# The Cancer Genomics Cloud powered by Seven Bridges: a secure and scalable cloud-based platform to access, share and analyze multi-omics datasets

**May 14, 2021**

**Sai Lakshmi Subramanian**
Program Manager, Seven Bridges
sai.subramanian@sbgenomics.com

CANCER GENOMICS CLOUD
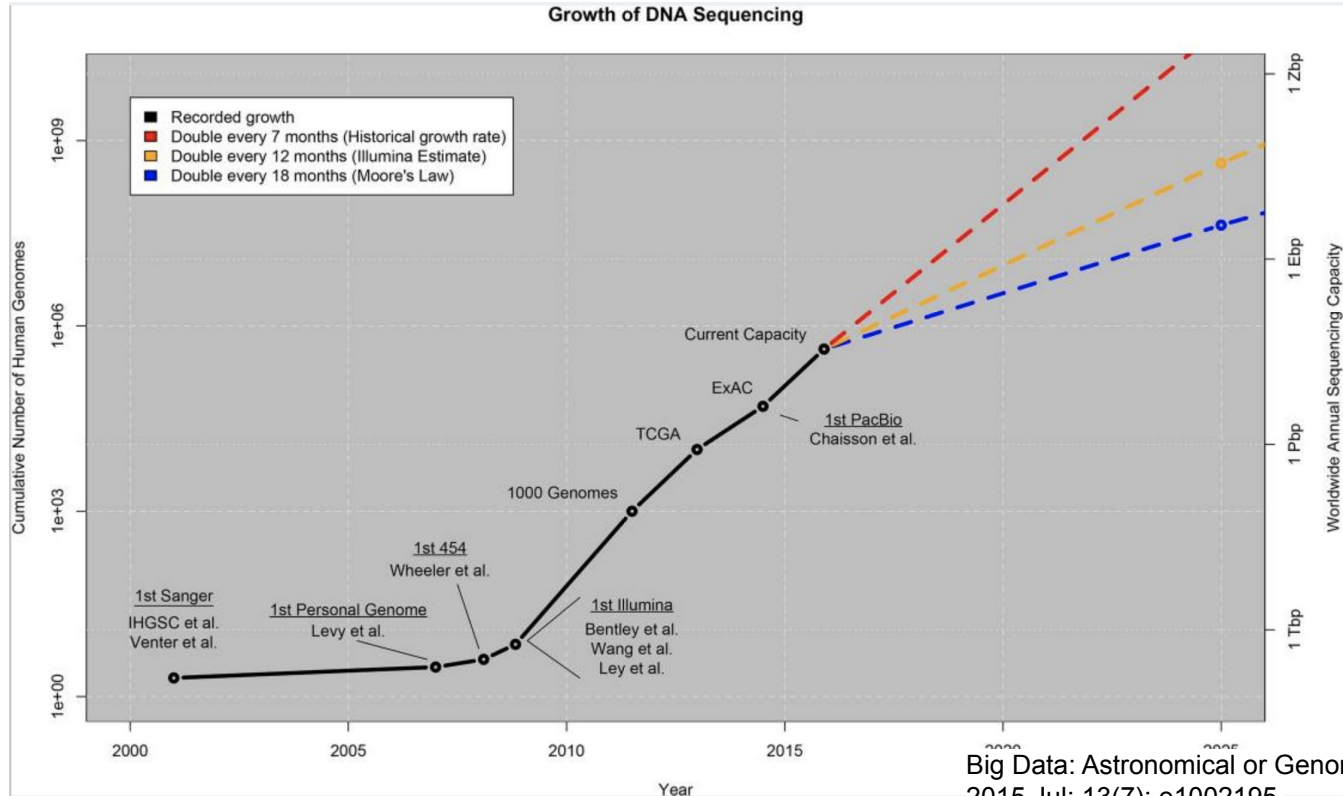SEVEN BRIDGES

# Agenda

- Background

- Access multi-omics datasets in the CGC

- Use Case: Exploring the human proteome - Analysis of CPTAC datasets

  - Features to control cloud costs

- Ongoing Projects - New workflows for novel science

- 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

# 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



https://datacommons.cancer.gov/cancer-research-data-commons
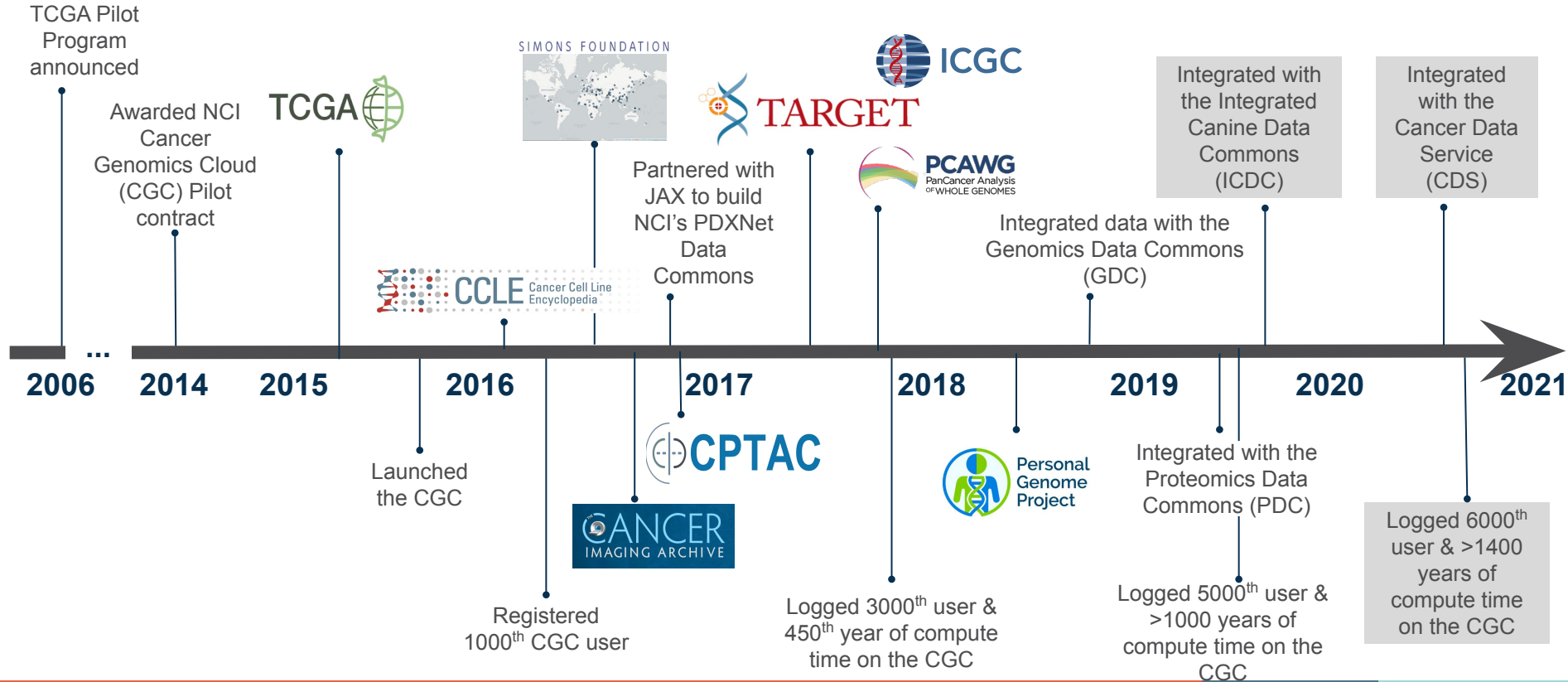
# Multi-omic data is critical for cancer research



*Cancer is a complex disease!*

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

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Growth of the Cancer Genomics Cloud Ecosystem



**2006** — TCGA Pilot Program announced

**2014** — Awarded NCI Cancer Genomics Cloud (CGC) Pilot contract

**2015** — TCGA

**2016** — CCLE Cancer Cell Line Encyclopedia; SIMONS FOUNDATION; Launched the CGC; Registered 1000th CGC user

**2017** — Partnered with JAX to build NCI's PDXNet Data Commons; TARGET; CPTAC; THE CANCER IMAGING ARCHIVE

**2018** — ICGC; PCAWG PanCancer Analysis of WHOLE GENOMES; Personal Genome Project; Logged 3000th user & 450th year of compute time on the CGC

**2019** — Integrated data with the Genomics Data Commons (GDC); Integrated with the Proteomics Data Commons (PDC); Logged 5000th user & >1000 years of compute time on the CGC

**2019–2020** — Integrated with the Integrated Canine Data Commons (ICDC)

**2020** — Logged 6000th user & >1400 years of compute time on the CGC

**2021** — Integrated with the Cancer Data Service (CDS)

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# The CGC democratizes complex analyses in a FAIR data ecosystem

