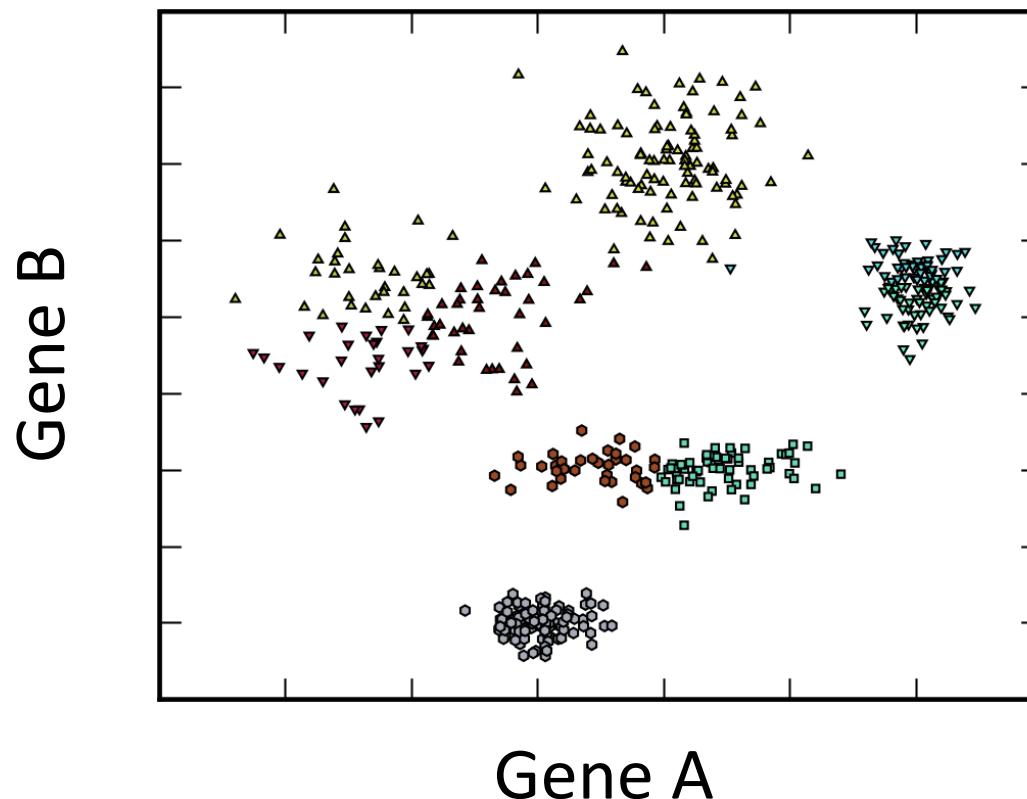


# Write on Board Before Class:

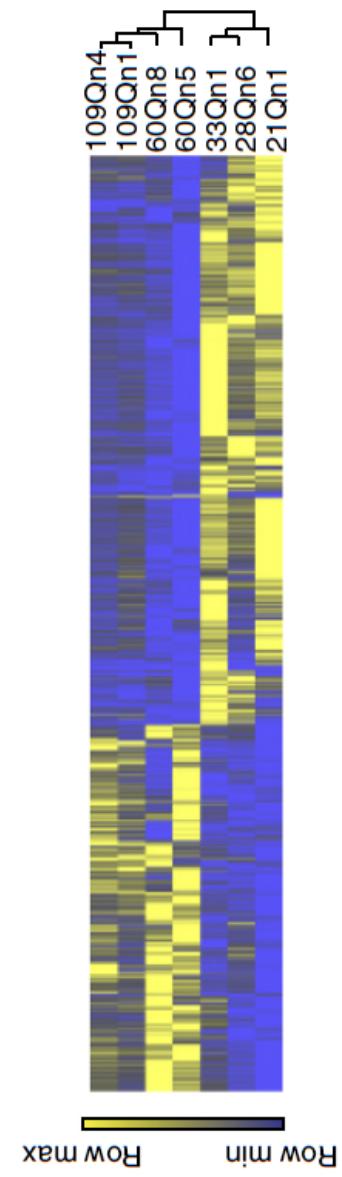
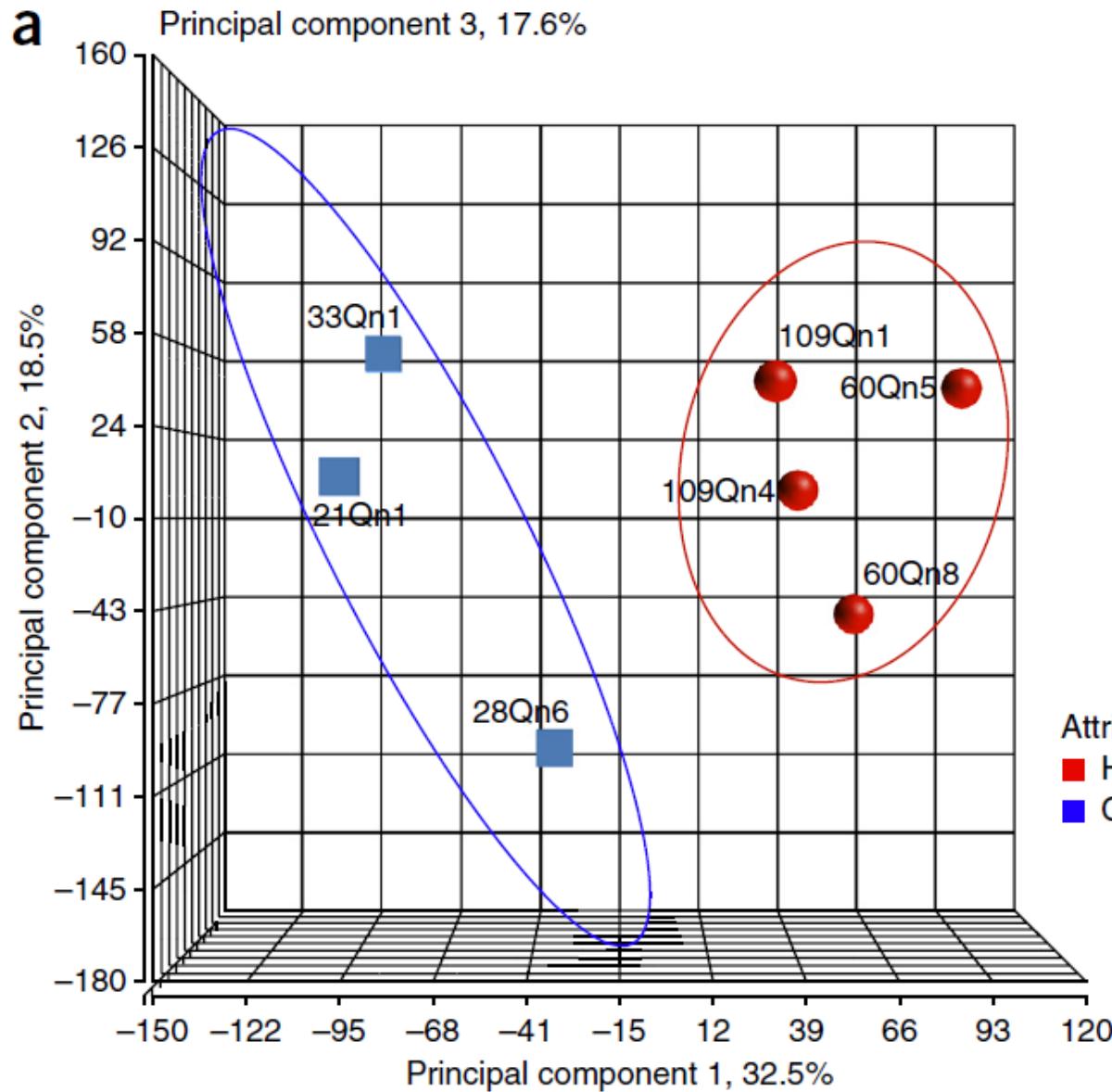
## Learning Objectives

- Describe the results of Principal Component Analysis (PCA)
- Understand the basis of an RNA-Seq experiment
- Describe the steps from “raw reads” to gene counts
- Calculate RPKM values
- Explain the role of DESeq2
- Interpret Gene Ontology
- Evaluate statistical significance of GO terms for sets of genes

Clustering works well for 2D data.  
But how could you visualize clusters in  
20,000D?



