The Cancer Genome Atlas

# The Cancer Genome Atlas:
## A Decade of Discovery

*Carolyn M. Hutter, Ph.D.*
*Acting Division Director*
*Division of Genome Sciences, NHGRI*

*National Advisory Council on Human Genome Research*
*September 11, 2017*

# The Cancer Genome Atlas (TCGA)

**THE CANCER GENOME ATLAS**
National Cancer Institute
National Human Genome Research Institute

**Mission:**

A comprehensive and coordinated effort to accelerate the understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing

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

https://cancergenome.nih.gov/

# A Decade of Discovery

## 2006: Pilot Phase of TCGA

*"The Complexity of Cancer Requires a Comprehensive Atlas of Genomic Alterations to Derive Medical Solutions"*
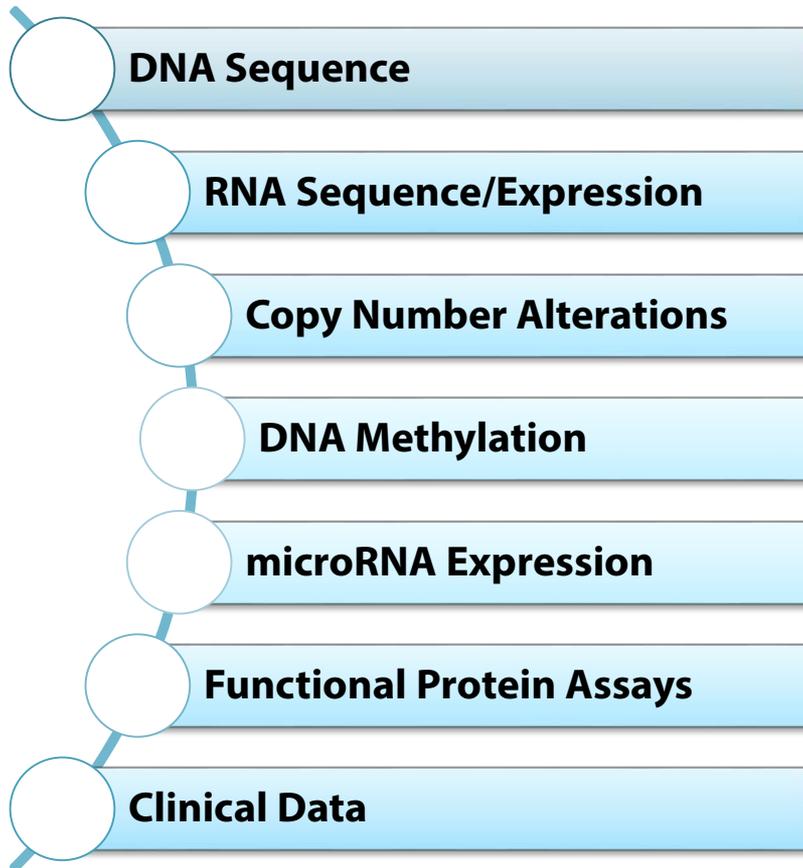*NCI/NHGRI Workshop*
*July, 2005*

National Cancer Institute

and

https://cancergenome.nih.gov/abouttcga/overview/history

# A Decade of Discovery

**Cost per Genome**

$100M

**2006: Pilot Phase of TCGA**

*"The Complexity of Cancer Requires a Comprehensive Atlas of Genomic Alterations to Derive Medical Solutions"*
*NCI/NHGRI Workshop*
*July, 2005*

tion

| 2001 | 2002 | 2003 | 2004 | 2005 | 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 |

# TCGA Data Size and Scope

## Data Types:

- DNA Sequence
- RNA Sequence/Expression
- Copy Number Alterations
- DNA Methylation
- microRNA Expression
- Functional Protein Assays
- Clinical Data

TCGA produced over

## 2.5

### PETABYTES

of data

https://cancergenome.nih.gov/

# TCGA Data as a Community Resource

A key goal of TCGA is to have the data made publicly and broadly available to the research community while protecting patient privacy

