

# RNASEQ: WHAT TO DO?

**Guest Lecture: Juli Petereit, PhD  
Bioinformatician and Data Scientist**

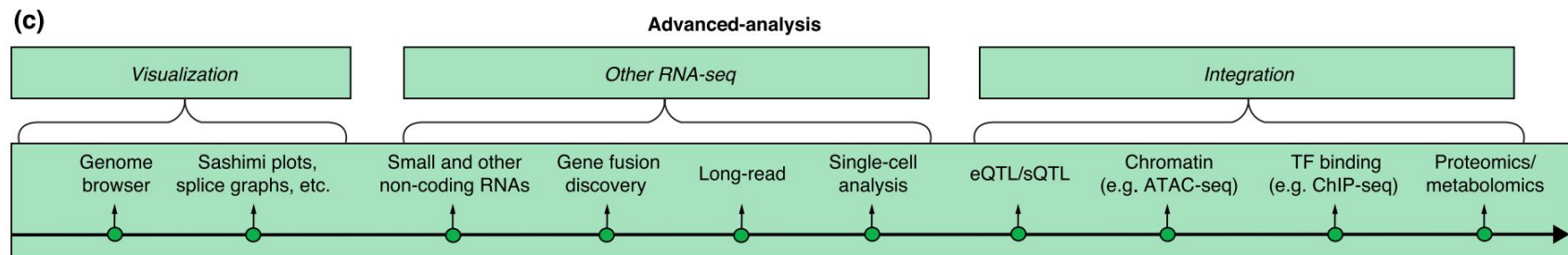
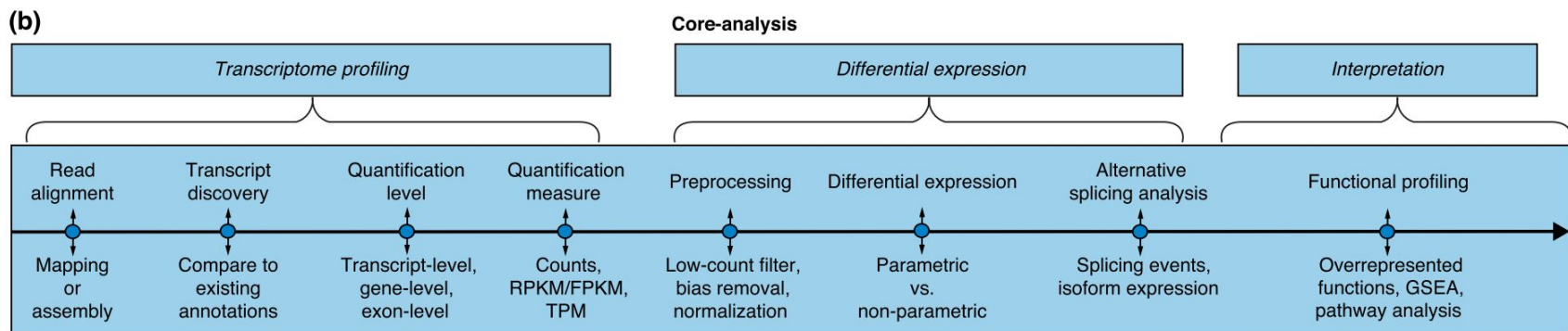
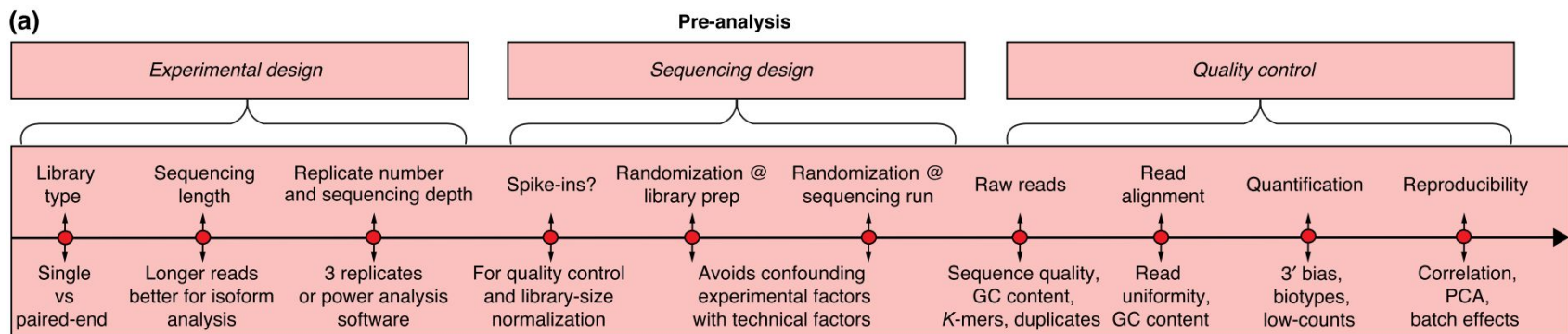
# WHAT IS RNASEQ?

