions

A user selects the files, apps, and parameters to create tasks

Completed tasks generate new files that are stored in the project

Also known as *workspaces* or *sandboxes*
Easily manage collaborators and permissions

**Create a project** ✕

Name

Purdue-Bioinformatics-Class

Project URL:

https://cgc.sbgenomics.com/u/sailakss/purdue-bioinformatics-class ✏

Billing Group

Pilot Funds (sailakss) ▾

Location ❓

AWS (us-east-1) ▾

Execution settings:
**Spot Instances** ❓          On ⬤

**Memoization (WorkReuse)** ❓     Off ◯

☑ This project will contain **CONTROLLED** Data. ❶

Cancel    Create

Projects are configurable, e.g.

- Customizable billing group - where costs should be attributed
- Cloud resources (AWS or GCP)
- **Spot** (or **preemptible**) instances
- Memoization - Intermediate file retention
- Using S3 or Glacier storage

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Collaborate and share results quickly and easily

# Billing groups

Clear advantages for collaboration and interoperability. Aligned to temporal dynamics of research funding.

Allow users to distribute costs appropriately per function, topic, lab, etc

Use different funding sources (e.g. R24, Pilot Funds, credit card)

SB can reimburse for task failure due to external factors

**Billing Group settings: Pilot Funds (sailakss)**

## Info

| Organization | **Seven Bridges Genomics** |
|---|---|
| Creator | sailakss |
| Primary contact | Seven Bridges Genomics |
| Address | One Broadway, 14th Fl. , Massachusetts, United States |

| **Remaining credits** | **$ 298.02** |
|---|---|
| **Pilot funds** | **$ 300.00** |
| **Total charges** | **$ 1.98** |

| **Analysis usage** | | **Storage usage** | |
|---|---|---|---|
| Analysis charges | $ 0.53 | Storage charges | $ 1.45 |
| Tasks | $ 0.13 | Active | $ 0.41 |
| Data Cruncher analyses | $ 0.40 | Downloaded | $ 0.00 |
| | | Storage deduction | $ 0.00 |

Instance limits

Total number of instances that can be run in parallel

Current usage: **0 of 80** ⓘ

# Multi-cloud implementation on the CGC

# Memoization allows use of previously computed results



**COMPLETED**  **Whole Genome Sequencing - BWA + GATK 4.0 (with Metrics)** ✏️

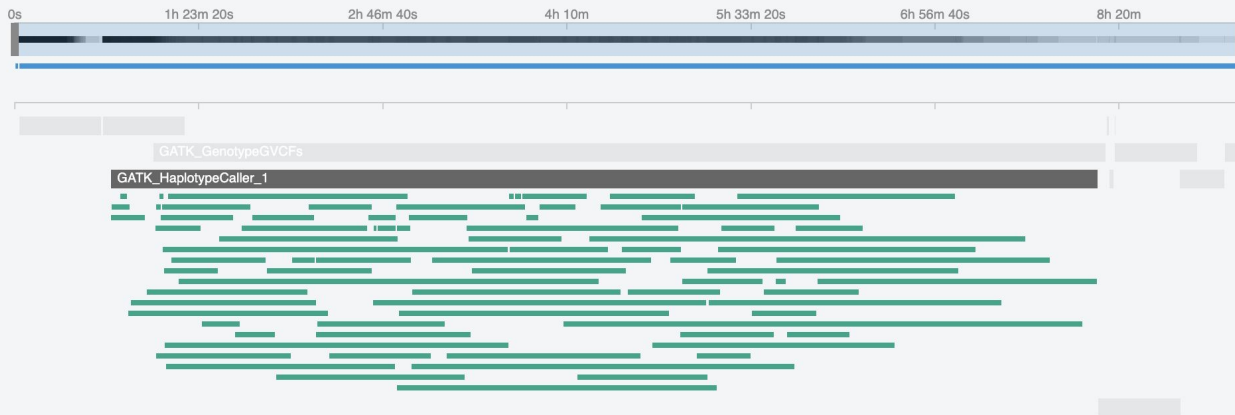Get support | View stats & logs | ► Edit and rerun

Executed on Aug. 23, 2020 22:12 by sinan.yavuz_demo

Preemptible Instances: **On** ⓘ | Memoization (WorkReuse): **On** ⓘ | Price: **$3.45** ⓘ | Duration: **9 hours, 14 minutes** ⓘ

▾ App: Whole Genome Sequencing - BWA + GATK 4.0 (with Metrics) - Revision: 0

ⓘ Precomputed outputs were used for some jobs. View task logs for more details.