**Cloud-based Environment for Collaborative Research and Bioinformatics Data Analysis**

- 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 |
|---|---|---|---|---|---|

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# 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

**The Jackson Laboratory**

**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

**PDXNet**

UCD · HCI · WashU · Wistar · JAX lab · MD Anderson · BCM

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

PDXNet · ⌂ Home · ☰ Resources · 📊 Analysis · 🔀 Metadata · ❓ Help · ℹ About · ✉ Contact

**PDXNet Portal**

## PDXNet Portal

Powered by Seven Bridges

The PDXNet Portal provides a way for researchers to learn about the PDX models, sequencing data (DNA and RNA), and PDX Minimum Information metadata tools generated by the network for public use.

The National Cancer Institute (NCI) launched the PDX (patient-derived xenografts) Development and Trial Centers Research Network (PDXNet) in September 2017 to accelerate translational research that uses PDX models and sequencing data. The PDXNet includes six PDX Development and Trial Centers (PDTCs) and the PDX Data Commons and Coordinating Center (PDCCCC). The two PDTCs added in 2018 focus exclusively on developing PDXs from minority patients. PDXNet also works closely with the NCI Patient-Derived Model Repository (PDMR) to ensure data are collected and provided in a standardized format.

Collectively, the PDTCs and the PDCCCC work together to test and advance multi-agent cancer treatments from PDX studies to human clinical trials. PDXNet is an inclusive consortium welcoming collaborations. Please contact us to discuss how we can work together to advance new cancer treatments.

