

PhD Course
Advanced Biostatistics

Lecture 6
Data Analysis in Transcriptomics

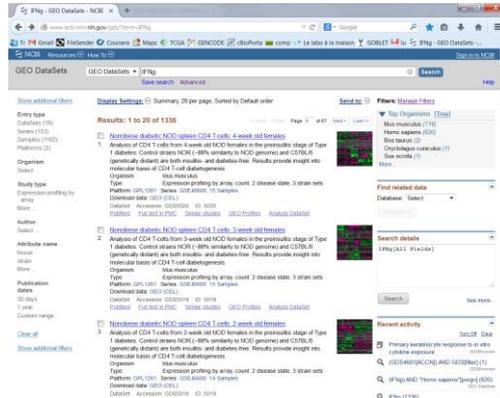
Peter Nazarov
petr.nazarov@lih.lu

31-05-2017

- ◆ **Microarray data analysis (L6.1)**
 - ◆ Principles
 - ◆ Pipeline for data analysis
 - ◆ Experiment description
 - ◆ APT import
 - ◆ QC, differential expression analysis
 - ◆ Differential expression analysis
- ◆ **RNA-seq data analysis (L6.2)**
- ◆ **Enrichment analysis (L6.1)**

Public Repositories

GEO: <http://www.ncbi.nlm.nih.gov/gds>

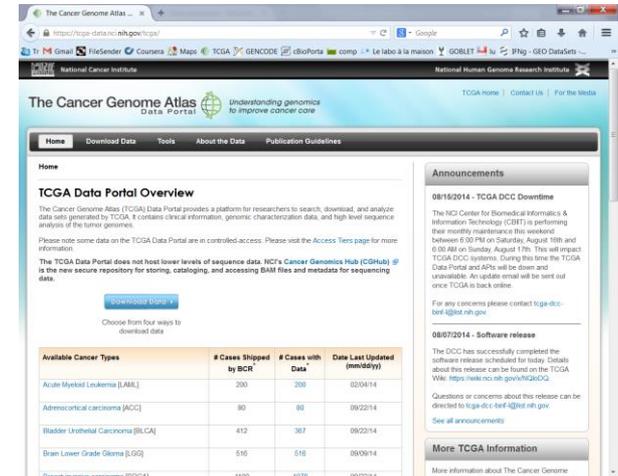


Browse Content

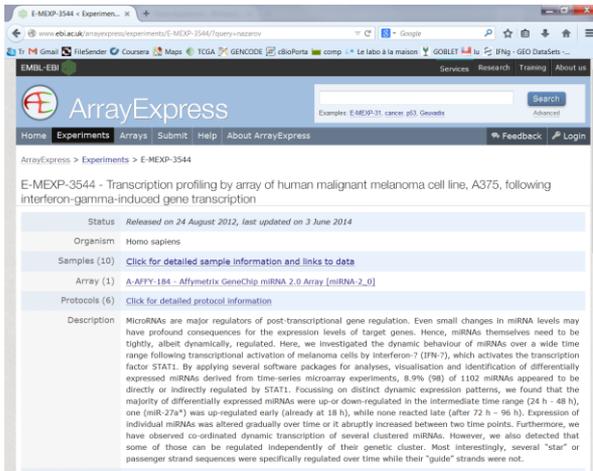
Repository Browser

DataSets:	3847
Series:	50810
Platforms:	13387
Samples:	1237318

TCGA: <https://tcga-data.nci.nih.gov/tcga/>



ArrayExpress: <http://www.ebi.ac.uk/arrayexpress/>

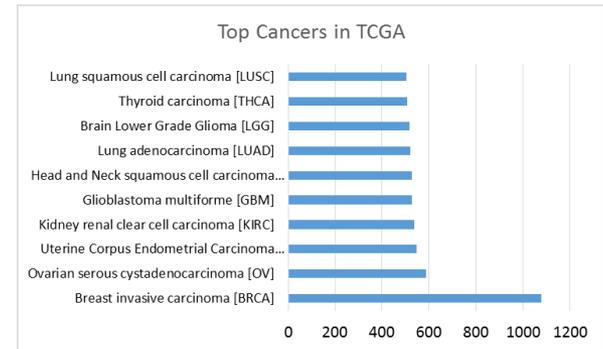


Data Content

Updated today at 06:00

- 52801 experiments
- 1555904 assays
- 24.99 TB of archived data

Sep 2014 – more than 10k patients



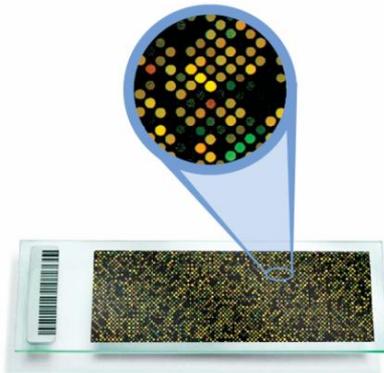
Analysis via:
<http://www.cbioportal.org/public-portal/>

Data for our course: <http://edu.sablab.net/transcript>

Types of Microarrays

Two-color Arrays (2C)

- ◆ Agilent full genome
- ◆ Thematic arrays



Pro

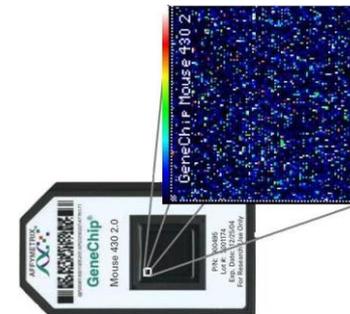
- ◆ Direct comparison
- ◆ Less sensitive to inaccuracies of spotting

Con

- ◆ Dye effects: need for “dye-swaps”
- ◆ Non-flexibility in analysis

One-color Arrays (1C)

