# TCGA Data Overview

## Xiaofei Zhou, Han Zhang and Meng Wang

The Ohio State University

## February 17, 2016

# Background

- Cancer is one of the biggest causes of deaths in the world.
- In 2012, about 8.2 million people globally died of cancer, accounting for 14.6 % all human deaths.
- In the US, one person dies from cancer every minute, that is, 1,500 deaths each day.
- Our goal is to control and conquer cancer.

# The Cancer Genome Atlas (TCGA)

- Aims at increasing the understanding of the molecular basis of cancer through genome analysis technologies.
- Is supported by NCI and NHGRI.
- We have already learned:
  - certain regions of the genome are linked with several types of cancers.
  - these regions are usually contain genes involved in the pathways of cell apoptosis.
  - signatures – specific changes in human genome – allow us to tell one type of cancer from another. These signatures help doctors with diagnosis, treatments and/or prognosis of cancer.

# TCGA Cancer Selection Criteria

TCGA selected cancers based on:

- overall public health impact
- poor prognosis
- availability of human tumor and matched-normal tissue samples that meet TCGA standards for
  - patient consent
  - quality
  - quantity

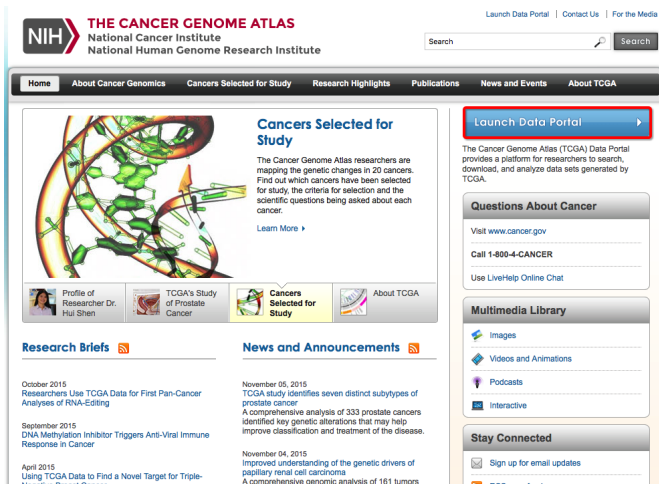There are 34 types of cancer currently available under TCGA program, click <u>this link</u> to see what they are.

# Accessing TCGA Data

TCGA data are available in two data repositories:

- The TCGA Data Portal:
  - provides access to almost all TCGA datasets, for example, SNP, copy number, methylation, expression, sequencing and so on.
- The Cancer Genomics Hub (CGHub):
  - links to TCGA primary sequence data.

# Accessing TCGA Data

To view more general information about the TCGA program, please visit `http://cancergenome.nih.gov/`

# Accessing TCGA Data

The tutorials under the <u>Download Data Section</u> in TCGA Data Portal can help you get familiar with methods of retrieving data.

# TCGA Wiki

If you encounter unfamiliar terms when you access the data, go to
`https://wiki.nci.nih.gov/display/TCGA/TCGA+Home`

# Accessing TCGA Data
## Quick View on Available Cancer Data

First click <u>here</u> to see available cancer types. Then click on the cancer you are interested in.

# Accessing TCGA Data
## TCGA Data Portal

We have 4 major methods to retrieve data through the TCGA Data Portal:

| Method | Functionalities | Use when: | Limitations |
|---|---|---|---|
| Data Matrix | • Allows you to select and download a subset of data for a particular cancer type<br>• Data are in tab-delimited format | You wish to download a subset of data (for example, just SNP data) for a particular cancer type. | Cannot be used to search for and download data across cancer types simultaneously |
| Bulk Download | Allows you to download full archives of data as uploaded by TCGA centers | You wish to download full archives of data. | Only one archive can be downloaded at a time. |
| HTTP Directories | Allows you to directly access the HTTP file system where the archives of data are stored. | You know the exact file or files you are searching for and the location of the files in the directory. | Only one archive can be downloaded at a time. |
| File Search | Allows users to filter and download data files using criteria such as Disease, Data Category, Data level and Access Tier. | Unlike the Data Matrix, the TCGA File Search allows cross-disease searching. | The File Search only provides the latest revision of each archive; older revisions are available through bulk download or HTTP access. |