🧑 PDXNet Models · 💿 PDTC Data · 📁 PDMR Data

**Data Summary**

| CONTRIBUTORS | FILES (PDTC/PDMR) | MODELS | CANCER TYPES |
|---|---|---|---|
| 👤 6 | 📄 2822 / 9492 | ⌛ 258 | ⏳ 33 |

PDXNet Models · Sequencing Files · Portal Update Timeline

**PDX Models by Contributor**

Adenocarcinoma - colon
Adenocarcinoma - pancreas
Adenocarcinoma - small intest.
Breast cancer, NOS
Cholangiocar.- intra/extrahepatic
Cystosarcoma phylloides - breast
Gastrointestinal stromal tumor
Invasive breast carcinoma
Liver/hepatobiliary cancer
Malig. periph. nerve sheath tum.
Neuroendocrine cancer, NOS
Non-small cell lung cancer, NOS
Papillary thyroid carcinoma
RCC, clear cell adenocarcinoma
Small cell lung cancer
Squamous cell lung carcinoma
Urothelial/bladder cancer, NOS

Model Count per Contributor

BCM · HCI · MDACC · UC Davis · WISTAR · WUSTL

## 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, Hyunsoo 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

Find a growing list of publications at: https://www.cancergenomicscloud.org/publications

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Participating in open standards groups helps make us more FAIR

# How do I get an account on the CGC?

CANCER GENOMICS CLOUD
SEVEN BRIDGES

- 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

CANCER GENOMICS CLOUD
SEVEN BRIDGES

**Log in**

eRA    Log in with eRA Commons

Log in with username and password

New to the CGC? Create an account

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# 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

# Use Case: Exploring the human proteome - Analysis of CPTAC datasets

# Proteogenome characterization of ccRCC

How do I analyze the spectra data of the proteome in clear cell renal cell carcinoma cases?

Kidney cancer is among the 10 most common cancers in both men and women
~73,000 new cases with >14,000 deaths in 2020
https://seer.cancer.gov/statfacts/html/kidrp.html



Cases by Primary Site

- Proteome and phosphoproteome data from the ccRCC tumors is available in PDC along with peptide spectrum analyses (PSMs) and protein summary reports from the CPTAC common data analysis pipeline (CDAP).
- Using the CGC, process high-throughput data-dependent acquisition (DDA) tandem mass spectrometry data acquired from peptides labeled with TMT tags.
- The multiplexed labeling, as applied in the CPTAC program, allows for differential quantitation across multiple tumor/normal samples.

Image Courtesy: https://pdc.cancer.gov/pdc/

# Typical 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

# Projects organize files, methods, and results

Project owner

Project owners can add collaborators to the project and define permissions

Project

Files | Apps | Tasks

File 1 | App 1 | Task 1
File 2 | App 2
… | …
File *n*

New file 1

…

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
Containers and Workflows Webinar Demo

Project URL:
https://cgc.sbgenomics.com/u/sailakss/containers-and-workflows-webinar-demo ✎

Billing Group
Pilot Funds (sailakss) ▾

Location ⓘ
AWS (us-east-1) ▾

Execution settings:
Spot Instances ⓘ                          On ⬤
Memoization (WorkReuse) ⓘ          Off ◯
Network Access ⓘ

**Block network access**
Executions within the project won't have network access

**Allow network access**
Executions will have unrestricted network access

☐ This project will contain **CONTROLLED** Data. ⓘ

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

# Multi-cloud implementation on the CGC

# Memoization allows use of previously computed results



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

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**

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

Volume

Project files

SB storage

Import

File_A

SB_file

* Alias

Output_X

Output_Y

Amazon/Google buckets

Export with API

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

User storage

SB resources

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

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)

**CANCER GENOMICS CLOUD**

1. User starts on PDC/ICDC/SRA (for CDS) portal to identify cohort of files

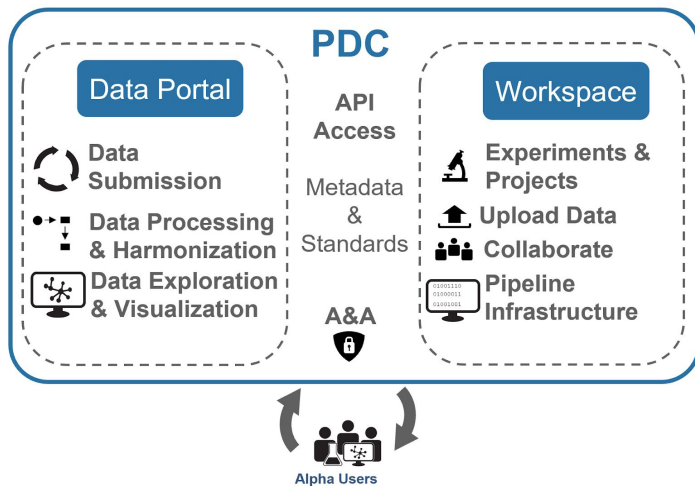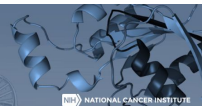2. User downloads **files manifest** of selected cohort

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

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Integration with the Proteomic Data Commons (PDC)