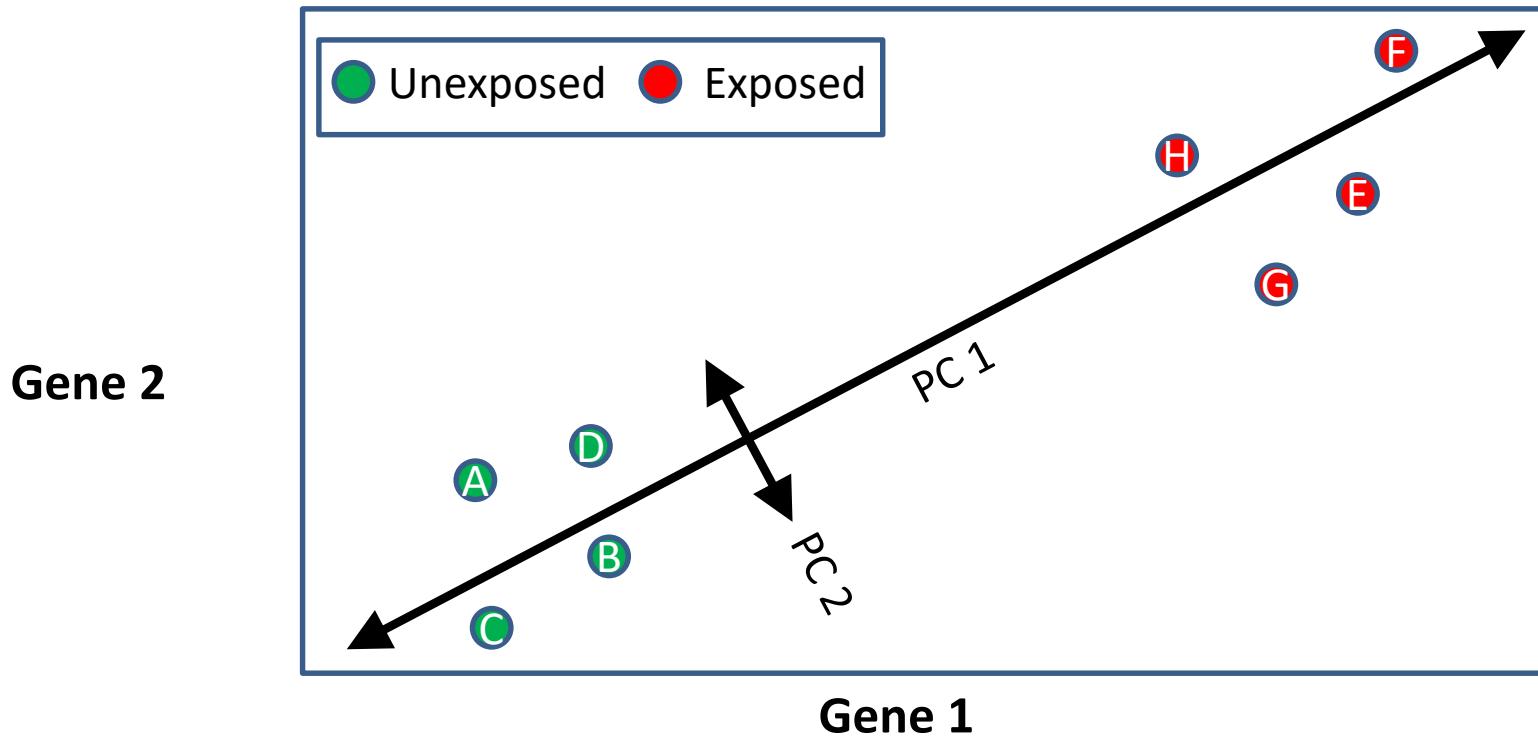
# The basics of PCA



# Principal Component Analysis

- Each sample is currently described by the expression of roughly 20,000 genes.
- Our goal:
  - to find a 2-D or 3-D way to present the data that captures the greatest variance
    - Obviously, I could select any two genes, but they might be the wrong ones.
    - Can we find “interesting” linear combinations of genes?

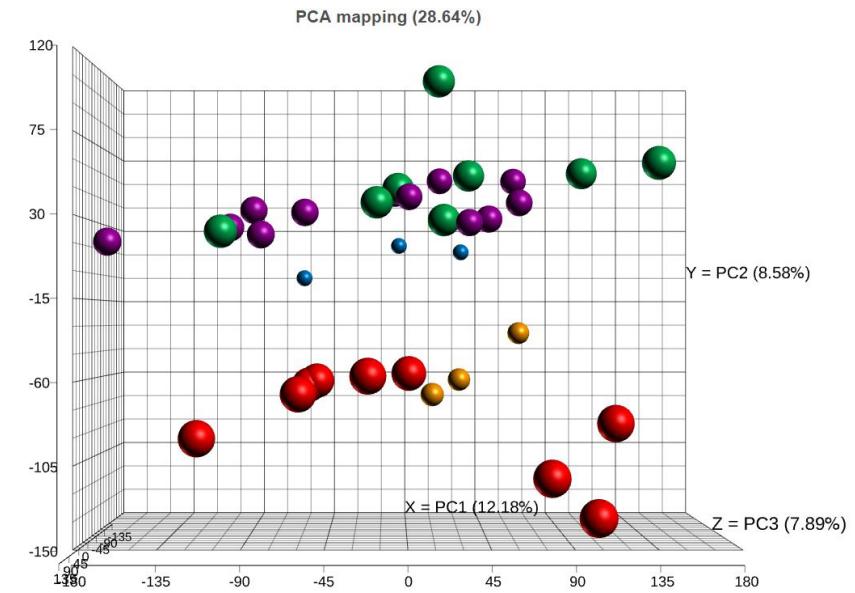
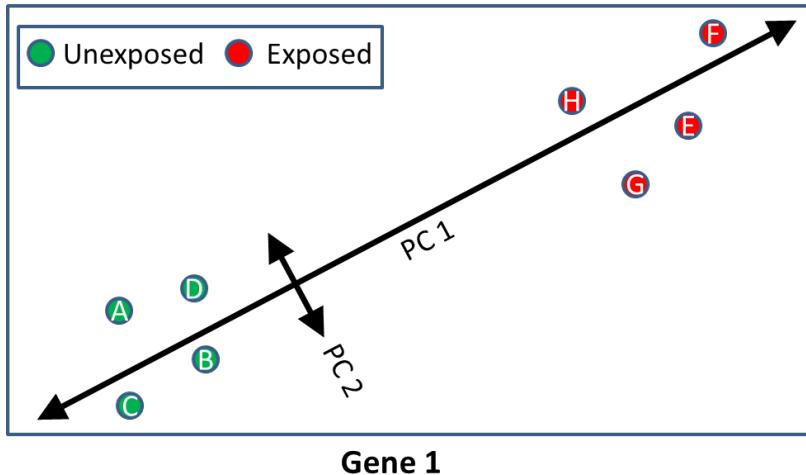
# Principal Component Analysis



Here's an example where one dimension is almost as good as two.

We can generalize this approach so that one dimension is almost as good as N, where N is large.

1. PCA finds useful linear combinations of thousands of variables.
2. There are as many PCs as there were dimensions in the original data.
3. The PCs are orthogonal.
4. Often, a few PCs will capture most of the variance.

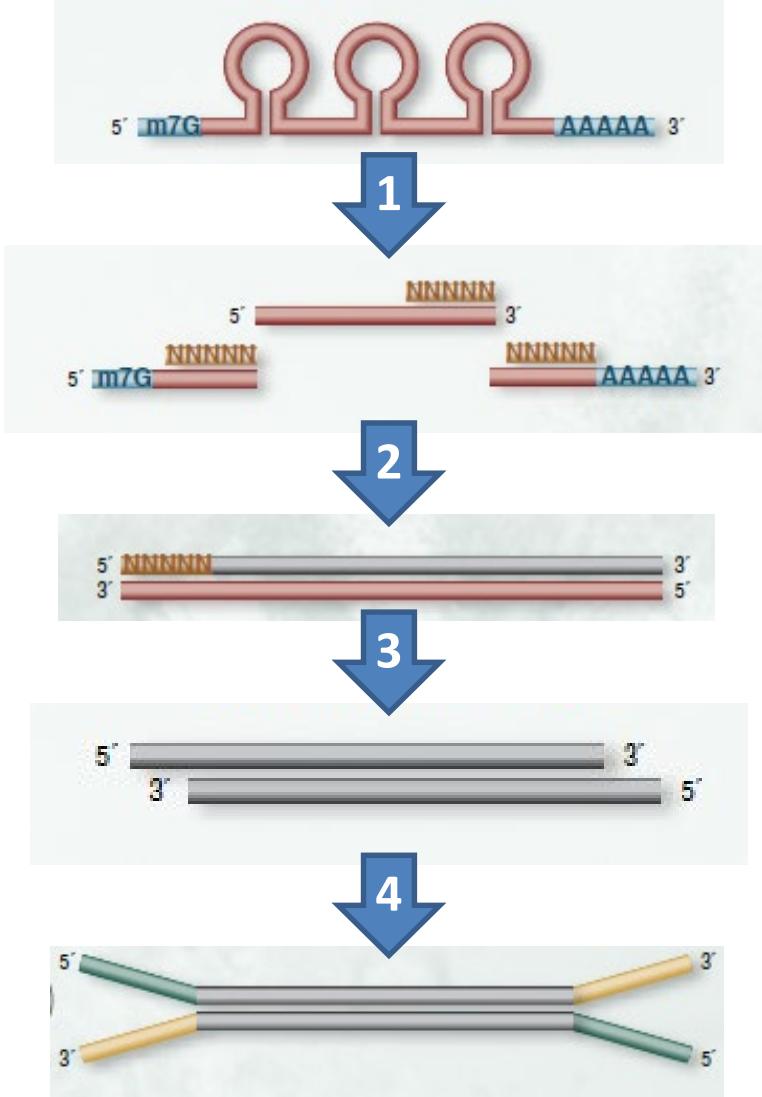


# Outline

- Overview of the steps of RNA-Seq
- Deriving expression levels from sequence data
- Gene Ontology
- Statistical significance

