

Toxicogenomics analysis

Giulia Callegaro, Hugo van Kessel, Steven Kunnen & Bob van de Water

EUROTOX 2023

CEC01 - A hands-on introduction to applied artificial intelligence in toxicology (CAAIT)



Outline

- Network analysis on omics data
- Small introduction on TempOseq and the data pipeline
- Practice!

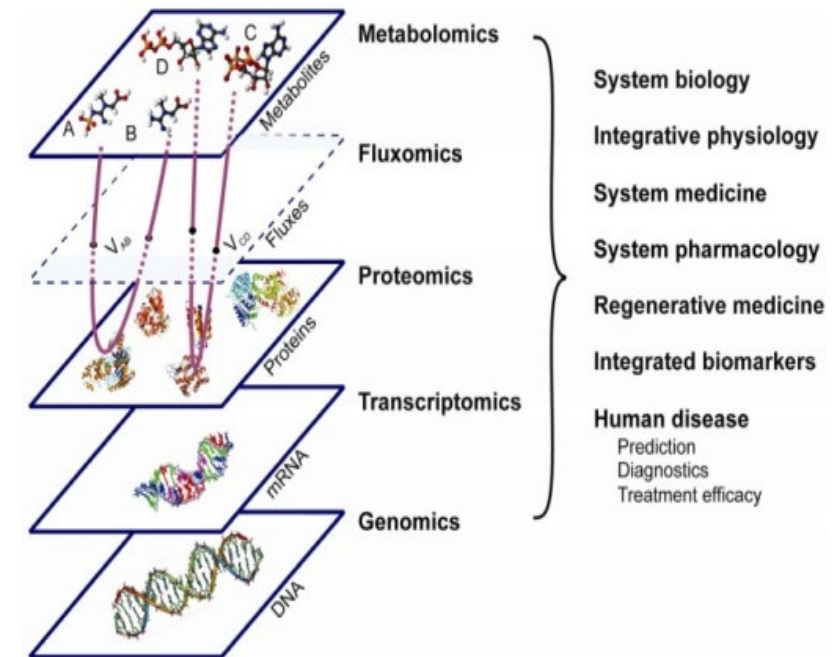
Gene network analysis

The existing model you will be using

Toxicogenomics

Systemic responses to substances

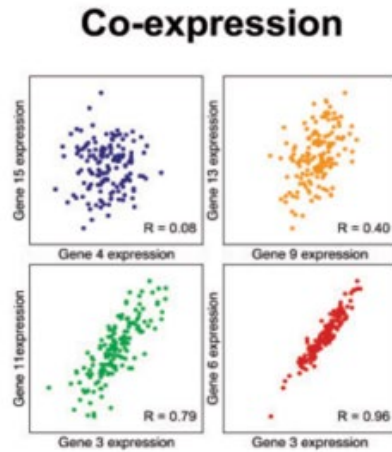
- **Genomics:** structure & function of the genome (e.g. allele frequencies/ types and drug sensitivity/ metabolism)
- **Epigenomics:** Reversible heritable changes in genome functions (often methylome status of genome or histone modifications)
- **Transcriptomics:** mRNA & ncRNA expression (e.g. gene activity/ involved pathways)
- **Proteomics:** Protein expression & activity (e.g. phosphorylation states, specific pathways)
- **Metabolomics:** Metabolic profiling (e.g. formed drug metabolites, or changes in glycolytic/ oxidative phosphorylation metabolites after drug exposure)



