# Gene Expression Data
## STOR 565

Andrew Nobel

January, 2021

# The Central Dogma



Process: Transcription
Purpose: RNA synthesis
Location: Nucleus

Process: Translation
Purpose: Protein synthesis
Location: Cytoplasm at a Ribosome

**DNA**          **RNA**          **Protein**

DNA contains the original codes for making the proteins that living cells need. mRNA is a copy of a gene located on the DNA molecule. mRNA will leave the nucleus of the cell and the ribosome will read its coding sequences and put the appropriate amino acids together.

# The Cancer Genome Atlas

Multi-Institution consortium supported by the National Institutes of Health.

**Goal:** "[T]o accelerate our understanding of the molecular basis of cancer through the application of genome analysis technologies, including large-scale genome sequencing"

Consortium collected tissue from thousands of tumors across numerous cancer types. Data derived from high-throughput technologies measuring

- ► Gene expression
- ► Micro-RNA
- ► DNA copy number
- ► Methylation

Home

## TCGA Data Portal Overview

The Cancer Genome Atlas (TCGA) Data Portal provides a platform for researchers to search, download, and analyze data sets generated by TCGA. It contains clinical information, genomic characterization data, and high level sequence analysis of the tumor genomes.

Please note some data on the TCGA Data Portal are in controlled-access. Please visit the Access Tiers page for more information.

The TCGA Data Portal does not host lower levels of sequence data. NCI's Cancer Genomics Hub (CGHub) is the new secure repository for storing, cataloging, and accessing BAM files and metadata for sequencing data.
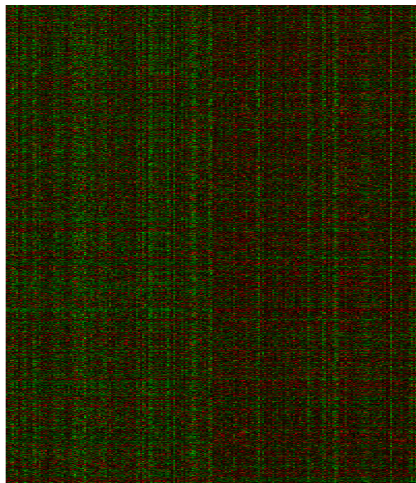
Download Data ›

Choose from four ways to download data

| Available Cancer Types | # Cases Shipped by BCR* | # Cases with Data* | Date Last Updated (mm/dd/yy) |
|---|---|---|---|
| Acute Myeloid Leukemia [LAML] | 200 | 200 | 04/29/15 |
| Adrenocortical carcinoma [ACC] | 80 | 80 | 08/27/15 |
| Bladder Urothelial Carcinoma [BLCA] | 412 | 412 | 08/27/15 |
| Brain Lower Grade Glioma [LGG] | 516 | 516 | 08/27/15 |
| Breast invasive carcinoma [BRCA] | 1100 | 1098 | 08/28/15 |
| Cervical squamous cell carcinoma and endocervical adenocarcinoma [CESC] | 308 | 308 | 08/21/15 |
| Cholangiocarcinoma [CHOL] | 36 | 36 | 08/21/15 |
| Colon adenocarcinoma [COAD] | 461 | 461 | 08/27/15 |