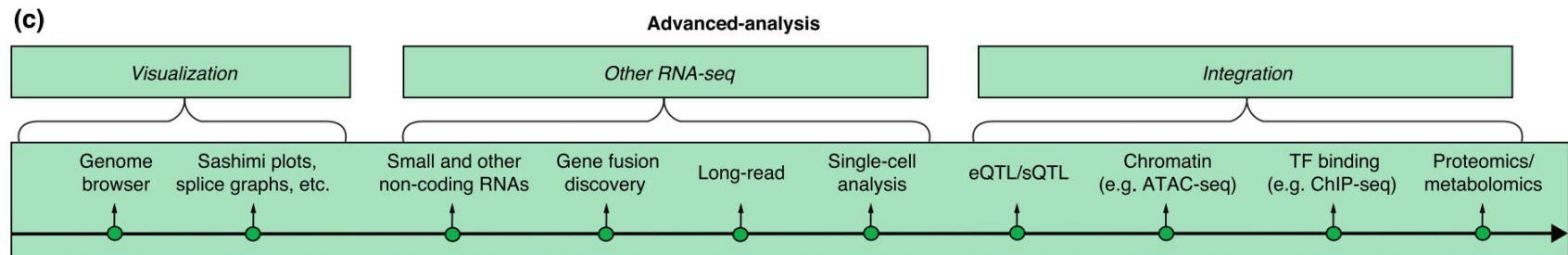
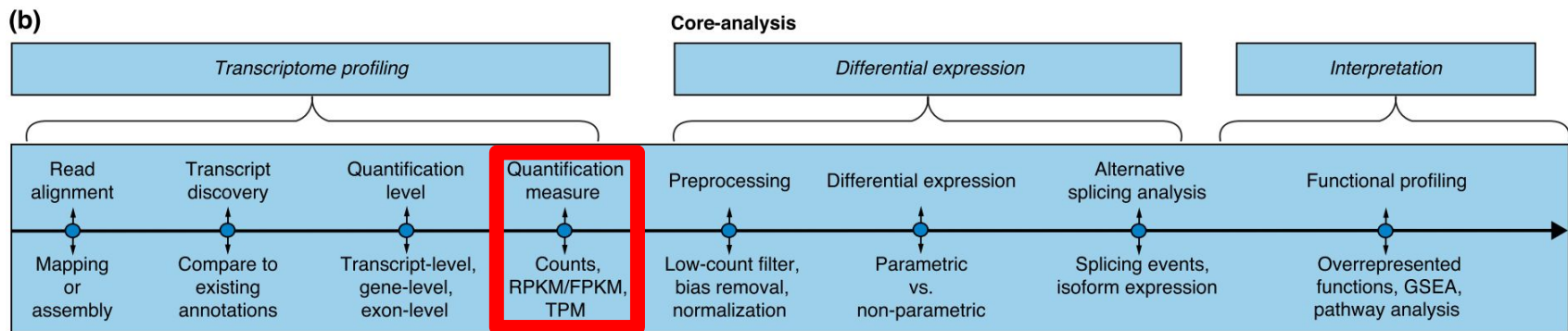
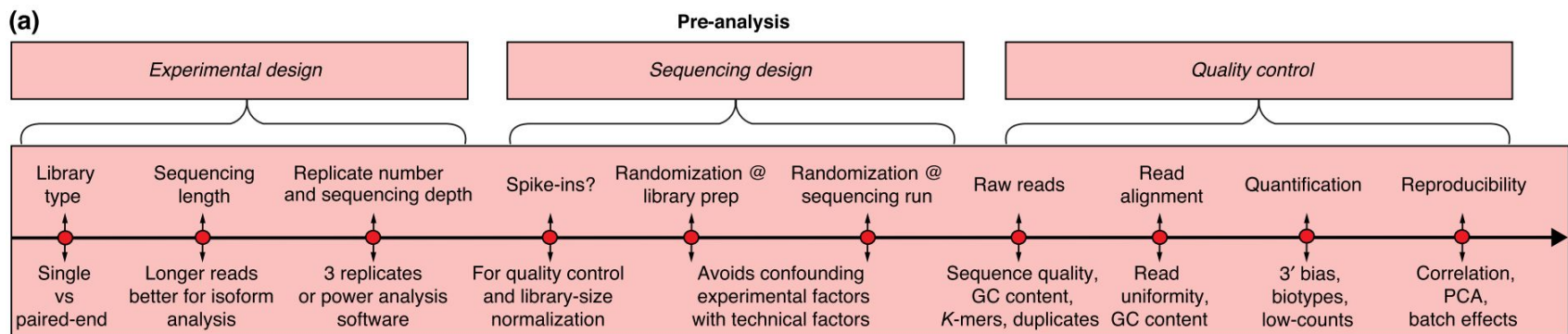
“RNA sequencing (RNA-Seq, RNAseq), also called whole transcriptome shotgun sequencing, uses next-generation sequencing (NGS) to reveal the presence and quantity of RNA in a biological sample at a given moment in time.” <https://en.wikipedia.org/wiki/RNA-Seq>

Most common application of RNAseq is to quantify gene and transcript expression.

So how do we go about it?







# RNASEQ, RNA-SEQ, NEXTGEN - WHICH ONE IS CORRECT?

Research community does not even agree on a common abbreviation: RNAseq, RNA-seq, RNA-Seq.

Thus it is not surprising that there are many variations of RNAseq protocols and analyses as well.

This is a challenge for new and newer users to appreciate all of the steps necessary to properly conduct an RNA-seq study.

I recommend to consult with as many people as possible, especially your bioinformatician, and don't try to tackle the problem alone.

# SO WHAT DO I HAVE TO DO? WHICH ANALYSIS PROCEDURE?

There is not one golden pipeline, -rule or -standard.

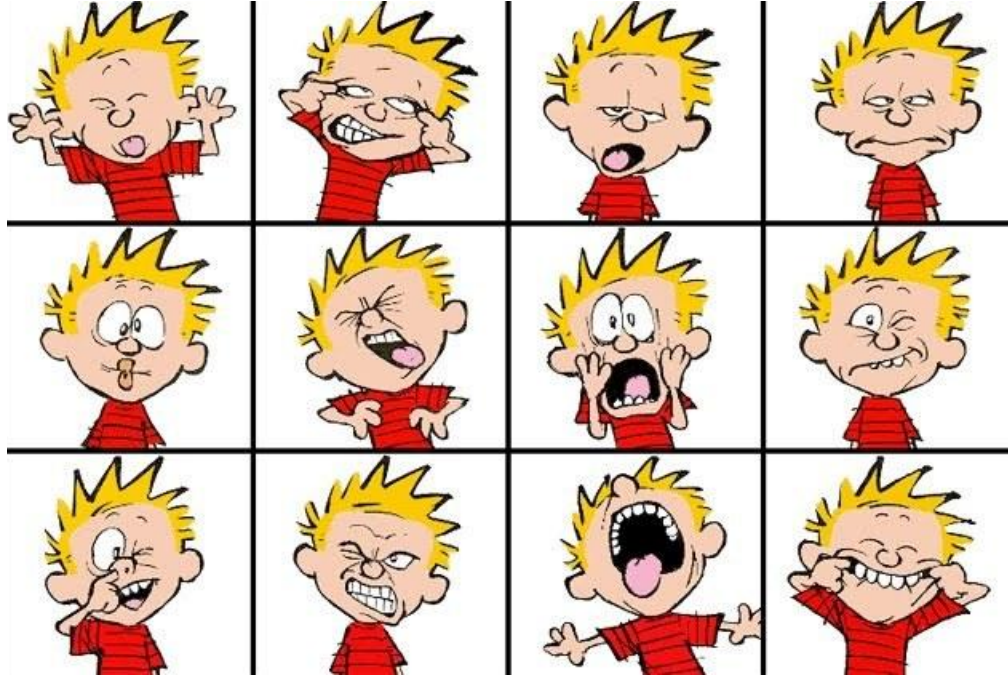
RNAseq analysis must be adjusted according to the scientific question/objective and organism being studied.

Individual experimental scenarios potentially require adjusted/different methods for transcript quantification, normalization, and other downstream analyses (e.g., differential gene expression)

What is the appropriate experimental setup and analysis strategy?

