

Emerging Approaches for Tumor Analyses In Epidemiological Studies

Session 2: Public Databases

November 9, 2022

Session Overview

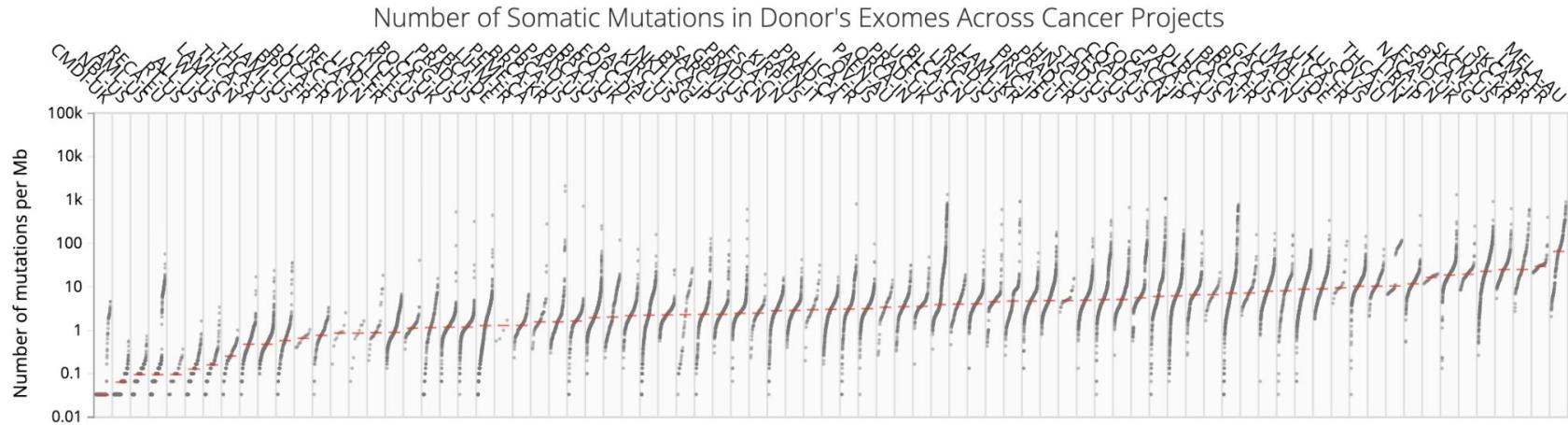
- Major Cancer Genomic Studies with Data Portal Available
- Common Genetics Data Resources
- Specialized Databases for Genomic Analyses
- Analytical Programming Packages for Cancer Genomic Datasets
- Cloud Resources with Available Cancer Genomic Datasets
- “Awesome” Bioinformatics Resources

Major Cancer Genomic Studies with Data Portal Available

WGS → WES → TS



WGS: International Cancer Genome Consortium (ICGC)



- Large proportion of TCGA samples included
- Controlled (sensitive genomic and clinical data) vs uncontrolled datasets
- Instant analysis with cloud computing
- Pan-Cancer Analysis of Whole Genomes (PCAWG) (~2700 donors)
- ICGC ARGO Data Platform

234,022 Files 16,307 Donors 1.73 PB

Data Release 28 March 27th, 2019

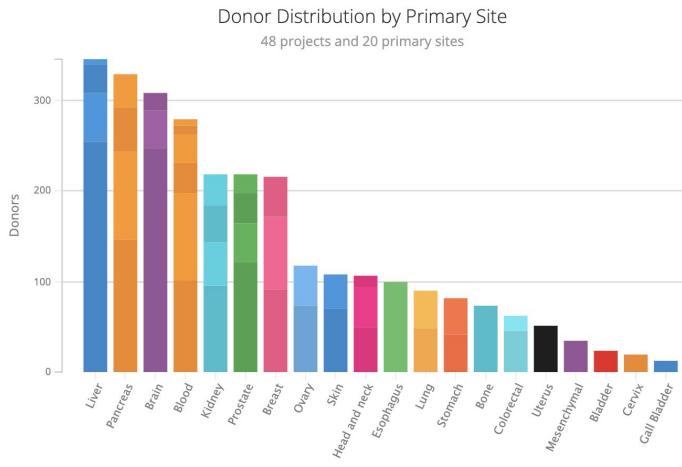
Cancer projects	86
Cancer primary sites	22
Donor with molecular data in DCC	22,330
Total Donors	24,289
Simple somatic mutations	81,782,588

Summary and Projects

DCC: ICGC Data Coordination Center

ARGO: Accelerating Research in Genomic Oncology

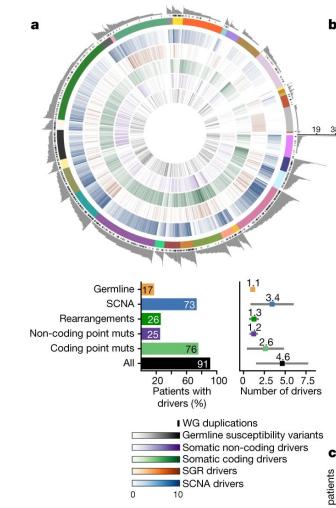
WGS:Pan-Cancer Analysis of Whole Genomes (PCAWG)



Data Type	# Donors	# Files	Format	Size
SGV	2,715	8,505	VCF	517.27 GB
StGV	2,715	5,668	VCF	7.29 GB
Aligned Reads	2,793	12,168	BAM	794.32 TB
Unaligned Reads	1	1	BAM	104.20 GB
Simple Somatic Mutations	2,715	25,501	VCF	189.99 GB
Copy Number Somatic Mutations	2,715	5,671	VCF	132.62 MB
Structural Somatic Mutations	2,715	14,195	VCF	1.61 GB

Available data as of Oct 25, 2021

2658 cancer whole genomes
with matched normal tissues
across 38 tumor types



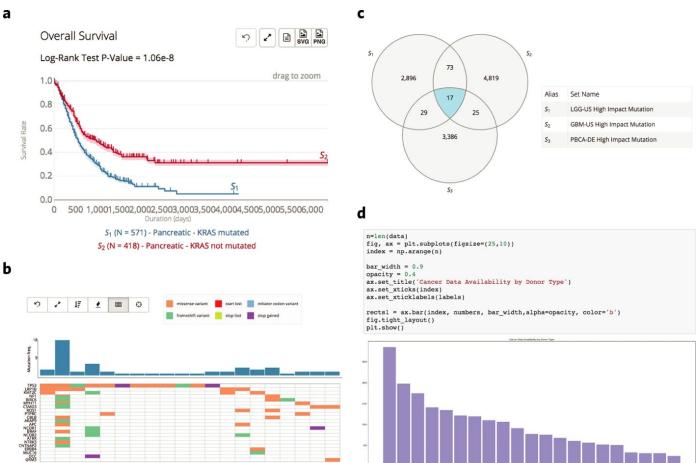
Major Publications

1. The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network. **Pan-cancer analysis of whole genomes**. Nature (2020).
2. Rheinbay, E. et al. **Analyses of non-coding somatic drivers in 2,693 cancer whole genomes**. Nature (2020).
3. PCAWG Transcriptome Core Group et al. **Genomic basis of RNA alterations in cancer**. Nature (2020).
4. Li, Y. et al. **Patterns of somatic structural variation in human cancer genomes**. Nature (2020).
5. Gerstung, M. et al. **The evolutionary history of 2,658 cancers**. Nature (2020).
6. Alexandrov, L. B. et al. **The Repertoire of Mutational Signatures in Human Cancer**. Nature (2020).
7. Phillips, M. et al. **Of Clouds and Genomic Data Protection**. Nature (2020).

ICGC Data Portal

The figure shows the ICGC Data Portal homepage. At the top, there are five navigation buttons: "Cancer Projects" (orange), "Advanced Search" (blue), "Data Analysis" (purple), "DCC Data Releases" (teal), and "Data Repositories" (green). Below these, a main title "Cancer genomics data sets visualization, analysis and download." is displayed. A search bar with placeholder text "e.g. BRAF, KRAS G12D, DO35100, MU7870, F1998, apoptosis, Cancer Gene Census, imatinib, GO:0016049" and a "Search" button is present. To the right, a box titled "Data Release 28" (March 27th, 2019) lists statistics: Cancer projects (86), Cancer primary sites (22), Donor with molecular data in DCC (22,330), Total Donors (24,289), and Simple somatic mutations (81,782,588). At the bottom, there are three buttons: "By donors", "By genes", and "By mutations".

The ICGC Data Portal Facet Search interface



Using the ICGC Data Portal

Repository		File Format	
<input type="checkbox"/> Collaboratory ...	121,467	<input type="checkbox"/> VCF	95,631
<input type="checkbox"/> GDC - Chicago	90,099	<input type="checkbox"/> FASTQ	63,543
<input type="checkbox"/> Azure - Toronto	77,486	<input type="checkbox"/> BAM	57,103
<input type="checkbox"/> AWS - Virginia	39,193	<input type="checkbox"/> BCR XML	17,715
<input type="checkbox"/> PDC - Chicago	22,200	<input type="checkbox"/> TGZ	23
<input type="checkbox"/> EGA - Hinxton	5,091	<input type="checkbox"/> XLSX	7
Select all	less	Select all	less
Data Type		Analysis Software	
<input type="checkbox"/> Unaligned Reads	69,131	<input type="checkbox"/> DKFZ/EMBL varia...	19,850
<input type="checkbox"/> SSM	61,591	<input type="checkbox"/> Broad variant call...	19,809
<input type="checkbox"/> Aligned Reads	51,538	<input type="checkbox"/> BWA with Mark D...	17,988
<input type="checkbox"/> StSM	14,195	<input type="checkbox"/> Sanger variant cal...	11,340
<input type="checkbox"/> Biospecimen Data	8,881	<input type="checkbox"/> STAR 2-Pass	9,253
<input type="checkbox"/> Clinical Data	8,841	<input type="checkbox"/> BWA-aligner	9,046
<input type="checkbox"/> SGV	8,505	<input type="checkbox"/> SomaticSniper An...	9,045
<input type="checkbox"/> CNSM	5,671	<input type="checkbox"/> MuTect2 Annotati...	9,037
<input type="checkbox"/> StGV	5,668	<input type="checkbox"/> MuSE Annotation	9,028
<input checked="" type="checkbox"/> No Data	1	<input type="checkbox"/> VarScan2 Annotat...	8,980
Select all	less	<input type="checkbox"/> BWA MEM	5,809
<input type="checkbox"/> WXS	110,698	<input type="checkbox"/> MUSE variant call ...	2,835
<input type="checkbox"/> WGS	83,530	<input type="checkbox"/> PCAWG SNV-MNV...	2,778
<input type="checkbox"/> RNA-Seq	12,840	<input type="checkbox"/> PCAWG InDel call...	2,778
<input type="checkbox"/> miRNA-Seq	9,046	<input type="checkbox"/> STAR	1,465
<input type="checkbox"/> Validation	100	<input type="checkbox"/> TopHat2	1,465
<input type="checkbox"/> Bisulfite-Seq	86	<input type="checkbox"/> Pilot50	150
<input checked="" type="checkbox"/> No Data	17,722	<input type="checkbox"/> Silver bullet	1
Select all	less	<input checked="" type="checkbox"/> No Data	93,365
Experimental Strategy		Select all	
<input type="checkbox"/> WXS	110,698	Select all	less
<input type="checkbox"/> WGS	83,530		
<input type="checkbox"/> RNA-Seq	12,840		
<input type="checkbox"/> miRNA-Seq	9,046		
<input type="checkbox"/> Validation	100		
<input type="checkbox"/> Bisulfite-Seq	86		
<input checked="" type="checkbox"/> No Data	17,722		
Select all	less		
Only Files in Study		Access	
<input type="checkbox"/> PCAWG	71,709	<input type="checkbox"/> Controlled	216,300
<input checked="" type="checkbox"/> None	162,313	<input type="checkbox"/> Open	17,722
Select all		Select all	

WGS:100,000 Genomes Project | Genomics England

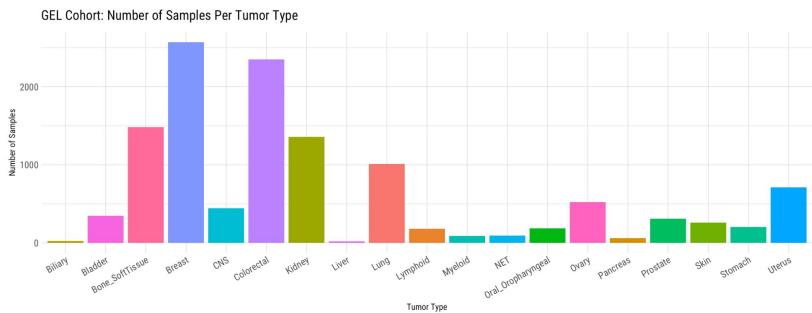
Genomics England's very first initiative – sequencing 100,000 genomes from around 85,000 NHS patients affected by rare disease or cancer – is leading to groundbreaking insights and continued findings into the role genomics can play in healthcare.

Highlighted Pan-Cancer findings:

Signatures of TOP1 transcription-associated mutagenesis in cancer and germline. *Nature*, 2022

Substitution mutational signatures in whole-genome-sequenced cancers in the UK population. *Science*, 2022

Nuclear-embedded mitochondrial DNA sequences in 66,083 human genomes. *Nature* 2022.