**2,775 Approved Data Access Requests for TCGA Data in 2007-2017**



**before 2012**

**2016 - forward**

https://twitter.com/ncigenomics?lang=en

The Cancer Genome Atlas

# Network Analysis: Tumor Specific Projects

## 32 Working Groups
## Defined by Tumor Type

- DNA Sequence
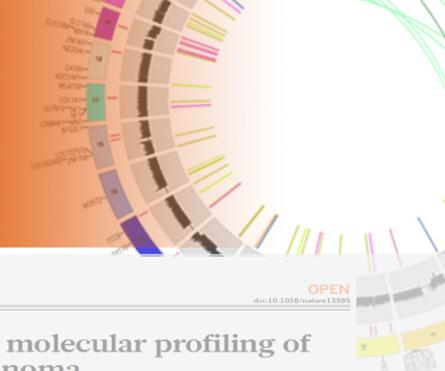- RNA Sequence/Expression
- Copy Number Alterations
- DNA Methylation
- microRNA Expression
- Functional Protein Assays
- Clinical Data

**Tumor Specific Working Groups**

GDACs

GSCs

GCCs

**Comprehensive Characterization of the Cancer Genomes**

**GSC:** Genome Sequencing Center
**GCC:** Genomic Characterization Center
**GDAC:** Genomic Data Analysis Center

**28 published**

**4 submitted**

# Key Results and Findings from TCGA

**MOLECULAR BASIS OF CANCER**

Improved our understanding of the genomic underpinnings of cancer

For example, a TCGA study found the basal-like subtype of breast cancer to be similar to the serous subtype of ovarian cancer on a 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.

**TUMOR SUBTYPES**

Revolutionized how cancer is classified

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

**THERAPEUTIC TARGETS**

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

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.

https://cancergenome.nih.gov/

# Network Analysis: Cross-Cancer Projects

## PanCan-12

Published Fall 2013



## PanCanAtlas

*In progress*

- TCGA Capstone

- Full data analysis

- Standards for data quality and reproducibility
  - Batch correction
  - Re-analysis

- Innovative "big data" distributed computing

# Pancancer Findings

Genomic analysis of multiple cancer types (aka "Pan-Cancer analysis) highlight commonalities and differences among various types and subtypes.



### 12 TCGA Pan-Cancer types

BLCA | BRCA | COAD | GBM | HSNC | KIRC | AML | LUAD | LUSC | OV | READ | UCEC

**Sequencing and Analysis**

**mutation rate**

Distribution of mutation rates across the twelve cancer types reveals interesting features, such as clusters in UCEC and COAD/READ that indicate factors other than age in the development of these tumors.

**mutation spectrum**

Environmental effects on cancer development can also be observed in mutation spectrum. For instance, lung tumors show higher proportions of C-to-A transversions – a signature of cigarette smoke exposure.

**mutated genes**

When grouped by mutation, we see that significant mutations fall into several distinct categories: transcription factors & regulators, histone modifiers, genome integrity, RTK signaling, cell cycle, and more.

**mutation relation**

We found 14 significant mutually exclusive pairs and 148 co-occurring pairs, We also identified a set consisting of TP53, PIK3CA, PIK3R1, SETD2, and WT1.

**clinical features**

We found TP53 to be significant, with mutations being associated with detrimental outcome through joint analysis of 12 tumor types. Mutations in BAP1 are correlated with detrimental outcome particularly in KIRC and UCEC.

**clonal architecture**

Mutations in TP53, DNMT3A, and PIK3CA play an initiation role in the tumorigenesis. Mutations in KRAS and/or NRAS largely play a progression role in the tumorigenesis of AML, BRCA, and UCEC.

Systematic analysis of the TCGA Pan-Cancer mutation dataset identifies SMGs, cancer-related cellular processes, and genes associated with clinical features and tumour progression.

# Pancancer Findings

Genomic analysis of multiple cancer types (aka "Pan-Cancer analysis") highlight commonalities and differences among various types and subtypes.