# Experimental Design for RNA-Seq

- Goal of RNA-Seq:
  - To measure the expression of all genes in a sample
- Sequencing machines are great but have limitations:
  - They work on DNA, not RNA
  - They are best for short fragments

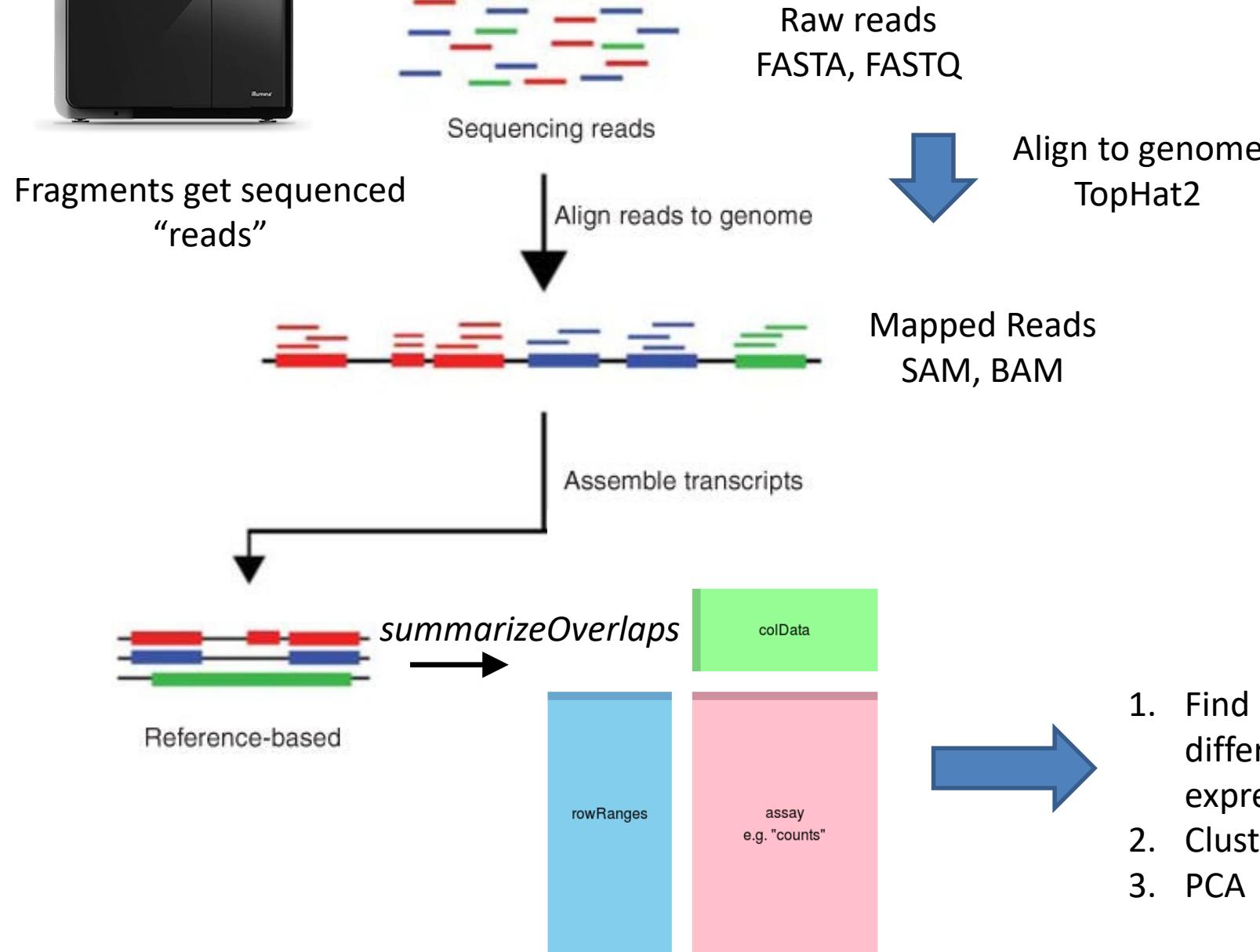


1. Fragment RNA and prime with random DNA primers
2. Synthesize second strand with Reverse Transcriptase
3. Remove RNA and synthesize second strand of DNA
4. Ligate adaptors for sequencing

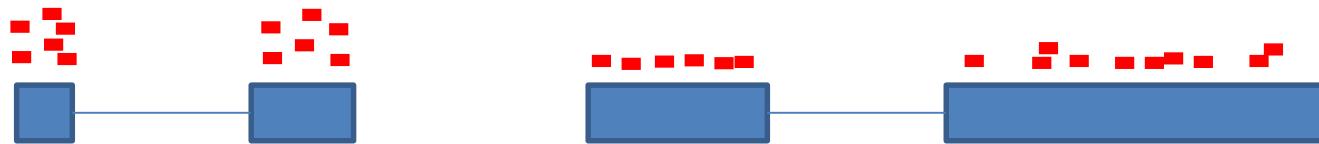
RNA	App	3' Adaptor	P5 Primer
DNA		5' Adaptor	P7 Primer
RT Primer			Barcode (BC)

# Outline

- Overview of the steps of RNA-Seq
- Deriving expression levels from sequence data
- Gene Ontology
- Statistical significance



# Raw counts are misleading



1. A long transcript with a low level of expression will still produce more sequence reads than a short, highly expressed transcript.
  2. An experiment that is sequenced more deeply will make all genes appear to be expressed at higher levels
- To correct for this, we use “Reads per Kilobase Million (RPKM)”

Gene	Length in KB	Replicate 1	Replicate 2	Replicate 3
A	2	1.0E6	1.2E6	3.0E6
B	4	2.0E6	2.5E6	6.0E6
C	10	0	0	1.0E5
Total reads			3.0E6	3.7E6
Reads/1,000,000			3	3.7

## Raw reads

1. Count the number of reads in each sample in millions.
2. Divide reads for a gene by the number of reads in the replicate (in millions)
3. Divide by gene length in kilobases

Reads per million	A	0.333	0.324	0.330
		B	0.667	0.676
C	0	0	0.011	
Reads per kilobase million	A	0.167	0.162	0.165
		B	0.167	0.169
RPKM	C	0.00	0.00	0.001

Gene	Length in KB	Replicate 1	Replicate 2	Replicate 3
A	2	1.0E6	1.2E6	3.0E6
B	4	2.0E6	2.5E6	6.0E6
C	10	0	0	1.0E5
Total reads		3.0E6	3.7E6	9.1E6
Reads/1,000,000		3	3.7	9.1