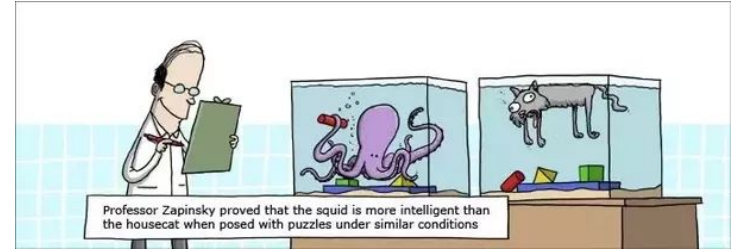


THIS IS HOW I FEEL MOST DAYS =)



# EXPERIMENTAL DESIGN

Define a good experimental design



- Define a precise research question.
- What is the underlying hypothesis?
- Which library type should be utilized?
- What sequencing depth and number of replicates are appropriate for the biological system under study?
- Does my protocol allow for adequate execution of the sequencing experiment itself?
- Are there any steps that possibly introduce unnecessary biases?



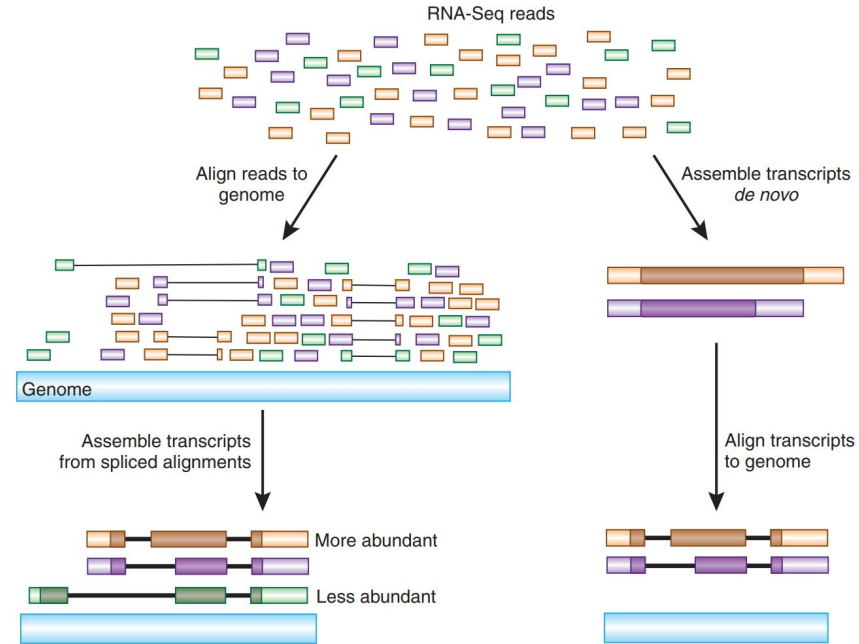
# ASSEMBLY VERSUS ALIGNMENT

For non-model organism without reference genome

- 1) Assembling reads de novo into contigs
- 2) Mapping these contigs onto the transcriptome

For model organism

Identify transcripts by mapping reads onto the genome



Nature Biotech (2010) 28, 421–423

# REPLICATES

Often the number of replicates is simply governed by availability of samples.

The number of samples should depend on:

- technical variability in the RNA-seq procedures
- biological variability of the system/organism under study
- desired statistical power (power analysis)

# TYPES OF SEQUENCING

Sequencing option depends on analysis goals

**Single-end (SE) reads:** cheaper, often sufficient for well-annotated organisms

**Paired-end (PE) reads:** preferable for de novo transcript discovery or isoform expression analysis

**Short reads:** could be good enough

**Long reads:** improve mappability and transcript identification, preferable for poorly annotated transcriptomes, faster alignment

For more: <http://www.anthonymbaldor.com/sequencing-run-types/>

# SEQUENCERS

Each technology has its own proprietary software, preparation kits, read lengths, quirks, etc. and therefore require different analyses procedures.

Illumina



PacBio



Oxford Nanopore



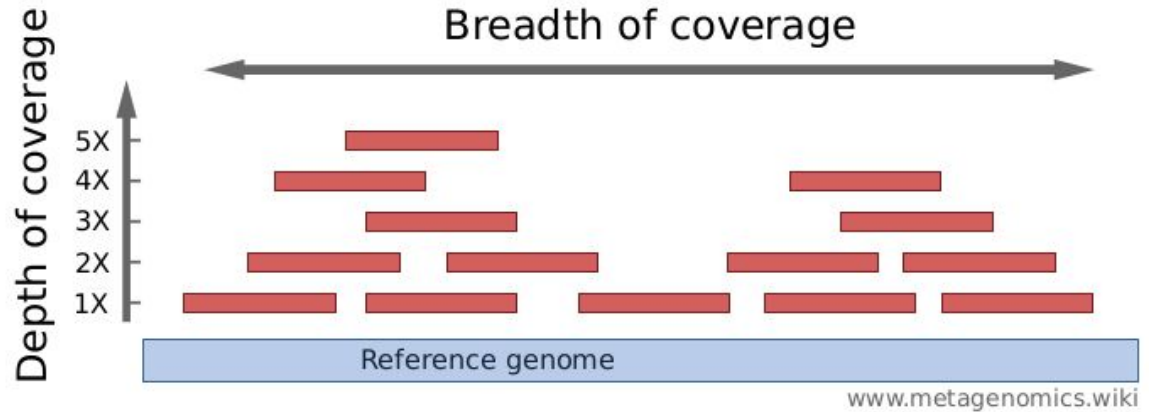
Ion Torrent



# SEQUENCING COVERAGE

Sequencing depth or library size is the number of sequenced reads for a given sample.

Deep sequencing improves quantification and identification, but not always is more necessarily better.



For more: Sims, D., *et al.* "Sequencing depth and coverage: key considerations in genomic analyses", Nature Reviews - Genetics(2014)

# QUALITY CONTROL, FOLLOWED BY QUALITY CONTROL

To present a high quality outcome, the results have to be reproducible and reliable.

Apply quality control checks at different stages of the analysis to ensure both:

- Raw reads
- Read Alignment
- Quantification
- Reproducibility

# ANALYSIS OF RNASEQ DATA

RNA-seq data undergoes multiple steps and thus several quality-control checkpoints need to be considered.

Raw read -> read alignment -> quantification -> counts -> expression measures

At each of these steps, specific checks should be applied to monitor the quality of the data.



# READ ALIGNMENT



An important mapping quality parameter is the percentage of mapped reads, which is a global indicator of the overall sequencing accuracy.

Other important parameters are the uniformity of read coverage on exons and the mapped strand.