Integrative analysis across cancer types identifies subgroups that are correlated, but not identical to, tissue-of-origin classifications of cancer.

http://news.ucsc.edu/2014/08/pan-cancer-study.html

# Pancancer: Large-scale genomics meets large-scale, reproducible analysis

## PanCan-12

Published Fall 2013
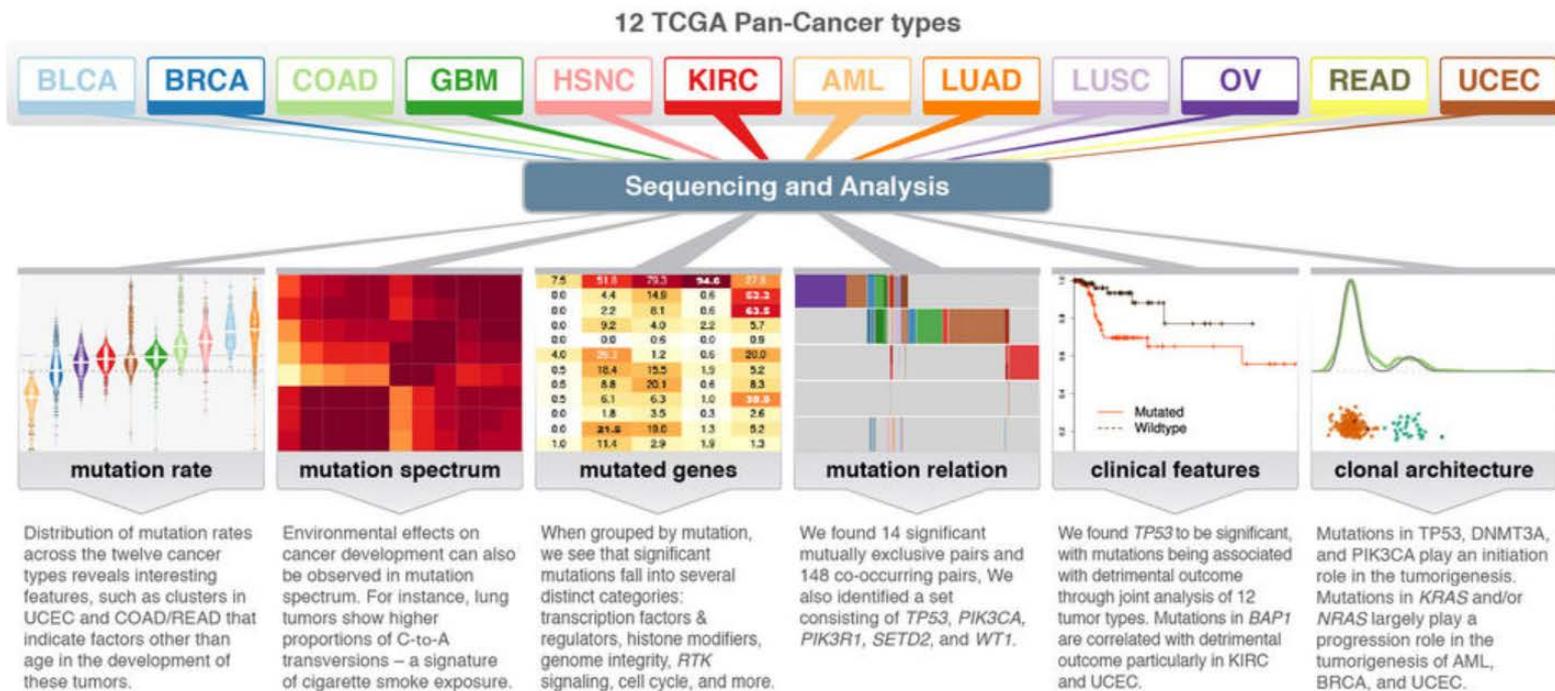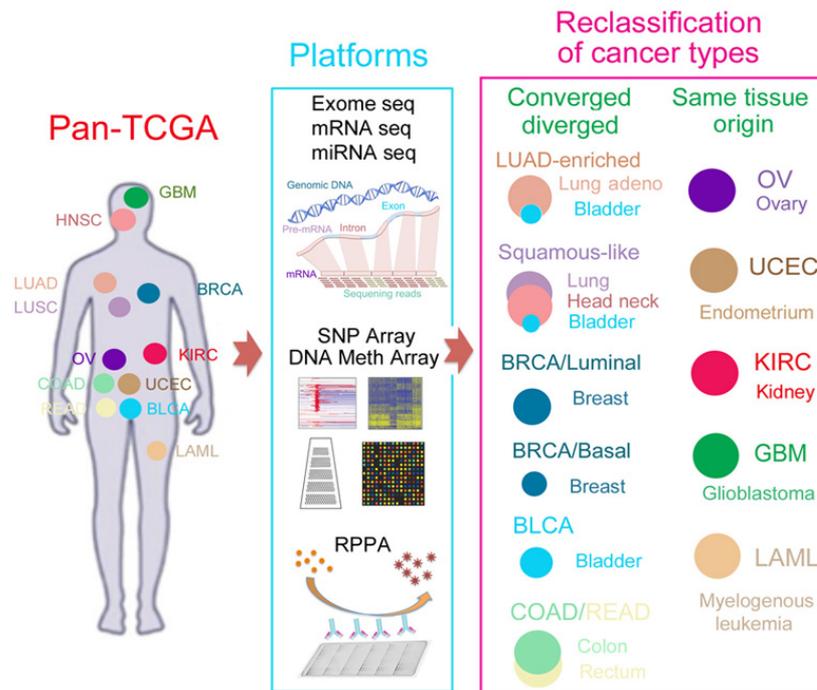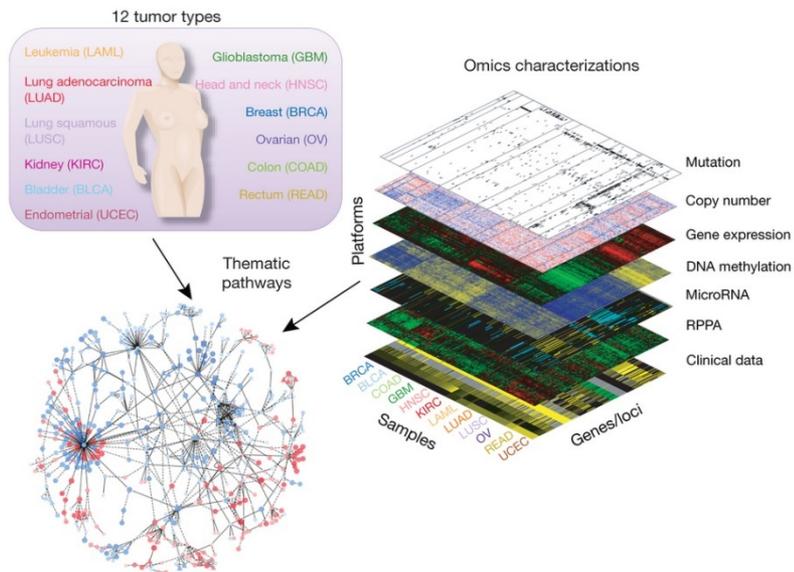