**Proteome Data Commons – democratize access to cancer-related proteomic datasets**

PDC

Data Portal
- Data Submission
- Data Processing & Harmonization
- Data Exploration & Visualization

API Access

Metadata & Standards

A&A

Workspace
- Experiments & Projects
- Upload Data
- Collaborate
- Pipeline Infrastructure

Alpha Users

Courtesy of Izumi HInkson

CGC integration with PDC enables both tool developers and researchers to take full advantage of rich data resources

Proteomics tools can run quickly and efficiently in the cloud, lowering barriers for usage

Login

NIH NATIONAL CANCER INSTITUTE
Proteomic Data Commons

e.g. BRCC3, 05BR003, kinase, PDC0001

NCI is pleased to release these data to the public. Some data are under an EMBARGO for publication and/or citation.

HOME    BROWSE    ANALYSIS    SUBMIT DATA    ABOUT

62 Studies   |   24 TB Data volume   |   81,275 Data files   |   > 357 M Spectra   |   > 1 M Peptides   |   15,007 Proteins

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# How to import data from PDC



**Researchers can also bring private data to co-analyze with public resources.**

# 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:<br>● Data Browser<br>● Desktop uploader<br>● Command line uploader<br>● Volumes | Tools, workflows, and software packages<br>● Public Apps Gallery<br>● Tools or workflows wrapped in CWL<br>● R packages<br>● Python libraries | Specify how an analysis will be run<br>● Task page<br>● Notebooks in RStudio or JupyterLab |

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# FragPipe: A complete proteomics pipeline

Resource

## Integrated Proteogenomic Characterization of Clear Cell Renal Cell Carcinoma

David J. Clark [1, 32], Saravana M. Dhanasekaran [2, 32], Francesca Petralia [3, 32], Jianbo Pan [1, 32], Xiaoyu Song [4, 5, 32], Yingwei Hu [1, 32], Felipe da Veiga Leprevost [2, 32], Boris Reva [3, 32], Tung-Shing M. Lih [1, 32], Hui-Yin Chang [2], Weiping Ma [3], Chen Huang [6], Christopher J. Ricketts [7], Lijun Chen [1], Azra Krek [3], Yize Li [8], Dmitry Rykunov [3], Qing Kay Li [1] ... Zhidong Tu

Show more ⌄

Resource

## Proteogenomic Characterization Reveals Therapeutic Vulnerabilities in Lung Adenocarcinoma

Michael A. Gillette [1, 2, 24, 27] ⚲ ✉, Shankha Satpathy [1, 24] ⚲ ✉, Song Cao [3, 25], Saravana M. Dhanasekaran [4, 25], Suhas V. Vasaikar [5, 25], Karsten Krug [1, 25], Francesca Petralia [6, 25], Yize Li [3], Wen-Wei Liang [3], Boris Reva [6], Azra Krek [6], Jiayi Ji [7], Xiaoyu Song [7], Wenke Liu [8], Runyu Hong [8], Lijun Yao [3], Lili Blumenberg [9], Sara R. Savage [10] ... Zhiao Shi
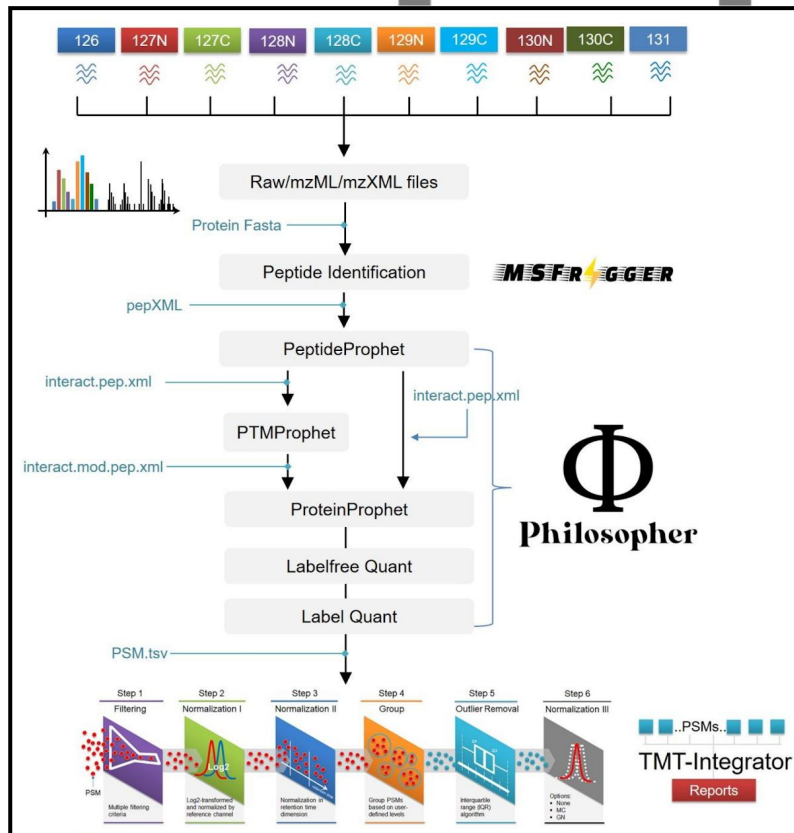
CANCER GENOMICS CLOUD
SEVEN BRIDGES