- ◆ Affymetrix GeneChip
- ◆ Affymetrix Exon
- ◆ Affymetrix mRNA



Pro

- ◆ Flexible analysis
- ◆ High level of standardization

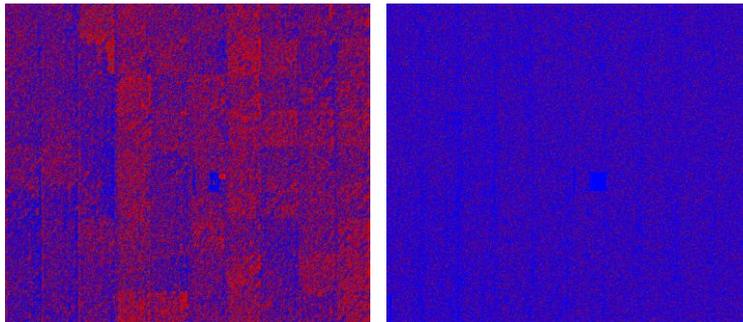
Con

- ◆ Price

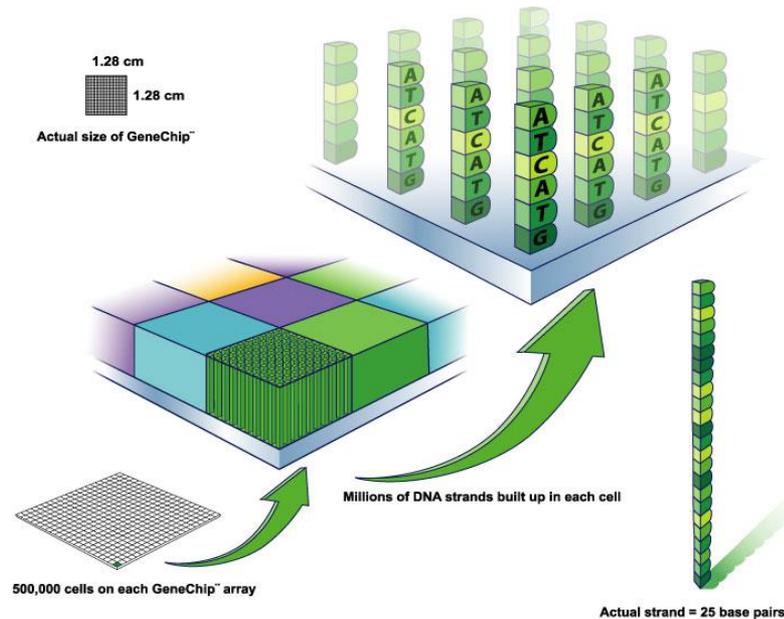
One-color Arrays

Raw

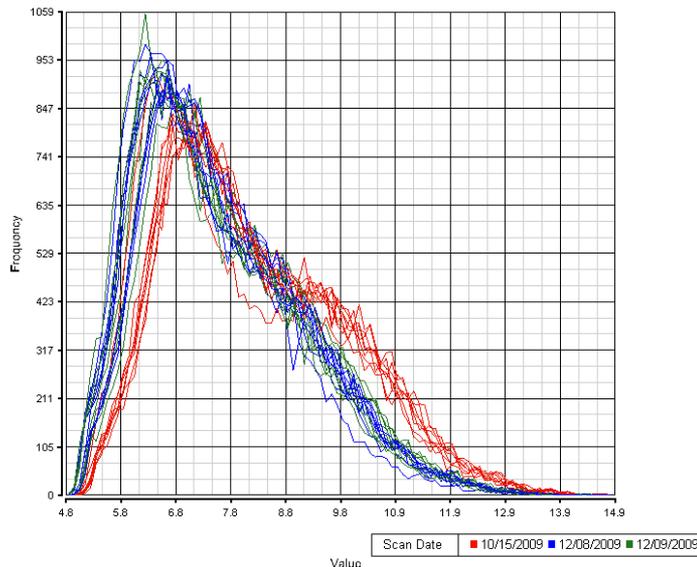
Normalized



High reproducibility and quality of spotting is required.
Affymetrix – “photolithography”-like technique



All Rows of 1



$$\text{LogIntensity} = \log_2(I)$$

Background is “removed” during normalization step

Filtering may help removing uninformative features

Affymetrix: Probes, Probesets and Transcript clusters

Probes

25-mer sequences targeted on a single region of transcriptome (hopefully)

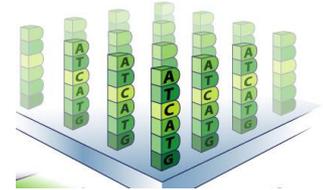
In old versions of Affy arrays (hgu95, hgu133, etc), there were:

PM – perfect match probes

MM – mismatch probes (having replacement in th 13th character)

This was done for background estimation.

But this approach is not used now!!



Probesets

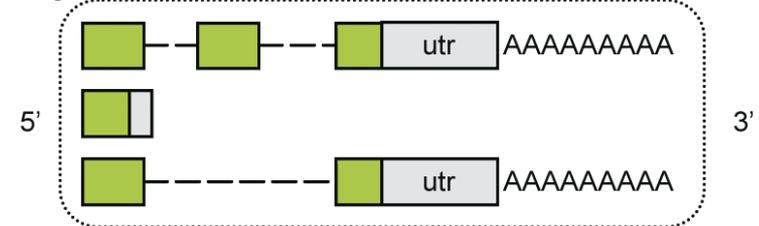
groups of closely located or overlapped probes (on average 4 probes)

3' IVT



Probeset

gene



5'

3'

Probesets

Exon



Exons

HuExon and HTA arrays allow measuring exon expression

Transcript clusters

For majority of features - synonymous to "genes". However, some distinct transcripts of genes are considered as different transcript clusters.

Okoniewski M, Comprehensive Analysis of Affymetrix Exon Arrays Using BioConductor, PLoS CompBio, 2008

Normalization of Affymetrix Arrays by RMA

Background
correction



Normalization
b/w arrays



Estimate
expression

Background and signal are strictly positive.
Noise is additive in log scale:

$$PM_{ij} = \underset{\text{exponential}}{S_{ijn}} + \overset{\text{normal}}{B_{ijn}}$$

Quantile **normalization** b/w arrays: makes distribution of probes the same across all arrays