## PanCanAtlas

*In progress*

- TCGA Capstone

- Full data analysis

- Standards for data quality and reproducibility

  – Batch correction

  – Re-analysis

- Innovative "big data" distributed computing

# TCGA MC3: Foundation for PanCanAtlas

**TCGA MC3 Project**



10K TCGA Exomes

GATK Preprocessing — BROAD INSTITUTE / UNIVERSITY OF CALIFORNIA SANTA CRUZ Genomics Institute

Coherent BAM collection — Cancer Genomics Hub

Consistent Mutation Calling across all TCGA exomes — BROAD INSTITUTE / BC Cancer Agency / UC SANTA CRUZ Genomics Institute / Baylor College of Medicine / Washington University in St.Louis / MD Anderson Cancer Network

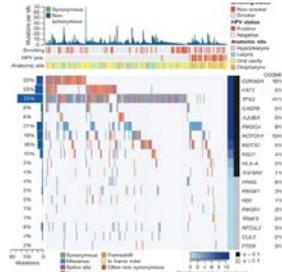Filtering and Annotations — COSMIC / dbGaP / Ensembl
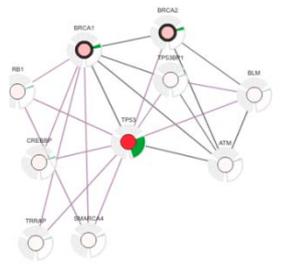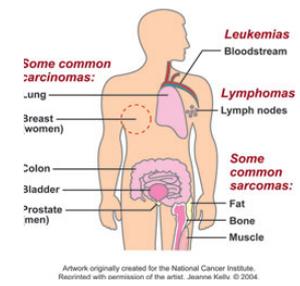
**Final MAF**