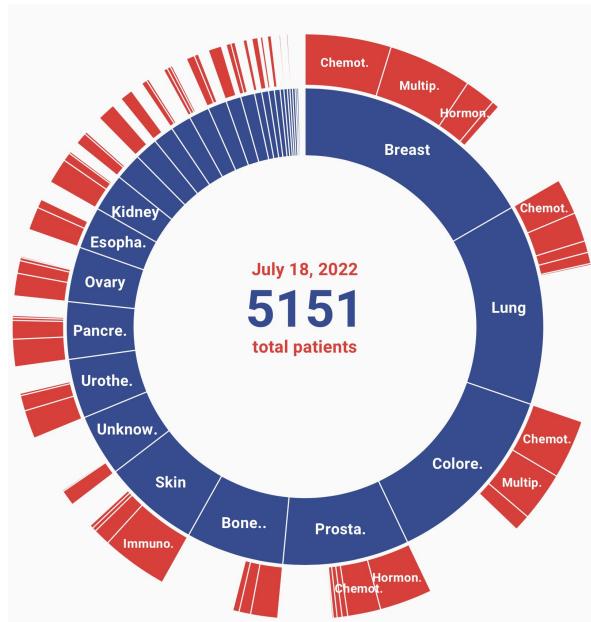
The GEL cohort:

15, 838 tumor-normal sample pairs

High-quality data derived from flash-frozen material, involving 12,222 GEL tumor samples from 11,585 individuals (several participants had synchronous or metachronous tumors).

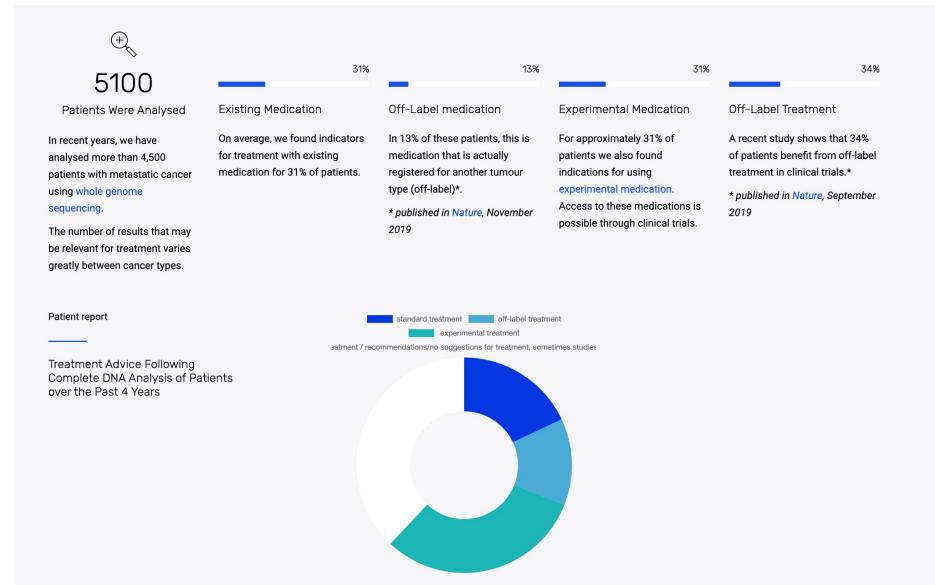
19 tumor types: skin, lung, stomach, colorectal, bladder, liver, uterus, ovary, biliary, kidney, pancreas, breast, prostate, bone and soft tissue, central nervous system (CNS), lymphoid, oropharyngeal, neuroendocrine tumors (NETs), and myeloid

WGS: Hartwig Medical Foundation



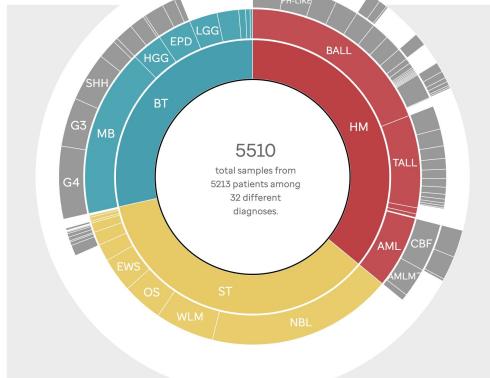
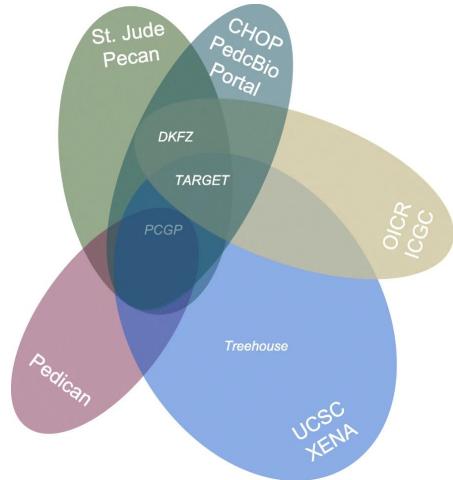
This is the largest database of **metastatic tumor data** obtained with WGS in the world. The uniqueness of the database consists of combining genetic data with treatment and treatment outcome data.

<https://database.hartwigmedicalfoundation.nl>



WGS: Pediatric Cancer Genome Project (PCGP) and Pediatric Cancer Databases (PeCan)

- The pediatric cancer genome project (PCGP) is a collaborative project created by St. Jude Children's Research Hospital and Washington University School of Medicine. The originally provided data portal "PCGP explore" was based on whole genome sequencing of pediatric tumors with the aim to cover the full spectrum of mutations in pediatric cancers [1]. PCGP is now part of St. Jude PeCan data portal.
- The Pan-Cancer Study of Childhood Cancers (PedPanCan) by the DKFZ includes various sources like ICGC Pedbrain Tumor, PCGP and from Heidelberg and others, and has been integrated into St. Jude PeCan.



5,510 Samples
5,213 Patients
32 Diagnoses
20,537 Genes
126,414 Mutations

St. Jude Cloud Genomics Platform

St. Jude Cloud Genomics Platform Sign in ☰

Select Data

Diagnoses Samples grouped by primary diagnosis **Publications** Samples grouped by publication **Studies** Datasets curated by St. Jude **Samples** All samples curated by St. Jude

[Download All Metadata](#) [Documentation](#)

Filter Selections

SEQUENCING TYPES Multiple RNA-Seq
 WES WGS

FILE TYPES BAM CNV
 Feature Counts gVCF
 Somatic VCF

SAMPLE TYPES Autopsy Cell Line
 Diagnosis Germline
 Metastasis Relapse
 Xenograft

[Share Selection](#)

Tumor	Paired Tumor-Normal	Germline	Search	Search Exact Phrase	
<input type="checkbox"/> Primary Diagnosis	Sequencing Types	File Types	Samples	Total Files	Total File Size
<input type="checkbox"/> Acoustic Neuroma	WGS WES RNA-Seq	BAM gVCF Feature Counts	3	23	188.48 GB
<input type="checkbox"/> Acute Leukemias of Ambiguous Lineage	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF Feature Counts	12	125	1.96 TB
<input type="checkbox"/> Acute Lymphoblastic Leukemia, NOS	WGS WES RNA-Seq	BAM gVCF Feature Counts	2	14	127.06 GB
<input type="checkbox"/> Acute Megakaryoblastic Leukemia	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF CNV Feature Counts	134	744	4.97 TB
<input type="checkbox"/> Acute Megakaryoblastic Leukemia, KMT2A rearrangement	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF Feature Counts	18	105	823.19 GB
<input type="checkbox"/> Acute Monoblastic/Monocytic Leukemia, KMT2A rearrangement	WGS WES RNA-Seq	BAM gVCF Feature Counts	2	22	352.16 GB
<input type="checkbox"/> Acute Myeloid Leukemia	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF Feature Counts	260	2,245	31.53 TB
<input type="checkbox"/> Acute Myeloid Leukemia, CEBPA alteration	WGS WES RNA-Seq	BAM gVCF Feature Counts	3	33	771.16 GB
<input type="checkbox"/> Acute Myeloid Leukemia, Core Binding Factor	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF CNV Feature Counts	249	1,444	11.76 TB
<input type="checkbox"/> Acute Myeloid Leukemia, KMT2A rearrangement	WGS WES RNA-Seq Multiple	BAM gVCF Somatic VCF Feature Counts	80	587	8.58 TB

Page 1 of 29 < >

St. Jude Cloud Visualization Community

WES: The Cancer Genome Atlas (TCGA)

Analysis platforms supporting TCGA dataset:

- Query based platforms: cBioPortal, FireBrowse, GEPIA, Genomic Data Commons, UCSC, TCGA Batch Effects Viewer, Tumor Map, TANRIC, SurvNet, Regulome Explorer, Xena
- Cloud-based platforms: Seven Bridges, Terra, DNAnexus etc.
- Programming-based packages: TCGAbiolinks, RTCGA etc.

NATIONAL CANCER INSTITUTE
THE CANCER GENOME ATLAS

TCGA BY THE NUMBERS

TCGA produced over

2.5
PETABYTES
of data

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

212,000
DVDs



TCGA data describes

33
DIFFERENT
TUMOR TYPES
...including
10
RARE
CANCERS

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

11,000
PATIENTS
...using

7
DIFFERENT
DATA TYPES


TCGA RESULTS & FINDINGS



MOLECULAR
BASIS OF
CANCER

Improved our
understanding of
the genomic underpinnings
of cancer



TUMOR
SUBTYPES

Revolutionized how
cancer is classified



THERAPEUTIC
TARGETS

Identified genomic
characteristics of tumors
that can be targeted with
currently available
therapies or used to help
with drug development

For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the same subtype found in other tissues at the molecular level, suggesting that despite arising from different tissues in the body, these subtypes may share a common path of development and respond to similar therapeutic strategies.

TCGA revolutionized how cancer is classified by identifying tumor subtypes with distinct sets of genomic alterations.*

TCGA's identification of targetable genomic alterations in lung squamous cell carcinoma led to NCI's Lung-MAP Trial, which will treat patients based on the specific genomic changes in their tumor.

THE TEAM



20
COLLABORATING
INSTITUTIONS

across the United States
and Canada

The Genomic Data Commons (GDC) houses TCGA and other NCI-generated data sets for scientists to access from anywhere. The GDC also has many expanded capabilities that will allow researchers to answer more clinically relevant questions with increased ease.



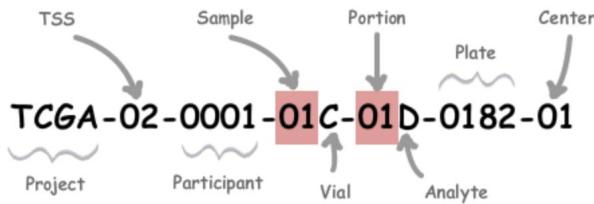
WHAT'S NEXT?

*TCGA's analysis of stomach cancer revealed that it is not a single disease, but a disease composed of four subtypes, including a new subtype characterized by infection with Epstein-Barr virus.

TCGA Cancers



TCGA Barcode



Cancer Type Studied	# Cases Characterized	Publication	Study Abbreviation
Acute Myeloid Leukemia	200 (200)	NEJM 2013	LAML
Adrenocortical Carcinoma	92 (91)	Cancer Cell 2016	ACC
Bladder Urothelial Carcinoma	412 (131)	Nature 2014, Cell 2017	BLCA
Breast Ductal Carcinoma	778 (430)	Nature 2012	BRCA/DCIS
Breast Lobular Carcinoma	201 (127)	Cell 2015	BRCA/LCIS
Cervical Carcinoma	307 (228)	Nature 2017	CESC
Cholangiocarcinoma	51 (38)	Cell Reports 2017	CHOL
Colorectal Adenocarcinoma	633 (276)	Nature 2012	COAD
Esophageal Carcinoma	185 (164)	Nature 2017	ESCA
Gastric Adenocarcinoma	443 (295)	Nature 2014	STAD
Glioblastoma Multiforme	617 (206)	Nature 2008, Cell 2013	GBM
Head and Neck Squamous Cell Carcinoma	528 (279)	Nature 2015	HNSC
Hepatocellular Carcinoma	377 (363)	Cell 2017	LIHC
Kidney Chromophobe Carcinoma	113 (66)	Cancer Cell 2014	KICH
Kidney Clear Cell Carcinoma	537 (446)	Nature 2013	KIRC
Kidney Papillary Cell Carcinoma	291 (161)	NEJM 2016	KIRP
Lower Grade Gioma	516 (293)	NEJM 2015	LGG
Lung Adenocarcinoma	585 (230)	Nature 2014, Nature Genetics 2016	LUAD
Lung Squamous Cell Carcinoma	504 (178)	Nature 2012, Nature Genetics 2016	LUSC
Mesothelioma	74 (87)	Cancer Discovery 2018	MESO
Ovarian Serous Adenocarcinoma	608 (489)	Nature 2011	OV
Pancreatic Ductal Adenocarcinoma	185 (150)	Cancer Cell 2017	PAAD
Paraganglioma & Pheochromocytoma	179 (173)	Cancer Cell 2017	PCPG
Prostate Adenocarcinoma	500 (333)	Cell 2015	PRAD
Sarcoma	261 (206)	Cell 2017	SARC
Skin Cutaneous Melanoma	470 (331)	Cell 2015	SKCM
Testicular Germ Cell Cancer	150 (137)	Cell Reports 2018	TGCT
Thymoma	124 (117)	Cancer Cell 2018	THYM
Thyroid Papillary Carcinoma	507 (496)	Cell 2014	THCA
Uterine Carcinosarcoma	57 (57)	Cancer Cell 2017	UCS
Uterine Corpus Endometrioid Carcinoma	560 (373)	Nature 2013	UCEC
Uveal Melanoma	80 (80)	Cancer Cell 2017	UVM