# Accessing TCGA Data
TCGA Data Portal – Retrieving Data Using Data Matrix

- Data Matrix allows user to search for a subset of data of a particular cancer type.
- Limitation:
  - cannot search across multiple cancer types simultaneously.
  - only returns the most recent data files.
- <u>TCGA Data Portal Data Matrix Access link</u>
- <u>TCGA Data Portal Data Matrix Tutorial link</u>

# Accessing TCGA Data

## TCGA Data Portal – Retrieving Data Using the Data Matrix



**Filter Settings**

**Select a disease:** GBM – Glioblastoma multiforme

**Data Type:**
- All
- CNV (CN Array)
- CNV (SNP Array)
- Clinical

**Batch Number:**
- All
- Batch 1
- Batch 2
- Batch 3

**Data Level:**
- Level 1
- Level 2
- Level 3

**Availability:**
- Available
- Pending
- Not Available

**Preservation:** Frozen  ⓘ Help

**Center/Platform:**
- All
- BCGSC (IlluminaHiSeq_miRNASeq)
- BCM (ABI)
- BI (ABI)

**Sample:**
ID Matches:
TCGA- [ -- ] [    ] [ -- ]  Remove

Add Row

**Paste Sample List:**

**Upload Sample List:**
Choose File  no file selected

**Access Tier:**
- All
- Protected
- Public

**Tumor/Normal:**
- Tumor - matched
- Tumor - unmatched
- Normal - matched
- Organ-Specific Control
- Cell Line Control

**Submitted Since (Date):** mm/dd/yyyy

**Submitted Up To (Date):** mm/dd/yyyy

☐ Only show samples with data available for all columns

Get web service URL for this filter    ✔ Apply   Clear

# Accessing TCGA Data
## TCGA Data Portal – Terminology

- **Data Type**: a label to categorize the many forms of platform data within TCGA Network.
- **Batch**: a set of related analytes from the same disease, or serial index.
- **Data Level:** 1 (Raw Data), 2 (Processed Data), 3 (Segmented or Interpreted Data) and 4 (Region of Interest Data).
- **Preservation**: the method used to preserve the sample after it has been removed from a participant.
- **Center**: a single or a collection of institutions and research centers that perform the same function within TCGA.
- **Platform**: a vendor-specific technology for assaying or sequencing.

# Accessing TCGA Data

TCGA Data Portal – Terminology

- **Sample**: use Biospecimen IDs to find data (See later slide).
- **Access Tier**: whether the data are open or control-access.
- **Tumor/Normal**
    - **Tumor - matched**: data for a tumor tissue for which matched normal tissue exists.
    - **Tumor - unmatched**: data for a tumor tissue for which there is no matched normal tissue.
    - **Normal - matched**: data for normal tissue for which matched tumor tissue exists.
    - **Organ-specific Control**: data for normal tissue from a participant who does not have cancer.
    - **Cell Line Control**: data for cell-line controls.

    Find the detailed definitions by using <u>TCGA Wiki</u> or <u>Data Selection Tutorial</u>.

# Accessing TCGA Data
## TCGA Data Portal – Data Types

To see all available data types under TCGA Data Portal, click <u>here</u>.

