# Expression data analysis with
# Chipster

Eija Korpelainen, Massimiliano Gentile
chipster@csc.fi

# Understanding data analysis - why?

➢ **Bioinformaticians might not always be available when needed**

➢ **Biologists know their own experiments best**
- Biology involved (e.g. genes, pathways, etc)
- Potential batch effects etc

➢ **Allows you to design experiments better**

- Enough replicates, reads etc → less money wasted

➢ **Allows you to discuss more easily with bioinformaticians**

# What will I learn?

- ➢ **How to operate the Chipster software**
- ➢ **How to analyze microarray data**
  - • Central concepts
  - • Analysis workflow
  - • What happens in the different analysis steps
- ➢ **How to analyze RNA-seq data**
  - • Short introduction to analysis workflow and central concepts

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - Array level
  - Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - Parametric and non-parametric tests
  - Linear modeling
  - Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Introduction to Chipster

# Chipster

➢ **Provides an easy access to over 370 analysis tools**
  - No programming or command line experience required

➢ **Free, open source software**

➢ **What can I do with Chipster?**
  - analyze and integrate high-throughput data
  - visualize data efficiently
  - share analysis sessions
  - save and share automatic workflows

# Analysis tools

➢ **170 NGS tools for**
- RNA-seq
- miRNA-seq
- exome/genome-seq
- ChIP-seq
- FAIRE/DNase-seq
- CNA-seq
- 16S rRNA sequencing
- Single cell RNA-seq

➢ **140 microarray tools for**
- gene expression
- miRNA expression
- protein expression
- aCGH
- SNP
- integration of different data

➢ **60 tools for sequence analysis**
- BLAST, EMBOSS, MAFFT
- Phylip

# Chipster

**Open source platform for data analysis**

## Welcome to Chipster

Chipster is a user-friendly software for analyzing high-throughput data such as NGS and microarrays. It contains over 360 analysis tools and a large collection of reference genomes. Users can save and share automatic analysis workflows, and visualize data interactively using for example the built-in genome browser. Chipster's client software uses Java Web Start to install itself automatically, and it connects to computing servers for the actual analysis. Chipster is open source, and the server environment is available as a virtual machine image free of charge. If you would like to use Chipster running on CSC's server, you need a user account.

### Launch Chipster v3.13

...or launch with more memory: 3 GB or 6 GB

*If you have trouble launching Chipster, read this*

**News and resources:**

- 13.6.2018 Version 3.13 released
- 17.4.2018 Course materials available!
- 17.4.2018 RNA-seq tutorial for differential expression analysis
- 19.8.2014 RNA-seq data analysis guidebook with Chipster instructions
- News archive

**Training:**

- 19.9.2018 Single cell RNA-seq data analysis, CSC
- 4.-5.9.2018 RNA-seq data analysis, University of Oulu
- 9.2.2018 Single cell RNA-seq data analysis using Chipster, CSC
- 16.1.2018 Webinar: VirusDetect pipeline

File   Edit   View   Workflow   Help

**Datasets**

- two-sample.tsv
- column-value-filter.tsv
- hc.tre
- kmeans.pdf
- kmeans.tsv
- extract.tsv
- seqs.txt.wee
- seqs.html
- annotations.tsv
- annotations.html
- cpdb-pathways.html
- cpdb-pathways.tsv
- cpdb-genes.tsv

**Analysis tools**

Microarrays   NGS   Misc          ✓   Show parameters      Run ▶

- Normalisation
- Quality control
- Preprocessing
- Statistics
- Clustering
- Annotation
- Pathways
- Promoter analysis
- Copy number aberrations
- Visualisation
- Utilities

One sample tests
Two groups tests
ROTS
SAM
Several groups tests
Linear modelling
Linear modelling using user-defined design ma
Test proportions
Correlate with phenodata
Correlate miRNA with target expression
Time series
Association analysis

Tests for comparing the mean gene expression of two groups. LPE only works, if the whole normalized data is used, i.e., the data should not be filtered. Other than empiricalBayes might be slow, if run on unfiltered data.