# QUANTIFICATION CHECKS

- GC content
- Gene length biases (relevant for normalization)
- Biotype composition (quality of the RNA purification step)
- Sample similarity and variance

R packages (e.g., **NOISeq**, **EDA-Seq**) provide useful plots for quality control of count data

**HTSeq-count** or **featureCounts** combine raw counts of mapped reads and quantify counts at the gene level.

# CONTROL ON ENTIRE DATASET

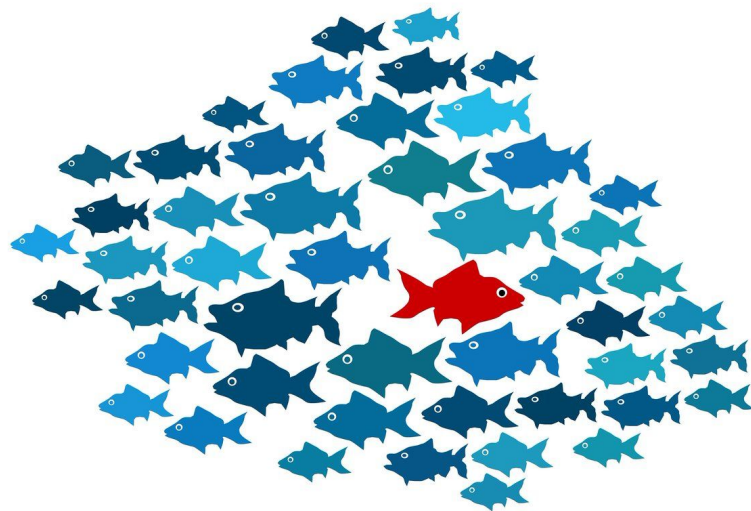
Check for global quality

Are the (technical and biological) replicates reproducible?

Is the distance relatively small and correlation relatively high among replicates?

Are there outlying replicates? Why are they outliers?

Is there a separation between experimental samples?

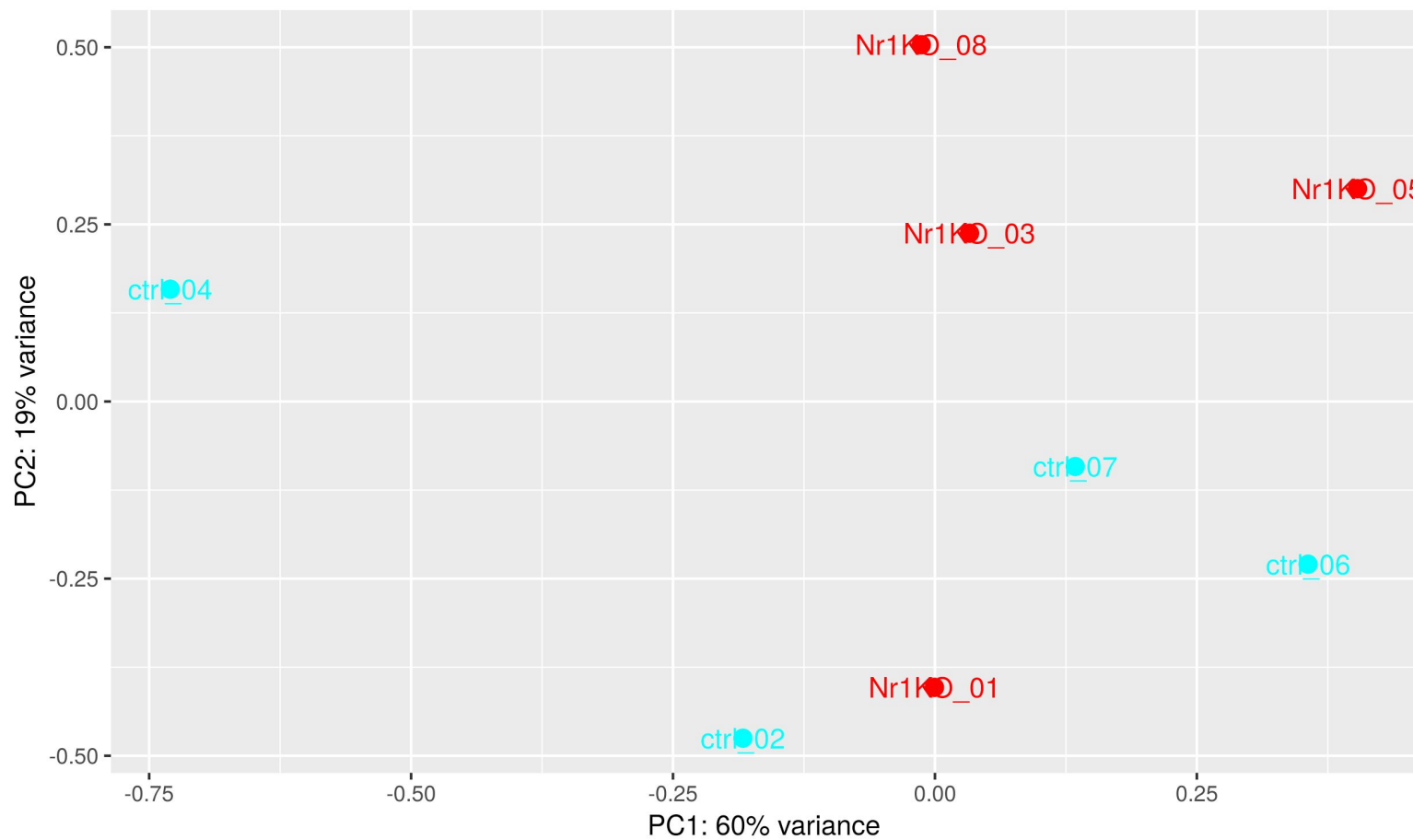


# CONTROL ON ENTIRE DATASET

There is no clear standard to control for variation with replicates. The amount of variation depends on the heterogeneity of the experimental system, organism, and experimental design.

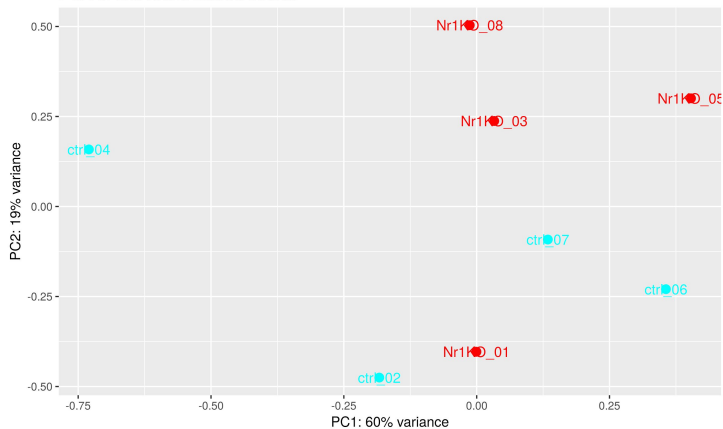
Principal component analysis (PCA) is a statistical procedure to investigate the relationships among samples.