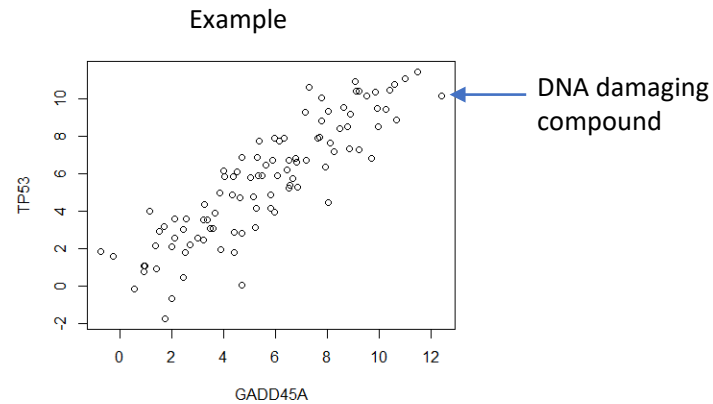
“Toxicogenomics was identified in TT21C as a transformative approach that was expected to play a pivotal role in **identifying** the toxicity pathways and cellular responses associated with exposure to environmental agents.”

Gene co-expression networks

Group genes with a tendency to co-activate across a group of samples



Co-expression == Pairwise
gene correlation across all
the conditions

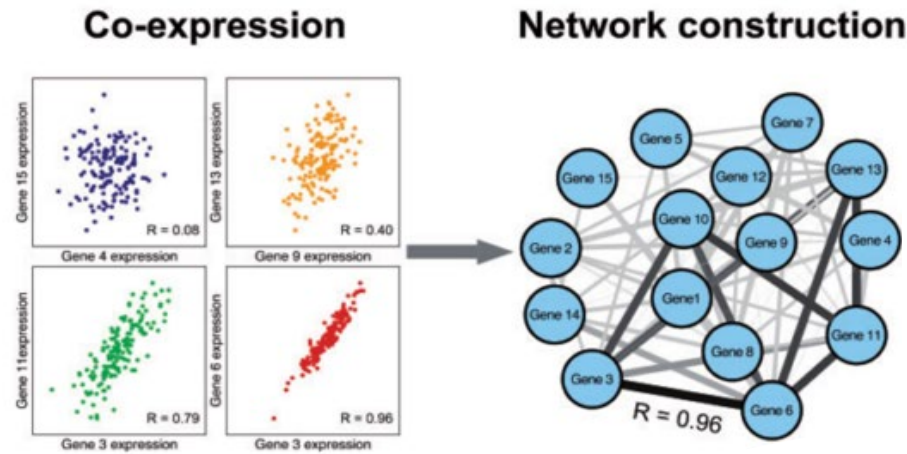


When TP53 is highly expressed,
so is GADD45.