TCGA Molecular Characterization Platforms

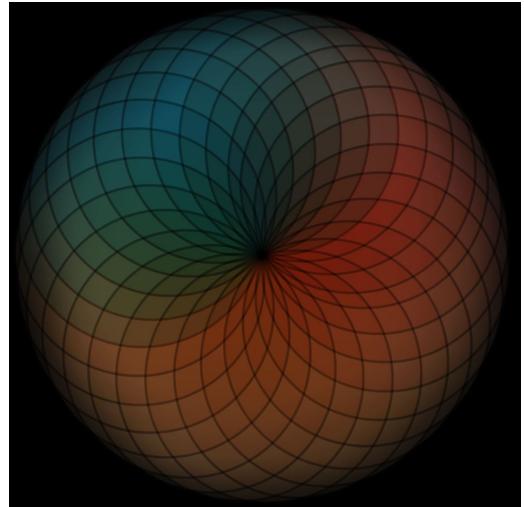
Center	TCGA Platform Code	DCC PlatformName	Instrument Support Materials	Sequence Download	TCGA ADF Download
Broad Institute of MIT and Harvard	ABI	Applied Biosystems Sequence Data	3730/3730xl DNA Analyzers	Primers	N/A
McDonnell Genome Institute at Washington University	ABI	Applied Biosystems Sequence Data	3730/3730xl DNA Analyzers	Primers	N/A
Human Genome Sequencing Center at Baylor College of Medicine	ABI	Applied Biosystems Sequence Data	3730/3730xl DNA Analyzers	Primers	N/A
University of North Carolina	AgilentG4502A_07_1	Agilent 244K Custom Gene Expression G4502A-07-1	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
University of North Carolina	AgilentG4502A_07_2	Agilent 244K Custom Gene Expression G4502A-07-2	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
University of North Carolina	AgilentG4502A_07_3	Agilent 244K Custom Gene Expression G4502A-07-3	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
Memorial Sloan Kettering Cancer Center	CGH-1x1M_G4447A	Agilent SurePrint G3 Human CGH Microarray Kit 1x1M	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
Broad Institute of MIT and Harvard	Genome_Wide_SNP_6	Affymetrix Genome-Wide Human SNP Array 6.0	Genome-Wide Human SNP Array 6.0	FASTA	TCGA ADF
McDonnell Genome Institute at Washington University	Genome_Wide_SNP_6	Affymetrix Genome-Wide Human SNP Array 6.0	Genome-Wide Human SNP Array 6.0	FASTA	TCGA ADF
University of North Carolina	H-miRNA_8x15K	Agilent 8 x 15K Human miRNA-specific microarray	Human mouse and rat miRNA Microarray	FASTA	TCGA ADF
University of North Carolina	H-miRNA_8x15Kv2	Agilent Human miRNA Microarray Rel12.0	Human mouse and rat miRNA Microarray	FASTA	TCGA ADF
Memorial Sloan Kettering Cancer Center	HG-CGH-244A	Agilent Human Genome CGH Microarray 244A	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
Harvard Medical School	HG-CGH-244A	Agilent Human Genome CGH Microarray 244A	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
Harvard Medical School	HG-CGH-415K_G4124A	Agilent Human Genome CGH Custom Microarray 2x415K	SurePrint G3 CGH+SNP Microarray	FASTA	TCGA ADF
McDonnell Genome Institute at Washington University	HG-U133_Plus_2	Affymetrix Human Genome U133 Plus 2.0 Array	Human Genome U133 Array	FASTA	TCGA ADF
Broad Institute of MIT and Harvard	HT_HG-U133A	Affymetrix HT Human Genome U133 Array Plate Set	Human Genome U133 Array	TSV	TCGA ADF
Berkeley Lab	HuEx-1_0-st-v2	Affymetrix Human Exon 1.0 ST Array	Exon 1.0 ST Array	FASTA	TCGA ADF
HudsonAlpha Institute for Biotechnology	Human1MDuo	Illumina Human1M-Duo BeadChip	Infinium HD BeadChip DNA	FASTA	TCGA ADF
HudsonAlpha Institute for Biotechnology	HumanHap550	Illumina 550K Infinium HumanHap550 SNP Chip	Infinium BeadChip DNA	FASTA	TCGA ADF
Johns Hopkins-USC Epigenome Center	HumanMethylation27	Illumina Infinium Human DNA Methylation 27	Infinium Human DNA Methylation 27	See ADF File	TCGA ADF
Johns Hopkins-USC Epigenome Center	HumanMethylation450	Illumina Infinium Human DNA Methylation 450	Infinium Human DNA Methylation 450	See ADF File	TCGA ADF
Johns Hopkins-USC Epigenome Center	IlluminaDNAMethylation_OMA002_CPI	Illumina DNA Methylation OMA002 Cancer Panel I	GoldenGate DNA Methylation	FASTA	TCGA ADF
Johns Hopkins-USC Epigenome Center	IlluminaDNAMethylation_OMA003_CPI	Illumina DNA Methylation OMA003 Cancer Panel I	GoldenGate DNA Methylation	See ADF File	TCGA ADF
McDonnell Genome Institute at Washington University	IlluminaGA_DNASeq	Illumina Genome Analyzer DNA Sequencing	Genome Analyzer IIx	N/A	N/A
Broad Institute of MIT and Harvard	IlluminaGA_DNASeq	Illumina Genome Analyzer DNA Sequencing	Genome Analyzer IIx	N/A	N/A
Human Genome Sequencing Center at Baylor College of Medicine	IlluminaGA_DNASeq	Illumina Genome Analyzer DNA Sequencing	Genome Analyzer IIx	N/A	N/A
University of California Santa Cruz	IlluminaGA_DNASeq	Illumina Genome Analyzer DNA Sequencing	Genome Analyzer IIx	N/A	N/A
University of North Carolina	IlluminaGA_DNASeq	Illumina Genome Analyzer DNA Sequencing	Genome Analyzer IIx	N/A	N/A
BC Cancer Agency	IlluminaGA_miRNASeq	Illumina Genome Analyzer miRNA Sequencing	Genome Analyzer IIx	N/A	N/A
Harvard Medical School	IlluminaGA_mRNA_DGE	Illumina Genome Analyzer mRNA Digital Gene Expression	Genome Analyzer IIx	N/A	N/A
BC Cancer Agency	IlluminaGA_RNASEq	Illumina Genome Analyzer RNA Sequencing	Genome Analyzer IIx	N/A	N/A
University of North Carolina	IlluminaGA_RNASeqV2	Illumina Genome Analyzer RNA Sequencing Version 2 analysis	Genome Analyzer IIx	N/A	N/A
Harvard Medical School	IlluminaHiSeq_DNASeq	Illumina HiSeq for Copy Number Variation	HiSeq 2000	N/A	N/A
BC Cancer Agency	IlluminaHiSeq_miRNASeq	Illumina HiSeq 2000 miRNA Sequencing	HiSeq 2001	N/A	N/A
University of North Carolina	IlluminaHiSeq_RNASeq	Illumina HiSeq 2000 RNA Sequencing	HiSeq 2002	N/A	N/A
University of North Carolina	IlluminaHiSeq_RNASeqV2	Illumina HiSeq 2000 RNA Sequencing Version 2 analysis	HiSeq 2003	N/A	N/A
University of North Carolina	IlluminaHiSeq_TotalRNASeqV2	Illumina HiSeq 2000 Total RNA Sequencing Version 2 analysis	HiSeq 2004	N/A	N/A
Johns Hopkins-USC Epigenome Center	IlluminaHiSeq_WGBS	Illumina HiSeq 2000 Bisulfite-converted DNA Sequencing	HiSeq 2005	N/A	N/A
MD Anderson Cancer Center	MDA_RPPA_Core	M.D. Anderson Reverse Phase Protein Array Core	MD Anderson RPPA Core Facility	N/A	N/A
Nationwide Children's Hospital	microsat_i	Microsatellite Instability Analysis	SOLID 3 Plus	N/A	N/A
University of California Santa Cruz	SOLID_DNASeq	ABI SOLID DNA System Sequencing	SOLID 3 Plus	N/A	N/A
Broad Institute of MIT and Harvard	SOLID_DNASeq	ABI SOLID DNA System Sequencing	SOLID 3 Plus	N/A	N/A
Human Genome Sequencing Center at Baylor College of Medicine	SOLID_DNASeq	ABI SOLID DNA System Sequencing	SOLID 3 Plus	N/A	N/A

TCGA Findings

TCGA outcomes & impact

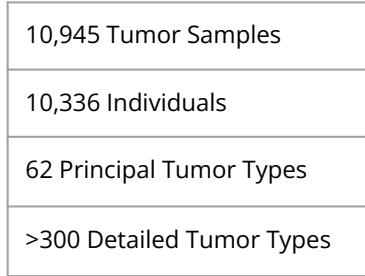
TCGA's Pan-Cancer Atlas

TCGA Research Network Publications

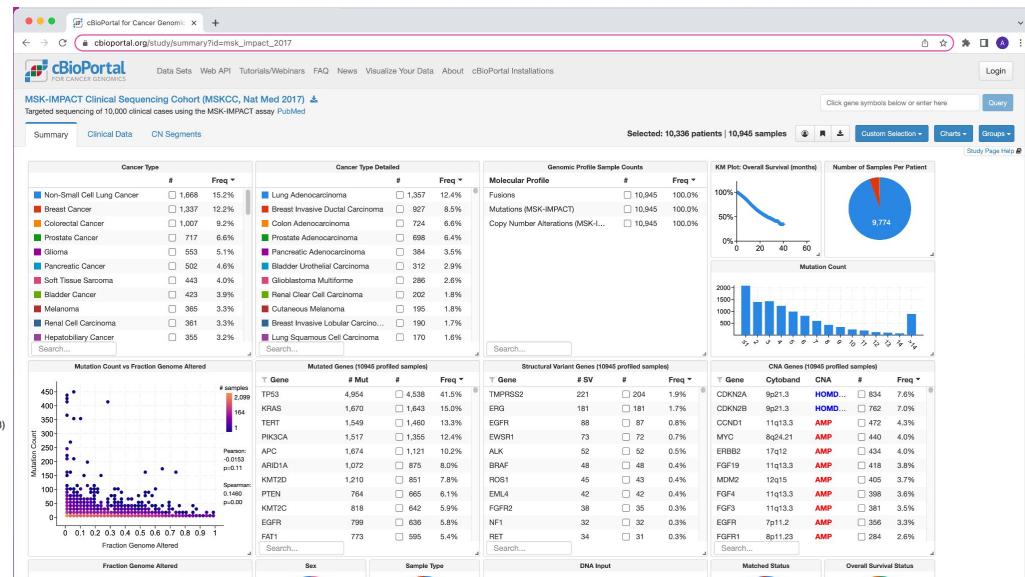
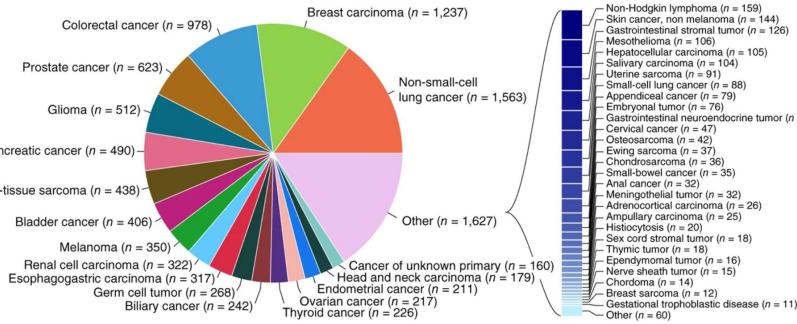


TS:MSK-IMPACT

- Hybridization capture-based NGS panel capable of detecting protein-coding mutations, copy number alterations (CNAs), and selected promoter mutations and structural rearrangements in **341 (and, more recently, 410) cancer-associated genes**
- Full data set publicly available through cBioPortal for Cancer Genomics: <http://cbioportal.org/msk-impact>



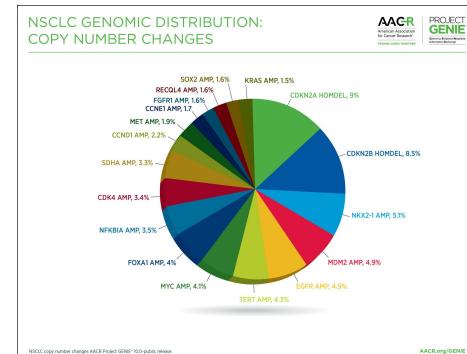
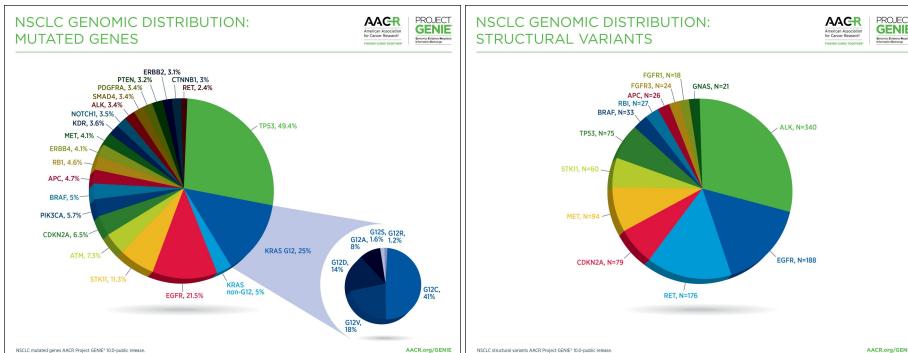
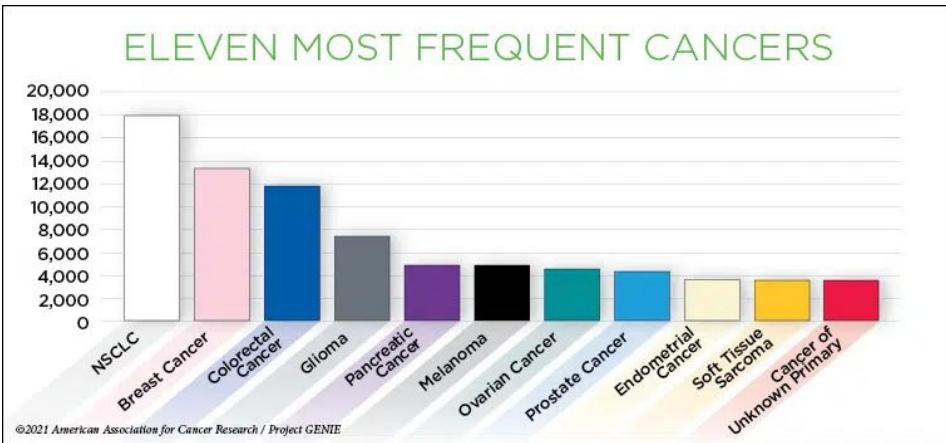
a



TS: Genomics Evidence Neoplasia Information Exchange (GENIE)

The AACR Project GENIE is an international data-sharing consortium focused on generating an evidence base for precision cancer medicine by integrating clinical-grade cancer genomic data with clinical outcome data for tens of thousands of cancer patients treated at multiple institutions worldwide.

