Name: _Kailin Liu_

Date: _09/11/2022_

# TCGA Website Scavenger Hunt

## TCGA (Home Page):

The Cancer Genome Atlas (TCGA), founded in December of 2005, is a cancer genomics program hosted by the _National Cancer Institute_ and the National Human Genome Research Institute. The publicly available data from this project includes ___genomic___, epigenomic, _transcriptomic_, and proteomic data. This data was collected from 20,000 different samples that span 33 different cancer types, including breast cancer, which we will be focusing on this semester.

## Program History:

Describe one outcome or impact of TCGA: ___deepened understanding of cancer through molecular characterizations (ex. aberrations in DNA sequence → functional consequences)___

Briefly skim the "Timeline & Milestones" page. When did TCGA publish their paper on breast cancer? _October 2012_

Because TCGA is a public dataset, and one of the first of its kind, they faced some initial concerns regarding the ethics of releasing health data to the public. Choose one of the papers in the "Ethics & Policies" section to skim. What is one way that your paper addresses these privacy concerns? ___The "Suggested Informed Consent Language" document outlines how to ask, educate, and clarify patients about the collection of their data, how it will be used, and additional contacts.___

## TCGA Cancers Selected for Study:

List three criteria used to select which cancers to study: _Public health impact, quality standards, poor prognosis_

Open the breast ductal carcinoma page and read TCGA's provided background. List one interesting fact you found: _there are 4 subtypes: HER2-enriched, Luminal A, Luminal B, & Basal-like_ _____

## Publications by TCGA:

TCGA published (at least) one paper on each of their studied cancer types. These papers, called marker papers, include an early analysis of the data, including any molecular characterizations that were performed. Read the abstract of the 2012 breast ductal carcinoma cancer paper. List any genes you come across (these may be good starting points for your future analyses of this cancer):
___TP53, PIK3CA, GATA3, MAP3K1, HER2, EGFR___ _____

## Using TCGA:

Go to the Genomic Data Commons (GDC) Data Portal via the link on TCGA home. This portal lets you view TCGA's data in a visual way. Let's explore this website. According to the Data Portal Summary, there are _72_ projects in the GDC data portal. Now click on the "Projects" tab. Notice that not all projects in this data portal are TCGA-affiliated, though TCGA does make up _33_ of the projects included.

## Using TCGA (Continued)

Under the "Program" tab, select just TCGA studies. According to the graph at the top of the page, _TP53_ is the most mutated gene in TCGA projects, affecting approximately _35_ % of cases.

Return to the GDC Portal home page. Now click the breast image in the diagram to the right of the page. This directs you to the "Exploration" tab and automatically selects all primary sites associated with breast cancers. Now select TCGA as the program, and TCGA-BRCA as the as the project. This is the data we will be focusing on this semester.

The table on this page shows each patient along with their data. Feel free to explore the data files by clicking on any of the links provided.

Now explore the Cases, Genes, Mutations, and OncoGrid tabs above the pie charts. What is one takeaway from the plots provided here: __There are a lot of tags to differentiate between each patient / data point (ex. type of mutation, impact, gene, etc.)__

As you can see, the GDC portal provides an overwhelming amount of information. Feel free to continue to explore it on your own time!

## Discussion:

Think through the following questions, and record your answers below:
1. What is the goal of TCGA?

____The main goals of TCGA include providing more readily available data for computational analysis, providing tools, creating a database of a wide range of cancers, bringing together researchers, understanding cancers, advancing therapy, and improving clinical treatments.____

2. What are some ways that we use TCGA's data for our own cancer research? (Think about the types of data available and brainstorm some research questions that can be proposed given that data.)

____We can use TCGA's data for building statistical analyses thanks to the large amount of data points. Although we are focused on breast cancer, TCGA holds data from many subtypes of the cancer including rare ones. Questions: How does expression of a missense mutation affect the transcription & translation of other genes/proteins in its signalling pathway?____

3. What are the benefits and drawbacks of TCGA or other large publicly available datasets?

____Benefits: wide array, pre-existing categorizations/filters, built in tools and visualizers
Drawbacks: missing info for certain details, takes out the humanistic correlation between data and patient____