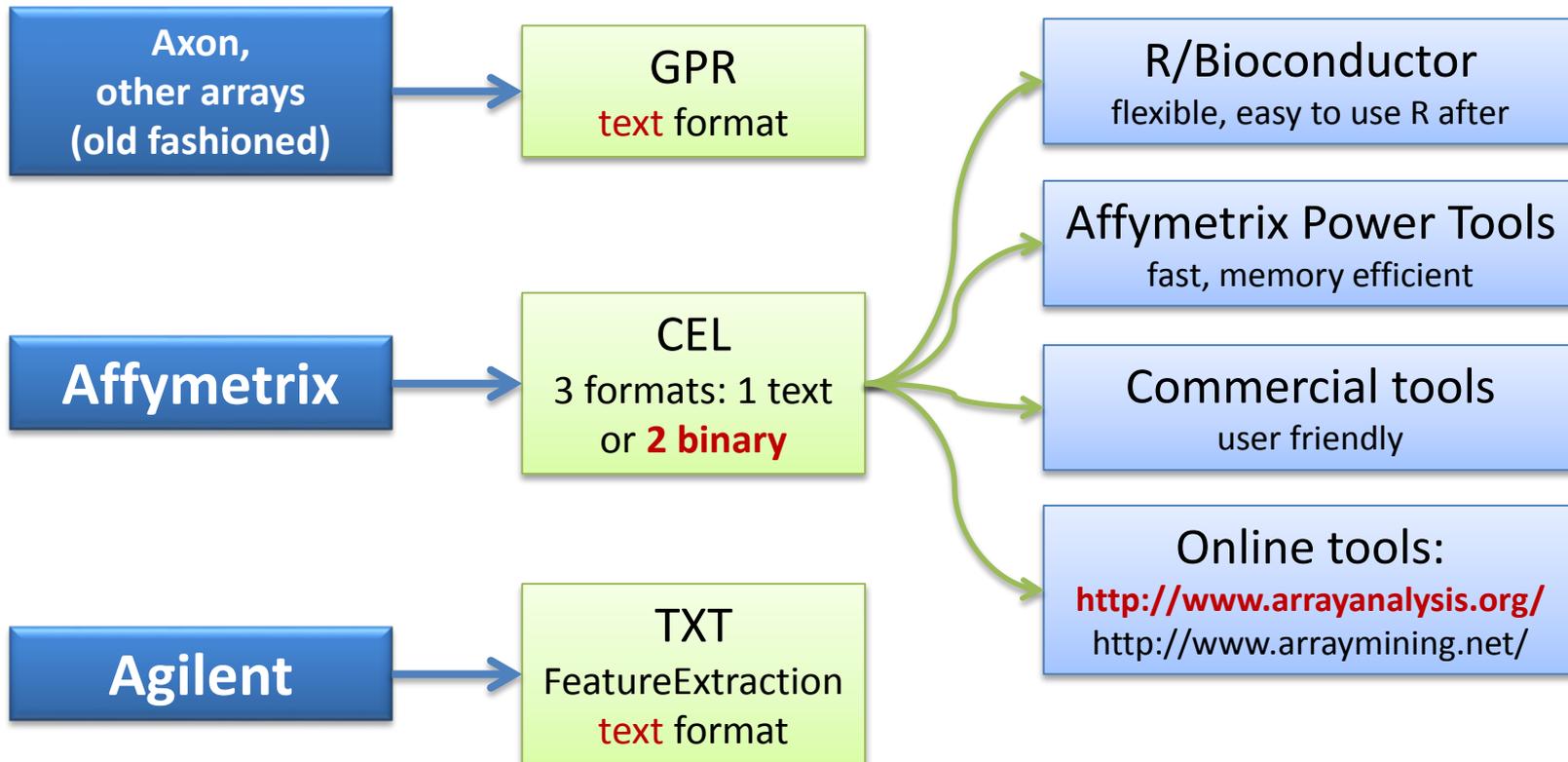
Probeset expression is estimated from a linear model:

$$Y_{ijn} = \underset{\text{observed}}{\mu_{in}} + \underset{\text{probe affinity}}{\alpha_{jn}} + \underset{\text{error with 0 mean}}{\varepsilon_{ijn}}$$

