Gene Expression (humans) hands-on (Prepared for BIFX550 in 2019; may be some menu options in NCBI might be outdated; if you happen to find any issues, please let me know; ravichandran@hood.edu )

Gene Expression (GE) varies from cell to cell.

Every cell has a complete set of genes

In a cell not all genes are expressed at a time.

GE varies during developmental stages, disease (or GE can cause disease), environmental state and final cells location

Gene expression (gene two copies; only one copy if active (methyl transferase is used for silencing the other)

Experiments to identify the Gene expression:
Older methods (not part of NCBI DB but commonly seen in many publications as images)

a.  Northern_blot Detect RNA in a sample   (http://en.wikipedia.org/wiki/Northern_blot )
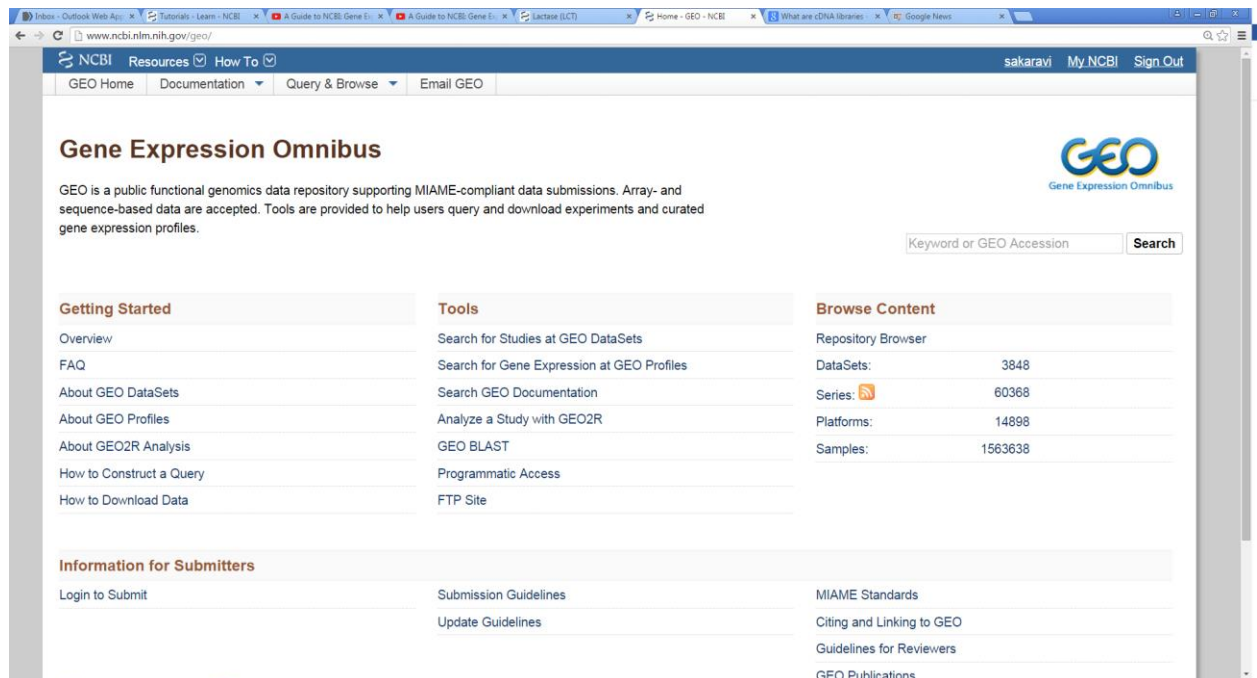b.  Western blot: Detect Protein in a sample ( http://en.wikipedia.org/wiki/Western_blot )

Expressed Sequences :

1.  Primary
    o  GenBank and EST (bulk)
        ▪  200-500 bases;
        ▪  mRNA (tissues/cells) → cDNA
        ▪  Why this is done? To identify genes; To calculate the abundance of proteins
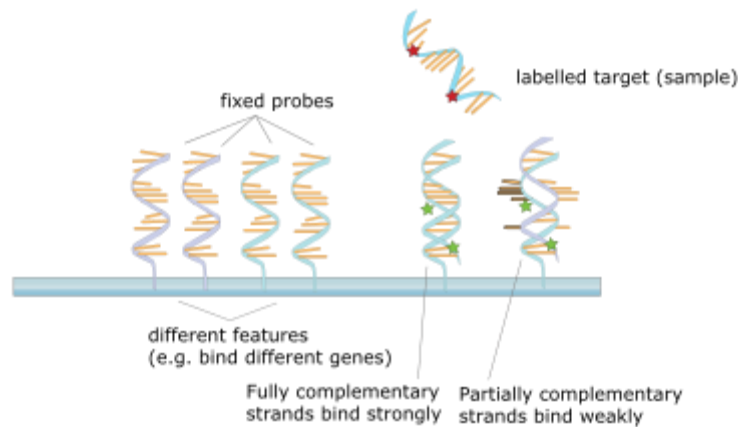    o  Newer: Sequence Read Archive (SRA)