Search apps 🔍

| 0s | 1h 23m 20s | 2h 46m 40s | 4h 10m | 5h 33m 20s | 6h 56m 40s | 8h 20m |

GATK_GenotypeGVCFs

GATK_HaplotypeCaller_1

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# User Flow

| Create a Project | Find datasets of interest | Bring/Build tools or workflows | Analyze |
|---|---|---|---|

**Create a Project**

Organizational unit within the CGC

**Find datasets of interest**

Many ways to find and bring in data:
- Data Browser
- Desktop uploader
- Command line uploader
- Volumes

**Bring/Build tools or workflows**

Tools, workflows, and software packages
- Public Apps Gallery
- Tools or workflows wrapped in CWL
- R packages
- Python libraries

**Analyze**

Specify how an analysis will be run
- Task page
- Notebooks in RStudio or JupyterLab

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Different options to bring data



* Public files
* Case Explorer & Data Browser
* Projects (that you are a member of)
* FTP/HTTP (signed URLs)
* Data tools
  - Command Line Uploader
  - Desktop Uploader
  - SBFS: Seven Bridges File System
  - API upload
* Volumes
* Import from manifest: ICDC/PDC

# Find open access TCGA data with Data Browser

# Easily connect cloud volumes



Control read/write permissions with IAM

Amazon/Google buckets

User storage

Volume

Import

* Alias

Project files

File_A

SB_file

Output_X

Output_Y

Export with API

SB resources

* Denotes files that reside in the bucket connected by the volume

SB storage

Outputs go to SB storage by default and can be exported by using API

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Enabling multi-omic research on the CGC through integrating with the PDC, ICDC, CDS

**NATIONAL CANCER INSTITUTE**
**Proteomic Data Commons**

**NIH** **NATIONAL CANCER INSTITUTE**
**Integrated Canine Data Commons**

**Cancer Data Service** (CDS)

1. User starts on PDC/ICDC/SRA (for CDS) portal to identify cohort of files
2. User downloads **files manifest** of selected cohort

**CANCER GENOMICS CLOUD**

1. User moves to CGC, creates a project
   a. Files → Add files → Import from a manifest
2. User prompted to upload the manifest from the PDC/ICDC/CDS
3. Data files from PDC/ICDC/CDS copied into user's project
4. Additional metadata accessed via Data Cruncher notebook

Links to doc pages to import data from: PDC, ICDC, CDS

# User Flow

**Create a Project**

Organizational unit within the CGC

**Find datasets of interest**

Many ways to find and bring in data:
- Data Browser
- Desktop uploader
- Command line uploader
- Volumes

**Bring/Build tools or workflows**

Tools, workflows, and software packages
- Public Apps Gallery
- Tools or workflows wrapped in CWL
- R packages
- Python libraries

**Analyze**

Specify how an analysis will be run
- Task page
- Notebooks in RStudio or JupyterLab

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Find the tools you need in the Public Apps Gallery

A curated collection of ~**500** bioinformatics tools & workflows

- ○ Optimized for speed & cost in the cloud
- ○ Fully parameterized & customizable
- ○ Accessible via the GUI & API



CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Bring Your Pipelines to the Platform with Web Composer

- An intuitive and flexible software development kit for developing and porting custom tools to the platform

- Conformance with community standards to ensure pipeline portability & reproducibility



docker

COMMON
WORKFLOW
LANGUAGE

Rabix
[Reproducible Analysis for Bioinformatics]