PCA of first round filtered counts

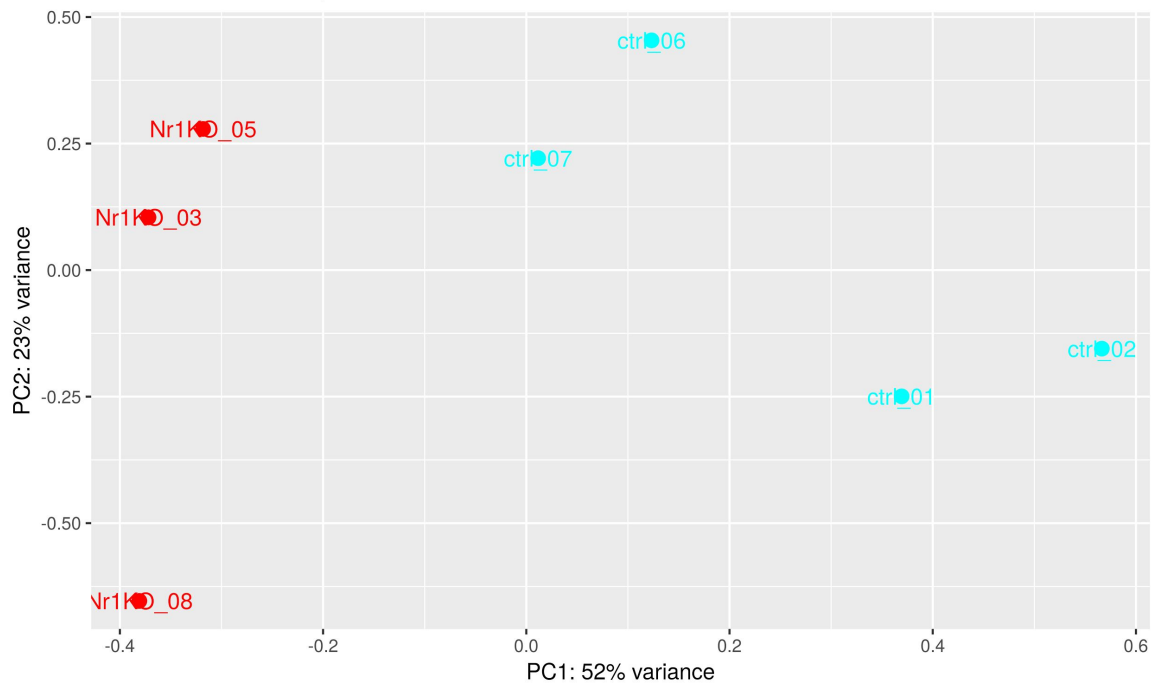


# HOURS LATER ...

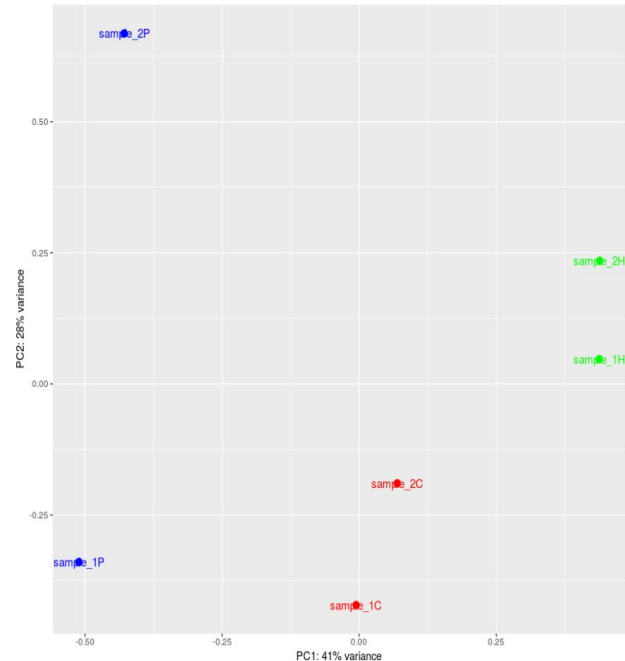
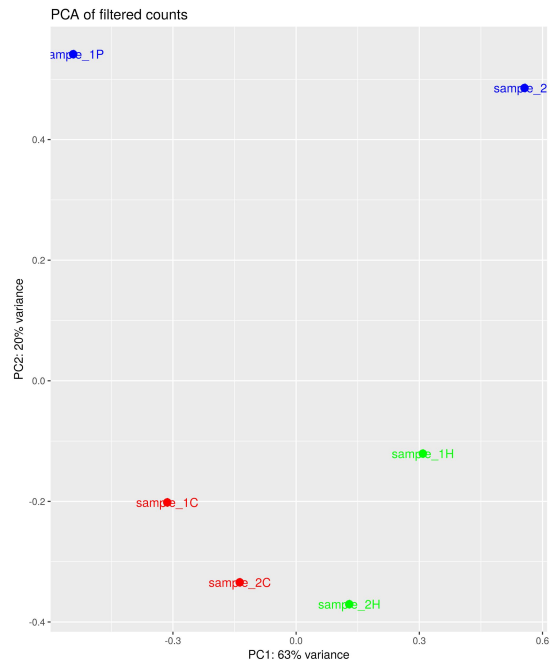
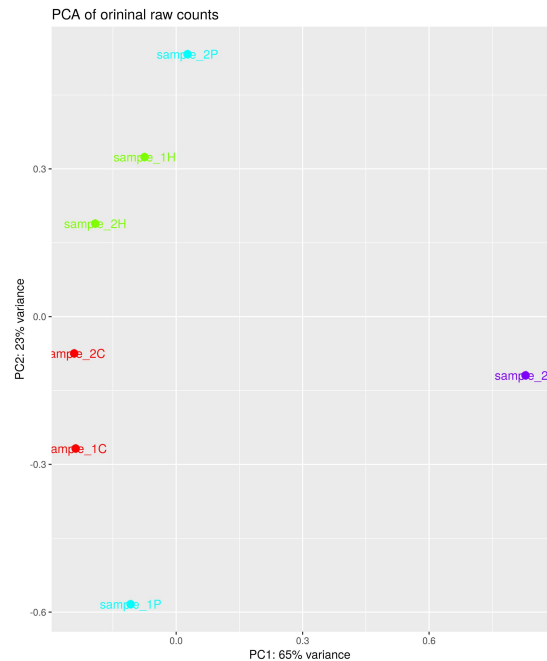
PCA of first round filtered counts