When TP53 is lowly expressed,
so is GADD45

Gene co-expression networks

Group genes with a tendency to co-activate across a group of samples

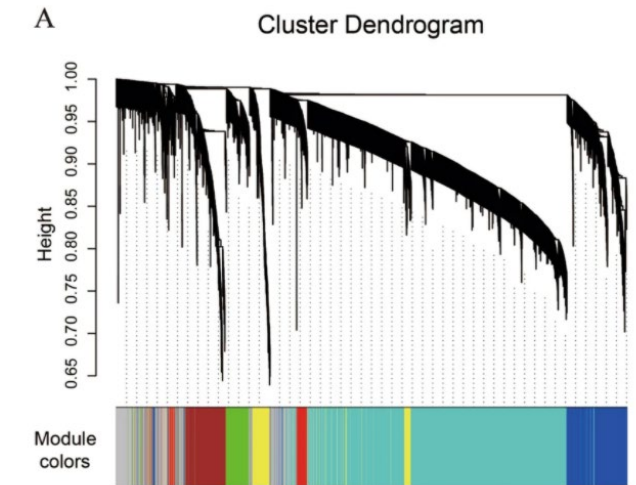
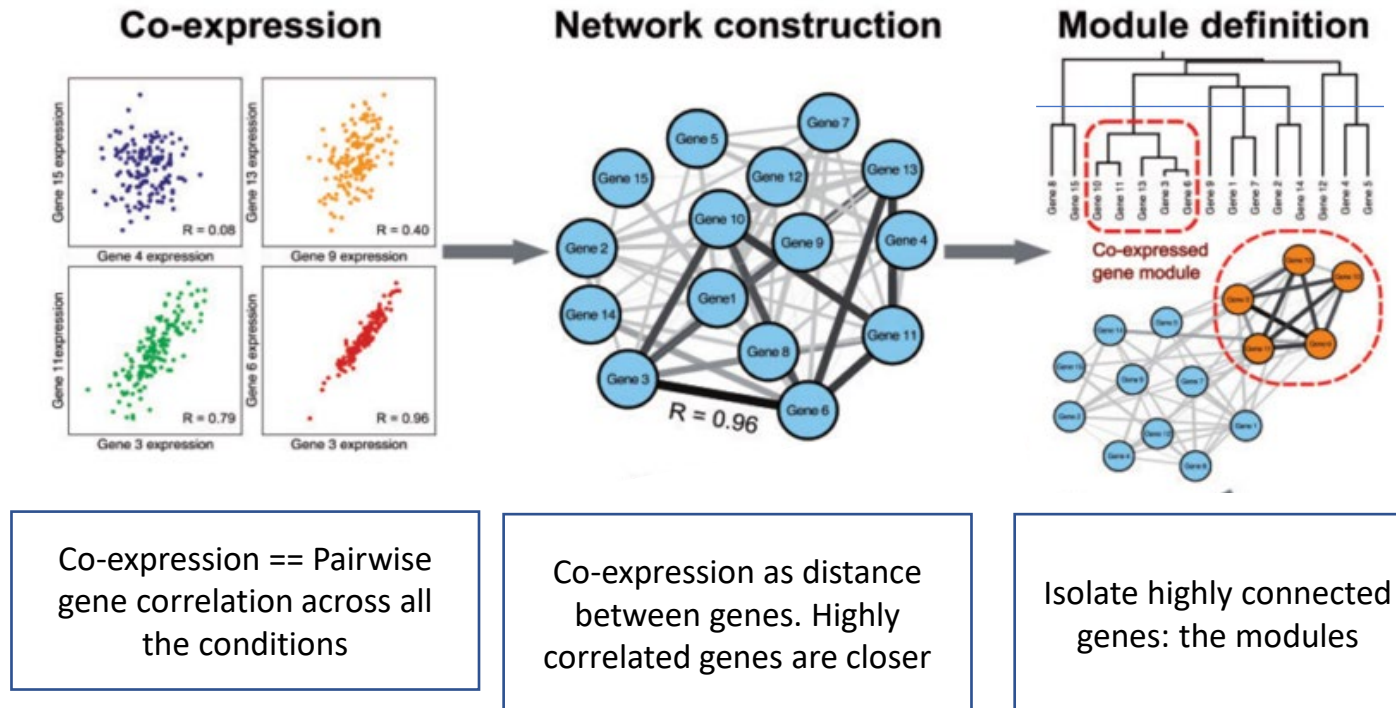