# FragPipe Proteomics Pipeline

- Developed by Nesvizhskii lab at the University of Michigan
- **FragPipe**: a complete proteomics pipeline for comprehensive analysis of proteomics data which is powered by
  - **MSFragger**, an ultrafast peptide identification tool for mass spec-based proteomics
  - **Philosopher** toolkit, for post-processing MSFragger results
  - **TMT-Integrator**, a tools for integrating channel abundances from multiple TMT or iTRAQ-labeled samples and generating reports

https://github.com/Nesvilab/FragPipe

# 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
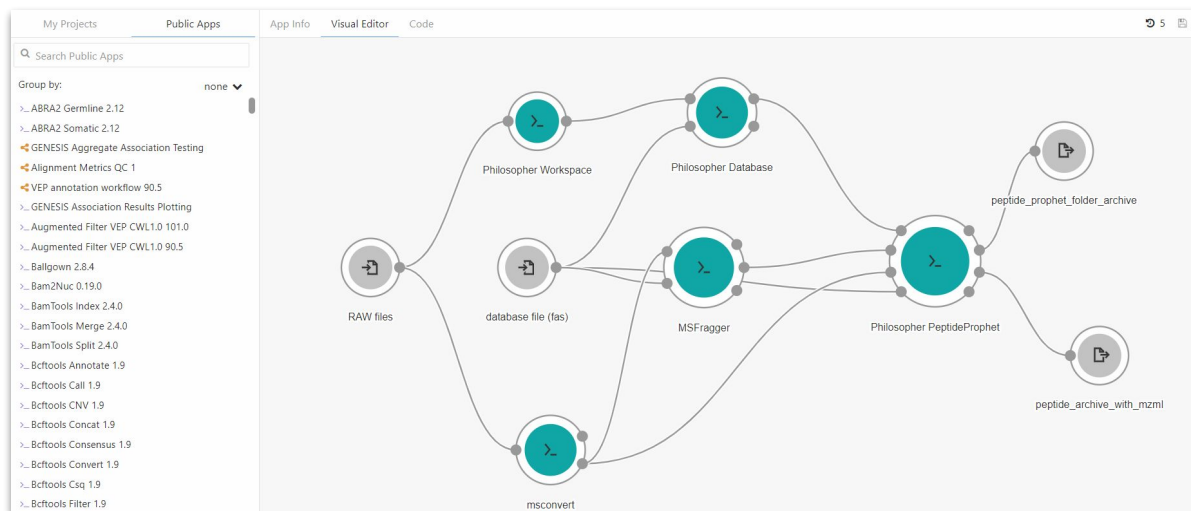- ○ Accessible via the GUI & API

https://cgc.sbgenomics.com/public/apps

# Bring Your Own Tools & Tailor new Pipelines in 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

# 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: <ul><li>Data Browser</li><li>Desktop uploader</li><li>Command line uploader</li><li>Volumes</li></ul> | Tools, workflows, and software packages <ul><li>Public Apps Gallery</li><li>Tools or workflows wrapped in CWL</li><li>R packages</li><li>Python libraries</li></ul> | Specify how an analysis will be run <ul><li>Task page</li><li>Notebooks in RStudio or JupyterLab</li></ul> |