Reads per million	A	0.333	0.324	0.330
	B	0.667	0.676	0.659
	C	0	0	0.011

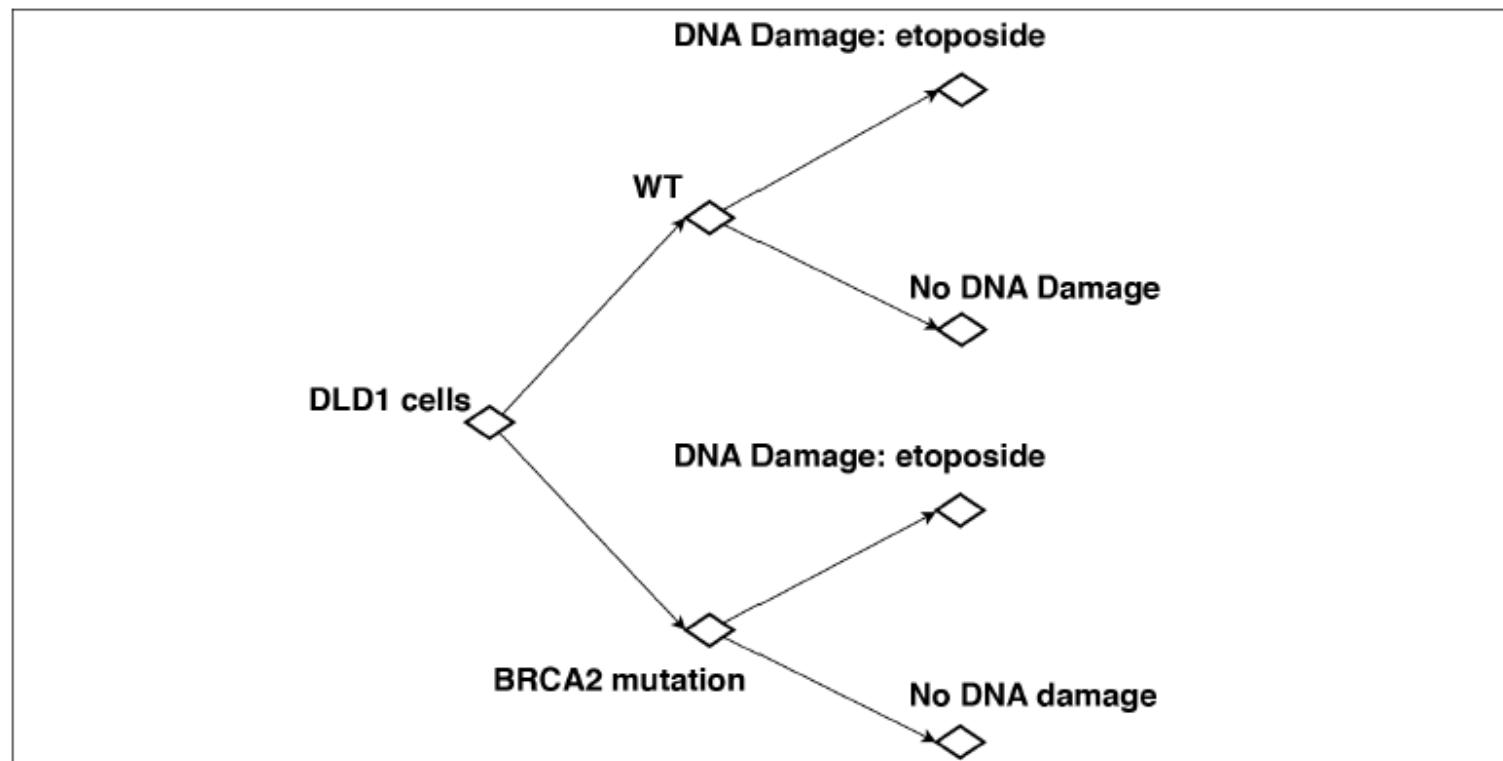
This step corrects for sequencing depth. Note that numbers are now more consistent across replicates

RPKM		Replicate 1	Replicate 2	Replicate 3
	A	0.167	0.162	0.165
	B	0.167	0.169	0.165
	C	0.00	0.00	0.001

This step corrects for gene length. Note that genes A and B have similar RPKMs but very different raw read counts.

# Differential expression

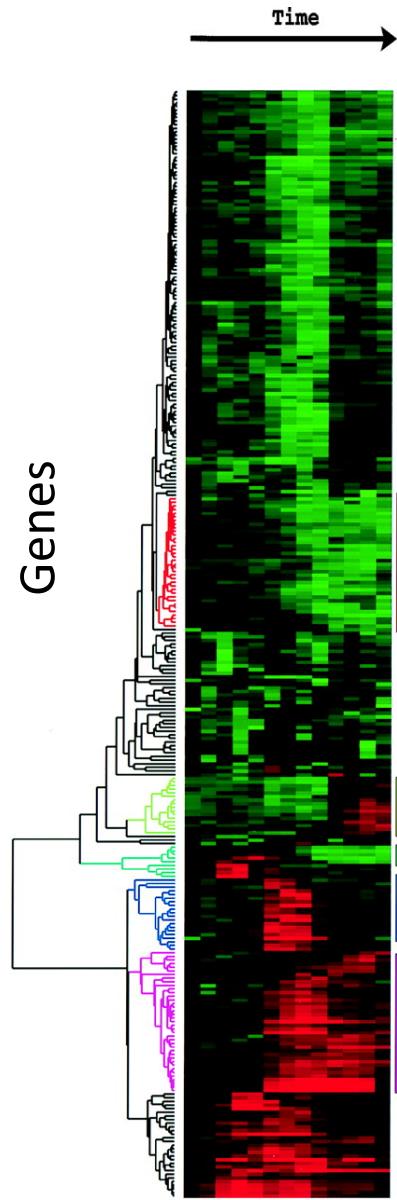
**DESeq2:** tests whether a difference in gene expression is a response to a change in condition vs. a random fluctuation



# Do your data make sense?

- Technical replicates should be very similar ( $R^2 > .9$ )
- Biological replicates should cluster together

# Interpreting your results



How did they figure out what  
the clusters of genes did?

- (A) cholesterol biosynthesis
- (B) the cell cycle
- (C) the immediate-early response
- (D) signaling and angiogenesis
- (E) wound healing and tissue remodeling

# Outline

- Overview of the steps of RNA-Seq
- Deriving expression levels from sequence data
- Gene Ontology
- Statistical significance

# Biological Insights

- What types of genes are being differentially expressed?

The screenshot shows the homepage of the Gene Ontology website. At the top left is a circular diagram illustrating biological processes. Next to it is the text "the Gene Ontology". In the center is the URL "http://www.geneontology.org". To the right is a search bar with the placeholder "Search" and a dropdown menu set to "gene or protein name" with a "go!" button. Below the header is a navigation bar with links for "Downloads", "Tools", "Documentation", "Projects", "About", and "Contact".

Controlled vocabulary to describe genes:

- Biological process
  - signal transduction; glucose transport
- Cellular component
  - nucleus; ribosome; protein dimer
- Molecular function
  - binding; transporter

<input type="checkbox"/> Gene/product	Gene/product name	Annotation qualifier	GO class (direct)
<input type="checkbox"/>	<a href="#">BRCA2</a>	Breast cancer type 2 susceptibility protein	mitotic cytokinesis
<input type="checkbox"/>	<a href="#">BRCA2</a>	Breast cancer type 2 susceptibility protein	telomere maintenance via recombination
<input type="checkbox"/>	<a href="#">BRCA2</a>	Breast cancer type 2 susceptibility protein	double-strand break repair via homologous recombination
<input type="checkbox"/>	<a href="#">BRCA2</a>	Breast cancer type 2 susceptibility protein	double-strand break repair via homologous recombination

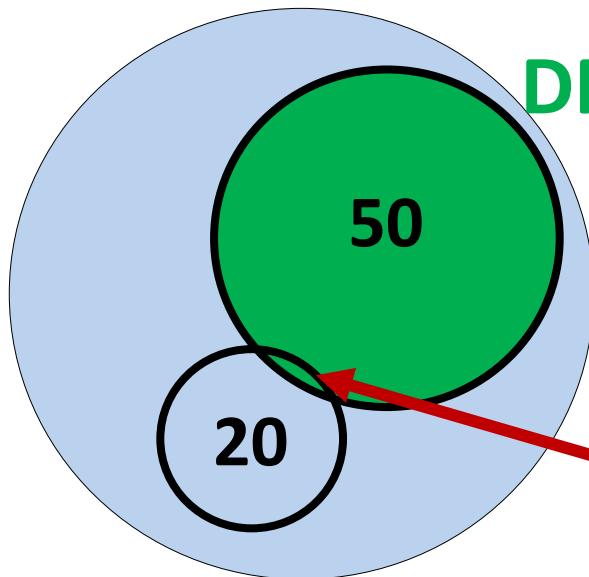
# Outline

- Overview of the steps of RNA-Seq
- Deriving expression levels from sequence data
- Gene Ontology
- Statistical significance

# Statistical Significance

- Your startup just developed a new drug, but related compounds cause cancer
- You want to know if it's safe
- Your idea: test it on cell lines and see what genes change in expression
- You find that it activates some genes involved in DNA Repair
- Could it be causing DNA damage?