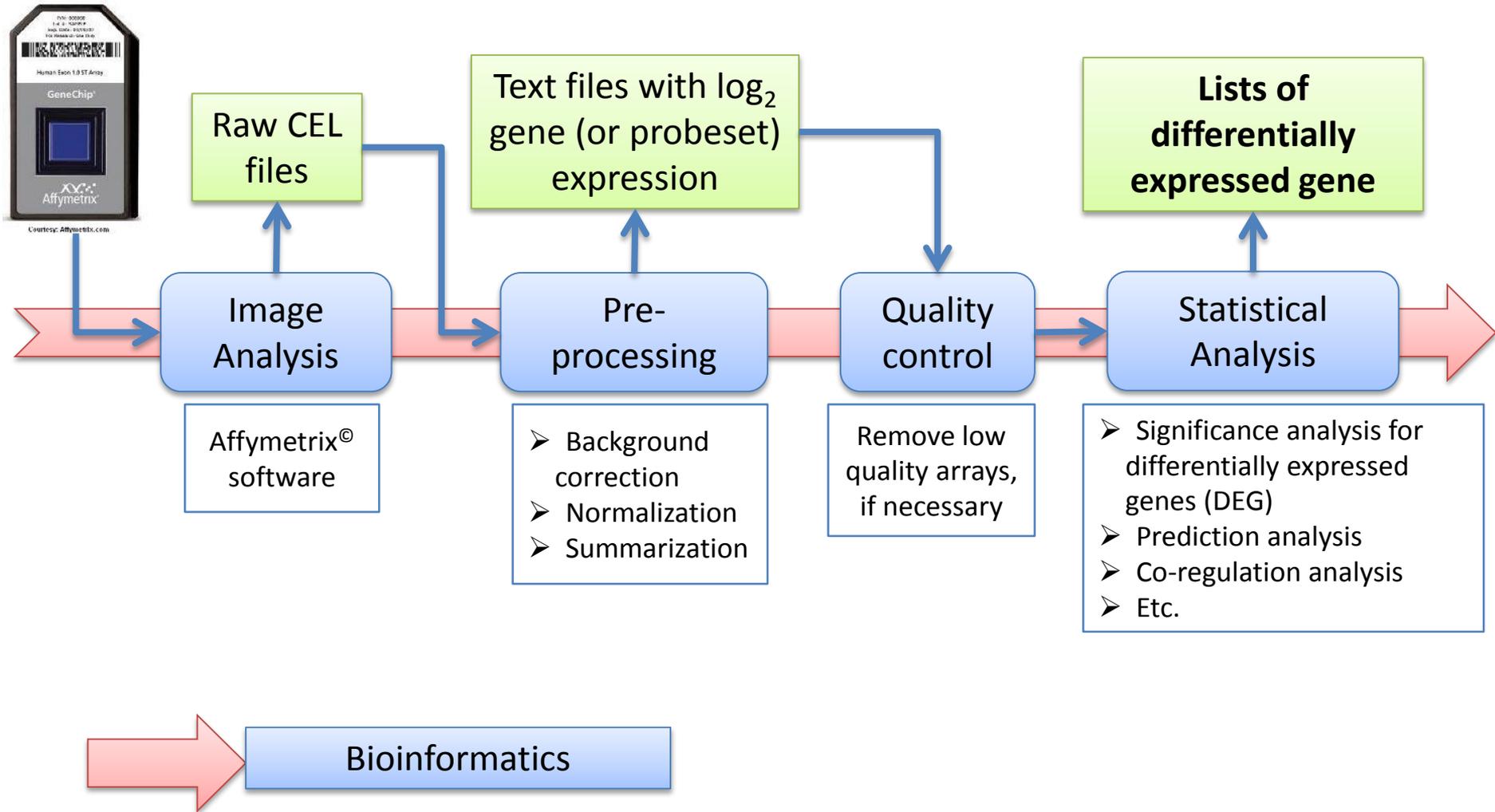
i -- array
j -- probe
n -- probeset

“Median polish” helps avoid outliers effect

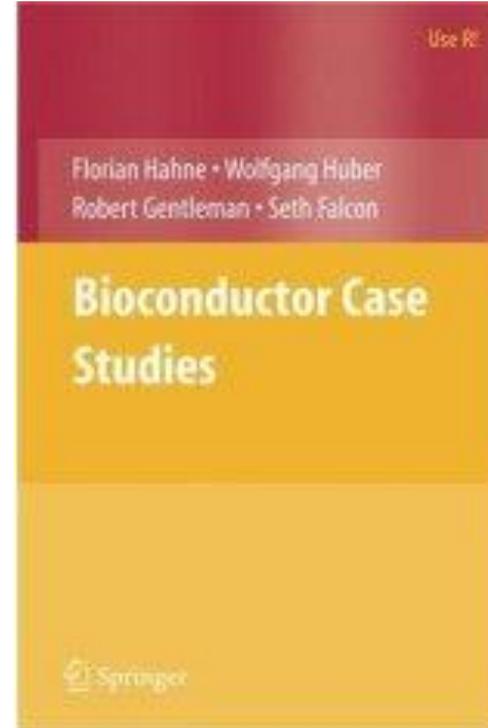
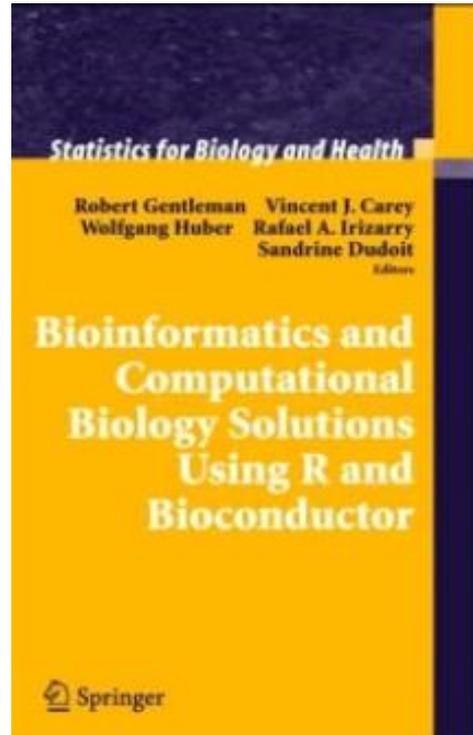
File Formats



Analysis Pipeline



R / Bioconductor



Import: Affymetrix Power Tools

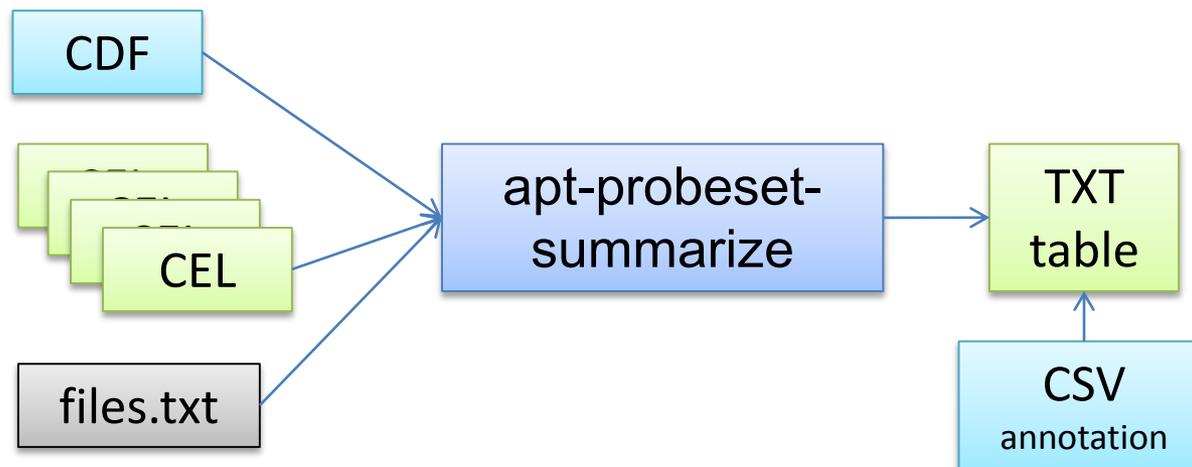
apt-probeset-summarize is a program for doing background subtraction, normalization and summarizing probe sets from Affymetrix expression microarrays. It implements analysis algorithms such as [RMA](#), [Plier](#), and DABG (detected above background).

The main features of **apt-probeset-summarize** not common in other implementations are: Quantile normalization using a subset (sketch) of the data which results in much smaller memory usage.

<http://www.affymetrix.com/support/developer/powertools/changelog/apt-probeset-summarize.html>

apt-probeset-summarize

```
-a rma-sketch -d chip.cdf -o output-dir --cel-files files.txt
```



Import: R

oligo is a Bioconductor package for preprocessing oligonucleotide microarrays. It currently supports chips produced by Affymetrix and NimbleGen and uses les provided by these manufacturers in their native format. The package provides a united framework for preprocessing and uses the data representation established by the Bioconductor project, which simplifies the interface with other packages in the project.

<https://www.bioconductor.org/packages/release/bioc/html/oligo.html>

```
library(oligo)
```

```
celFiles = list.celfiles(full.names=TRUE)
```

```
## import
```

```
rawData = read.celfiles(celFiles)
```

```
## Bg correction, normalization and summarization
```

```
normData = rma(rawData)
```