The first public data release was available to the global community in January 2017; our current release, **GENIE 11.0-public**, was released in January 2022. The registry now contains over **136,000 sequenced samples** from over **121,000 patients**, making the AACR Project GENIE registry among the largest fully public cancer genomic data sets released to date.



Data can be access via [cBioPortal](#) or download data directly from [Sage Bionetworks](#).

Major Cancer Genomic Studies with Data Portal Available

Other Data Portals



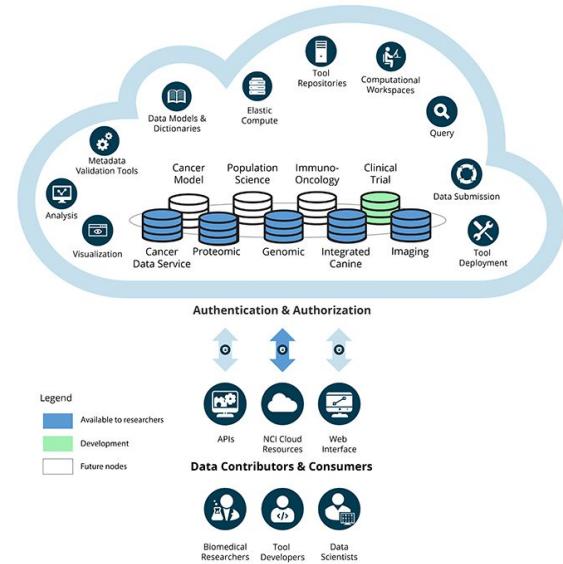
Cancer Research Data Commons (CRDC) at NCI

Key Datasets

Dataset Name	Description	Available Resources
The Cancer Genome Atlas (TCGA)	A collaboration between NCI and the National Human Genome Research Institute (NHGRI) that has characterized tumor and normal tissues from 11,000 patients, covering 33 cancer types	GDC, Broad, SB, ISB, IDC
Therapeutically Applicable Research to Generate Effective Treatments (TARGET)	A consortium of extramural and NCI investigators working to characterize and understand hard-to-treat childhood cancers and translate findings into the clinic.	GDC, Broad, SB, ISB
Clinical Proteomic Tumor Analysis Consortium (CPTAC)	A national effort to accelerate the understanding of the molecular basis of cancer through the application of large-scale proteome and genome analysis, or proteogenomics.	GDC, PDC, Broad, ISB, SB, IDC
Human Cancer Model Initiative (HCFI)	An international consortium that is generating novel, next-generation, tumor-derived culture models complete with genomic and clinical data.	GDC, SB, ISB
Cancer Genome Characterization Initiatives (CGCI)	An initiative examining genomes, exomes, and transcriptomes of various types of adult and pediatric cancers.	GDC, SB, ISB
Foundation Medicine (FM) ↗	Targeted sequencing data from ~18,000 adult patients generated by the Foundation Medicine Inc., molecular information company seeking to match patients with personalized treatment plans.	GDC, SB, ISB
Multiple Myeloma Research Foundation (MMRF) ↗	Data from nearly 1,000 patients with extensive molecular and clinical data, including longitudinal information collected over the course of disease for many patients.	GDC, SB, ISB
Genomics Evidence Neoplasia Information Exchange (GENIE) ↗	Over 44,000 cases from the international pan-cancer registry continuing to be collected by the American Association for Cancer Research (AACR) initiative.	GDC, SB, ISB
International Cancer Proteogenomic Consortium (ICPC)	An international consortium that brings together more than a dozen countries to study the application of proteogenomic analysis in predicting cancer treatment success and to share data and results with researchers worldwide, hastening progress for patients.	PDC, SB, ISB
Children's Brain Tumor Tissue Consortium (CBTTC) ↗	A collaborative research consortia focused on identifying therapies for children with brain tumors	PDC, SB, ISB
Genetics and Epidemiology of Colorectal Cancer Consortium (GECCO) ↗	A research collaboration to detect colorectal cancer susceptibility loci using genome-wide sequencing	CDS, SB
Comparative molecular life history of spontaneous canine and human gliomas (GLIOMA01) ↗	Characterization of the genomic and transcriptomic landscape of canine glioma to enable cross-species comparative genomic analysis of sporadic glioma	ICDC, SB

CRDC is a cloud-based data science infrastructure that connects data sets with analytics tools to allow users to share, integrate, analyze, and visualize cancer research data to drive scientific discovery.

NCI Cancer Research Data Commons (CRDC)



Genomic Data Commons (GDC)

Proteomic Data Commons (PDC)

Integrated Canine Data Commons (ICD)

Imaging Data Commons (IDC)

Cancer Data Commons (CDS)

NCI Cloud Resources

GDC Data Portal

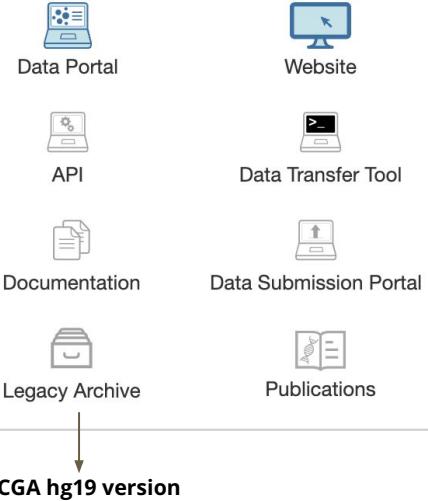
The screenshot shows the GDC Data Portal homepage. At the top, there's a navigation bar with links for Home, Projects, Exploration, Analysis, Repository, Quick Search, Manage Sets, Login, Cart, and GDC Apps. A red box highlights the "GDC Apps" link. Below the navigation is a search bar with placeholder text "e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2". To the right is a "Data Portal Summary" section showing: PROJECTS 70, PRIMARY SITES 67, CASES 85,416, GENES 21,754, FILES 819,832, and MUTATIONS 2,670,227. A large central area features a human body silhouette with colored regions representing different cancer types. To the right is a bar chart titled "Cases by Major Primary Site" with the following data:

Major Primary Site	Cases (in thousands)
Adrenal Gland	0.1
Bile Duct	0.1
Bladder	0.2
Brain	0.5
Bone Marrow	1.0
Cervix	0.2
Esophagus	0.5
Eye	0.1
Head and Neck	0.2
Lung	1.5
Lymph Nodes	0.1
Nervous Tissue	0.1
Ovary	0.2
Pancreas	0.1
Prostate	0.2
Soft Tissue	0.1
Stomach	0.1
Thyroid	0.1
Uterus	0.1

Below the summary is a section titled "GDC Applications" with icons for Data Portal, Website, API, Data Transfer Tool, Documentation, Data Submission Portal, Legacy Archive, and Publications. At the bottom, there's a footer with links for Site Home, Policies, Accessibility, FOIA, Support, U.S. Department of Health and Human Services, National Institutes of Health, National Cancer Institute, USA.gov, NIH - Turning Discovery Into Health, UI v1.28.0 @ 0798511, API v3.0.0 @ d0ab653b, Data Release 32.0 - March 29, 2022.

Bioinformatics pipelines:

- Bioinformatics Pipeline: DNA-Seq Analysis
- Bioinformatics Pipeline: mRNA Analysis
- Bioinformatics Pipeline: miRNA Analysis
- Bioinformatics Pipeline: Copy Number Variation Analysis
- Bioinformatics Pipeline: Methylation Analysis Pipeline
- Bioinformatics Pipeline: Protein Expression



Analysis:

Set Operations

Cohort Comparison

Clinical Data Analysis

Visualization:

Cases/Genes/Mutations/Oncogrid

cBioPortal for Cancer Genomics

Major features:

- Host multi-omics and clinical data from >340 cancer genomic studies and >30 tissue sites.
- Support many genomic visualization and analyses, including mutational distribution, oncoplot, mutual exclusivity analysis, co-expression, group comparison analysis, survival analysis, integrative analysis, etc.
- All the data are harmonized in the same format and can be directly downloaded from the web or [Datahub](#).
- Provide R client for accessing the study datasets.
- Support OQL & Expression for query.
- Provide link to share the query.
- Support local installation.
- Web-based API (Application Programming Interface) available.

Exclusively for tumor data !!

What data is in cBioPortal?

Data sources



THE CANCER GENOME ATLAS

National Cancer Institute
National Human Genome Research Institute



International
Cancer Genome
Consortium



CCLE Cancer Cell Line Encyclopedia



PROJECTGENIE
Genomics Evidence Neoplasia Information Exchange

- Clinical Data:
- Treatments
 - Survival
 - Etc

- omic data:
- Mutations
 - Fusions
 - Copy number
 - mRNA expression
 - Protein levels
 - DNA Methylation

Background biological data
(e.g. networks, 3D protein structure)

 **cBioPortal**
for Cancer Genomics

<https://genie.cbiportal.org/>

Curated effect & therapy implications



Precision Oncology Knowledge Base



CLINICAL INTERPRETATIONS OF

VARIANTS IN CANCER



MY CANCER GENOME®
GENETICALLY INFORMED CANCER MEDICINE

Predicted functional effect

PolyPhen-2

 mutationassessor.org
functional impact of protein mutations



Variant recurrence



Catalogue of somatic mutations



Download Datasets

luad_tcga_pan_can_atlas_2018.tar.gz



Data Sets Web API R/MATLAB Tutorials/Webinars FAQ News Visualize Your Data About cBioPortal Installations

case_lists	
data_armlevel_cna.txt	
data_clinical_patient.txt	
data_clinical_sample.txt	
data_clinical_supp_hypoxia.txt	
data_cna_hg19.seg	
data_cna.txt	
data_fusions.txt	
data_log2_cna.txt	
data_microbiome.txt	
data_mrna_seq_v2_rsem_zscores_ref_all_samples.txt	
data_mrna_seq_v2_rsem_zscores_ref_diploid_samples.txt	
data_mrna_seq_v2_rsem_zscores_ref_normal_samples.txt	
data_mrna_seq_v2_rsem.txt	
data_mutations.txt	
data_normals_RNA_Seq_v2_mRNA_median_Zscores.txt	
data_normals_RNA_Seq_v2_mRNA_median.txt	
data_rppa_zscores.txt	
data_rppa.txt	
LICENSE	
meta_armlevel_cna.txt	
meta_clinical_patient.txt	
meta_clinical_sample.txt	
meta_cna_hg19_seg.txt	
meta_cna.txt	
meta_fusions.txt	
meta_log2_cna.txt	
meta_microbiome.txt	
meta_mrna_seq_v2_rsem_zscores_ref_all_samples.txt	
meta_mrna_seq_v2_rsem_zscores_ref_diploid_samples.txt	
meta_mrna_seq_v2_rsem_zscores_ref_normal_samples.txt	
meta_mrna_seq_v2_rsem.txt	
meta_mutations.txt	
meta_rppa_zscores.txt	
meta_rppa.txt	
meta_study.txt	
normals	
data_mrna_seq_v2_rsem_normal_samples_zscores_ref_normal_samples.txt	
data_mrna_seq_v2_rsem_normal_samples.txt	
meta_mrna_seq_v2_rsem_normal_samples_zscores_ref_normal_samples.txt	
meta_mrna_seq_v2_rsem_normal_samples.txt	

Datasets

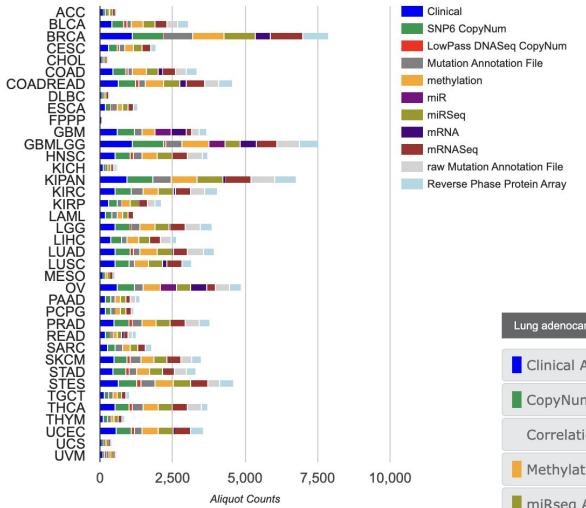
The table below lists the number of available samples per cancer study and data type. It also provides links to download the data for each study. For alternative ways of downloading, see the [Download Documentation](#).