Genome (100)



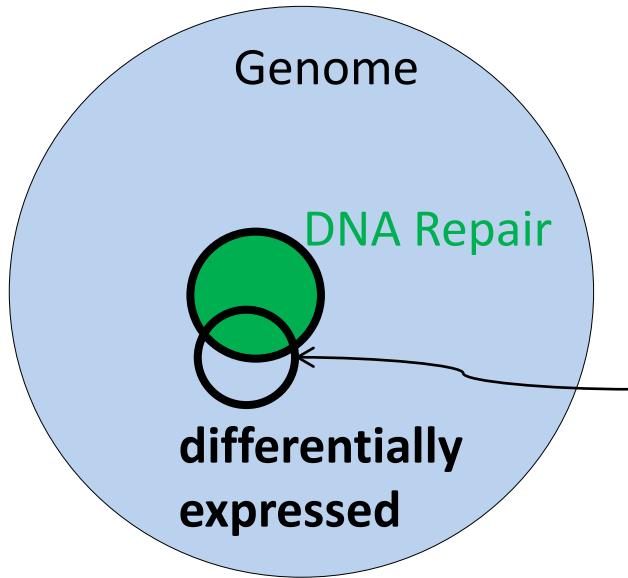
Differentially  
expressed

DNA Repair

One gene differentially expressed gene is related to DNA repair.

Should I worry that our drug causes DNA damage?

# Statistical significance

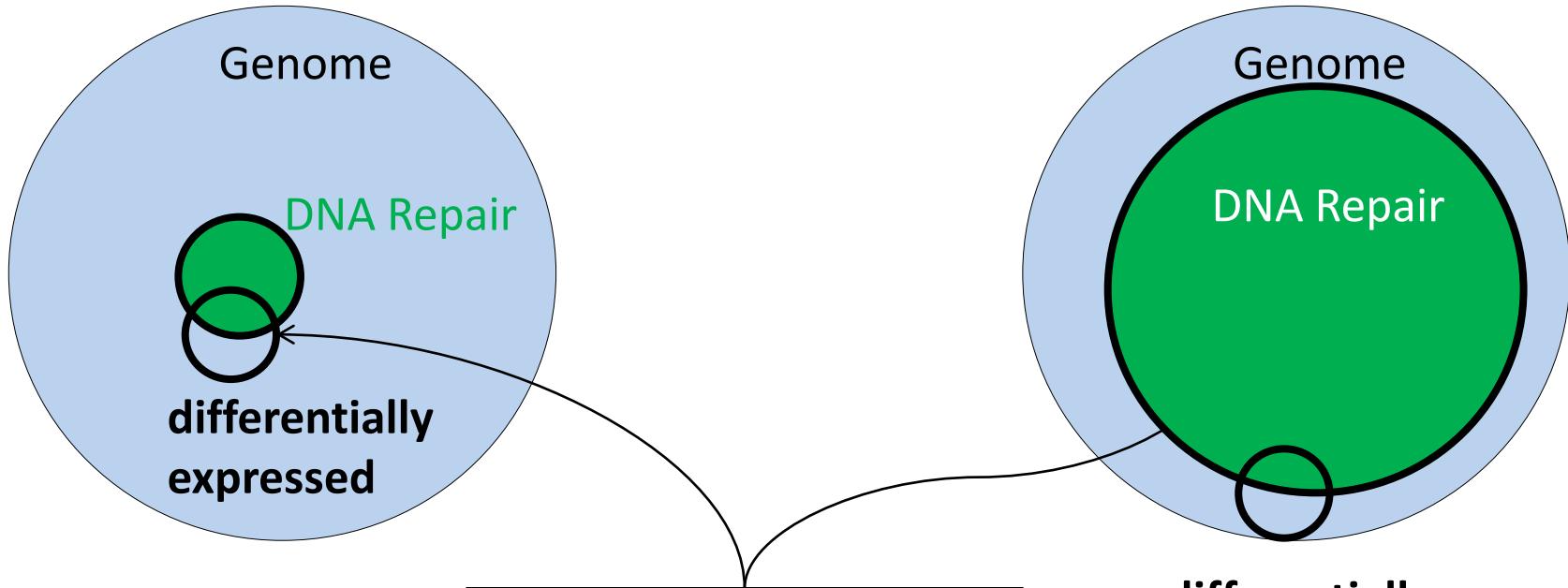


If I get many **more** repair genes than I would expect by chance, I need to find out if my drug is causing DNA damage.

In other words: are the differentially expressed genes **enriched** for ones involved in DNA repair?

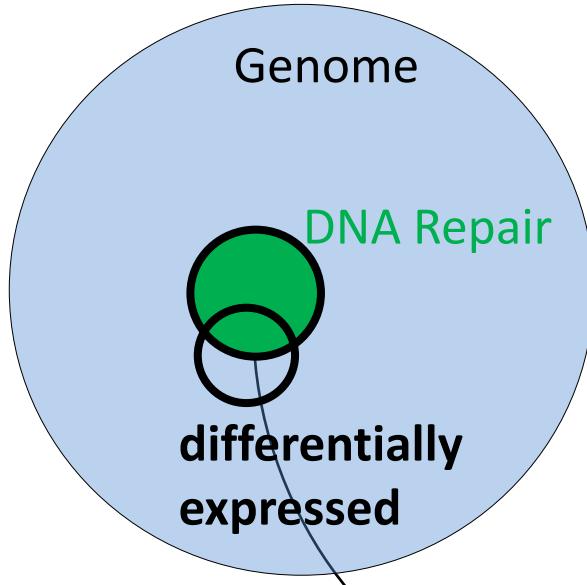
# Statistical significance

The significance depends on the size of the lists.

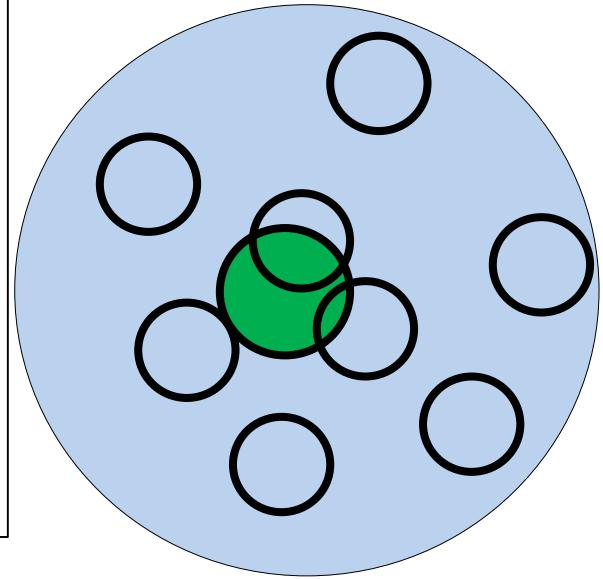


If the two lists had  
nothing in common,  
could we still get this  
degree of overlap?

# Statistical significance

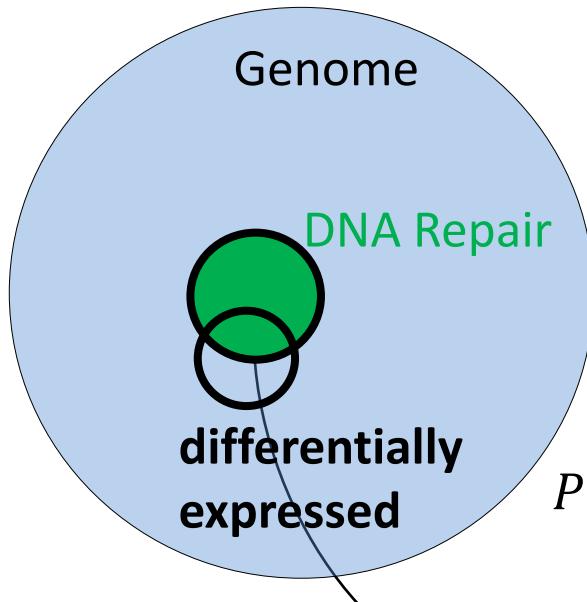


**Empirical approach:**  
Find the distribution of observed “green genes” by random sampling



Is this overlap significant?