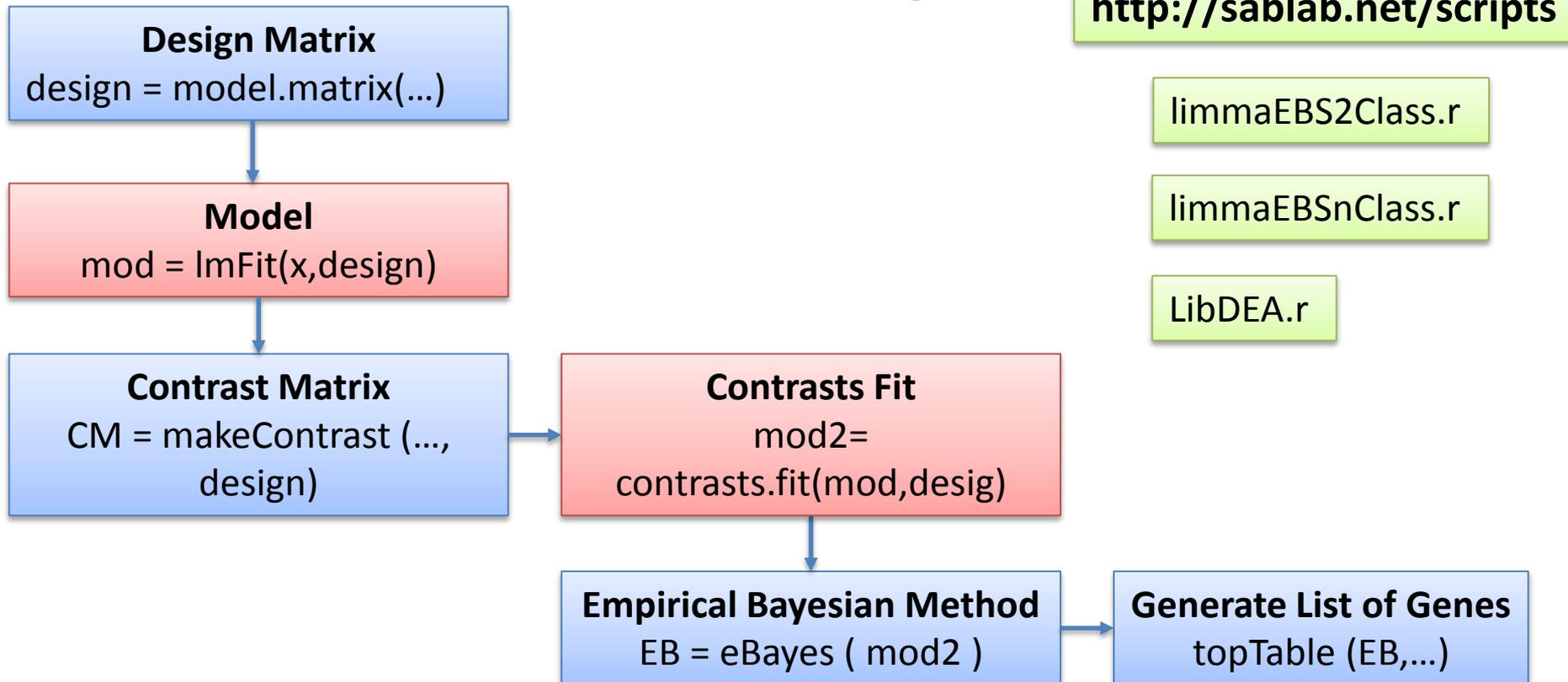
```
## get expression matrix
```

```
X = exprs(eset)
```

<http://edu.sablab.net/data/gz/>

Differential Expression Analysis

Factorial design



Please go through the code at:

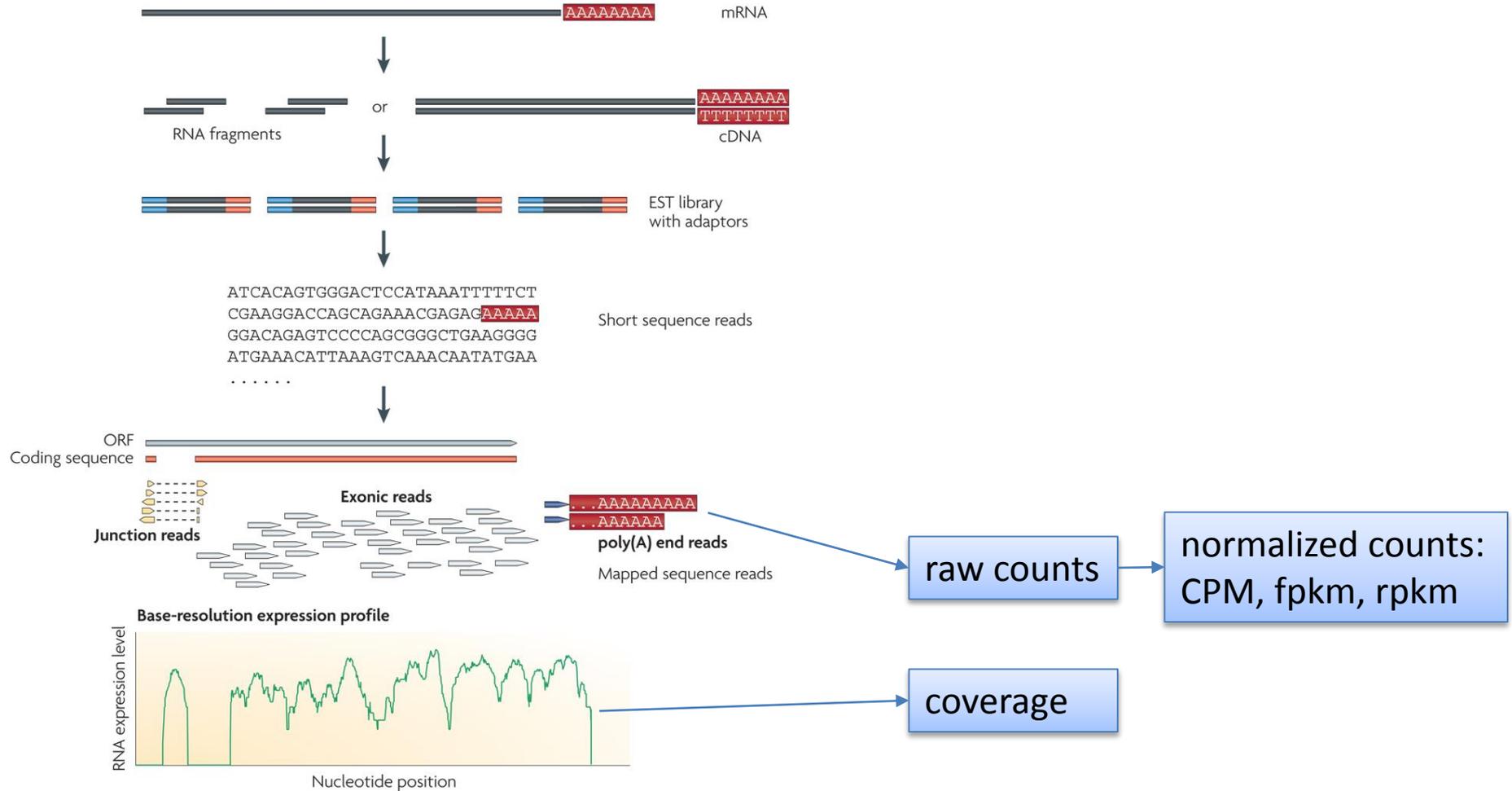
<http://edu.sablab.net/abs2017/scripts6.html>