Co-expression == Pairwise gene correlation across all the conditions

Co-expression as distance between genes. Highly correlated genes are closer

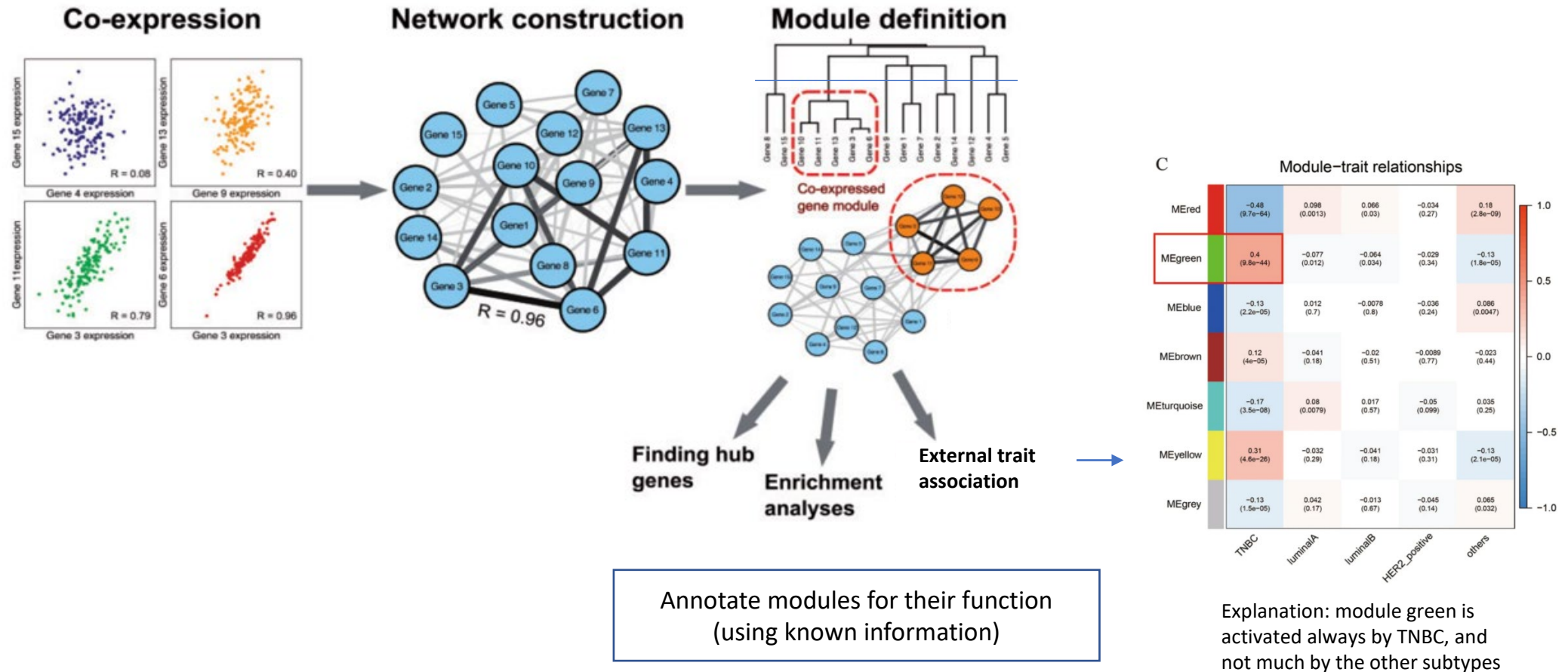
Gene co-expression networks

Group genes with a tendency to co-activate across a group of samples



Gene co-expression networks

Group genes with a tendency to co-activate across a group of samples