# Statistical significance



## Analytical Approach:

The probability of getting **exactly** this amount of overlap for two randomly chosen sets of genes of the same size is given by the hypergeometric distribution:

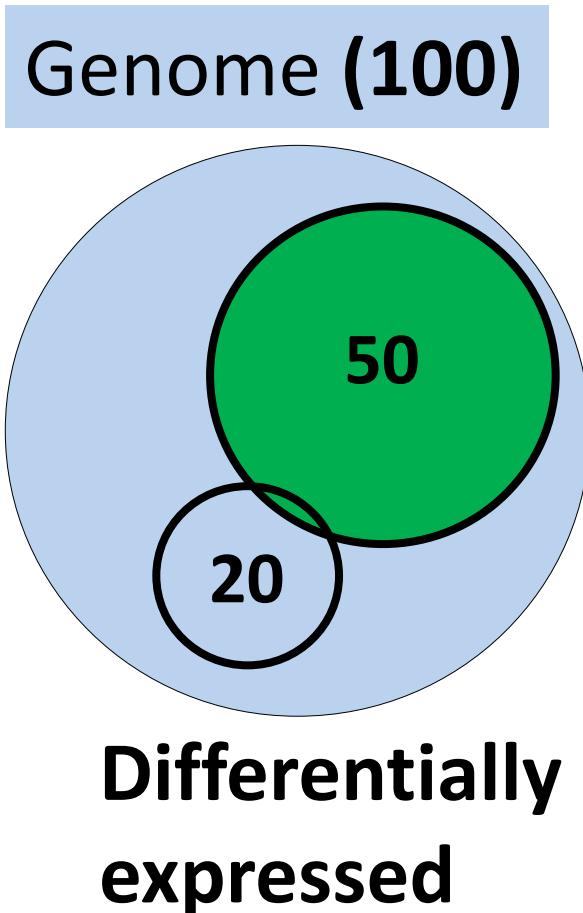
$$P(Overlap) = \frac{\binom{DNA\ repair}{Overlap} \binom{Genome - DNA\ repair}{DiffExp - Overlap}}{\binom{Genome}{DiffExp}}$$

Is this overlap significant?

Recall that  $\binom{n}{k}$  ("n choose k") is the binomial coefficient.

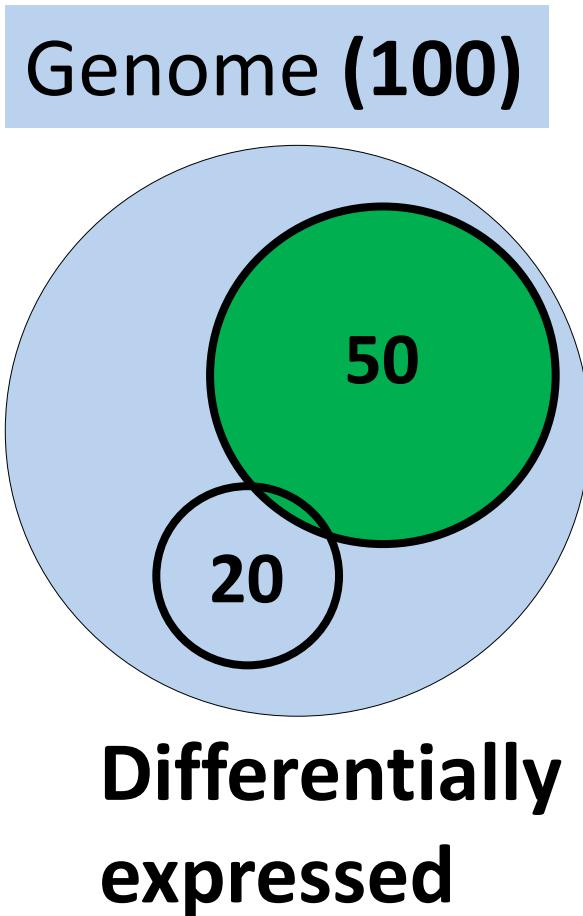
= the number of ways to choose k items from a set of n.

# How you might use the HG test:



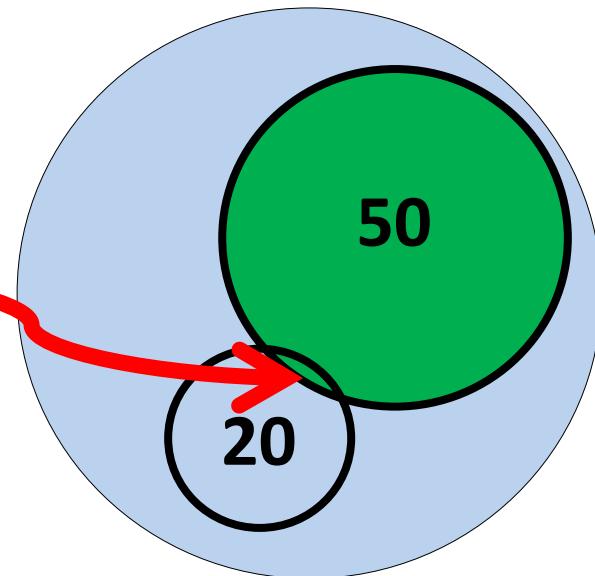
- Your startup just developed a new drug, but related compounds cause cancer
- You want to know if it's safe
- Your idea: test it on cell lines and see what genes change in expression
- You find that it activates some genes involved in DNA Repair
- Could it be causing DNA damage?

# Statistical significance



- Usually, we wish to test if a term is “**enriched**” in our data.
- But the hypergeometric gives the probability of getting **exactly** this amount of overlap for two randomly chosen sets of genes of the same size.
- Using the CDF, we can ask if we see **more** of a term than we would expect under the null model.

One gene overlaps.  
Should I worry that our  
drug causes DNA damage?

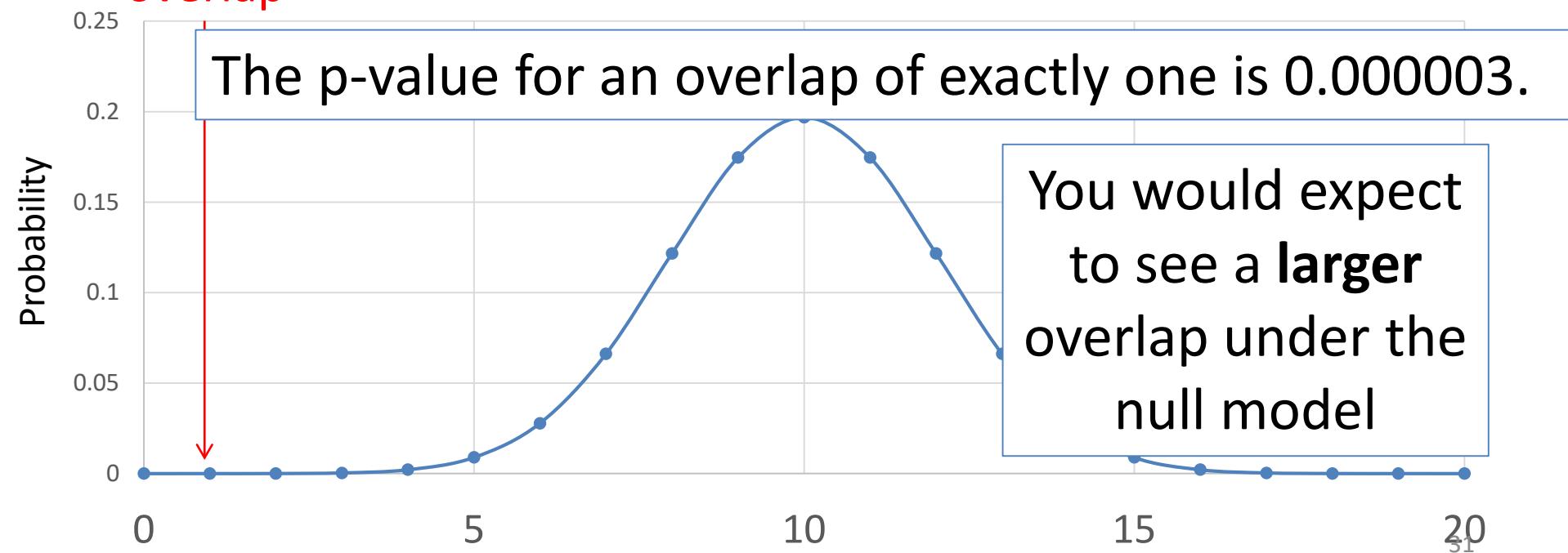


Observed  
overlap

Hypergeometric Distribution

The p-value for an overlap of exactly one is 0.000003.