Section 6.1

Do Exercises 6.1

RNA-seq

Next Generation Sequencing: RNA-Seq



Wang Z et al. RNA-Seq: a revolutionary tool for transcriptomics. **Nat Rev Genet.** 2009

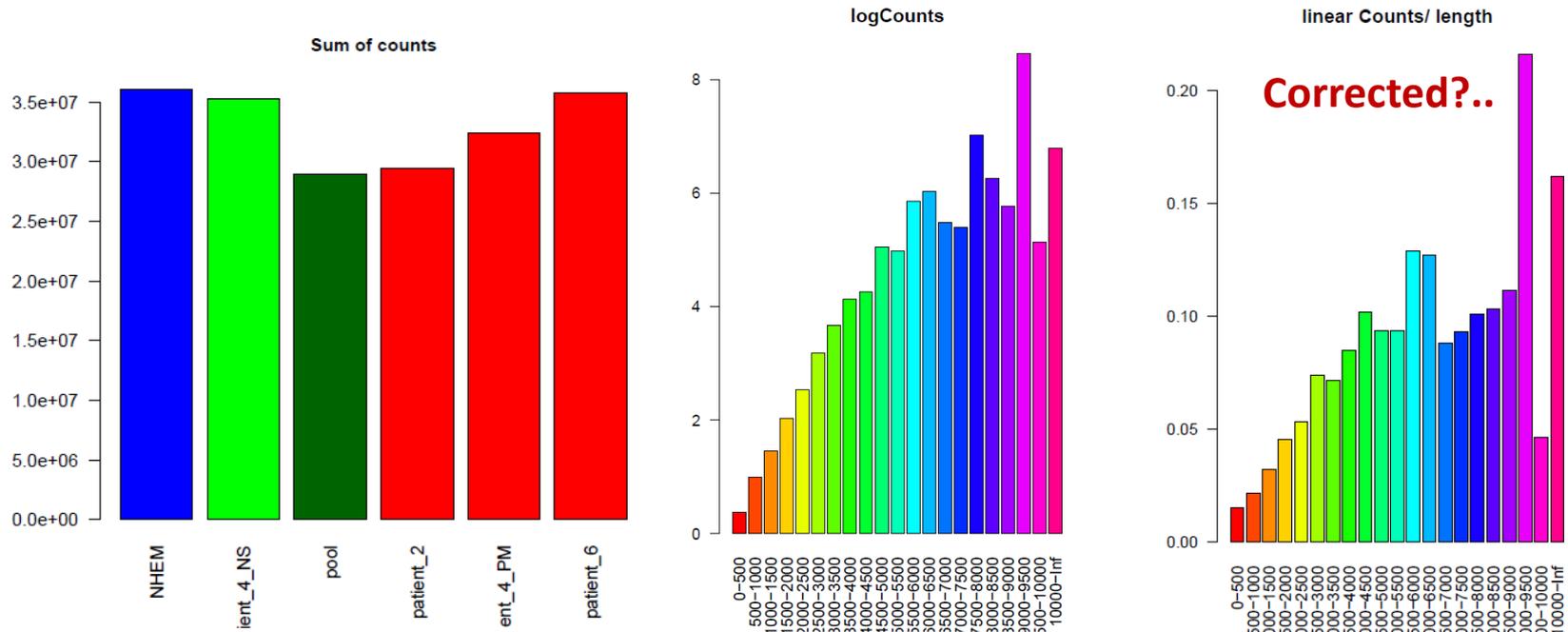
Normalization

Problems:

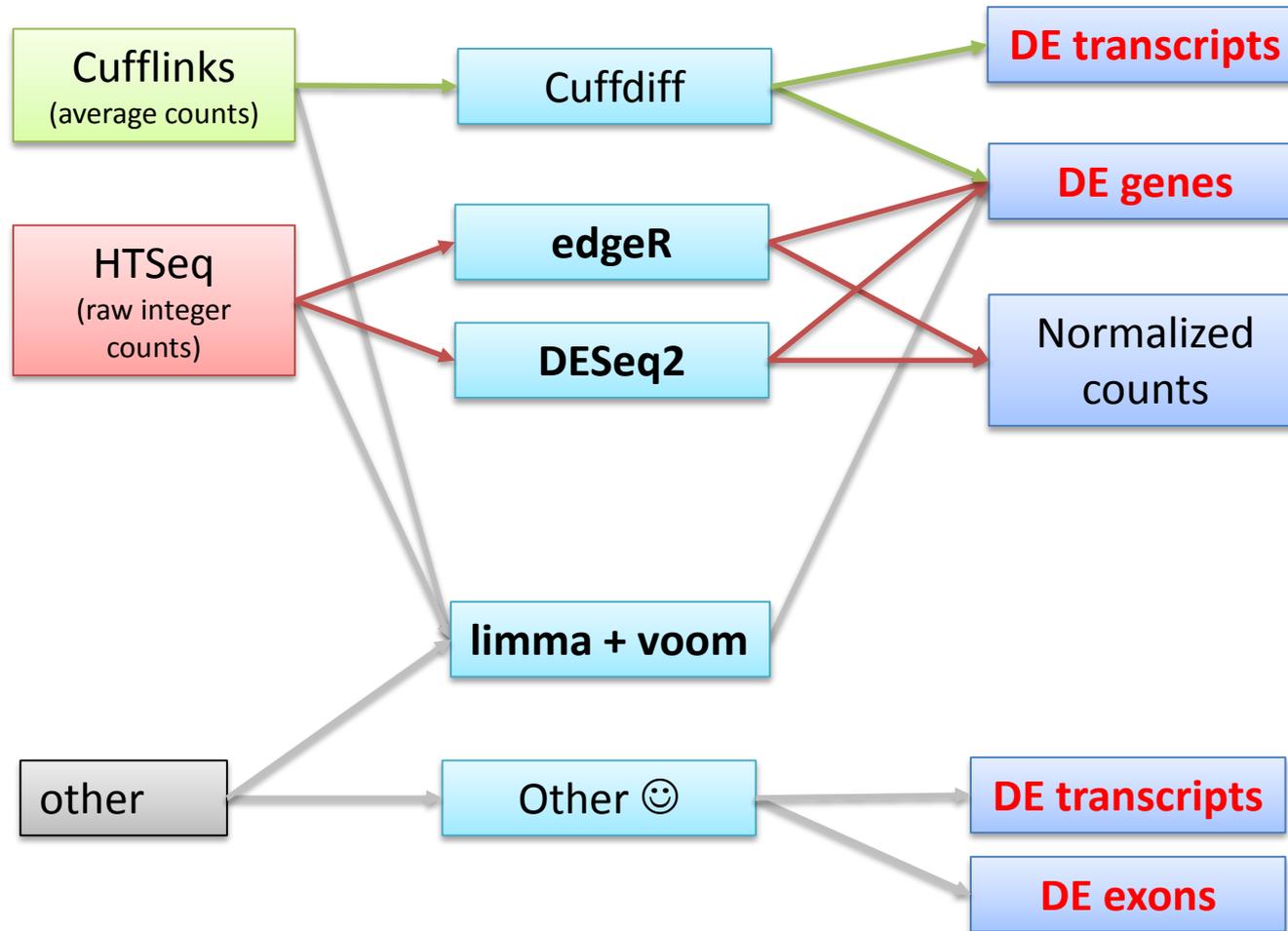
- ◆ Libraries has different size (different number of reads from samples)
- ◆ Long transcripts produce more reads