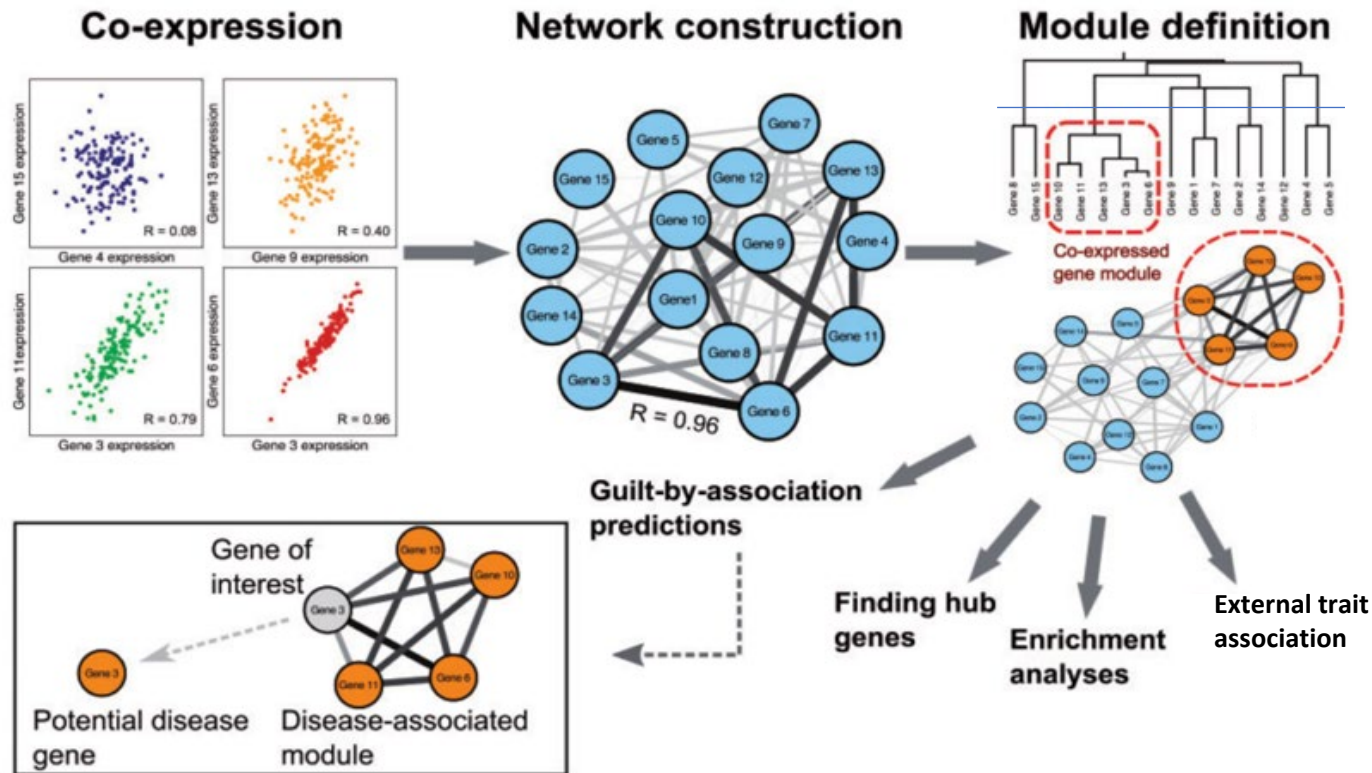


Annotate modules for their function
(using known information)

but not all genes are annotated..

Gene co-expression networks

Group genes with a tendency to co-activate across a group of samples



library (WGCNA)

The TXG-MAPr tool

<https://txg-mapr.eu/login/>



Credentials:

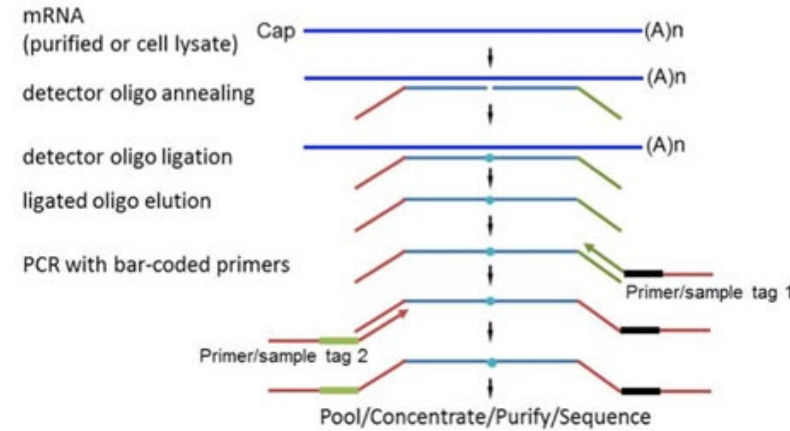
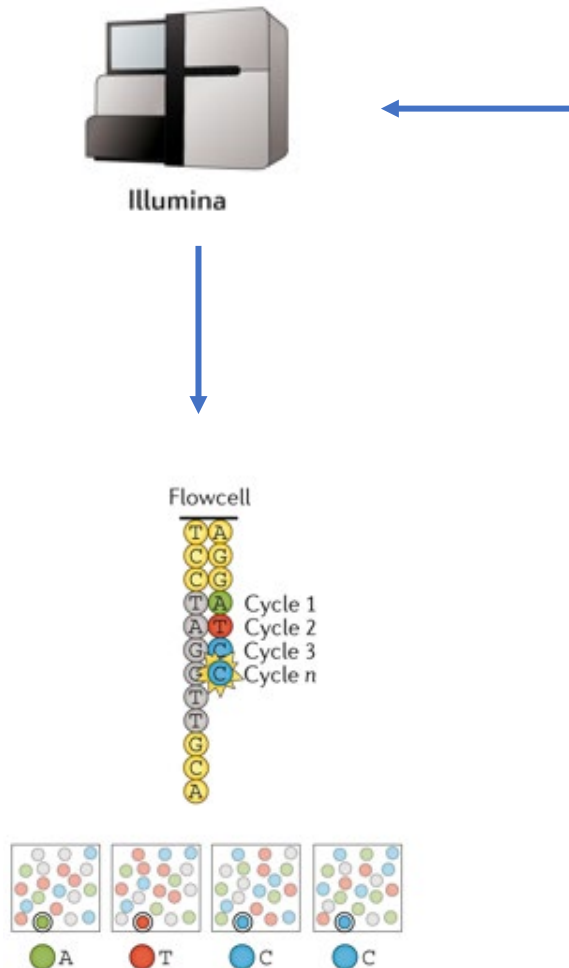
- username: eurotox2023
- password: txg_mapr01

Unknown genes are associated to existing information!