More help      Show tool sourcecode

**Workflow**

🔍 🔍 ☑ Fit

**Visualisation**

🖵 Maximise   🗗 Detach   ✕ Close

two-sample.tsv

472 kB, Wed Sep 03 06:56:07 EEST 2014

(Click here to add your notes)

Analysis history

**Statistics / Two groups tests**

| | |
|---|---|
| Column | group |
| Pairing | EMPTY |
| Test | empirical Bayes |
| p-value adjustment method | BH |
| p-value threshold | 0.01 |
| Show NA | no |

Spreadsheet

Heatmap

Expression profile

Volcano plot

Scatterplot

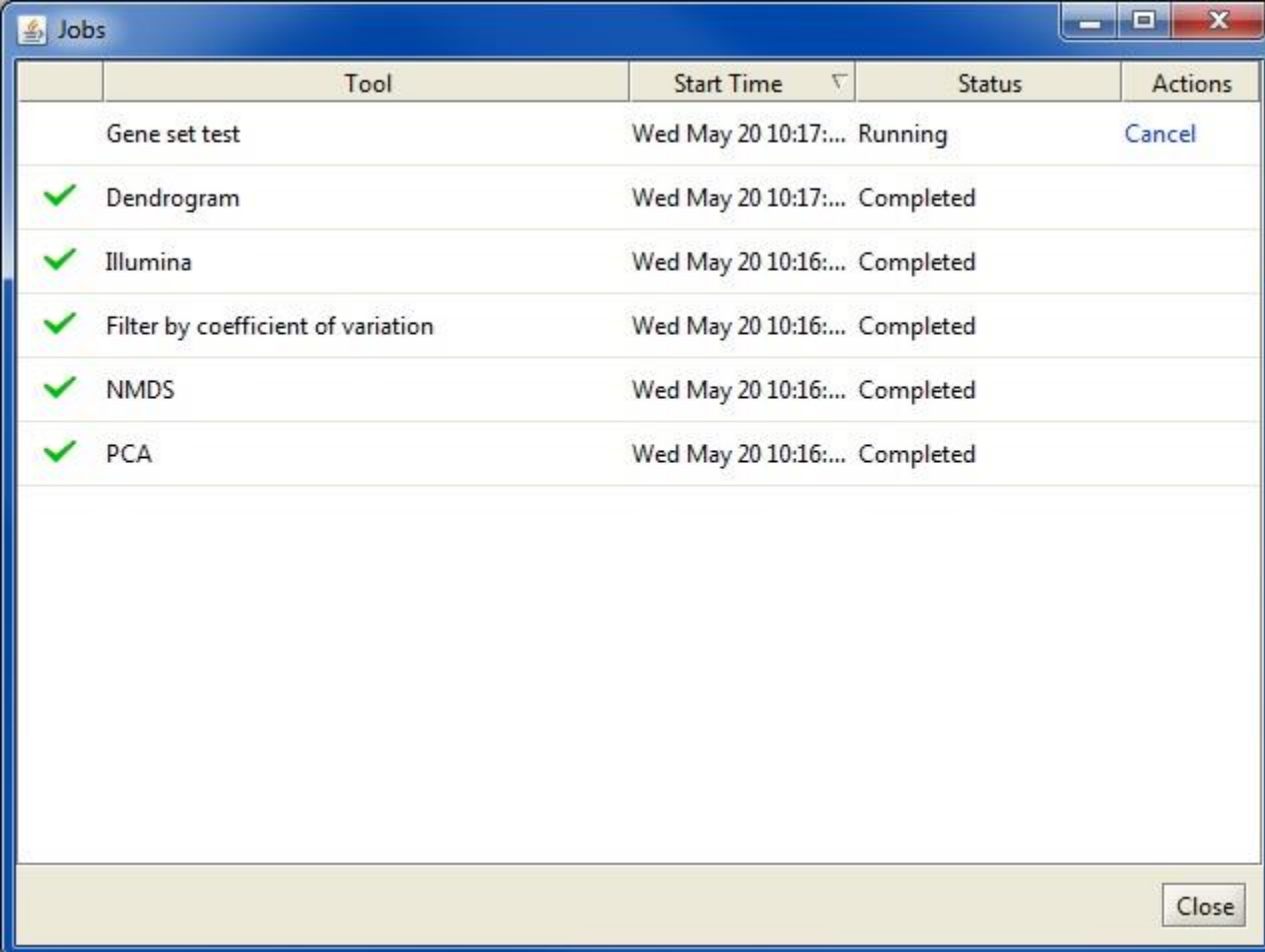3D Scatterplot

Histogram

Open in external web browser

Connected to chipster.csc.fi          View jobs   0 jobs running   Used memory 118M / 870M