Name	Reference	All	Mutations	CNA	RNA-Seq
Acinar Cell Carcinoma of the Pancreas (UHL, J Pathol 2014)	Jia et al. J Pathol 2014	23	23	0	0
Acral Melanoma (TGEN, Genome Res 2017)	Liang et al. Genome Res 2017	38	38	38	36
Acute Lymphoblastic Leukemia (St Jude, Nat Genet 2015)	Anderson et al. Nat Genet 2015	93	93	0	0
Acute Lymphoblastic Leukemia (St Jude, Nat Genet 2016)	Zhang et al. Nat Genet 2016	73	73	0	0
Acute Myeloid Leukemia (OHsu, Nature 2018)	Tyner et al. Nature 2018	672	622	0	451
Acute Myeloid Leukemia (TGCA, Firehose Legacy)		200	197	191	173
Acute Myeloid Leukemia (TGCA, NEJM 2019)		200	200	191	173
Acute Myeloid Leukemia (TGCA, PanCancer Atlas)		200	200	191	173
Acute myeloid leukemia or myelodysplastic syndromes (WashU, J. 2017)		136	136	0	0
Adenoid Cystic Carcinoma (FMI, Am J Surg Pathl, 2014)	Ross et al. Am J Surg Pathl 2014	28	28	28	0
Adenoid Cystic Carcinoma (JHU, Cancer Prev Res 2016)	Retig et al. Cancer Prev Res 2016	25	25	0	0
Adenoid Cystic Carcinoma (MDA, Clin Cancer Res 2015)	Mitani et al. Clin Cancer Res 2015	102	65	0	0
Adenoid Cystic Carcinoma (MGH, Nat Gen 2016)	Drier et al. Nature Genetics 2016	10	10	0	0
Adenoid Cystic Carcinoma (MSKCC, Nat Genet 2013)	Ho et al. Nat Genet 2013	60	60	60	0
Adenoid Cystic Carcinoma (Sanger/MDA, JCI 2013)	Stephens et al. JCI 2013	24	24	0	0
Adenoid Cystic Carcinoma of the Breast (MSKCC, J. Pathol. 2015)	Martelotto et al. J Pathol 2015	12	12	12	0
Adenocarcinoma Project (J Clin Invest 2019)	Allen et al. J Clin Invest 2019	1049	1049	928	0
Adrenocortical Carcinoma (TGCA, Firehose Legacy)		92	90	90	79
Adrenocortical Carcinoma (TGCA, PanCancer Atlas)		92	91	89	78
Adult Soft Tissue Sarcoma (TGCA, Cell 2017)		206	206	206	206
Amputary Carcinoma (Baylor College of Medicine, Cell Reports 2018)	Gingras et al. Cell Rep 2016	160	160	0	0
Anaplastic Oligodendroglioma and Anaplastic Oligoastrocytoma (MSKCC, Neuro Oncol 2017)	Thomas et al. Neuro Oncol 2017	22	22	22	0
Basal Cell Carcinoma (UNIGE, Nat Genet 2016)	Bonilla et al. Nat Genet 2016	293	293	0	0
Bladder Cancer (MSKTCGA, 2020)		476	474	442	296
Bladder Cancer (MSKCC, Eur Urol 2014)	Kim et al. Eur Urol 2015	109	109	109	0
Bladder Cancer (MSKCC, J Clin Oncol 2013)	Iyer et al. J Clin Oncol 2013	97	97	97	0
Bladder Cancer (MSKCC, Nat Genet 2016)	Al-Ahmadi et al. Nat Genet 2016	34	34	33	0
Bladder Cancer (TGCA, Cell 2017)	Robertson et al. Cell 2017	413	412	408	408
Bladder Urothelial Carcinoma (BGI, Nat Genet 2013)	Guo et al. Nat Genet 2013	99	99	0	0
Bladder Urothelial Carcinoma (DFCI/MSKCC, Cancer Discov 2014)	Van Allen et al. Cancer Discov 2014	50	50	0	0
Bladder Urothelial Carcinoma (TGCA, Firehose Legacy)		413	130	408	408
Bladder Urothelial Carcinoma (TGCA, Nature 2014)		131	130	128	129
Bladder Urothelial Carcinoma (TGCA, PanCancer Atlas)		411	410	408	407
Brain Lower Grade Glioma (TGCA, Firehose Legacy)		530	286	513	530
Brain Lower Grade Glioma (TGCA, PanCancer Atlas)		514	514	511	514
Brain Tumor PDXs (Mayo Clinic, 2018)		97	83	83	66
Breast Cancer (HTAN, Cell Rep Med 2022)	Johnson et al. Cell Rep Med. 2022	5	5	5	0
Breast Cancer (METABRIC, Nature 2012 & Nat Commun 2016)	Pereira et al. Nat Commun 2016, Rueda et al. Nature 2019, Curtis et al. Nature 2012	2509	2509	2173	1904
Breast Cancer (MSK, Cancer Cell 2018)	Razavi et al. Cancer Cell 2018	1918	1918	1918	0
Breast Cancer (MSK, Clinical Cancer Res 2020)		60	60	0	0
Breast Cancer (NCI, NCI-60, NCI-80, NCI-96)	Quinn et al. Mol Cancer Ther 2009	444	444	444	444

FIREBROWSE (Broad GDAC)

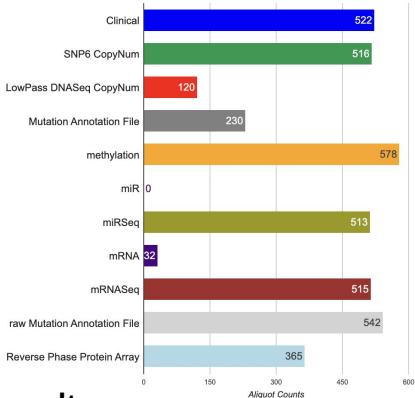
Dataset: TCGA

TCGA data version 2016_01_28

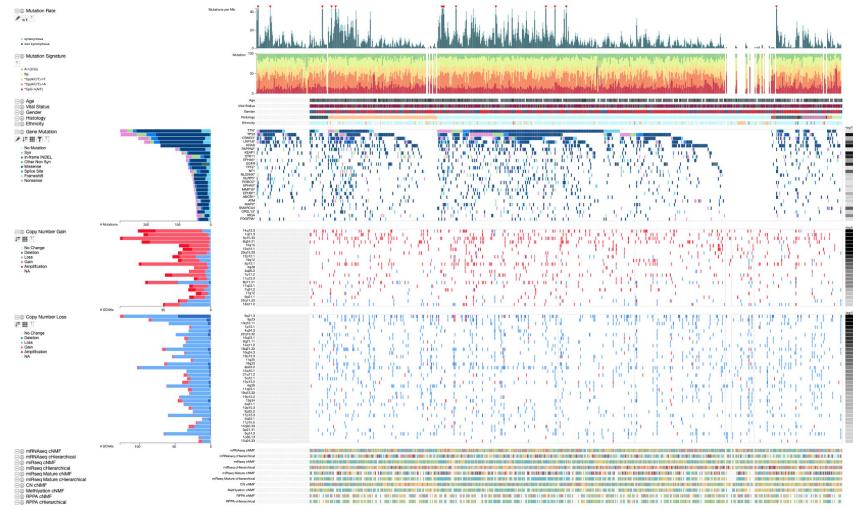


Major genomic analysis results

TCGA data version 2016_01_28 for LUAD



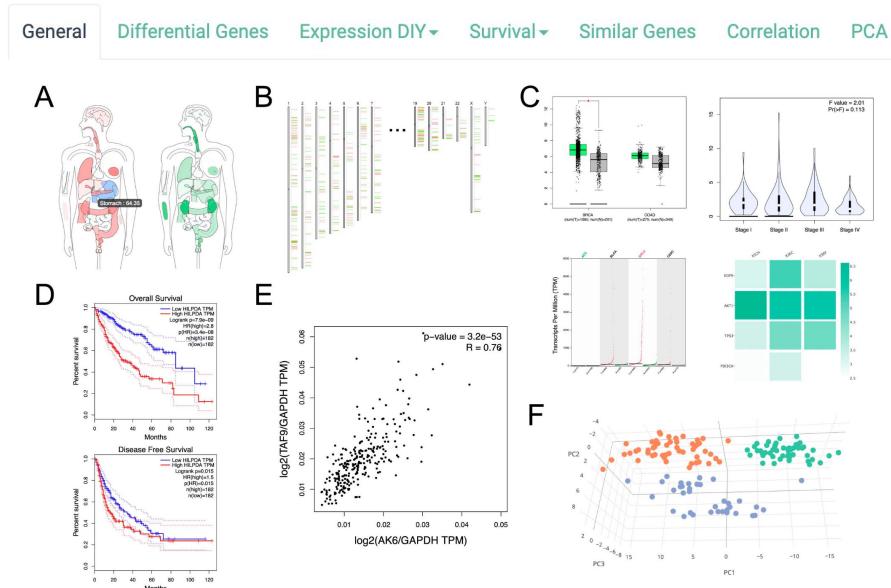
iCoMut for FireBrowse



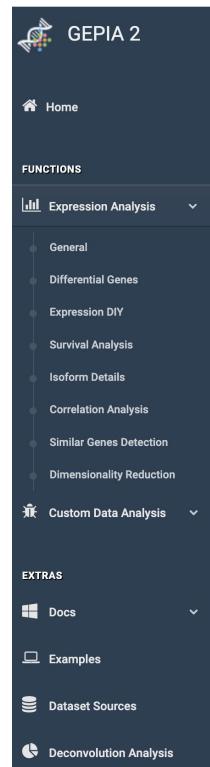
Gene Expression Profiling Interactive Analysis (GEPIA)

Datasets: **TCGA** and **GTEX** data

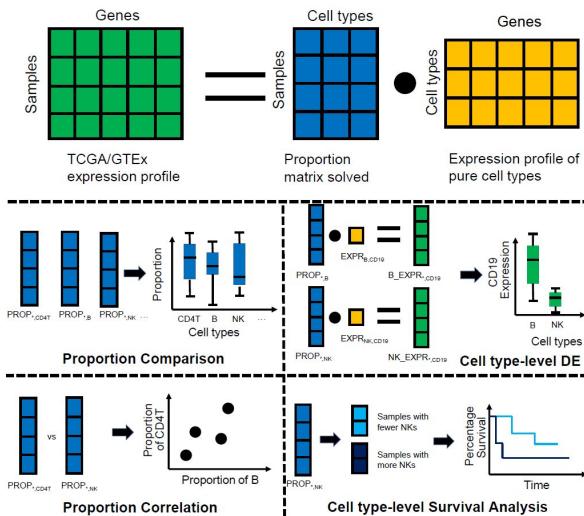
- Single Gene Analysis
- Cancer Type Analysis
- Multiple Gene Analysis



[GEPIA2](#)



[GEPIA2021](#)

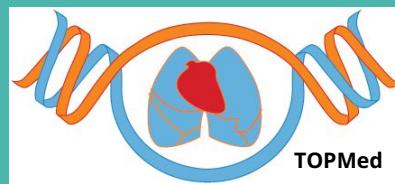
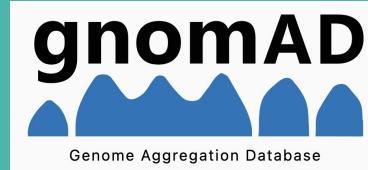


a standalone extension with multiple deconvolution based analysis for GEPIA. We deconvolute each sample tool in TCGA/GTEX with the bioinformatics tools **CIBERSORT**, **EPIC** and **quantiTseq**. Based on the inferred cell proportions in each bulk-RNA sample, we can then perform multiple downstream analysis:

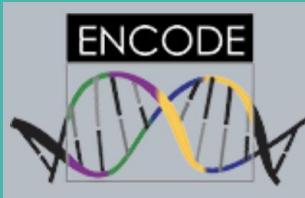
Summary of Data Portals

Data Portals	Studies	Data type	Interactive analysis	Integrative analysis	Visualization	Access Control	Recommendation
GEPIA	TCGA; GTEx	Expression data	Yes	No	Limited	No	Tumor-Normal expression analysis
GTEx Portal	GTEx	Gene expression and genotyping data	Limited	Limited	Yes	Raw data and most analyzed data	Germline analysis
ICGC Portal	PCAWG and other ICGC Data	Multi-Omics Data	Limited	No	No	Raw data and most analyzed data	Data download
GDC	>20 Studies (e.g. TCGA, GENIE)	Multi-Omics Data	Limited	No	Limited	Raw data and most analyzed data	Data download
Firebrowse	TCGA	Multi-Omics Data	No	Yes	Limited	No	TCGA deep analysis
St. Jude Cloud	PeCan	Multi-Omics Data	No	No	Yes	Yes	Data download and visualization
cBioPortal	>340 Studies (e.g. TCGA, IMPACT, GENIE, PCAW, ICGC)	Multi-Omics Data	Yes	Yes	Yes	No	Integrative analysis

Other Genetic Data Resources



Cell Model Passports
A Hub for Preclinical Cancer Models



Genotype-Tissue Expression (GTEx)

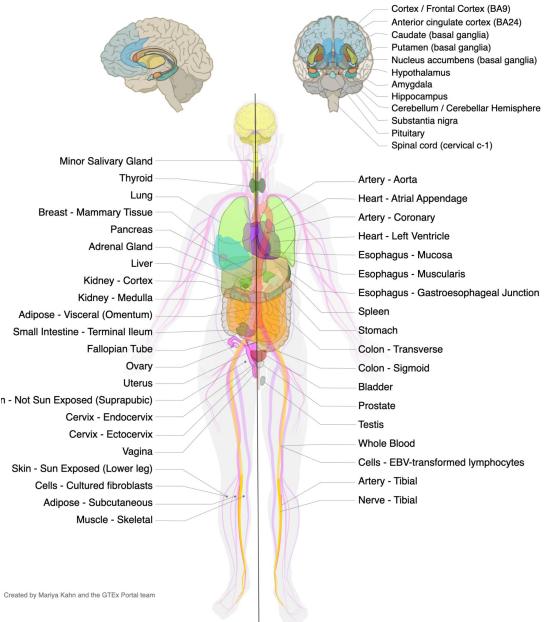
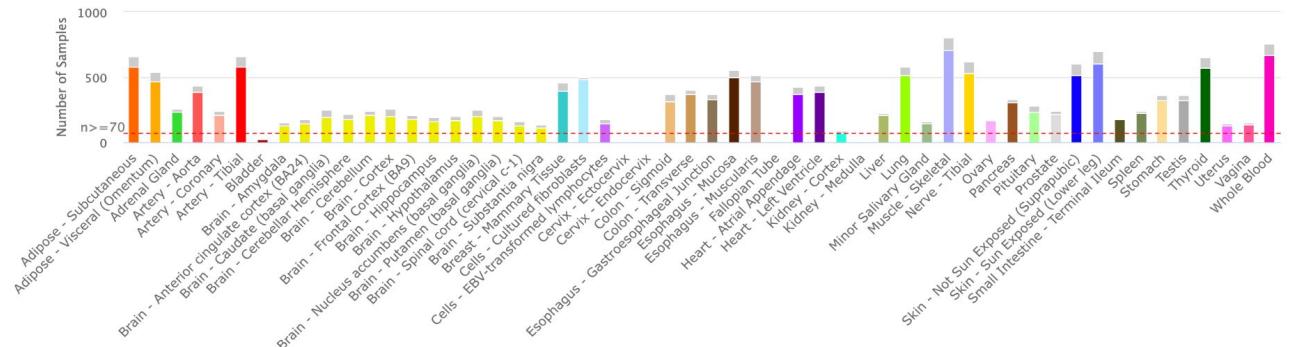
Public resource of tissue-specific gene expression