# Screenshot of Expression Data

**Heat map:** Means of displaying a data matrix with numerical entries

- ▶ positive entries are red

- ▶ negative entries are green

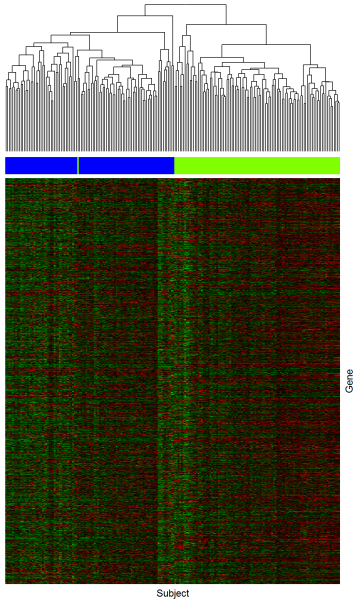- ▶ entries close to zero are black

# Example: Gene Expression Data from Breast Cancer



Gene

Subject

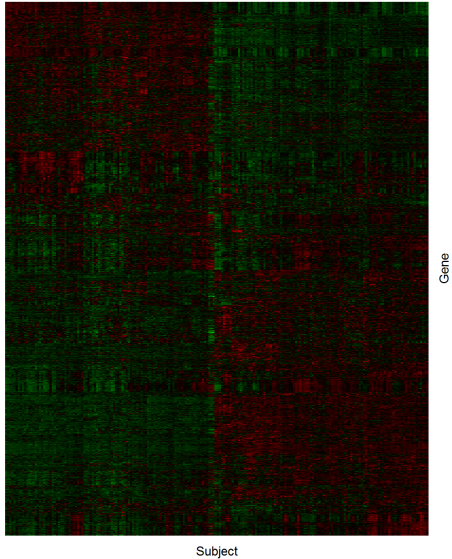Heat map of gene expression data from
The Cancer Genome Atlas (TCGA)

- Samples: $n = 200$ tumors from two
  breast cancer subtypes
    - 100 Luminal A tumors
    - 100 Basal tumors

- Variables: $p = 11,000$ genes (post-filtering)

# Clustering Samples of TCGA Data



Colors: Luminal A and Basal
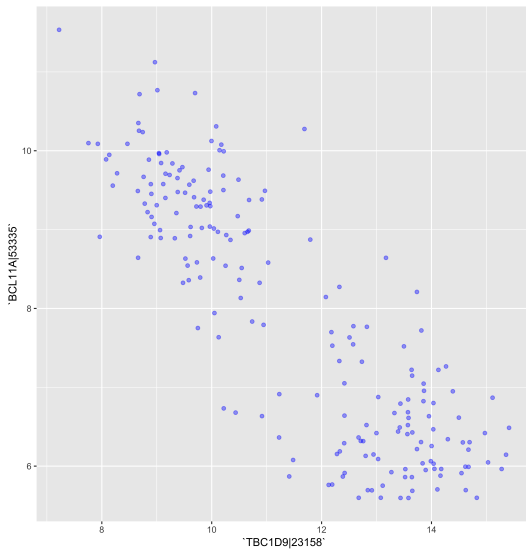
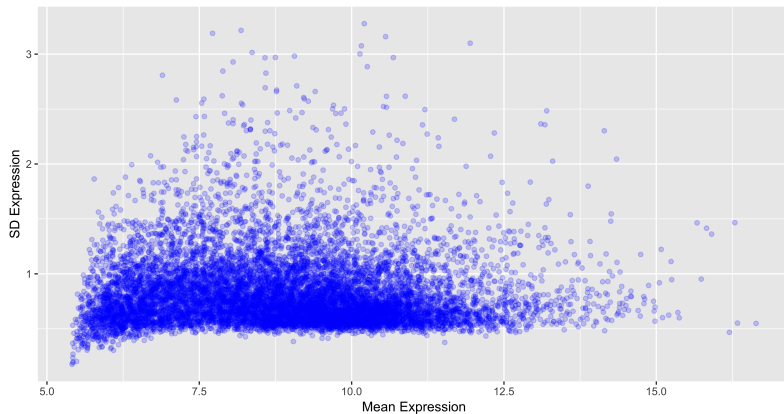# Clustering Rows and Columns



Gene

Subject

# Two Uncorrelated Genes

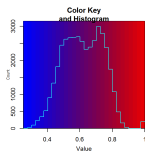# Two Positively Correlated Genes

# Two Negatively Correlated Genes

# Scatterplot of SD (expression) for Two Subtypes

# Heatmap: Correlation Matrix of Samples ($n \times n$)

# Heatmap: Correlation Matrix of Genes ($p \times p$)