Practice: generate your data to load into the model

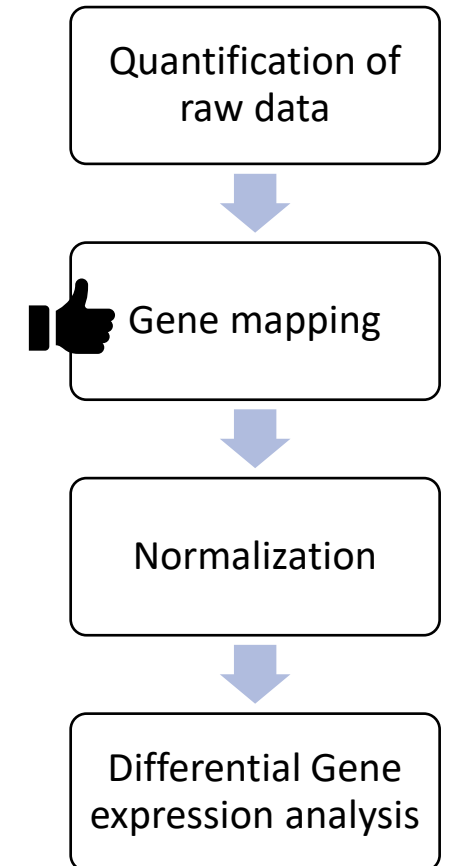
Transcriptomic data processing

Transcriptomics – TempOseq RNA sequencing



Principle: hybridization with complementary probes AND Sequencing by synthesis

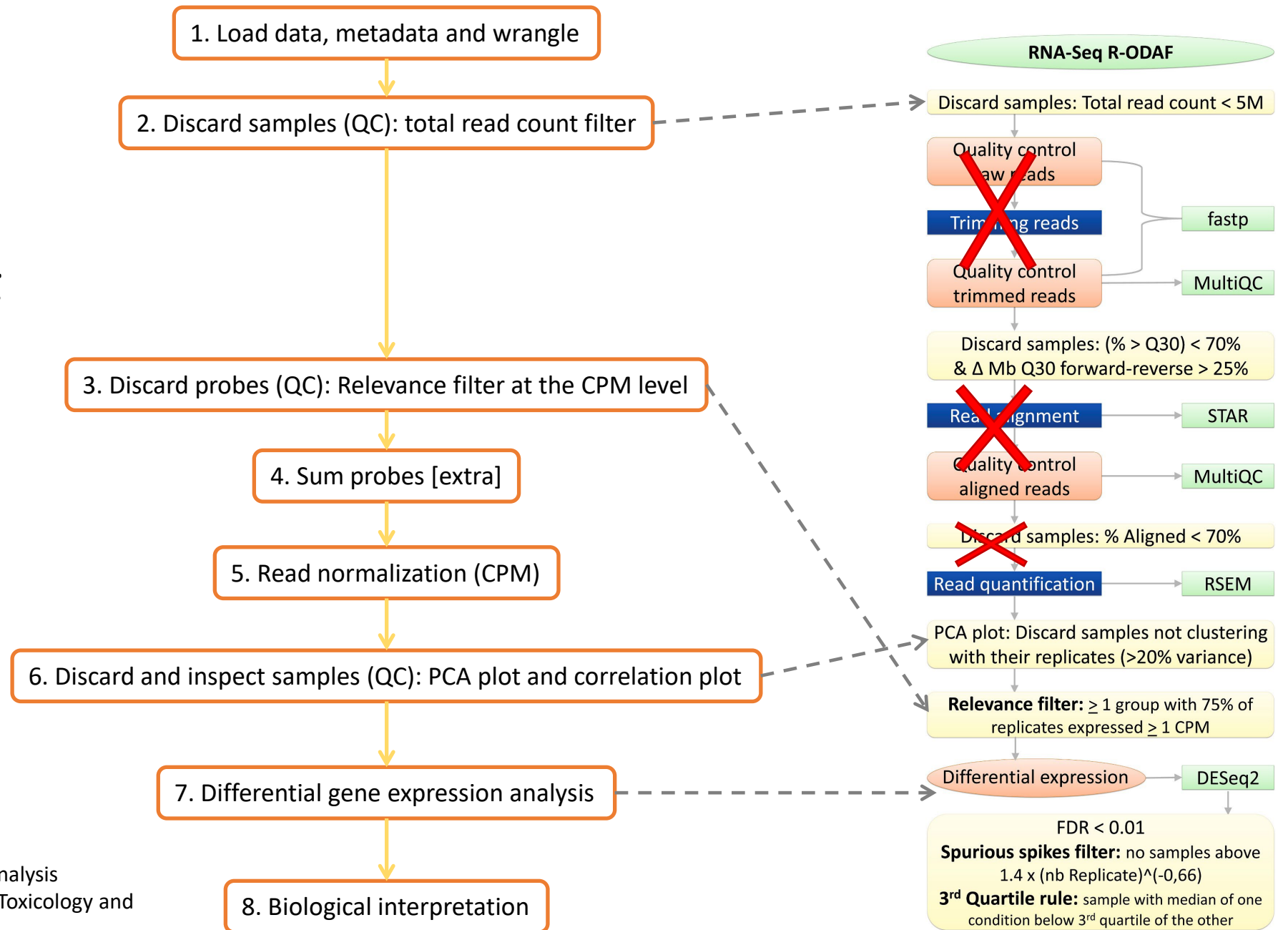
- Prebuilt **set** of oligos:
 - Complementary to mRNA of interest
 - Including a primer
 - Including a sample tag
- Only oligos annealing with mRNA in the samples are retained
- Amplification
- Sequencing of amplified oligos
 - “we know what we sequence” concept
- Easy to define target space: suitable for high- throughput