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ▸ mRNA Sequencing | | | | | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| ▸ Total RNA Sequencing | | | | | | | |

| Array-based Expression | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Data Subtype** | **Cancer Types Applicable** | **Data Type Name** | **Level 1** | **Level 2** | **Level 3** | **Important Metadata** | **How to Retrieve Data Files** |
| Gene | BRCA, COAD, GBM, KIRC, KIRP, LAML, LGG, LUAD, LUSC, OV, READ, UCEC | Expression - Gene | Raw signals per probe for each participant's tumor sample | Normalized signals per probe or probe set for each participant's tumor sample | Expression calls for genes, per sample | Experimental protocol, including calculation methods, is included in the **MAGE-TAB** archive | Data Matrix & Bulk Download: Select 'Expression-Gene" For Data Type |
| | | | File type: tab-delimited (.txt) | File type: tab-delimited (.txt) | File type: tab-delimited (.txt) | Probe information is contained in the | File Search: Select 'Other' for Data Category |

Note that non-genomic info are available under clinical data. (see later slides)

Also note that primary DNA and RNA sequence data are only available through CGHub, not TCGA Data Portal.

# Accessing TCGA Data

The TCGA Data Portal - more useful links

Data Download Page

- Data Matrix Tutorial
- Bulk Download Tutorial
- HTTP Download Tutorial
- File Search Download Tutorial

# Accessing TCGA Data
Understanding Biospecimen IDs

Biospecimen IDs exist in 2 forms:

- Universally Unique Identifier (UUID)
- Barcode

**UUID:** is NOT human readable.
*e.g.* ebf3e73f-41a0-4ca5-b608-fe1c629e16de

# Accessing TCGA Data

Understanding Biospecimen IDs

**Barcode:** consists of many identifiers and is human readable.
*e.g.*



- **Project:** Project name, such as TCGA
- **TSS:** A Tissue Source Site collects samples and clinical metadata
- **Participant:** someone who contributes one or more samples to a TSS
- **Sample:** Tumor types range from 01 to 09; normal types range from 10 to 19

# Accessing TCGA Data

Understanding Biospecimen IDs

To see how to retrieve a sample or a set of samples from Biospecimen IDs, click here.

To see how to get one type of Biospecimen ID from another type of Biospecimen ID, click here.

To see how to find the corresponding participants or aliquots from Biospecimen IDs, click here.

# Access the TCGA Data
Clinical Data

TCGA clinical data include:

- Clinical information about the participant
- Information about the how participant samples (biospecimens) were processed

Note that clinical data contains a lot of non-genomic information such as the participant's age, gender, and race.

Clinical data are contained in four types of files:

- **Biotab** Files: contain clinical and biospecimen information for a set of participants.
- **XML** Files: Each TCGA participant has a separate clinical and biospecimen XML file.
- **svs** Files: Tissue slide images.
- **pdf** Files: Original pathology reports.
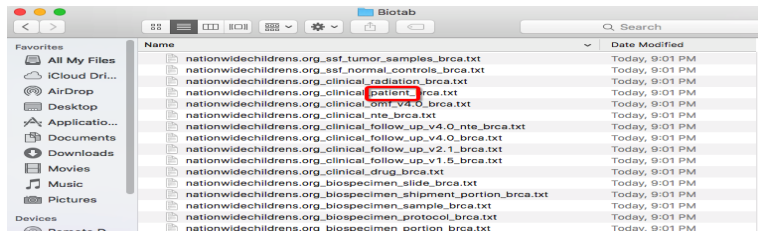
# Access the TCGA Data
Clinical Data

Note that clinical data are retrieved differently from other data types.

Data Download Link

- Clinical Data:  Data Matrix Tutorial
- Clinical Data:  Bulk Download Tutorial
- Clinical Data:  File Search Tutorial
- Image and Pathology Reports Retrieval Tutorial

# Accessing TCGA Data

Clinical Data - Obtaining non-genomic information

# Accessing TCGA Data
Understanding the Downloaded Data Files

- **Metadata:** the summary data about the downloaded data. Often describe biospecimen-related elements.
- **idf** file: provides general information about the investigation and experiment.
- **sdrf** file: encapsulates a succession of processes applied to samples and the multiple states it takes on as a result of the processes.

# Thank you!