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# FragPipe Workflows

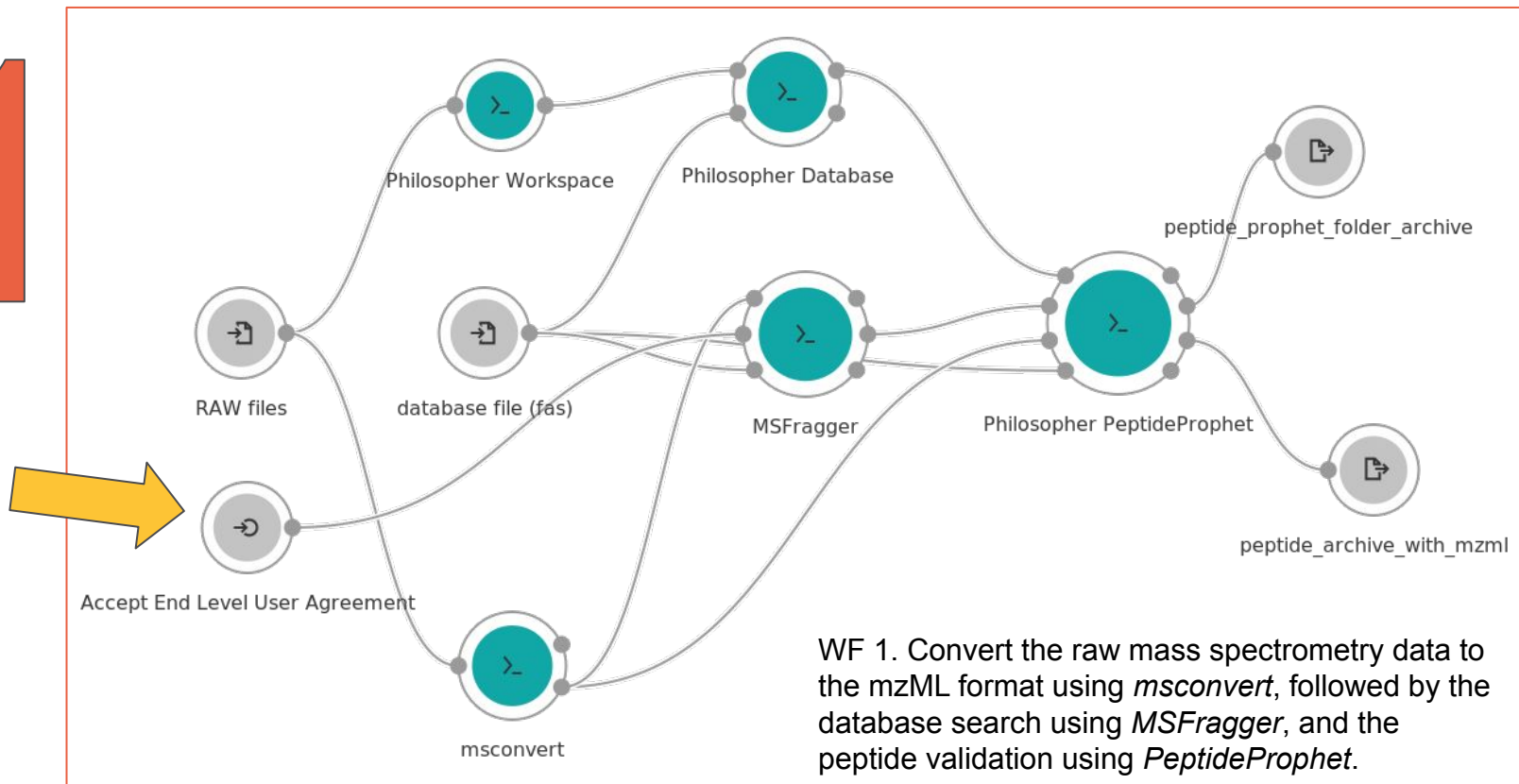| Search names and description 🔍 | Category: All ▾ | Toolkit: All ▾ | CWL Version: All ▾ | Status: Available ▾ |

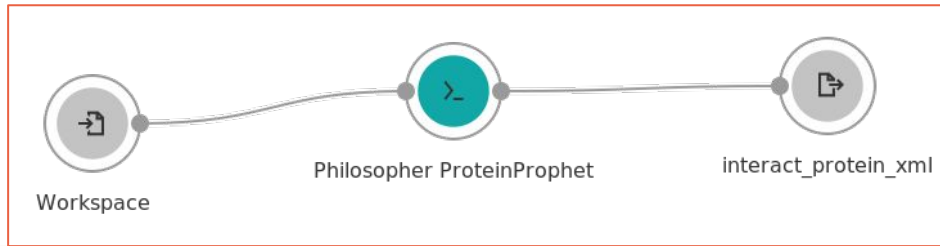| ▲ **Name** | **Type** |
|---|---|
| Ⓦ 01 FragPipe: Convert - Identify - Peptide Prophet   The first step of the workflow consists of converting the raw mass spectrometry data to t… | Workflow |
| Ⓦ 02 FragPipe: Protein Prophet   This workflow step takes the PeptideProphet output files from the first step containing the peptide validation and… | Workflow |
| Ⓦ 03 FragPipe: Filter - Quant - Report   This workflow takes the PeptideProphet, and the ProteinProphet output files, and applies a stringent Fal… | Workflow |
| Ⓦ 04 FragPipe: TMT Integrator and QC   This workflow step executes TMT-Integrator using the report tables generated by Philosopher. The pro… | Workflow |

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# FragPipe - 4 Workflows



WF 1. Convert the raw mass spectrometry data to the mzML format using *msconvert*, followed by the database search using *MSFragger*, and the peptide validation using *PeptideProphet*.

CANCER GENOMICS CLOUD
SEVEN BRIDGES

**2**

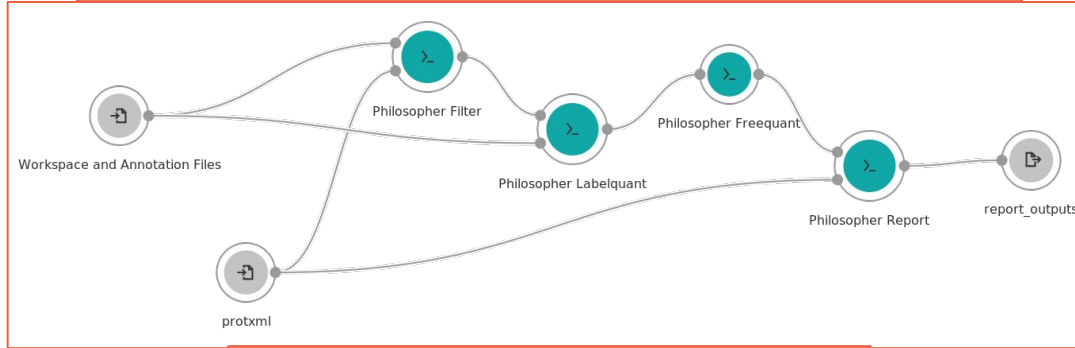**Workspace** — Philosopher ProteinProphet — interact_protein_xml
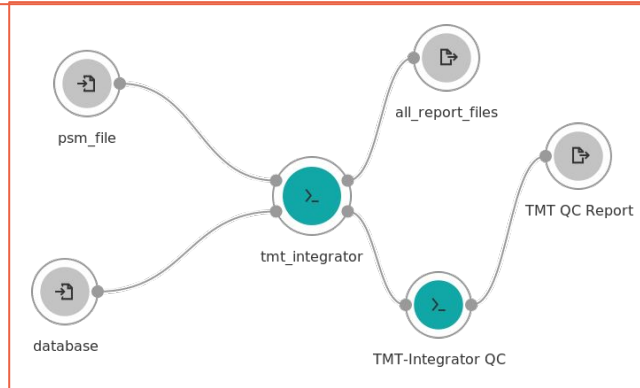
WF 2. Takes the PeptideProphet output files from the first step containing the peptide validation and calculates the protein inference using ProteinProphet.