Solutions (?) :

- ◆ Accounting for library size during analysis (standard) or direct correction for it
- ◆ Correction for transcript size (but which transcript is expressed?)



Differential Expression Analysis



Differential Expression Analysis (edgeR)

<https://sablabs.net/scripts>

LibDEA.r

Differential Expression Analysis (DESeq2)

<https://sablabs.net/scripts>

LibDEA.r

Enrichment Analysis

1. Category Enrichment Analysis

Fisher's exact test: based on hypergeometrical distributions

Hypergeometrical: distribution of objects taken from a "box", without putting them back

$$P = 1 - \sum_{i=0}^{k-1} \frac{\binom{M}{i} \binom{N-M}{n-i}}{\binom{N}{n}}$$

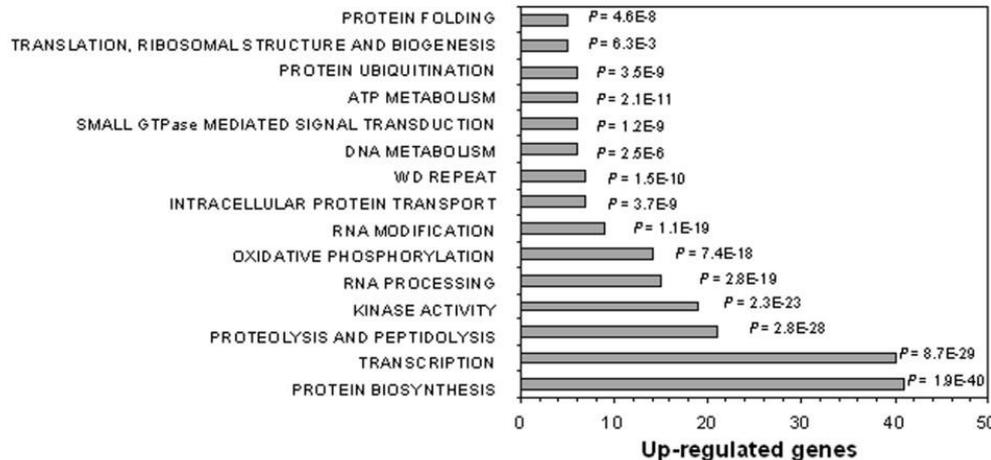
N: total number of genes

M: total number of genes annotated with this term

n: number of genes in the list

k: number of genes in the list annotated with this term

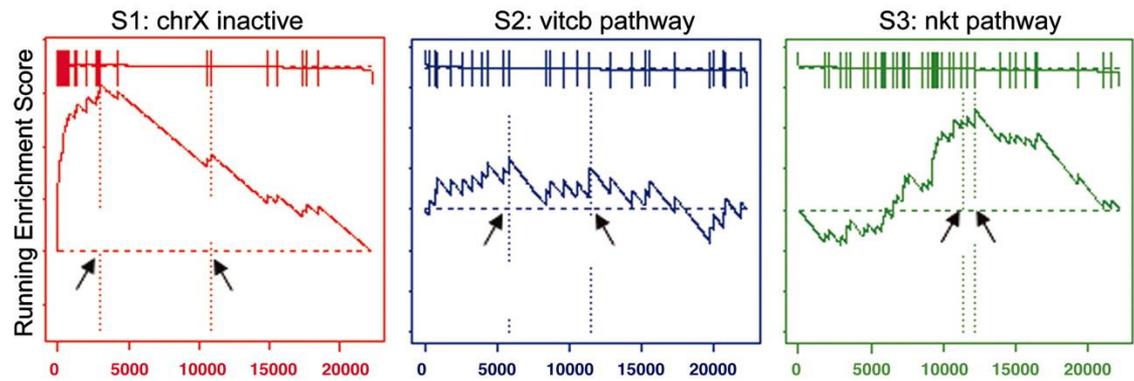
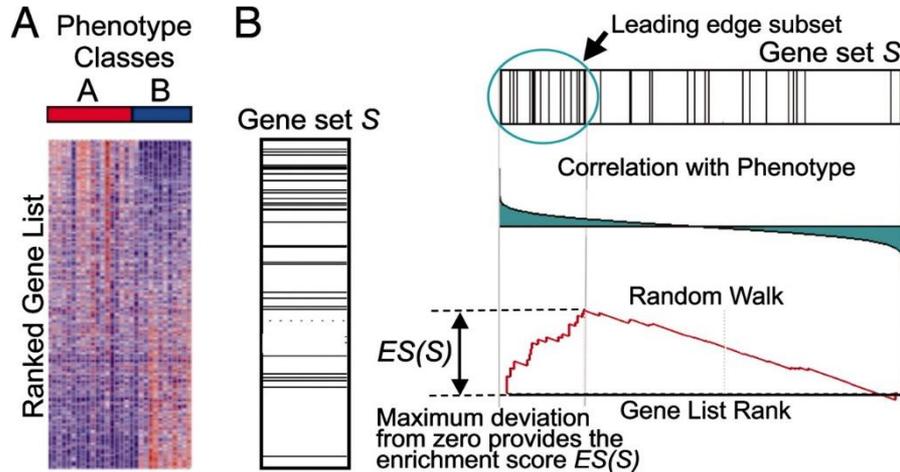
$$C_k^n = C_n^k = \binom{n}{k} = \frac{n!}{k!(n-k)!}$$



Okamoto et al. Cancer Cell International 2007 7:11 doi:10.1186/1475-2867-7-11

2. Gene Set Enrichment Analysis (GSEA)

Is direction of genes in a category random?