Definitions

- Read depths: (aimed) number of reads per sample
- Total read count: the actual number of reads per sample
- Count per million (CMP): one way to correct for different sequencing depths $\rightarrow (\text{read}/\text{total_reads}) * 10^6$

Operational pipeline and comparison with R-ODAF



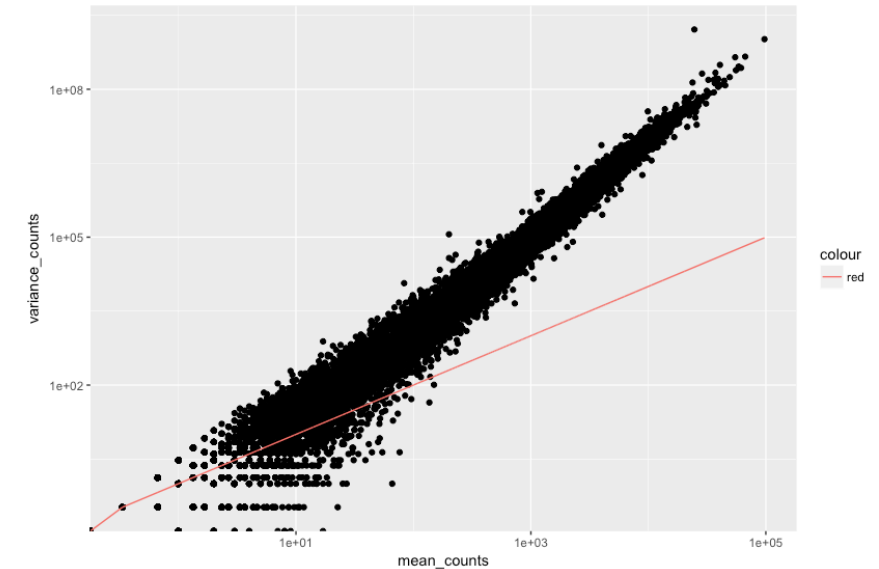
Differential gene expression analysis

Reads counts cannot be approximated with a normal distribution:

- They are integer counts
- low number of counts associated with a large proportion of genes
- long right tail due to the lack of any upper limit for expression

Count data is modeled using the negative binomial distribution (why not Poisson? Because mean < variance)

- For each gene, a generalized linear model (GLM) is fitted, modelling the read counts K_{ij} as following a negative binomial distribution



raw count for gene i , sample j

The mean is taken as “normalized counts” scaled by a normalization factor

one dispersion per gene

$$K_{ij} \sim \text{NB}(s_{ij}q_{ij}, \alpha_i)$$

Deseq2 in practice

- Create DESeqDataSet object

```
deseq_object = DESeqDataSetFromMatrix(countData = as.data.frame(countdata_raw_fsampl_fprobe_sumprobe),  
                                       colData = metadata %>% mutate(mean_id = as.factor(mean_id)),  
                                       design = ~ mean_id,  
                                       tidy = TRUE)
```

- Define size factors (CMP)

```
sizeFactors(deseq_object) = colSums(column_to_rownames(countdata_raw_fsampl_fprobe_sumprobe, var = "gene_symbol"))/1E6
```

- Run Deseq model

```
deseq_object = DESeq(deseq_object)
```

- (for each contrast) extract the results

```
results(loop_deseq_object, contrast = c("mean_id", contrast$mean_id_treatment[i], contrast$mean_id_control[i])) %>%
```




Universiteit
Leiden
The Netherlands

LACDR

We are hiring!

Currently online:

**Postdoctoral researcher in Toxicogenomics
2.0: bringing interpretation and
FAIRification to the next level**



Soon online:

**PhD researcher in Systems Toxicology:
connecting gene co-expression networks to
cellular injury and clinical pathology**

**Bioinformatician – (Postdoc or PhD
researcher) - in Translational Quantitative
Gene Network Analysis for Systems
Toxicology-based Human Chemical Safety
Assessment**



**RISK
HUNT3R**