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Mutation Annotation Format (*maftools*)



| **Visualization** |
| Oncoplot (*oncoplot*) |
| Oncostrip (*oncostrip*) |
| Compare two cohorts (*coOncoplot, forestPlot*) |
| Lollipop plot (*lollipopPlot*) |
| TiTv plot (*titv, plotTiTv*) |
| Rainfall plot (*rainfallPlot*) |
| Genecloud (*geneCloud*) |
| GISTIC plots (*gisticBubblePlot, gisticChromPlot, gisticOncoPlot*) |
| APOBEC and Signature plots (*plotApobecDiff, plotSignatures*) |
| MAF summary (*plotmafSummary*) |

**Input MAF**
*read.maf*
*readGistic\**

*MAF object*

**Set operation**
Subset (*subsetMaf*)

*MAF object*

**Variant Annotations**
Variant annotations via oncotator API (*oncotate*)
Annovar output to MAF conversion (*annovarToMaf*)
ICGC simple somatic to MAF (*icgcSimpleMutationToMAF*)

| **Analysis** |
| Driver gene detection (*oncodrive*) |
| Mutual exclusive and co-occuring events (*somaticInteractions*) |
| Differentially mutated genes (*mafCompare*) |
| De-novo Mutational Signature analysis (*trinucleotideMatrix, extractSignatures*) |
| APOBEC enrichment estimation (*trinucleotideMatrix*) |
| Pan Cancer comparison (*pancanComparision*) |
| Survival analysis (*mafSurvival*) |
| Heterogeneity estimation (*inferHeterogeneity, math.score*) |
| Pfam domain summarization (*pfamDomains*) |
| MutSig gene symbol correction (*prepareMutSig*) |
| Enrichment Analysis (*clinicalEnrichment, signatureEnrichment*) |

Mayakonda A, Lin D, Assenov Y, Plass C, Koeffler PH (2018). "Maftools: efficient and comprehensive analysis of somatic variants in cancer." *Genome Research*. doi: 10.1101/gr.239244.118.

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Maftools Workflow



1. Dockerize individual tools
2. Wrap each tool in CWL
3. Connect tools into a workflow and set parameters

# User Flow

| Create a Project | Find datasets of interest | Bring/Build tools or workflows | Analyze |
|---|---|---|---|

**Organizational unit within the CGC**

**Many ways to find and bring in data:**
- Data Browser
- Desktop uploader
- Command line uploader
- Volumes

**Tools, workflows, and software packages**
- Public Apps Gallery
- Tools or workflows wrapped in CWL
- R packages
- Python libraries

**Specify how an analysis will be run**
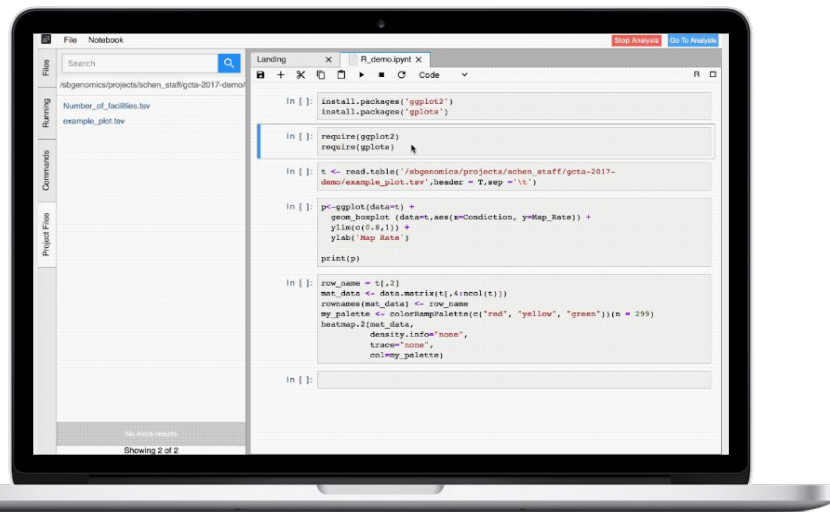- Task page
- Notebooks in RStudio or JupyterLab