A. Subramanian et al. PNAS 2005,102,43

Example: GO enrichment

<http://edu.sablab.net/transcript>

Strategy 1:

Take all DEG and use them in enrichment.

- Safe
- No additional assumptions
- Cannot distinguish \uparrow and \downarrow functions

Strategy 2:

Separate DEG to down- and up- regulated genes. Then perform independent enrichment by these 2 groups

- Can be biased (gene can be $\uparrow\downarrow$)
- Assume \uparrow gene \Rightarrow \uparrow function
- Can distinguish \uparrow and \downarrow functions

Enrichr

<http://amp.pharm.mssm.edu/Enrichr/enrich>

BioCompendium

<http://biocompendium.embl.de/>

LUSC Example

<http://edu.sablab.net/data/txt/lusc.zip>

<http://amp.pharm.mssm.edu/Enrichr/>

0. Prepare lists of DE genes...

1. Put up-regulated into **enrich**

3. Check: Down CMAP, Disease Signatures from GEO up,

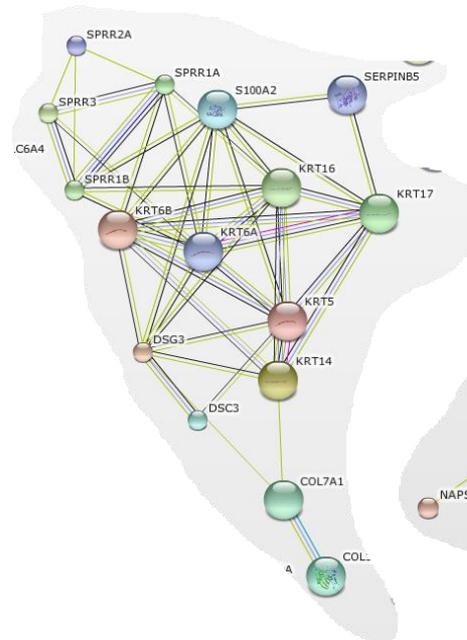
4. Try **biocompendium**

5. Put top 100 genes into String to see PP-interactions

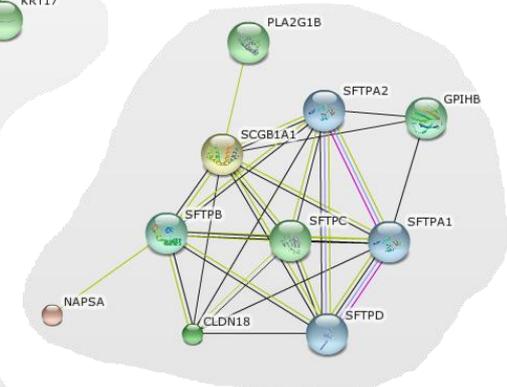
<http://biocompendium.embl.de/>

<http://string-db.org>

Up regulated



Down regulated



In R

```
#####
## enrichGOens - warpup for topGO package: enrichment analysis of GO-terms
## based on Ensembl IDs
#####
## genes - vector with list of ENSEMBL IDs (character)
## fdr - vector of FDR for each gene (numeric)
## fc - vector of logFC for each gene (numeric)
## thr.fdr - significance threshold for FDR (numeric)
## thr.fc - significance threshold for absolute logFC (numeric)
## db - name of GO database: "BP", "MF", "CC" (character)
## genome - R-package for genome annotation used. For human - 'org.Hs.eg.db' (character)
## do.sort - if TRUE - resulted functions sorted by p-value (logical)
## randomFraction - for testing only, the fraction of the genes to be randomized (numeric)
#####
## (c)GNU GPL P.Nazarov 2014. petr.nazarov[at]crp-sante.lu
#####
```

<https://sablabs.net/scripts>

enrichGOens.r

```
enrichGOens =
function (genes, fdr, fc, thr.fdr=0.05, thr.fc=0, db="BP", genome="org.Hs.eg.db", do.sort=TRUE,
        randomFraction=0) {
  ## load libraries
  if (!require(genome, character.only=TRUE)){
    cat("MESSAGE enrichGO: '", genome, "' package is not found. Installing...\n", sep="")
    source("http://bioconductor.org/biocLite.R")
    biocLite(genome)
    library(genome, character.only=TRUE)
  }
  if (!require("topGO")){
    cat("MESSAGE enrichGO: ' topGO ' package is not found. Installing...\n")
    source("http://bioconductor.org/biocLite.R")
    biocLite("topGO")
    library("topGO")
  }
  if (!exists("sortDataFrame")) source("http://sablabs.net/scripts/sortDataFrame.r")
  ## prepare gene categories and score
  myGO2genes <- annFUN.org(db, mapping = "org.Hs.eg.db", ID = "ensembl")
  score = (-log10(fdr)*abs(fc))
  names(score)=genes
  score[fdr>thr.fdr | abs(fc)<=thr.fc]=0

  ## add randomness if required, to test stability
  if (randomFraction>0){
    ## define remove probability: low score have more chances
    prob1 = 1/(1+score)
    prob1[is.na(prob1)]=0
    prob1[score == 0] = 0
    ## define add probability: high score has more chances
    prob2 = -log10(fdr)*abs(fc)
    prob2[is.na(prob2)]=0
    prob2[score > 0] = 0
    ## add and remove
    n=round(sum(score>0)*randomFraction)
    score[sample(1:length(genes),n,prob=prob2)]=1+rexp(n,1/mean(score[score>0]))
    score[sample(1:length(genes),n,prob=prob1)]=0
  }
  ## create topGOdata object
  SelectScore = function(sc){return(sc>0)} ## simple function for significance
  GOdata = new("topGOdata", ##constructor
    ontology = db,
    allGenes = score,
    geneSelectionFun = SelectScore,
    annot = annFUN.GO2genes,
    GO2genes = myGO2genes)

  ## run testing
  resFisher = runTest(GOdata, algorithm = "classic", statistic = "fisher")
  ## transform results into a table
  enrichRes = GenTable(GOdata, classicFisher = resFisher,
    ranksOf = "classicFisher", topNodes = length(resFisher$score))
  enrichRes$classicFisher[grep("<", enrichRes$classicFisher)] = "1e-31"
  enrichRes$classicFisher = as.double(enrichRes$classicFisher)
  enrichRes$FDR = p.adjust(enrichRes$classicFisher, "fdr")
  enrichRes$Score = -log10(enrichRes$FDR)
  ## by default sorted by p-value. If needed - sort by ID
  if (!do.sort) enrichRes = sortDataFrame(enrichRes, "GO.ID") ## remove sorting
  return(enrichRes)
}
```

Thank you for your attention !