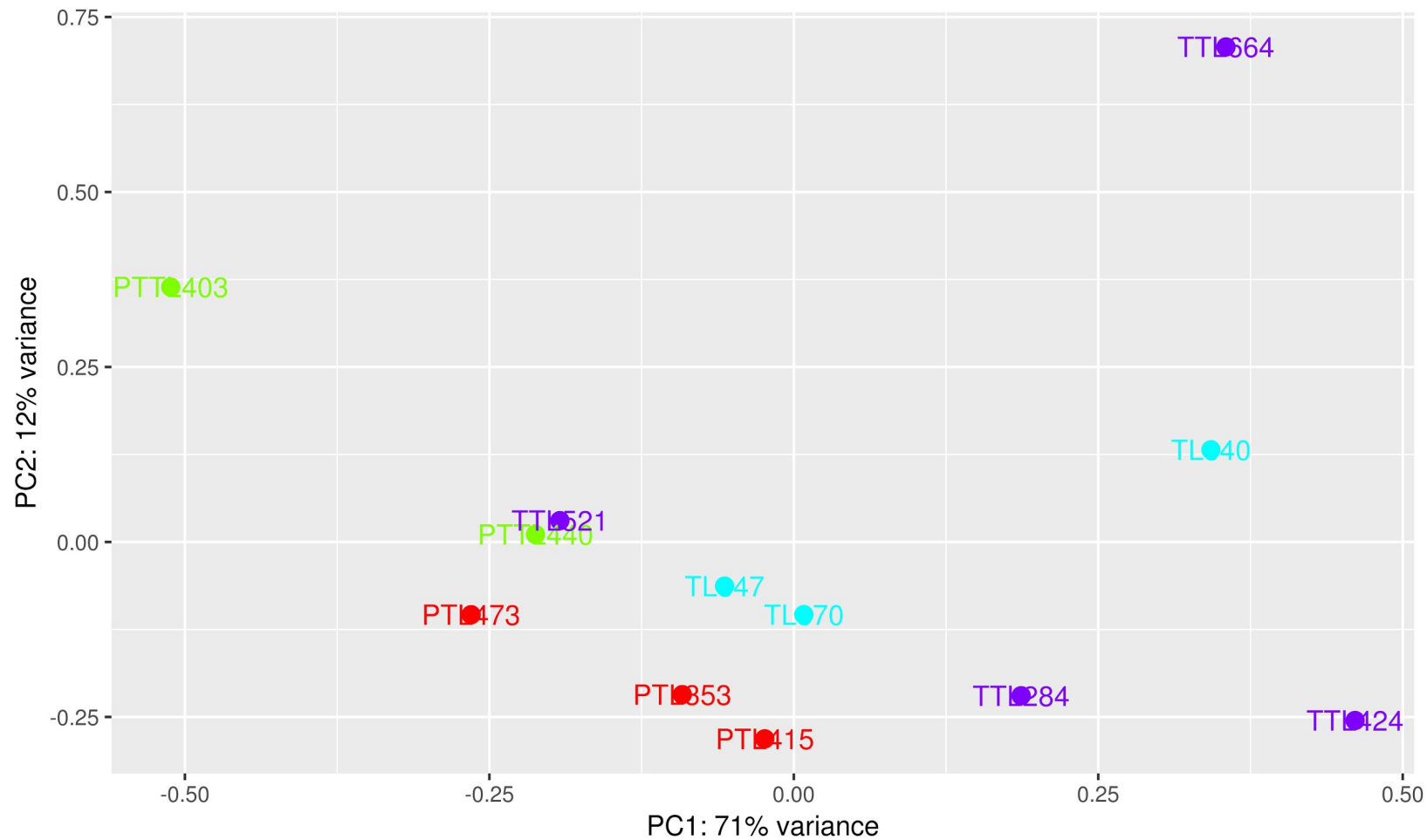
PCA of third filtered, normalized mod feature counts



# PREPROCESSING



# Normalized Counts

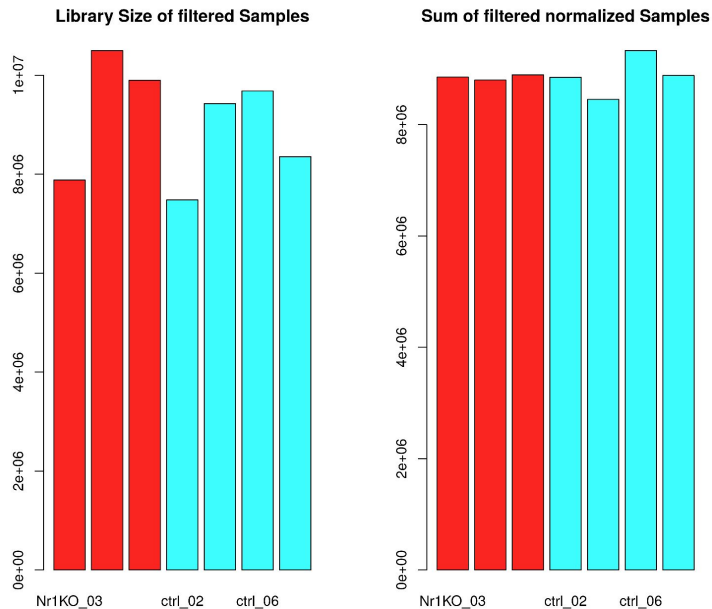


# RNASEQ DATA NORMALIZATION

Many different normalization methods: RPKM, FPKM, TPM, CPM, TMM, ...

TPMs (Transcripts Per Kilobase Million) are normalized for gene length and library size and denotes the proportion of counts within that sample of a gene.

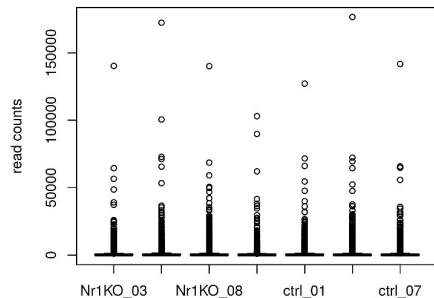
TPM values are used for clustering, heatmaps, and manuscripts, but NOT for differential gene expression analysis.



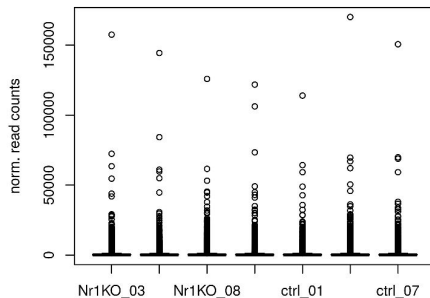


# RNASEQ DATA NORMALIZATION

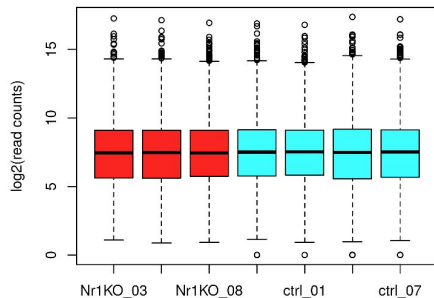
filtered read counts



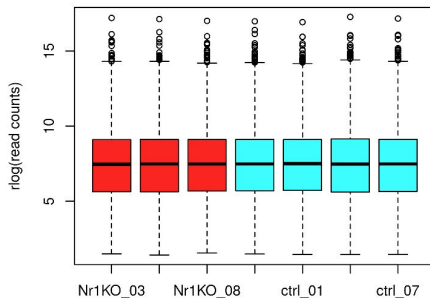
untransformed, normalized read counts



log2-transformed, normalized read counts



rlog-transformed, normalized read counts



Not all variation is bad.