# Powerful, collaborative, & reproducible interactive analysis

Users create interactive analysis sessions within a project - all files are available and over 50 instances can be used (*c3.xlarge* to *x1.32xlarge* on AWS)
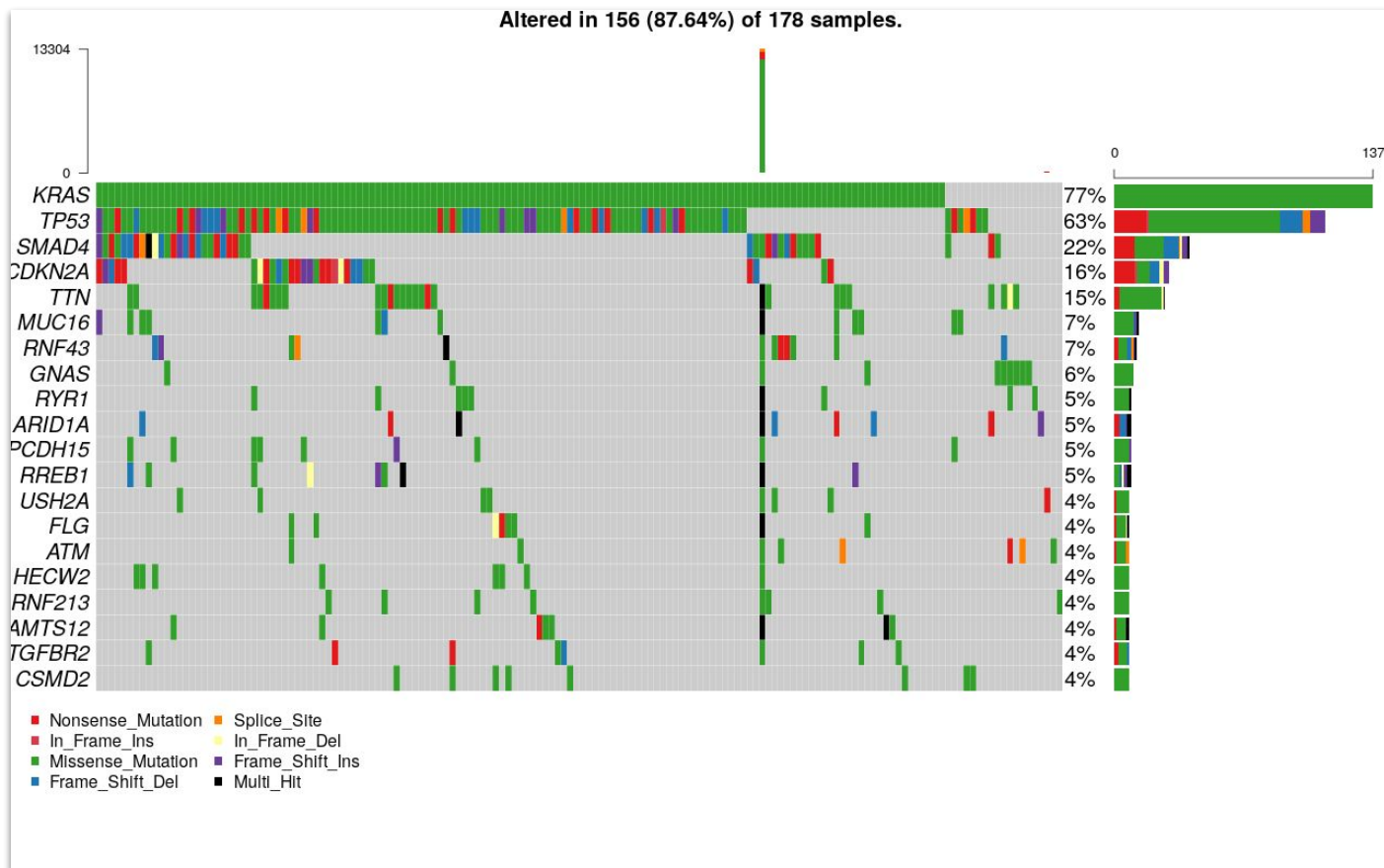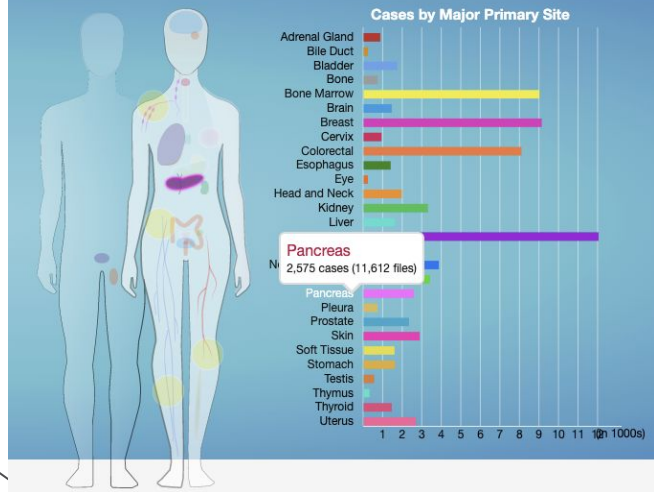
# PAAD Oncoplot