You would expect  
to see a **larger**  
overlap under the  
null model

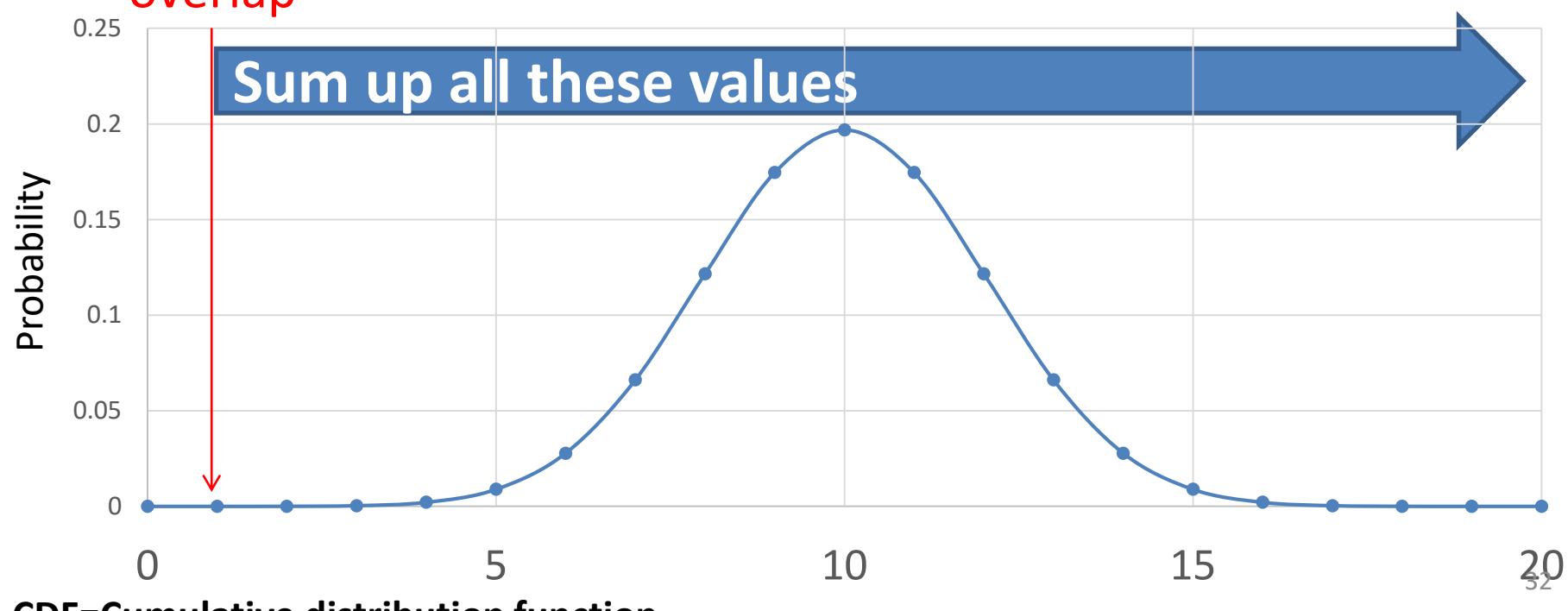


# The CDF helps us find enriched terms

$$CDF(Overlap) = \sum_{n=overlap}^{\text{Number of genes in DNA Repair}} \frac{\binom{DNA\ repair}{n} \binom{Genome - DNA\ repair}{DiffExp - n}}{\binom{Genome}{DiffExp}}$$

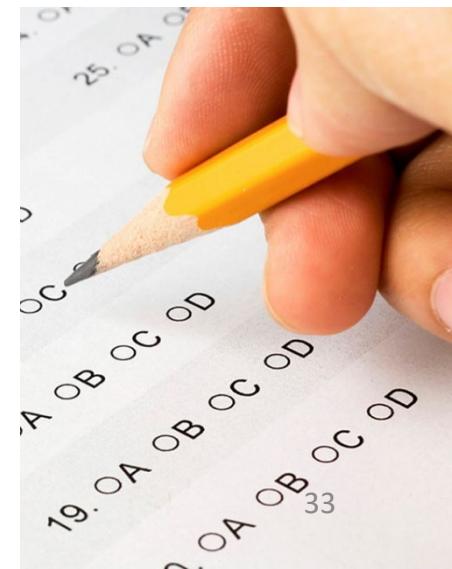
Observed overlap

Hypergeometric Distribution



# Testing Multiple Hypotheses

- Example: Filter GO terms using a p-value threshold of 0.01
- By definition, the null-hypothesis has a 1% probability of being correct for each test.
- There are roughly 30,000 terms in GO.
- At this level, we expect roughly 300 false positives!



# Multiple Hypotheses

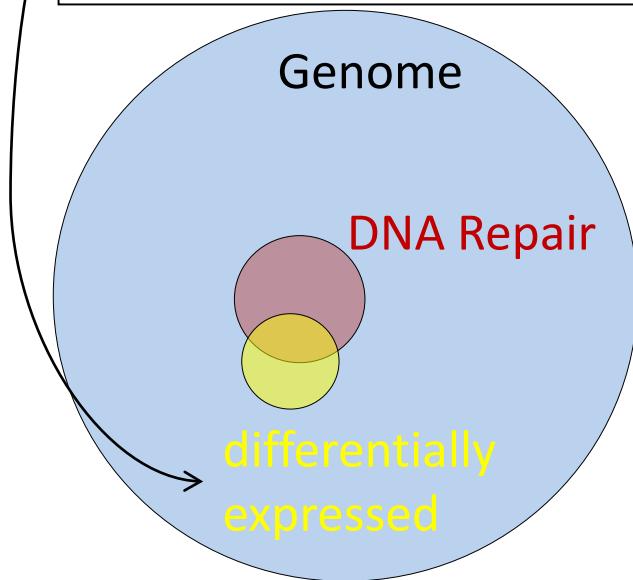
- A simple solution: require that the p-value be small enough to reduce the false positives to the desired level.
- This is called the Bonferroni correction.
- In our case, we would only accept terms with a

$$p \leq \frac{0.01}{30,000} = \frac{\text{desired threshold}}{\text{number of tests}}$$

- Since our tests are not all independent, this is very conservative, and will miss many true positives
- More sophisticated approaches exist, such as controlling the “false discovery rate”.

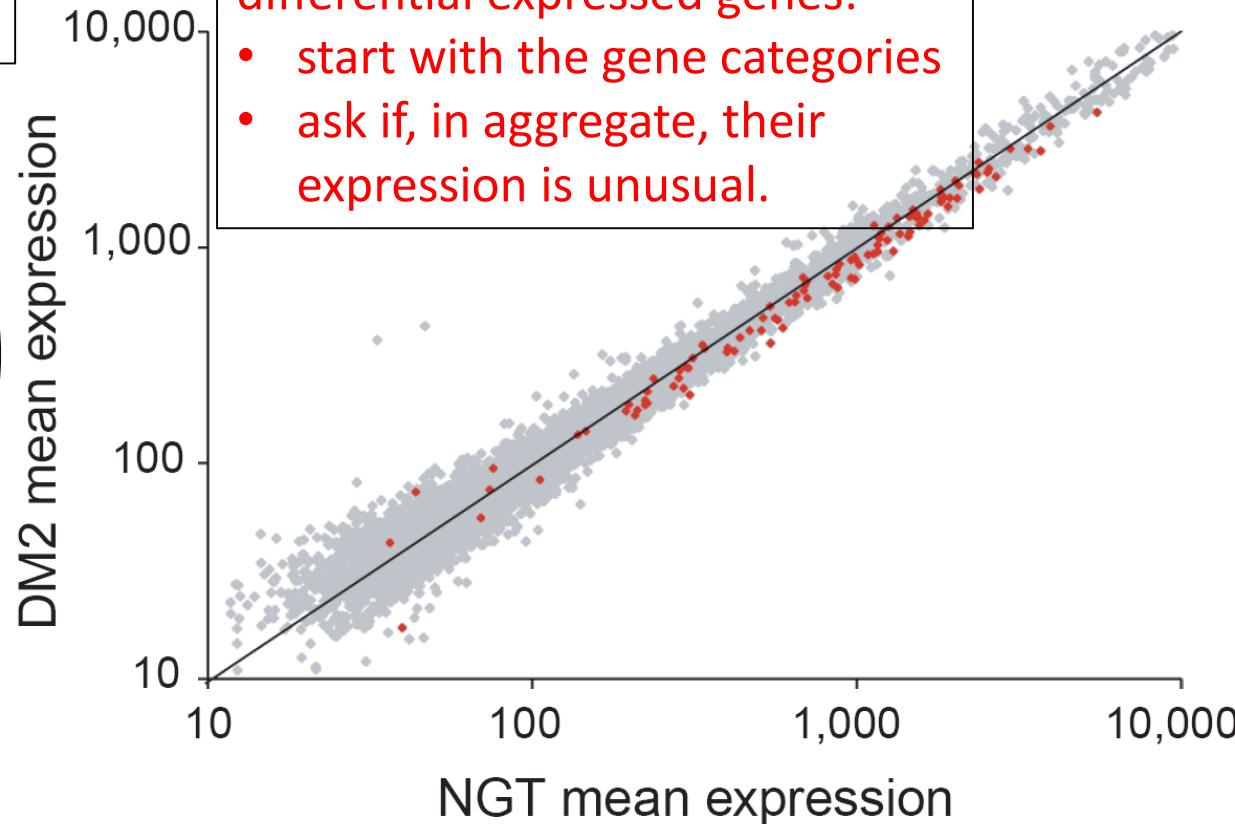
# Aggregate score statistics

My results depend on how I defined “differentially expressed”



Instead of starting with differential expressed genes:

- start with the gene categories
- ask if, in aggregate, their expression is unusual.