Be careful transforming data,  
as biological variation might  
be massaged away!!

# YAY WE ARE FINALLY READY TO ANALYZE THE DATA

The most common application of RNAseq is to estimate gene and transcript expression and to quantify differential expression between tissues/conditions at the genome level.

Downstream analyses are somewhat similar to microarray data, but don't forget about the data distribution.

Microarray data often transform into a normal distribution, whereas RNAseq has a negative binomial distribution.

# DIFFERENTIAL GENE EXPRESSION ANALYSIS

The most popular R packages are: **edgeR** and **DESeq2**

Both applications require the raw filtered counts as their input.

Correcting for gene length is not necessary when comparing changes in gene expression within the same gene across samples.

Recommendation: At least three biological replicates.

The **NOISeq** R package includes diagnostic plots to identify origins of biases in RNA-seq data.

# WHOLE SYSTEM NETWORKS

Very current topic in life sciences.

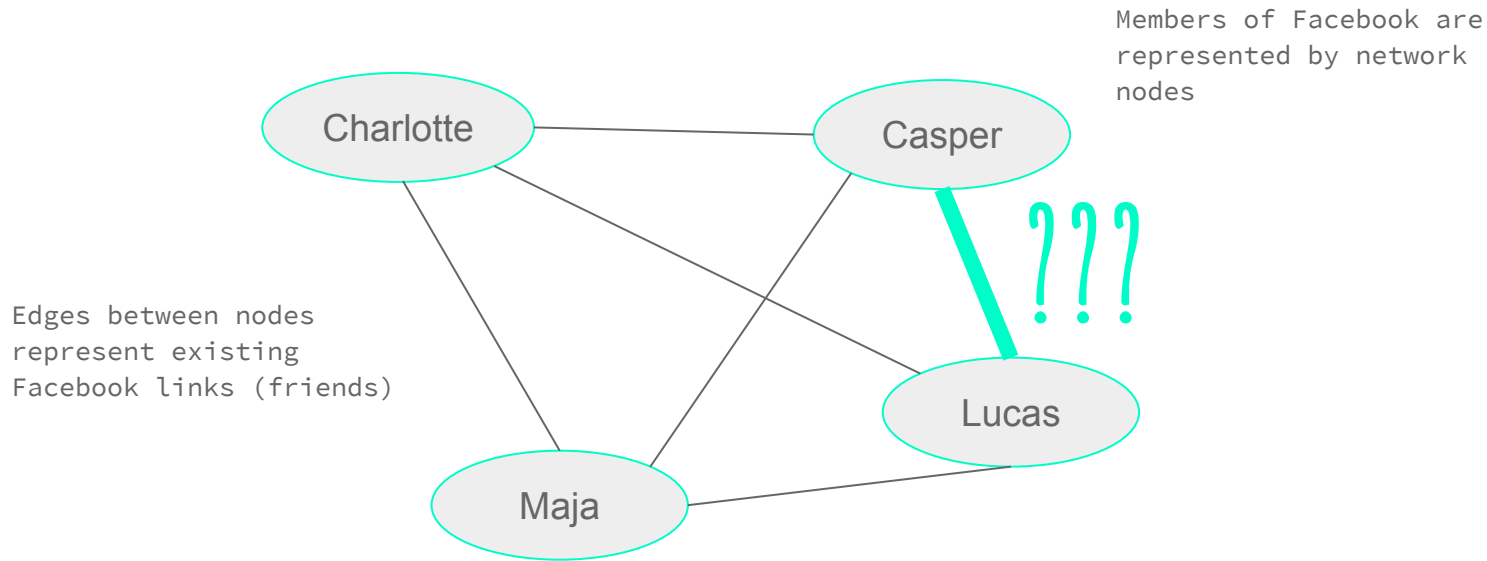
Network analysis is a versatile tool and can present a variety of interactions and organizations at a systems level.

The study of node-edge graphs:

- Node = elements in the system
- Edges = a type of relationship/association/interaction between elements

# CONCEPTUAL EXAMPLE

How does Facebook recommend people you might know?



# BIOLOGICAL NETWORKS - A NODE-EDGE GRAPH



Nodes are elements in the system, molecular entities:

- genes, DNA, RNA, peptides, proteins, metabolites, transcription factors, other biomolecules

Edges are relationships/associations between elements:

- interaction, regulation, transformation, activation, similar in expression behavior or value, sequence similarity

Can reveal organizational principles and structure of biological systems at the cellular level.

# NETWORKS IN BIOLOGY

1970, Kauffman was the first to use Boolean networks as a biological network modeling paradigm.

Use in biology has increased exponentially.

Established role in research: Systems Biology, Network Biology, Network Medicine.

Have been used for structural and functional studies of genes, proteins, and metabolites.

Network model construction is specific to research question and dataset.

Figure 2 : Yeast protein interaction network.

From: Network biology: understanding the cell's functional organization

# BIOLOGICAL NETWORK MODELS

Protein-Protein Interaction (PPI) Networks

Co-expression Networks

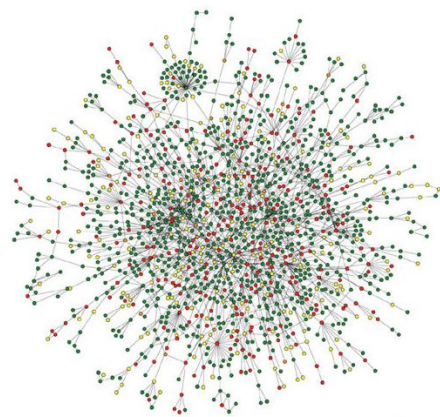
Metabolic Networks

Regulatory Networks

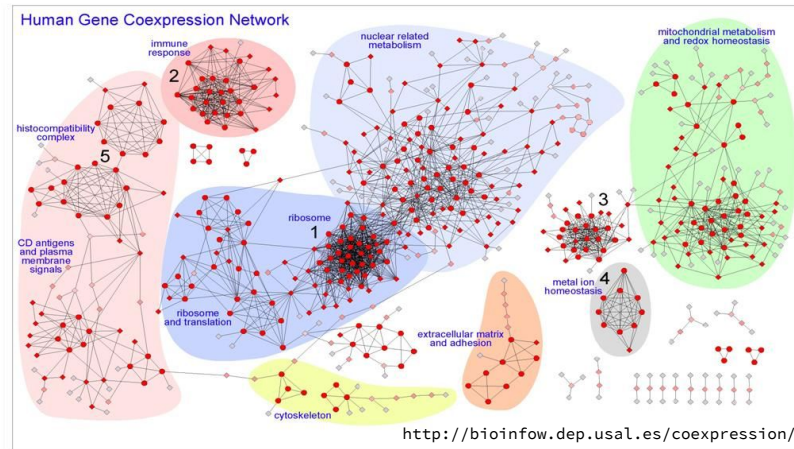
Signaling Network

Evolutionary Tree of Life

...