**Top 3 mutated genes**

- KRAS
- TP53
- SMAD4



Altered in 156 (87.64%) of 178 samples.

Nonsense_Mutation
In_Frame_Ins
Missense_Mutation
Frame_Shift_Del
Splice_Site
In_Frame_Del
Frame_Shift_Ins
Multi_Hit

CANCER GENOMICS CLOUD
SEVEN BRIDGES

Image courtesy: https://portal.gdc.cancer.gov/

# Want to learn more?

- Learn how to perform cloud based loading of single cell data, quality control, normalization, PCA and clustering and biomarker identification.
- Using open data
- The workflow makes tables and an html report
- Also learn to use cloud based RStudio to dive deeper into the data



**Hands-on Demo on 04/08 (Thursday)**

# Use Case 2: microRNA biogenesis in cancer

# Using the CGC to understand microRNA biogenesis in cancer

***Collaborative Project program to advance your research***

- Submit a proposal for up to **$10,000** in cloud credits to cgc@sevenbridges.com
- Get additional access to our CGC team and bioinformatics support
- Projects have resulted in dozens of papers, many users submit multiple papers from one project
- We encourage applications from students and postdocs

Xavier Bofill de Ros - Research Fellow
Gu Lab, Center for Cancer Research NCI

**NIH** NATIONAL CANCER INSTITUTE
Center for Cancer Research

| 2017 | 2018 | 2019 | 2020 |

- **I Collaborative application**
- Joint efforts on QuagmiR development
- TCGA data analysis with QuagmiR

**Bioinformatics**

- QuagmiR publication
- QuagmiR tool available on CGC

- **II Collaborative application**
- Subsequent research publication

**Cell Reports**

- Third research publication
- Currently working with Case Explorer and other

**nature COMMUNICATIONS**
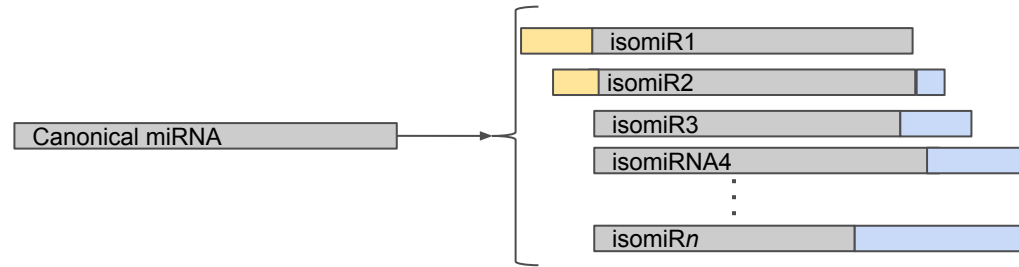
CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Case study - microRNAs & isomiRs