# Aggregate score statistics

<http://www.broadinstitute.org/gsea/>

The screenshot shows the GSEA (Gene Set Enrichment Analysis) website. At the top, there's a navigation bar with links for "GSEA Home", "Downloads", "Molecular Signatures Database", "Documentation", and "Contact". The "GSEA Home" link is highlighted. On the right side of the header is the "BROAD INSTITUTE" logo with a search bar. Below the header, there's a large diagram illustrating the GSEA process: "Molecular Profile Data" (represented by a scatter plot) and a "Gene Set Database" (represented by a cylinder) both feed into a central "Run GSEA" box. The output of "Run GSEA" is "Enriched Sets", shown as a bar chart.

## Overview

**Gene Set Enrichment Analysis** (GSEA) is a computational method that determines whether *a priori* defined set of genes shows statistically significant, concordant differences between two biological states (e.g. phenotypes).

## What's New

02/19/10: We have a new release of GSEA 2.0.6 that fixes the FTP problems that have been experienced recently. Please discontinue use of older versions and use the new version instead.

12/10/09: Leading Edge Analysis now works correctly in Release GSEA 2.0.5. There are no changes to the algorithm or functionality.

12/07/2009: Release GSEA 2.0.5 of the GSEA java application is now available. The new release has been updated to work on Snow Leopard. There are no changes to the algorithm or functionality. This update requires Java 6 (on all platforms).

## Getting Started

A quick tutorial to get you up and running.

## Tools and Information

**Downloads:** Implementations of GSEA plus additional resources to analyze, annotate and interpret enrichment results.

**Molecular Signatures Database:** A collection of gene sets for use with GSEA software and tools for exploring them.

**Documentation:** Information on the GSEA software, the GSEA algorithm.

## Registration

Please register to download the GSEA software and view the MSigDB gene sets. After registering, you can log in at any time using your email address. Registration is free. Its only purpose is to help us track usage for reports to our funding agencies.

## Contributors

GSEA is maintained by the [GSEA team](#). Our thanks to our many contributors. Funded by: National Cancer Institute, National Institutes of Health, National Institute of General Medical Sciences.



## Citing GSEA

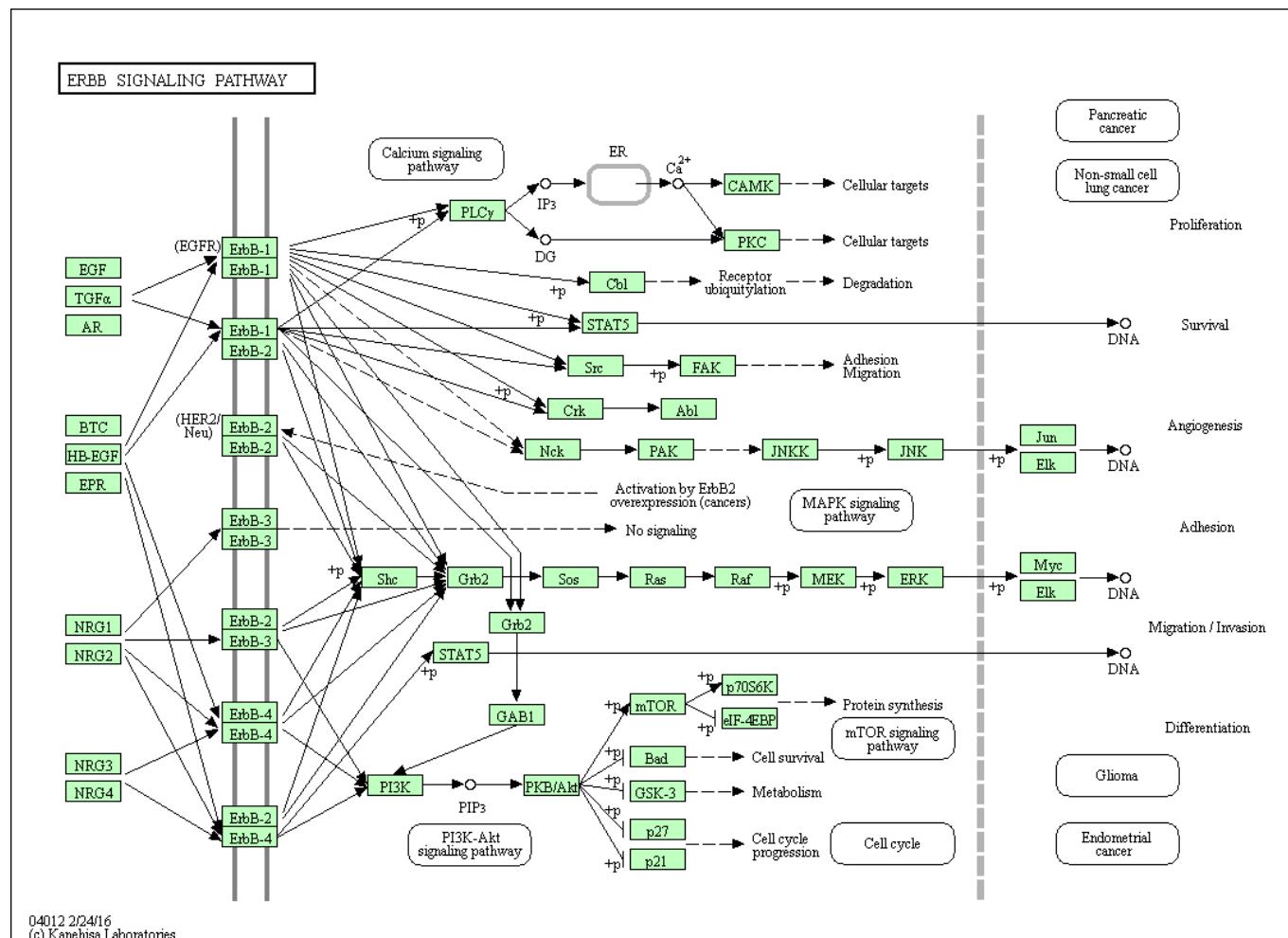
To cite your use of the GSEA software, please reference Subramanian, Tamayo, et al. (2005, PNAS 102, 15545-15550) and Mootha, Lindgren, et al. (2003, Nat Genet 34, 267-273).

- **Biological process**
- A biological process is not equivalent to a pathway.
  - Does not represent the dynamics or dependencies of a pathway.

# GO

BTC	NRAS
CDC37	NRG1
Cpne3	NRG2
CPNE3	NRG4
CUL5	PIK3CA
EGF	PIK3R1
EGFR	PRKCA
ERBB2	PTK6
ERBB3	PTPN12
ERBB4	PTPN18
ERBIN	Ptprr
EREG	PTPRR
GAB1	RPS27A
GRB2	SHC1
GRB7	SOS1
HBEGF	SRC
HRAS	STUB1
HSP90AA1	Symbol
KRAS	UBA52
MATK	UBB
Myoc	UBC
MYOC	

# KEGG Pathway



Select all

Clear all

Perform an action with this page's selected terms...

### Accession, Term

### Ontology

### Qualifier

### Evidence

<input type="checkbox"/> GO:0030520 : estrogen receptor signaling pathway	41 gene products <a href="#">view in tree</a>	biological process		NAS
<input type="checkbox"/> GO:0043526 : <b>neuroprotection</b>	67 gene products <a href="#">view in tree</a>	biological process		IEA With Ensembl:ENSRNOP00000026350
<input type="checkbox"/> GO:0048386 : positive regulation of retinoic acid receptor signaling pathway	9 gene products <a href="#">view in tree</a>	biological process		IDA
<input type="checkbox"/> GO:0045885 : positive regulation of survival gene product expression	56 gene products <a href="#">view in tree</a>	biological process		IEA With Ensembl:ENSRNOP00000026350
<input type="checkbox"/> GO:0006355 : regulation of transcription, DNA-dependent	16904 gene products <a href="#">view in tree</a>	biological process		NAS
<input type="checkbox"/> GO:0043627 : response to estrogen stimulus	354 gene products <a href="#">view in tree</a>	biological process		IEA With Ensembl:ENSRNOP00000026350
<input type="checkbox"/> GO:0007165 : signal transduction	18490 gene products <a href="#">view in tree</a>	biological process		TAS
				TAS

Not just the  
obvious categories

## GO Evidence Code Decision Tree