Nature Reviews | Genetics

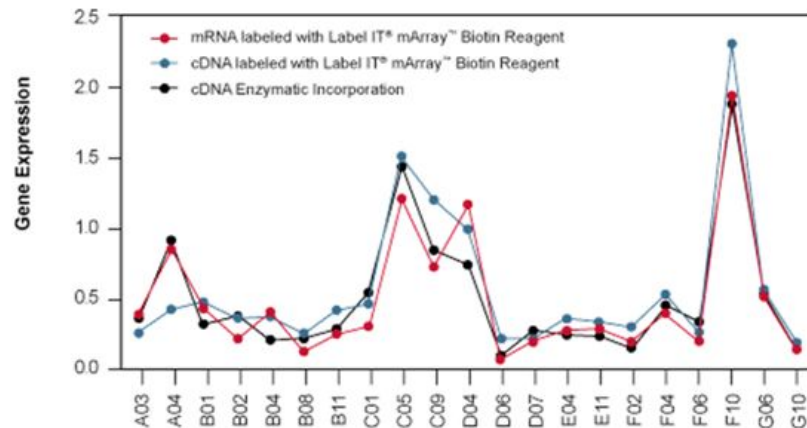




# CO-EXPRESSION NETWORKS

Measurement of expression levels of individual genes (or proteins, etc.) over a time-series or/and multiple conditions.

Genes with similar expression are hypothesized to have common functions.

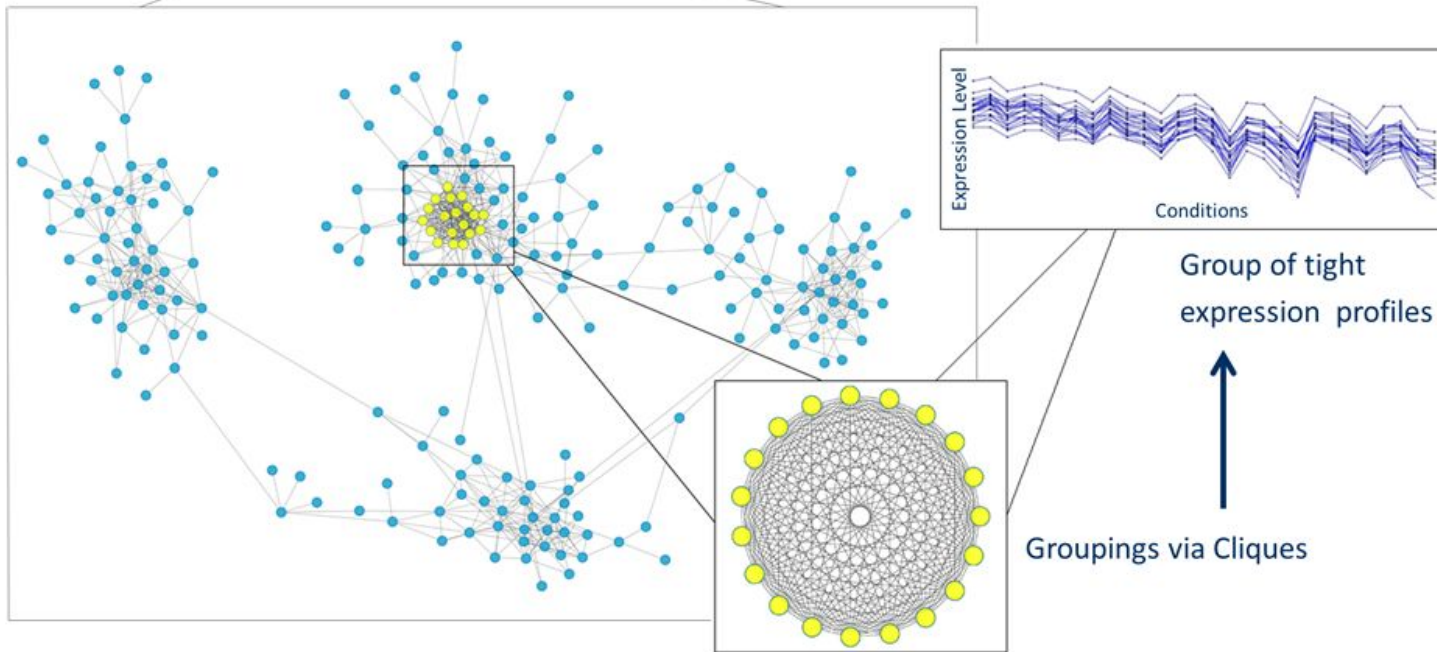


```
10_baseline_2_baseline_3_baseline_4_optim_1_optim_2_optim_3_optim_4_inhibition20_1_inhibition20_2_inhibition20_3_inhibition20_4
244901_at,0.01813439,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244902_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244903_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244904_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244905_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244906_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244907_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244908_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244909_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244910_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244911_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244912_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244913_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244914_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244915_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244916_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244917_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244918_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244919_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244920_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244921_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244922_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244923_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244924_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244925_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244926_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244927_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244928_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244929_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244930_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244931_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244932_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244933_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244934_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244935_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244936_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244937_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244938_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
244939_at,0.01711487,0.01544416,0.00708344,0.17004215,0.09440024,0.00807716,0.49374224,0.9970032,0.115444
```

Expression Data

Constructing gene co-expression networks

Network Model



# TUTORIAL

```
source("https://bioconductor.org/biocLite.R")  
biocLite("DESeq2")
```

```
install.packages("ggplot2")
```

```
install.packages("ggfortify")
```

See file

# REFERENCES

Conesa, A, *et al.* “A survey of best practices for RNA-seq data analysis”, Genome Biology(2016)

Barabási, A-L and Oltvai, ZN. “Network biology: understanding the cell's functional organization”, Nature Reviews - Genetics(2004)

Prieto, C, *et al.* “Human Gene Coexpression Landscape: Confident Network Derived from Tissue Transcriptomic Profiles” PLoS One(2008)

Petereit, J, *et al.* “petal: Co-expression network modelling in R”, BMC Systems Biology(1016)