Samples collected from **54** non-diseased tissue sites across over 900 individuals

Datasets include SNP array, WGS, WES, bulk RNA-Seq and snRNA-Seq

Available protected access data include:

- BAM files for RNA-Seq, Whole Exome Seq, and Whole Genome Seq
- Genotype Calls (.vcf) for WES and WGS
- Allele Specific Expression (ASE) tables
- All expression matrices from the Portal, including samples that did not pass the Analysis Freeze QC
- Expanded sample attributes
- Expanded subject phenotypes, including age and ethnicity



Created by Marily Kahn and the GTEx Portal team

V8 Release	# Tissues	# Donors	# Samples
Total	54	948	17382
With Genotype	54	838	15253
Has eQTL Analysis*	49	838	15201

* Number of samples with genotype >= 70

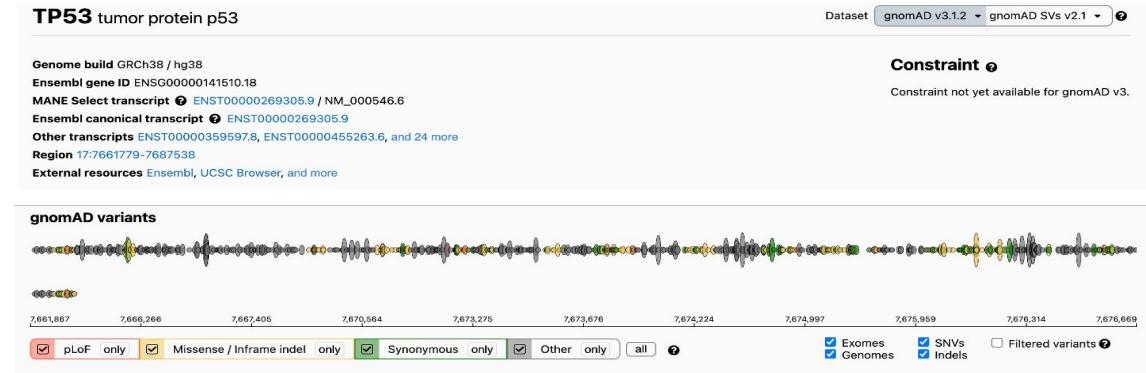
gnomAD: Genome Aggregation Database



gnomAD is a coalition of investigators seeking to aggregate and harmonize exome and genome sequencing data from a variety of large-scale sequencing projects, and to make summary data available for the wider scientific community.

Available access data include:

- Variants & Coverage &Constraint
- Multi-nucleotide variants (MNVs)
- Proportion expressed across transcripts
- Structural variants
- Loss-of-function curation results
- Variant co-occurrence
- Linkage disequilibrium
- Ancestry classification



Different Version:

- gnomAD v2.1 data set contains data from **125,748 exomes** and **15,708 whole genomes**, all mapped to the **GRCh37/hg19** reference sequence.
- gnomAD v3.1 data set contains **76,156 whole genomes**, all mapped to the **GRCh38** reference sequence.
- gnomAD v3.1 contains a substantially larger number of African American samples than v2.1 and provides allele frequencies in the Amish population for the first time. gnomAD v3.1 also has a fully genotyped callset available from the Human Genome Diversity Project and 1000 Genomes Project, representing > 60 distinct populations.

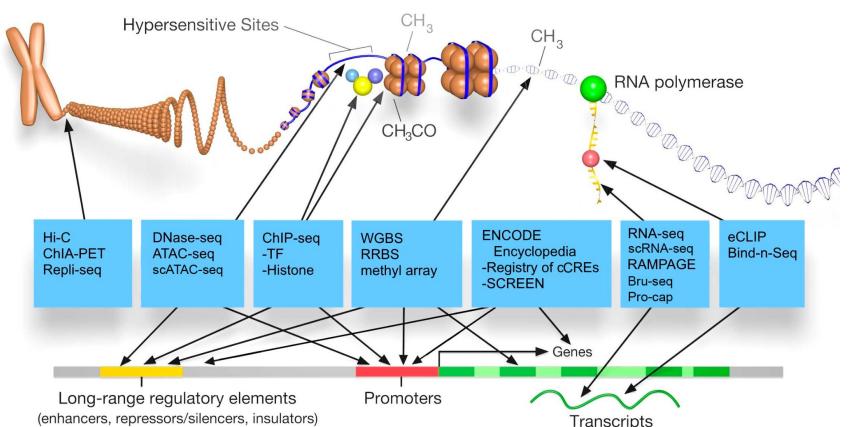
EnCODE and Portal

ENCODE Encyclopedia Version 5:

The ENCODE Consortium not only produces high-quality data, but also analyzes the data in an integrative fashion. The ENCODE Encyclopedia organizes the most salient analysis products into annotations and provides tools to search and visualize them. The

Encyclopedia has two levels of annotations:

- Integrative-level annotations integrate multiple types of experimental data and ground level annotations.
 - Ground-level annotations are derived directly from the experimental data, typically produced by uniform processing pipelines.



Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

Functional genomics

Functional characterization

Encyclopedia of elements

Ruth Ashkenazi's

BMTex

Deeply profiled cell lines

Computational and integrative products

Human donors

ENCORE

Southeast Stem Cell Consortium

Imputed experiments

Immune cells

Mouse development

Reference epigenome

Functional genomic series

Single-cell experiments

RNA-seq

Register search

Encyclopedia browser

ChIP-seq experiments

SSAY →

cell line	TF ChIP-seq	Histone ChIP-seq	Control ChIP-seq	Disease-seq	miRNA-seq	polyA+RNA-seq	shRNA-seq	ATAC-seq	CRISPR-seq	microRNA-seq	snRNA-seq	RNAmicroarray	WES	DNAmarray	Control ChIP-seq	ChIP-seq	RIBS	small RNAseq	ChIP-PET	long read sequencing	
K562	3057	904	885	253	194	138	205	563	206	18	415	42	32	110	21	91	236	229	103	112	145
HepG2	716	15	84	3	7	11	270	64	230	3	3	12	1	3	127	124	1	8	11	4	
A549	200	15	60	3	11	11	22	6	185	2	3	7	2	1	109	105	2	3	5	3	
GM12878	11	15	12	3	6	14	22	13	4	3	8	2	3	1	2	6	5	1	2	2	
HCT116	27	17	17	18	1	88	1	17	2	1	1	2	1	2	1	1	1	32	1	1	

tissue	middle frontal area	heart	adrenal gland	liver	stomach		
middle frontal area	50	188	61	56	120		
heart	10	10	16	7	43		
adrenal gland	8	40	13	23	11	8	59
liver	42	93	29	16	3	22	8
stomach	20	75	26	23	5	15	6

primary cell	CD4+ positive, alpha-beta T cell	CD4+ positive, alpha-beta T cell	CD4+ positive memory T cell	naive thymus-derived CD4+ positive, alpha-beta T cell	CD4+ positive, alpha-beta T cell	T cell
07	342	142	561	219	714	169
18	3	16	4	74	9	2
11	3	62	5	18	1	6
21	5	8	41	8		
			11	2	66	6
				3		
15	4	5	2	48	7	

whole organisms	1084	3	1073	272	97	105
whole organism	1084	3	1073	268	95	105
			4	12		

in vitro differentiated cells	dendrite cell	mesenchymal stem cell	motor neuron	endothelial cell	chondrocyte	
59	344	73	48	55	78	36
				41	6	
				18	11	7
				24	6	
				18	12	2
						11

cell-free sample	cell-free sample
cell-free sample	199

organoid	brain	nephron
7	4	3
Y	4	2
2	29	18

technical sample	technical sample
technical sample	7
	1

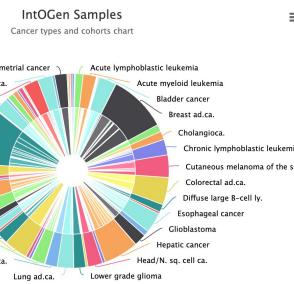
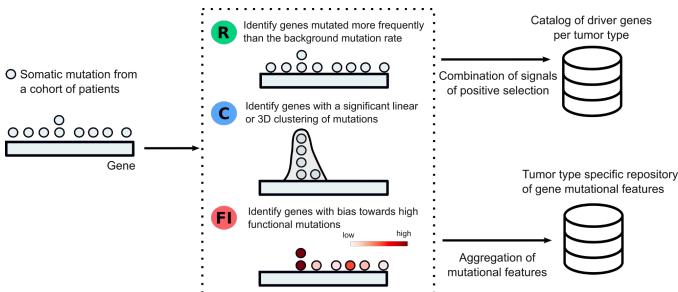
Specialized Databases for Genomic Analyses

Driver Gene Analysis

Integrative Onco Genomics (IntOGen)

Collects and analyzes somatic mutations in tumor genomes to identify cancer driver genes

66 cancer types
221 cohorts
28,076 samples
>200 million mutations
568 drivers



OncoKB

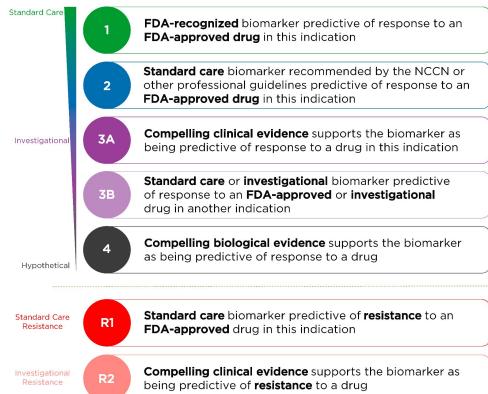
MSK's Precision Oncology Knowledge Base

Annotation of biological consequences and clinical implications of genetic variants in cancer

Database contains: 688 genes, 5729 alterations, 133 cancer types, 111 drugs

Potential driver genes in different categories. Each category has a series of levels:

- Therapeutic levels
- Diagnostic levels
- Prognostic levels
- FDA levels



Cosmic Cancer Gene Census

The Cancer Gene Census (CGC) is an ongoing effort to catalogue those genes which contain mutations that have been causally implicated in cancer and explain how dysfunction of these genes drives cancer.

Cancer Hotspots

A resource for statistically significant mutations in cancer identified in large scale cancer genomics data

Single residue and in-frame indel mutation hotspots in 24,592 tumor samples

Can download hotspot results and mutation data (MAF)

Show/Hide † Mouse over Variants and Samples values for more information

Search:

Gene	Residue	Type	Variants †	Q-value	Samples †
NRAS	Q61	single residue	R K L	0	422
PIK3CA	E545	single residue	K	0	633
IDH1	R132	single residue	H C	0	766
PIK3CA	H1047	single residue	R L	0	647
BRAF	V600	single residue	E	0	897
EGFR	L858	single residue	R	0	144
TP53	R175	single residue	H	0	416
KRAS	Q61	single residue	H R L K	0	190
KRAS	G13	single residue	D C	0	264
TP53	R248	single residue	O W	0	560
KRAS	G12	single residue	D V C R	0	2175
TP53	R273	single residue	C H L	0	609
PIK3CA	E542	single residue	K	0	372
AKT1	E17	single residue	K	1.15e-288	349
GNAS	R201	single residue	H C	7.47e-257	139
FGFR3	S249	single residue	C	3.49e-239	114
PIK3CA	N345	single residue	K I	3.87e-219	98
PTEN	R130	single residue	Q G *	1.92e-210	168
HRAS	Q61	single residue	R K L	8.15e-210	102
TP53	Y220	single residue	C	4.54e-208	152

Showing 1 to 20 of 1,165 mutations

Previous 1 2 3 4 5 ... 59

FIREBROWSE (Broad GDAC) for Somatic Analysis

CopyNumber Analyses

Aggregate AnalysisFeatures
 CopyNumber Clustering CNMF
 CopyNumber Clustering CNMF thresholded
CopyNumber Gistic2
 CopyNumberLowPass Gistic2

Correlate Clinical vs CopyNumber Arm
 Correlate Clinical vs CopyNumber Focal

Correlate CopyNumber vs mRNA

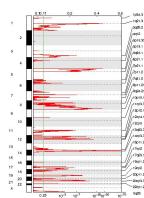
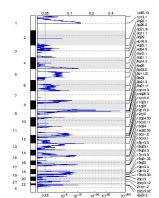
Correlate CopyNumber vs mRNASeq

Correlate molecularSubtype vs CopyNumber Arm

Correlate molecularSubtype vs CopyNumber Focal

Pathway Paradigm mRNA And Copy Number

Pathway Paradigm RNASeq And Copy Number



Mutation Analyses

Aggregate AnalysisFeatures
 Correlate Clinical vs Mutation
 Correlate Clinical vs Mutation APOBEC Categorical
 Correlate Clinical vs Mutation APOBEC Continuous
 Correlate Clinical vs MutationRate
 Correlate molecularSubtype vs Mutation
 Correlate mRNASeq vs Mutation APOBEC

Mutation APOBEC

Mutation Assessor

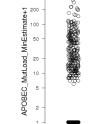
Mutation CHASM

MutSig2.0

MutSig2CV

MutSigCV

Pathway Overlaps MSigDB MutSig2CV



mRNAseq Analyses

Aggregate AnalysisFeatures
Correlate Clinical vs mRNASeq
 Correlate CopyNumber vs mRNASeq
 Correlate mRNASeq vs Mutation APOBEC