- microRNAs regulate gene expression
- Isoforms of miRNA (isomiRs) are correlated with cancer progression
- isomiRs very difficult to study because they are so heterogeneous
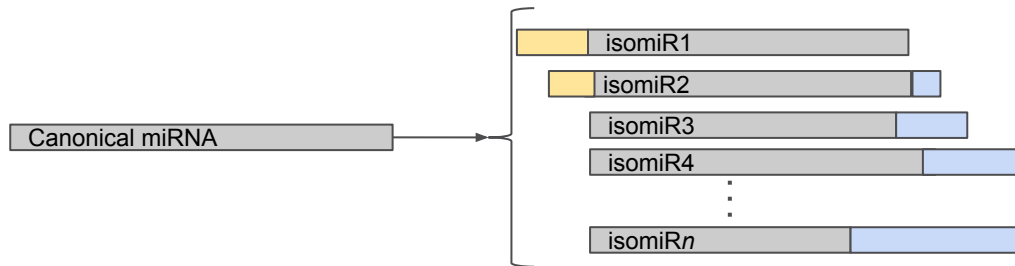


miRNA Pathway

# Case study: QuagmiR
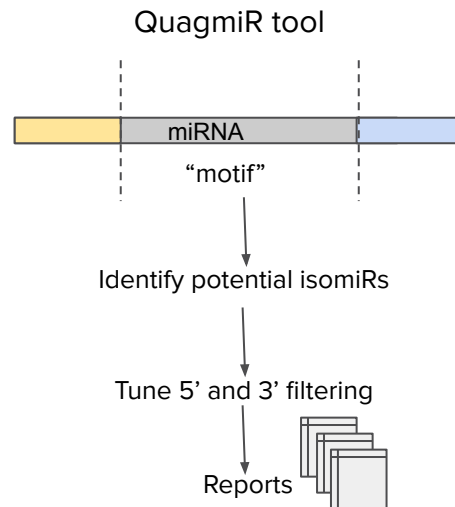
OXFORD

Sequence analysis

## QuagmiR: a cloud-based application for isomiR big data analytics

Xavier Bofill-De Ros[1], Kevin Chen[1], Susanna Chen[1], Nikola Tesic[2], Dusan Randjelovic[2], Nikola Skundric[2], Svetozar Nesic[2], Vojislav Varjacic[2], Elizabeth H. Williams[2], Raunaq Malhotra[2], Minjie Jiang[1] and Shuo Gu[1],*

[1]RNA Mediated Gene Regulation Section, RNA Biology Laboratory, Center for Cancer Research, National Cancer Institute, Frederick, MD, USA and [2]Seven Bridges Genomics Inc., Cambridge, MA, USA

- **QuagmiR**: a tool that pulls specific reads and aligns them against a consensus sequence in the middle of a miRNA, allowing mismatches on the ends to capture 3′ isomiRs
- **Initial idea:** reprocess all TCGA miRNAs with QuagmiR
- **The CGC enabled an efficient and highly scalable analysis, hence more research projects resulted from the initial one**
- Xavier used the fact that miRNA data tends to be smaller in size and leveraged CGC capabilities to easily analyze up to **70 samples per task across dozens of tasks**

Canonical miRNA

isomiR1
isomiR2
isomiR3
isomiR4
⋮
isomiR*n*

QuagmiR tool

miRNA
"motif"

Identify potential isomiRs

Tune 5' and 3' filtering

Reports

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Quantify and visualize isomiR differences

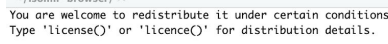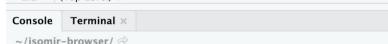## QuagmiR on the Seven Bridges Cancer Genomics Cloud (CGC)

Xavier edited this page on Jul 19, 2018 · 10 revisions

---

### Create a **CGC** account and project

1. Create an account on the CGC.
2. Create a project.

### Import data into your project.

3. Import data into your project using one of the following approaches:

- Upload your own data using the CGC Uploader (recommended) or other available tools.