## Goals
- Produce uniform mutation calls across all TCGA exome data
- 500 terabytes BAM exomes
- Multiple variant calling tools and QC filters
- Reproducible tools and workflows
- Modern computing paradigm
  - Docker Containers
  - CWL/WDL workflows
  - Leverage clouds and clusters
  - Global Alliance (GA4GH) Standards

Ellrott et al "Automating Somatic Mutation calling for Ten Thousand Tumor Exomes" (Submitted)

15

The Cancer Genome Atlas

Oncogenic Processes

Cell-of-origin

Pathways

https://www.cancer.gov/about-nci/organization/ccg/blog/2017/tcga-pancan-atlas

The Cancer Genome Atlas

# PanCancer Analysis of Whole Genomes

The Cancer Genome Atlas

# Use of distributed computing with large-scale genomics data



**PCAWG** — PanCancer Analysis of WHOLE GENOMES

ICGC Project
**2583** WGS donors
**50M** somatic variants

Collaboratory, OICR
AWS Ireland
Sanger, Hinxton
DKFZ, Heidelberg
EBI, London
NCI Cluster, UCSC
Microsoft Azure
AWS Virginia & Seven Bridges
iDASH, UCSD
BSC, Barcelona
ETRI, Seoul
UTokyo
PDC, UChicago

HPC
Academic clouds
Commercial clouds

**23** months actual
(**6 months theoretical**)

Compute centers (C), GNOS repositories (G), and S3-compatible data storage (S)

18

PCAWG-Tech: Christina Yung & Junjun Zhang; slide courtesy of Lincoln Stein

The Cancer Genome Atlas

# Impact of TCGA

- Comprehensive Data Resource

- Forward-looking Data Sharing

- Transformational Scientific Advances

- Innovative Pipelines and Approaches

## The TCGA Legacy: Multi-Omic Studies in Cancer

**27-29 September, 2018**

**Washington, DC**

June 15, 2018
Abstract Submission

August 10, 2018
Early Registration

The Cancer Genome Atlas

# Acknowledgements

**Current TCGA Project Team**
*NHGRI*
Heidi Sofia
Melpi Kasapi

*NCI*
Jean-Claude Zenklusen
Samantha Caeser-Johnson
John Demchok
Ina Felau
Martin Ferguson
Roy Tarnuzzer
Zhining Wang
Liming Yang

*TCGA DAC*
Vivian Ota Wang
Jeff Kim

**Past TCGA Project Team Members**

**TCGA Research Participants**

**TCGA Network**



British Columbia Cancer Agency
*Vancouver, Canada*

Broad Institute
*Cambridge, Massachusetts*

Institute for Systems Biology
*Seattle, Washington*

Brigham & Women's Hospital and Harvard Medical School
*Boston, Massachusetts*

Oregon Health and Science University
*Portland, Oregon*

Washington University School of Medicine
*St. Louis, Missouri*

The Research Institute at Nationwide Children's Hospital
*Columbus, Ohio*

Memorial Sloan-Kettering Cancer Center
*New York, New York*

Buck Institute for Research on Aging
*Novato, California*

Johns Hopkins University
*Baltimore, Maryland*

University of California, Santa Cruz
*Santa Cruz, California*

University of Texas M.D. Anderson Cancer Center
*Houston, Texas*

SRA International
*Arlington, Virginia*

University of Southern California
*Los Angeles, California*

Baylor College of Medicine
*Houston, Texas*

University of North Carolina at Chapel Hill
*Chapel Hill, North Carolina*

Biospecimen Core Resource (BCR)
Genome Characterization Center (GCC)
Genome Sequencing Center (GSC)
Genome Data Analysis Center (GDAC)
Cancer Genomics Hub (CGHub)
Data Coordinating Center (DCC)