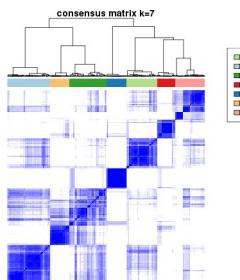
miRseq FindDirectTargets

mRNASeq Clustering CNMF

mRNASeq Clustering Consensus Plus

Pathway Paradigm RNASeq

Pathway Paradigm RNASeq And Copy Number



Correlations Analyses

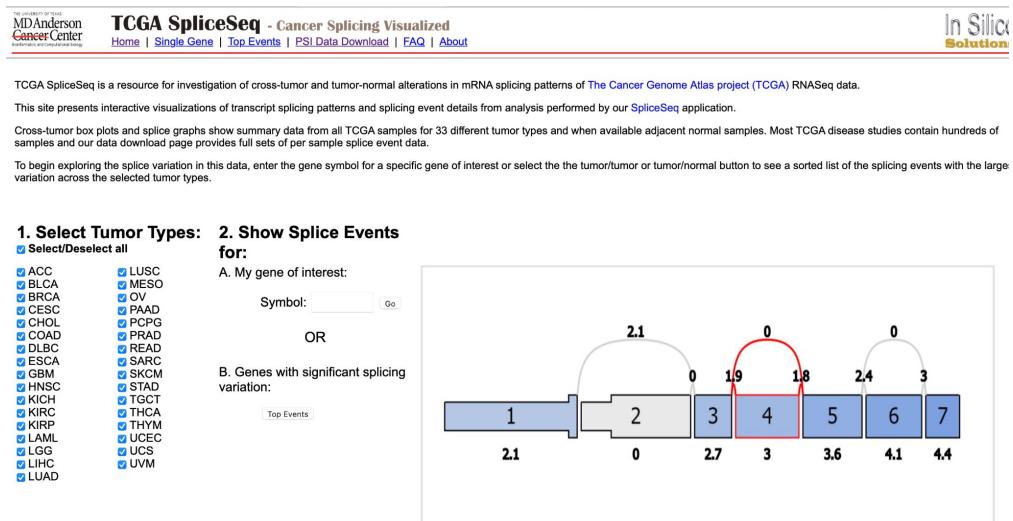
Correlate Clinical vs CopyNumber Arm
 Correlate Clinical vs CopyNumber Focal
 Correlate Clinical vs Methylation
 Correlate Clinical vs miRseq
 Correlate Clinical vs Molecular Subtypes
 Correlate Clinical vs mRNA
 Correlate Clinical vs mRNASeq
 Correlate Clinical vs Mutation
 Correlate Clinical vs Mutation APOBEC Categorical
 Correlate Clinical vs Mutation APOBEC Continuous
 Correlate Clinical vs MutationRate
 Correlate Clinical vs RPPA
 Correlate CopyNumber vs mRNA
 Correlate CopyNumber vs mRNASeq
 Correlate Methylation vs mRNA
 Correlate molecularSubtype vs CopyNumber Arm
 Correlate molecularSubtype vs CopyNumber Focal
 Correlate molecularSubtype vs Mutation
 Correlate mRNASeq vs Mutation APOBEC

And more...

Alternative Splicing

Oncosplicing

Oncosplicing is a database to systematically study clinically relevant alternative splicing in 33 TCGA cancers and 31 GTEx tissues.



TCGA SpliceSeq

TCGA SpliceSeq is a resource for investigation of cross-tumor and tumor-normal alterations in mRNA splicing patterns of **The Cancer Genome Atlas project (TCGA)** RNASeq data.

PCAWG Data Portal for WGS Based Somatic Analysis

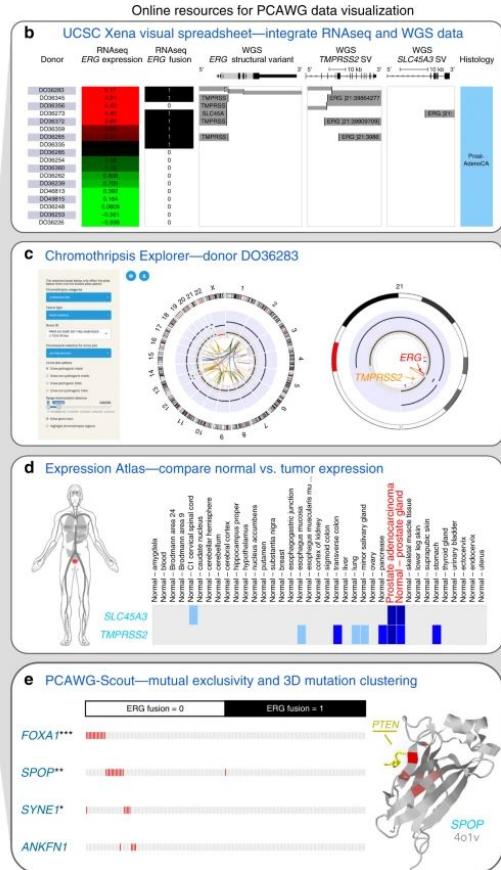
a

Search and download > 70,000 primary BAMs and VCFs

ICGC data repository

PCAWG AWGs primary results (mutation, CNV, expression, etc.)

Search and download primary results files



PCAWG Data Portal

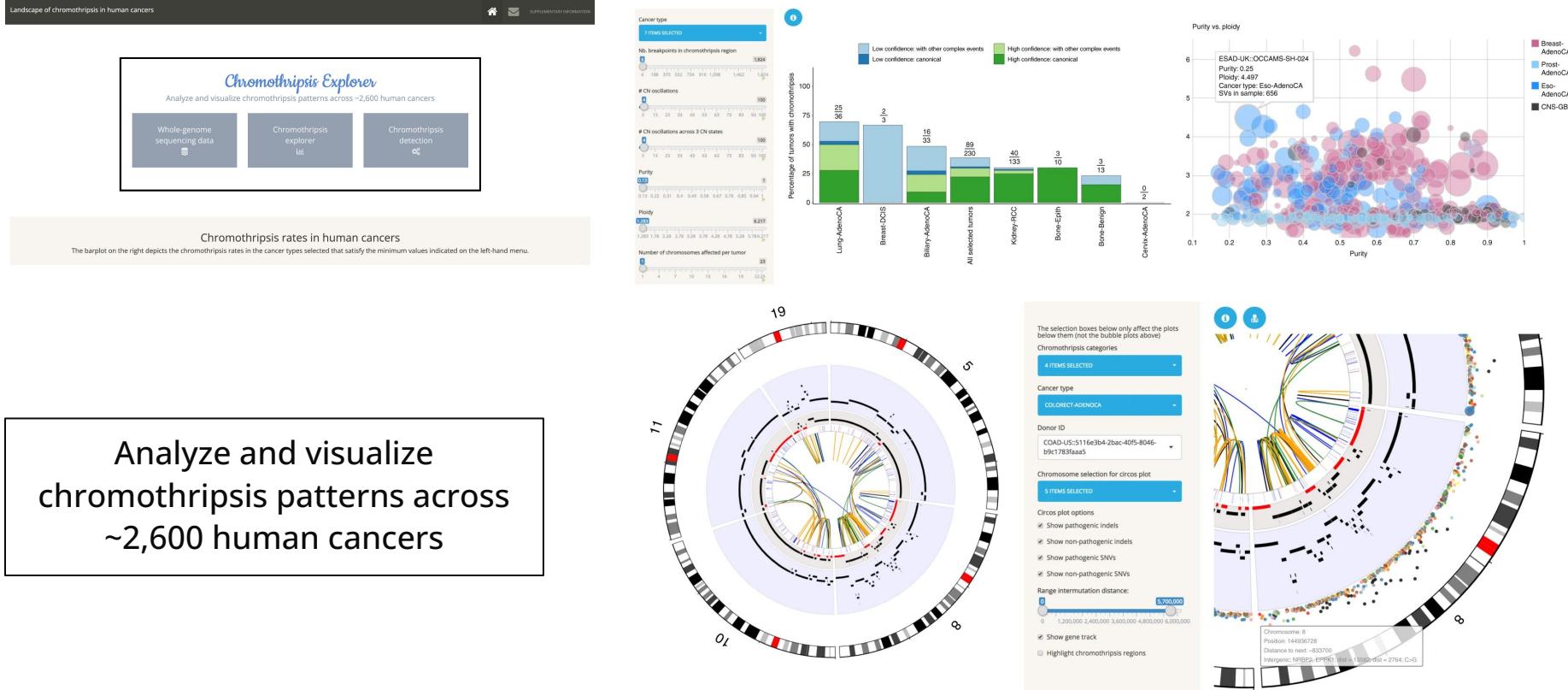
DCC / PCAWG / Filter by file name...

Name

- README.md
- APOBEC_mutagenesis
- benchmarking_data
- cell_lines
- clinical_and_histology
- consensus_cnv
- consensus_snv_indel
- consensus_sv
- data_releases
- donors_and_biospecimens
- driver_mutations
- drivers
- evolution_and_heterogeneity
- germline_variations
- Hartwig
- hla_and_neoantigen
- msi
- mutational_signatures
- networks
- pathogen_analysis

A user guide for the online exploration
and visualization of PCAWG data

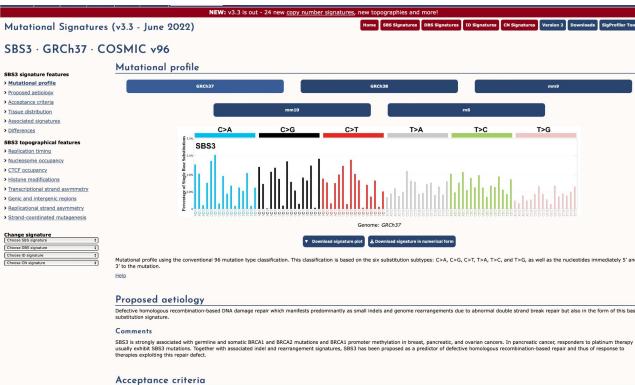
Chromothripsis Explorer



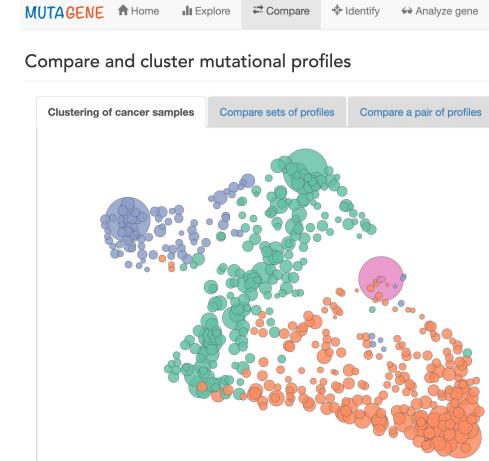
Mutational Signature Analysis

MUTAGENE

COSMIC Mutational Signatures



Signal



- Support different mutational profiles (SBS/ID/DBS/CNV)
- Detail annotations (proposed aetiology and acceptance criteria, tissue distribution, signature associations, version difference)
- Topographies associations (Replication timing, Nucleosome occupancy, CTCF occupancy, Histone modifications, Transcriptional strand asymmetry, Genic and intergenic regions, Replication strand asymmetry, strand-coordinated mutagenesis)

- Explore different source of signatures (cancer, gene edits, environmental mutagenesis)
- Support different mutational profiles (SBS/ID/DBS/SV)
- Support signature data analyses
- Include largest cancer genomic studies, allow to discovery the rare mutational signatures in cancer
- Enhanced signature and study visualizations with API available

- Explore context-dependent mutational profiles and signatures
- Compare and cluster mutational profiles
- Identify mutational processes based on NMF algorithm
- Compares observed mutational frequencies to expected background mutability to identify potential drivers

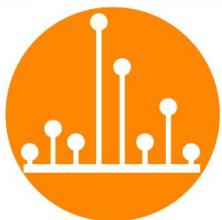
mSigPortal: Integrative mutational signature portal for cancer genomic studies

MSIGPORTAL

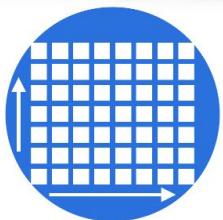
Integrative mutational signature portal for cancer genomic studies



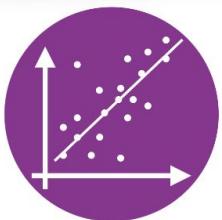
Signature Catalog



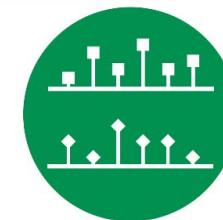
Signature Visualization



Signature Exploration



Signature Association



Signature Extraction

Coming soon!