    High-Throughput Studies:
    Are routinely used to estimate gene Expression, Genome Variation, Epigenomics (DNA methylation, Histone modification etc). Data from these studies are submitted into Gene Expression Omnibus (GEO)

- o GEO (Microarray)
    - GEO Datasets (RNA-Seq, ChiP Seq, ChIP-ChIP)
    - GEO Profiles (Global expression of protein coding genes-probe-based microarray expt)
    - What is a microarray?
        - Plate with spots. Each spot includes a gene PROBE (in the form of short piece of DNA attached to the plate) that will bind with a short piece of cDNA from a specific gene
        - mRNA that are harvested from a sample (cell/tissue) are transcribed into cDNA. These cDNAs are also labelled into red or green fluorescent dye (treated/control)
        - Light intensity for each color is measured for each spot. Statistics methods are used for processing.

Public domain figure: https://upload.wikimedia.org/wikipedia/commons/thumb/a/a8/NA_hybrid.svg/800px-NA_hybrid.svg.png

fixed probes

labelled target (sample)

different features
(e.g. bind different genes)

Fully complementary          Partially complementary
strands bind strongly         strands bind weakly

We need to understand how GEO handles data submission and data curation.

a. Investigator carries out an experiment.
b. He submits the details of his experiment by submitting a single series records.
    a. General information about the experiment (how and other details)
c. Samples are separately entered (sample records; ex 2 controls (replicates) and 2 treated (replicates); in this case 4 sample records)
d. Finally, the investigator submits a platform record which provides the details of the Chips used in the experiments.
    a. Chips are also called arrays and in the GEO context they are called Platforms

All these ends up in the GEO DataSets database and gets assigned a number or ID (ex. GSE33253)

Link for GEO DataSets is www.ncbi.nlm.nih.gov/gds

Note the left side menu shows that there are one **Series** record, four **Samples** record and one **Platforms** record. Look back to see the definitions of these records. Click on each one to explore the details.

**Series**    Accession is GSE33253
**Sample**  Accession IDs are:  GSM822870-GSM822873 (4 samples)
**Platform** Accession: GPL1261 ID: 100001261

Note that there are no reference sequences for GEO. For example, in NCBI, the redundant GenBank sequences are reduced to one RefSeq sequence. The same with other DBs (PubChem etc.). But, NCBI groups the Series entry, their samples, the platform used for the series and groups them into one curated DataSet Entry (called **GDS**). For example, let us look at a curated and non-curated dataset display within GEO DataSets. The left-side entry is for curated and right-side dataset is yet to be curated. GDS entries can be analyzed using Curated DataSet Browser (http://www.ncbi.nlm.nih.gov/sites/GDSbrowser/ )

Please note that when you look at this page later, you will see different number of hits

Show additional filters

**Entry type**
DataSets (1)
Series (1)
Samples (6)
Platforms (1)

**Organism**
Select ...

**Study type**
Expression profiling by array
More

Show additional filters

**Entry type**
Series (1)
Samples (4)
Platforms (1)

**Organism**
Select ...

**Study type**
Expression profiling by array
More ...

## Uncurated dataset



## Curated dataset

Expression Profiles will give the individual Gene profiles associated with the study. The other options are self-explanatory.



Follow the gene profiles to look at the profile of MAP Kinase (ERK)

We can use GEO2R to carry out some analysis on the non-curated datasets. Let us do that now. Curated datasets can also be used in GEO2R to carry out your own analysis. The big help of having a curated dataset is the links to other parts of the Database.

Other expression data

**RNA-Seq (Sequence Read Archive)**

**RNA-Seq can be used just like ESTs.  The problem is these are raw reads.**

The Sequence Read Archive (SRA) stores raw sequence data from "next-generation" sequencing technologies including 454, IonTorrent, Illumina, SOLiD, Helicos and Complete Genomics.

RNA-seq data is often used to identify exon-intron boundary.

**Biosystems can also be used to learn about Gene Expression or regulation or epigenetics etc.**

## Genomic Testing Registry (GTR)