IsomiRs by sample type - hsa-let-7c-3p

How to use QuagmiR in less than 2 minutes
123 views • Oct 5, 2018

# Support and Resources

## CGC Monthly Webinar Series

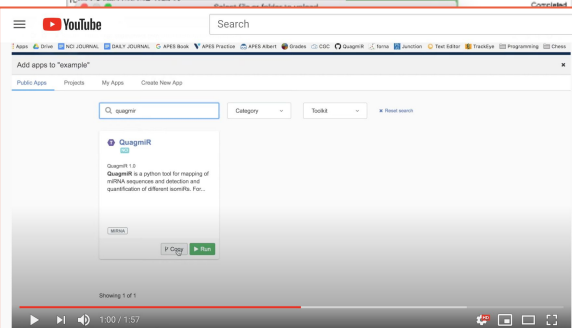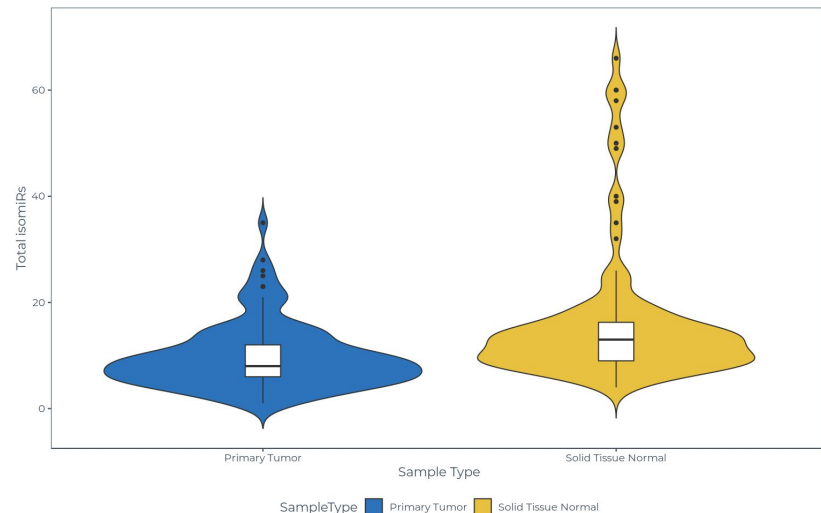Learn about CGC platform features that you can use in your projects.

Variety of research and technical topics in the field of cancer research using the CGC

**Resources:** Upcoming webinar info, slides and recordings are available at:
https://www.cancergenomicscloud.org/webinars

**Save the date/time:** 4th Wednesday of each month at 2pm ET

## CGC Knowledge Center

https://docs.cancergenomicscloud.org/

Contact CGC Support: cgc@sevenbridges.com

Office Hours: Every week on Thursdays

https://www.cancergenomicscloud.org/officehours

HOME    ABOUT    RESOURCES    POLICIES    KNOWLEDGE CENTER    HELP    LOGIN

WEBINARS          QUICK START TUTORIAL
PUBLICATIONS      COMPREHENSIVE USER GUIDE
RELEASE NOTES     TROUBLESHOOTING
PUBLIC APPS       OFFICE HOURS
EVENTS            CONTACT US

*Learn from cancer ge*

# FASTER

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# In Summary

**Data Access**
Immediately access petabytes of *Open and Controlled* TCGA, CPTAC, TCIA, and other omics datasets
Bring your own private cohorts alongside public data.

**Collaborate on the cloud**
Collaborate with other researchers around the world in a secure workspace
Access to high-throughput, cost-effective cloud computing resources and storage on demand and at cost.

**Interactive Analysis**
The ability to perform custom, interactive analysis and visualization on the platform using Python, RStudio.

**Tools and Workflows**
Standard bioinformatics pipelines
Bring your own analysis tools directly to the platform
Connect multiple tools together using our interactive custom workflow builder

6000 users

>80 countries

1,600,000+ computational tasks

1400+ years of total compute time

66,800+ workflows

500+ public apps

**Support & Resources**
Access comprehensive online documentation and training resources; Technical support from a team of >200 expert scientists, bioinformaticians, and engineers.

CANCER GENOMICS CLOUD
SEVEN BRIDGES

Questions?