**3**

Workspace and Annotation Files — Philosopher Filter — Philosopher Labelquant — Philosopher Freequant — Philosopher Report — report_outputs

protxml

WF 3. takes the PeptideProphet, and the ProteinProphet output files, and applies a stringent False Discovery Rate (FDR) filtering. Peptide and proteins are filtered individually at 1% FDR. The high-quality PSMs, peptides, and proteins are then quantified using a label-free algorithm that uses the apex peak intensity as a measurement. Finally, the isobaric tags are quantified and annotated with the correct sample labels.

**4**

psm_file — tmt_integrator — all_report_files

database — TMT-Integrator QC — TMT QC Report

WF 4. Executes TMT-Integrator using the report tables generated by Philosopher. The program applies a series of statistical filters, and high-quality thresholds to filter the data. Summary report tables are created containing peptides, proteins, genes, and phosphosites (only for phospho-enriched data sets).

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Features to control cloud costs

# Monitor analysis at each step in **macro** or micro levels

# Monitor analysis at each step in macro or **micro** levels

# Scale Up with batching by PDC file metadata

# Review Outputs without downloading

# Metrics of FragPipe Cloud vs. Local

Analysis of ccRCC whole cell lysate samples @Michigan on local server (non-cloud)

**Computation**:
- local server
- 56 cores (Xeon(R) CPU, 2.60GHz)
- 500GB RAM

**Time for analysis**:

27 hours, and a total of 16 GB RAM

Analysis of ccRCC samples on the CGC

**Time for accessing the data:**

Less than a minute; no need for downloads

**Time for analysis on CGC**:
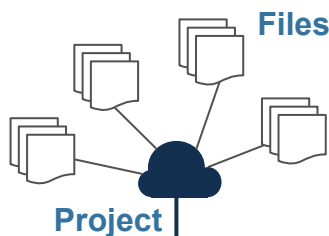
4.5 hours total. Up to 23 aws nodes in parallel @ 8 vCPU 16GB RAM, 48 aws nodes cumulative

**Cost of analysis on CGC**:

$0.49 per TMT plex, ~$12 for total analysis

# 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)

# Multi-omic data is critical for cancer research



110 Proteogenomic Characterized Renal Cell Carcinoma Cases

Transcriptomics
Immune-based subtypes

Genomics
Chromosome translocations

Histopathology

Phosphoproteomics
Kinase inhibition targets

Proteomics
Protein-specific metabolic alterations

*Cancer is a complex disease!*

Many research questions go beyond genomics data!

Different modalities should be examined to comprehensively understanding the full picture of a research question

***As a researcher, it is essential to focus on the data. The CGC provides all tools for a multi omics analysis, so you can spend more time interpreting the results, and not configuring programs.***

Clark et al., Integrated Proteogenomic Characterization of Clear Cell Renal Cell Carcinoma, Cell. 2019 Oct 31; 179(4): 964–983.e31. DOI:https://doi.org/10.1016/j.cell.2019.10.007

# Ongoing Projects

Adding new resources and capabilities to CGC
to support broader areas of cancer research

# Epigenetic mechanisms for transcription regulation in cancer

**Overall aim:**

- Enable secondary analysis of epigenomics data on CGC, focusing on ATACseq and ChIPseq experiments
- Develop multi-omics analysis for data coming from different omics experiments (RNAseq, ATACseq, ChIPseq, WGBS, RRBS, proteomics)
- Analyze data from several publicly available repositories in order to characterize epigenetic markers in cancer.

**Deliverables (publicly available on CGC soon):**

- Workflows for ATACseq and ChIPseq analysis, based on ENCODE's specification
- Data Cruncher Interactive Analysis combining ATACseq, ChIPseq, WBGS, RNAseq and proteomics data
- Workflow for multi-omics analysis as a one-click solution

CANCER GENOMICS CLOUD
SEVEN BRIDGES

# Multi-omics and imaging machine learning analysis

**Overall aim:**

Create a predictive model from both genomics and image features combined with the available clinical data in order to **predict therapy response among patients with cancer**.