All existing human and mouse signatures based on different genome builds and algorithm versions

Allows identification of signature features at sample level and discovery of new signatures

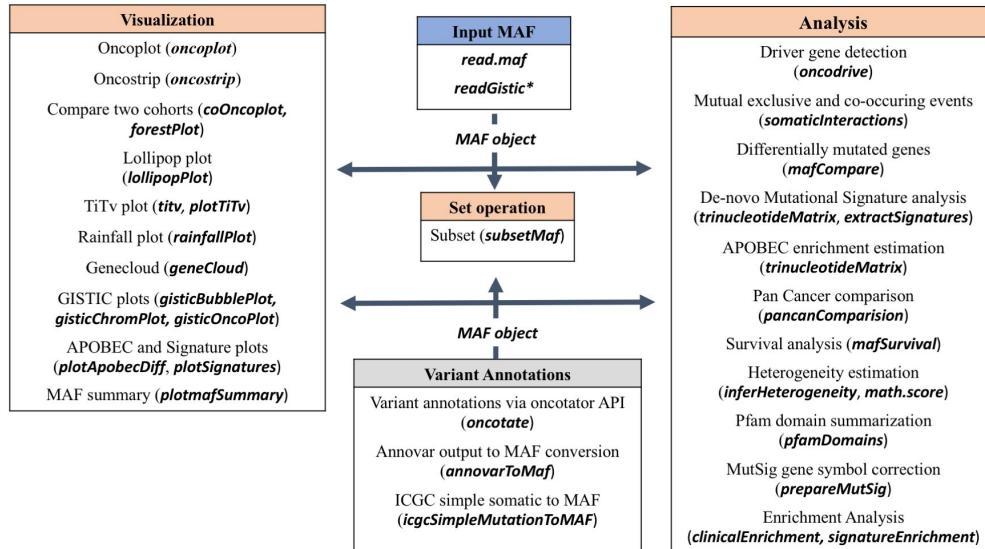
Explore etiological factors associated with signature at sample level

Analyze signature association with other genomic features and clinical data

Extract and compare mutational signatures using state-of-the-art algorithms

Analytical Programming Packages for Cancer Genomic Datasets

MAFtools: Summarize, Analyze and Visualize MAF Files



[**TCGAmutations**](#) - An R data package for TCGA somatic mutations

Data source:
[Broad firehose](#)
[TCGA MC3](#)

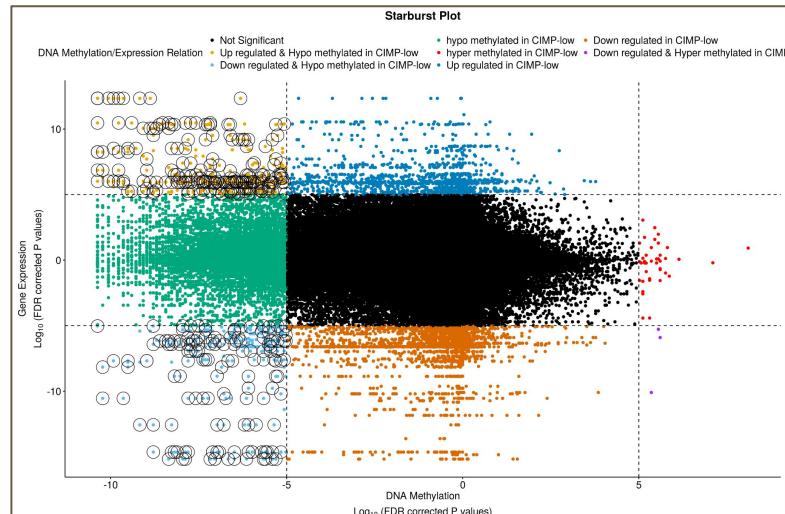
(Multi-Center Mutation Calling in Multiple Cancers)

Most of MAFs from different cancer genomic studies can also be download from NCI GDC

TCGAbiolinks

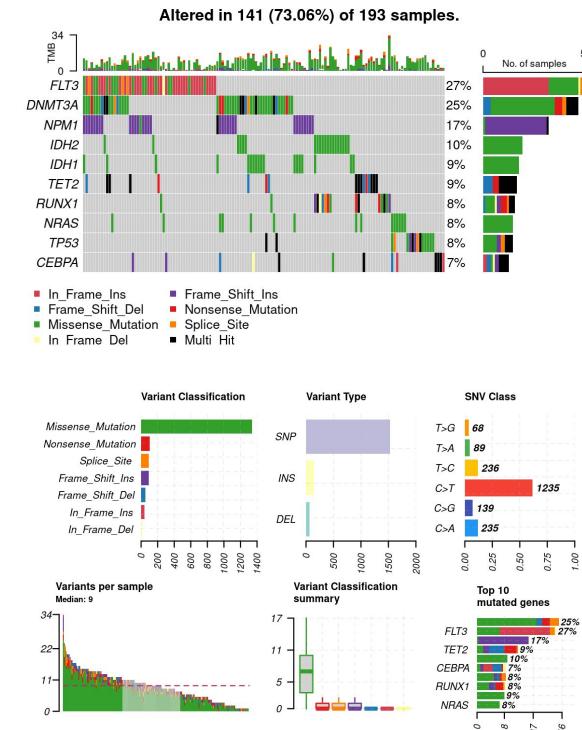
- | | | |
|----------------------|--------------------------|---|
| HTML | R Script | 1. Introduction |
| HTML | R Script | 10. Classifiers |
| HTML | R Script | 10. TCGAbiolinks_Extension |
| HTML | R Script | 11. Stemness score |
| HTML | R Script | 2. Searching GDC database |
| HTML | R Script | 3. Downloading and preparing files for analysis |
| HTML | R Script | 4. Clinical data |
| HTML | R Script | 5. Mutation data |
| HTML | R Script | 6. Compilation of TCGA molecular subtypes |
| HTML | R Script | 7. Analyzing and visualizing TCGA data |
| HTML | R Script | 8. Case Studies |
| HTML | R Script | 9. Graphical User Interface (GUI) |
| PDF | | Reference Manual |
| Text | | NEWS |

- Case study 1: Pan Cancer downstream analysis BRCA
- Case study 2: Pan Cancer downstream analysis LGG
- Case study 3: Integration of methylation and expression for ACC
- Case study 4: ELMER pipeline - KIRC



TCGA Workflow: Analyze cancer genomics and epigenomics data using Bioconductor packages

- Studies: TCGA, ENCODE, Roadmap
- Data access
 - TCGA data
 - NCI Genomic Data Commons (GDC) (level 1 to 3 data)
 - Via TCGAbiolinks
 - GDC Legacy Archive
 - Via TCGAbiolinks
 - Broad Institute's GDAC Fire (level 3 and 4 data)
 - Via RTCGAToolbox
 - ROADMAP via AnnotationHub
 - ENCODE via ENCODExplorer
- Data analysis and visualization
 - maftools :
 - Summarize, Analyze and Visualize MAF Files
 - ELMER
 - Inferring Regulatory Element Landscapes and Transcription Factor Networks Using Cancer Methylomes



cBioPortalData

- cBio Cancer Genomics Portal
 - <https://www.cbioportal.org/>
 - Platform for exploratory and interactive visualization, analysis and download of large-scale cancer genomic data sets.
 - Data sets
 - Public data (TCGA, ICGC, published sequencing studies)
 - Private instances
 - Visualize your own data
 - Open source
- cBioPortalData R package: Obtain data from the cBioPortal API using R
 - Identifying available studies
 - Download studies via cBioDataPack or cBioPortalData

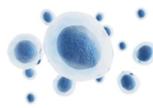
Cloud Resources with Available Cancer Genomic Datasets



- Developed at the Broad Institute
- Uses Google Cloud Platform as compute and storage infrastructure
- Workflows need to be written in WDL (Workflow Description Language)
- Lots of common workflows available in Dockstore (<https://dockstore.org/>) and firecloud (firecloud.org)
- Easily shareable workspaces

Terra supports researchers in many biomedical disciplines

Cancer Genomics



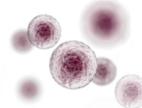
The Van Allen Lab is using Terra to advance clinical oncology through cancer genomics.

Single-Cell Transcriptomics



The Shalek Lab is using Terra to improve the scalability, accessibility, and reproducibility of single-cell analysis.

Medical and Population Genetics



The Natarajan Lab is using Terra to study genetic factors of heart diseases.

[Read More](#)

Infectious Diseases



The Broad's Viral Genomics Group is using Terra to advance genomic epidemiology and surveillance of viral pathogens.

Datasets included in Terra:

- Human Reference Datasets
- [Human Cell Atlas](#)
- FireCloud resources
(GTEX/TARGET/TCGA)

How to access?



Seven Bridges Genomics

- Backend is AWS
- One of the 3 available cloud resources for analysis if the data is in NCI Cancer Research Data Commons (CRDC)

Datasets included:

- Cancer Genomics Cloud
- Gabriella Miller Kids First Data Center
- Trans-Omics for Precision Medicine (TOPMed), Genotype-Tissue Expression (GTEx)
- Model Organism Databases (MODs) datasets
- ICGC and PCAWG
- Blood Profiling Atlas in Cancer (BloodPAC)



Figure 3. Overview of the Seven Bridges Platform. A Cloud-based data and compute infrastructure underlies the discovery layer, which is built around features to streamline data management, search, and analysis. An application programming interface and collaboration features ensure flexibility for users. Data security and regulatory compliance controls operate at all levels.

https://www.sevenbridges.com/wp-content/uploads/2016/11/WP_Scalable_Web.pdf

DNAnexus

- Backend is AWS and Azure
- Enables easy data sharing
- Provides easy to use tools, APIs and visualization
- Data accessible with DNAnexus:
 - St Jude Cloud Genomics Platform
 - UKBiobank
 - ICGC Pan-Cancer dataset
 - Hosted on AWS S3
 - DNAnexus provides ICGC Data Fetcher
 - TCGA
 - Hosted on AWS S3



[Log In to St. Jude Cloud](#)

St. Jude Cloud data is vended through the DNAnexus interface. If you have a DNAnexus employee, please log in with your St. Jude credentials.



“Awesome” bioinformatics resources

“Awesome” bioinformatics resources related to cancer genomic study

Awesome genomics

Cancer Data Science's go to place for excellent genomics tools and packages.

Awesome multi-omics

A community-maintained list of software packages for multi-omics data analysis.

Awesome cancer variant databases

A community-maintained repository of cancer clinical knowledge bases and databases focused on cancer and normal variants

Awesome cancer evolution

Papers for studying cancer evolution

Awesome genome visualization

A list of interesting genome visualizers, genome browsers, or genome-browser-like implementations. [New website](#).

Awesome Clonality

A curated list of awesome resources on clonality and tumor heterogeneity.

Awesome expression browser

A curated list of software and resources for exploring and visualizing (browsing) expression data, but not only limited to that.

Awesome microbes

List of resources, including software packages (and the people developing these methods) for microbiome (16S), metagenomics (WGS, Shot-gun sequencing), and pathogen identification/detection/characterization.

Awesome bioinformatics benchmarks

A curated list of bioinformatics benchmarking papers and resources.

Awesome single cell

List of software packages (and the people developing these methods) for single-cell data analysis, including RNA-seq, ATAC-seq, etc.

Awesome bioinformatics

A curated list of awesome Bioinformatics software, resources, and libraries. Mostly command line based, and free or open-source.

Awesome ChIP-Seq

A curated list of ChIP-Seq analysis

Awesome mutational signature resource in mSigPortal

Original Research Papers Including Specific Mutational Signatures in mSigPortal

Cancer Type	Experimental Strategy	Year	Journal	Title
Human germline	WGS	2021	Science	Population sequencing data reveal a compendium of mutational processes in the human germ line
Glioma	WGS & WES	2021	Nature Genetics	Radiotherapy is associated with a deletion signature that contributes to poor outcomes in patients with cancer
iPSC	WGS	2021	Nature Cancer	A systematic CRISPR screen defines mutational mechanisms underpinning signatures caused by replication errors and endogenous DNA damage
Hematopoietic stem and progenitor cells (HSPCs)	WGS	2021	Cell Stem Cell	Antiviral treatment causes a unique mutational signature in cancers of transplantation recipients
Melanoma	circle-damage-seq	2021	Science Advances	The major mechanism of melanoma mutations is based on deamination of cytosine in pyrimidine dimers as determined by circle damage sequencing
Colorectal Cancer	WES	2021	Cancer Discovery	Discovery and features of an alkylating signature in colorectal cancer
Skin Cancer	WGS	2021		Pre-mutagenic and mutagenic changes imprinted on the genomes of mammalian cells after irradiation with a nail polish dryer
PanCancer	WGS	2021		Mutational impact and signature of ionizing radiation
PanCancer	WGS & WES	2021	Brief Bioinform	Comprehensive analysis reveals distinct mutational signature and its mechanistic insights of alcohol consumption in human cancers
Liver Cancer	WGS & WES	2021	Hepatology	Mutational Signature Analysis Reveals Widespread Contribution of Pyrrolizidine Alkaloid Exposure to Human Liver Cancer

10 Showing 1 to 10 of 38 entries

Review Papers Focused on Mutational Signatures

Year	Journal	Title
2021	Nature Reviews Cancer	Mutational signatures: emerging concepts, caveats and clinical applications
2021	DNA Repair	Significance and limitations of the use of next-generation sequencing technologies for detecting mutational signatures
2020	Nature Genetics	Are carcinogens direct mutagens?
2019	Briefings in Bioinformatics	Computational approaches for discovery of mutational signatures in cancer
2019	Nature Reviews Genetics	Switching APOBEC mutation signatures
2019	Cell	Local Determinants of the Mutational Landscape of the Human Genome
2018	Trends in Cancer	Mutation Signatures Depend on Epigenomic Contexts
2018	Journal of the National Cancer Institute	Biomarkers for Homologous Recombination Deficiency in Cancer
2018	Nature Communications	The therapeutic significance of mutational signatures from DNA repair deficiency in cancer
2018	Briefings in Bioinformatics	Mutational signatures and mutable motifs in cancer

Search

Columns

10 Showing 1 to 10 of 23 entries

Name	Method	Year	Journal	Title
RepairSig	Non-additive model	2021	Cell Systems	RepairSig: Deconvolution of DNA damage and repair contributions to the mutational landscape of cancer
TensorSignatures	Tensor factorisation framework	2021	Nature Communications	Learning mutational signatures and their multidimensional genomic properties with TensorSignatures
SparseSignatures	LASSO regularization	2021	PLOS Computational Biology	De novo mutational signature discovery in tumor genomes using SparseSignatures
Signal	NMF+KLD	2020	Nature Cancer	A practical framework and online tool for mutational signature analyses show inter-tissue variation and driver dependencies
CANCERSIGN	NMF	2020	Scientific Reports	CANCERSIGN: a user-friendly and robust tool for identification and classification of mutational signatures and patterns in cancer genomes
YAPSA	LCD	2020	Genes Chromosomes Cancer	Analysis of mutational signatures with yet another package for signature analysis
MutSignatures	NMF/fcnlms	2020	Scientific Reports	MutSignatures: an R package for extraction and analysis of cancer mutational signatures
iMutSig	Probabilistic Method	2020	F1000Research	iMutSig: a web application to identify the most similar mutational signature using shiny
Sigflow	Bayesian NMF	2020	Bioinformatics	Sigflow: an automated and comprehensive pipeline for cancer genome mutational signature analysis
pyCancerSig	NMF	2020	BMC Bioinformatics	pyCancerSig: subclassifying human cancer with comprehensive single nucleotide, structural and microsatellite mutational signature deconstruction from whole genome sequencing

10 Showing 1 to 10 of 32 entries

<< < > >>

Page: 1

THANKS FOR YOUR ATTENTION!

Questions?

Next: Practical session 2 (10:45am)