- Build machine learning tools for processing multi-omics and imaging data from dbGAP/TCGA dataset.
- Use existing deep learning algorithms and libraries and adapt them for execution on CGC platform.

**Deliverables (publicly available on CGC soon):**

- Data Cruncher Interactive Analysis
- Tools and workflows for deep learning models adapted for imaging data, utilizing GPU instances and Tensorflow Python deep learning library.

# Support and Resources

## CGC Monthly Webinar Series

- Learn about CGC platform features that you can use in your projects.
- 4th Wednesday of each month at 2pm ET
- Upcoming webinar info, slides and recordings are available at:
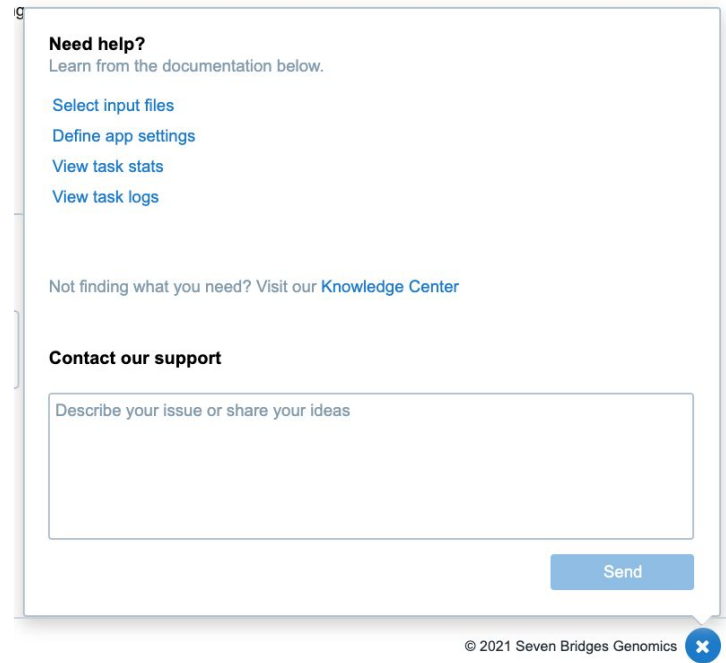  https://www.cancergenomicscloud.org/webinars

## CGC Knowledge Center

https://docs.cancergenomicscloud.org/

Contact CGC Support: cgc@sevenbridges.com

Office Hours: Every week on Thursdays

https://www.cancergenomicscloud.org/officehours

**Need help?**
Learn from the documentation below.

Select input files
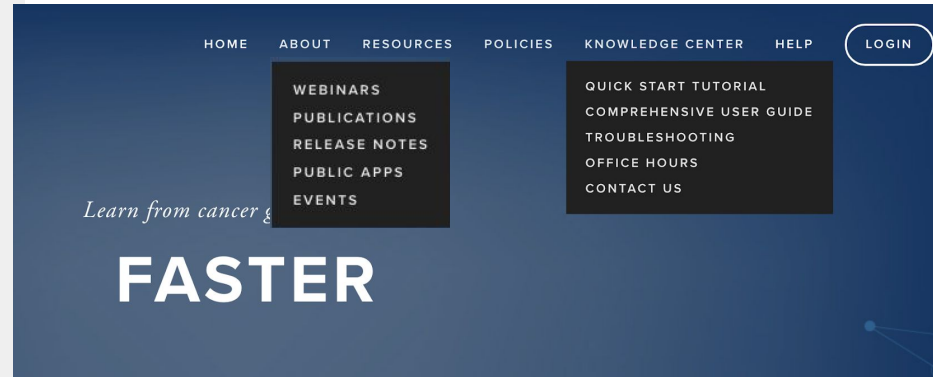Define app settings
View task stats
View task logs

Not finding what you need? Visit our **Knowledge Center**

**Contact our support**

Describe your issue or share your ideas

Send

© 2021 Seven Bridges Genomics

HOME    ABOUT    RESOURCES    POLICIES    KNOWLEDGE CENTER    HELP    LOGIN

WEBINARS
PUBLICATIONS
RELEASE NOTES
PUBLIC APPS
EVENTS

QUICK START TUTORIAL
COMPREHENSIVE USER GUIDE
TROUBLESHOOTING
OFFICE HOURS
CONTACT US

*Learn from cancer g*

# FASTER

# In Summary

**Data Access**
Immediately access petabytes of *Open and Controlled* TCGA, CPTAC, 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

CANCER GENOMICS CLOUD
SEVEN BRIDGES

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

# Acknowledgements

**CGC Team**

***Manisha Ray***

Jelena Radenkovic

Milos Stanojevic

Milos Trboljevac

Marko Tosic

Ana Stelkic

***Dave Roberson***

***Sai Lakshmi Subramanian***

Jack DiGiovanna

Brandi Davis-Dusenbery

**&** The Global Seven Bridges Team

**UMichigan Team**

Alexey I. Nesvizhskii

***Felipe da Veiga Leprevost***

Hui-Yin Chang

Guo Ci Teo

Fengchao Yu

**PDC Team**

***Paul Rudnick***

Rajesh Thangudu

CANCER GENOMICS CLOUD
SEVEN BRIDGES

Questions?