# Mode of operation
Select: data → tool category → tool → run → visualize

# Job manager

➢ **You can run many analysis jobs at the same time**

➢ **Use Job manager to**

- view status
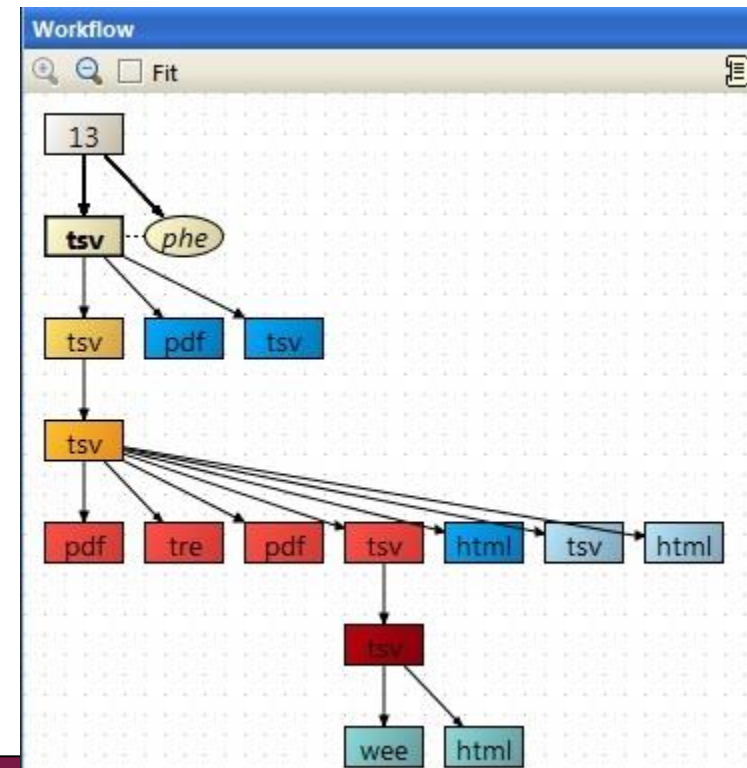- cancel jobs
- view time
- view parameters

# Workflow panel

➢ **Shows the relationships of the files**

➢ **You can move the boxes around, and zoom in and out.**

➢ **Several files can be selected by keeping the Ctrl key down**

➢ **Right clicking on the data file allows you to**

- <u>Save an individual result file ("Export")</u>
- Delete
- Link to another data file
- Save workflow

# Workflow – reusing and sharing your analysis pipeline

➢ **You can save your analysis steps as a reusable automatic "macro", which you can apply to another dataset**

➢ **When you save a workflow, all the analysis steps and their parameters are saved as a script file, which you can share with other users**

# Analysis history is saved automatically
-you can add tool source code to reports if needed

# Technical aspects

- ➤ **Client-server system**
  - Enough CPU and memory for large analysis jobs
  - Centralized maintenance
- ➤ **Easy to install**
  - Client uses Java Web Start
  - Server available as a virtual machine

# Analysis sessions

➢ **In order to continue your work later, <u>you have to save the analysis session.</u>**

- Session includes all the files, their relationships and metadata (what tool and parameters were used to produce each file).

➢ **Session is saved into a single .zip file on your computer.**

- In you can also save it on the server ("Save cloud session")

➢ **Session files allow you to continue the analysis on another computer, or share it with a colleague.**

➢ **You can have multiple analysis sessions saved separately, and combine them later if needed.**

# Problems? Send a support request
-request includes the error message and link to analysis session (optional)

# Two types of data visualizations

1. **Interactive visualizations produced by the client program**
   - Select the visualization method from visualization panel icons
   - Save by right clicking on the image

2. **Static images produced by the analysis tools on the server**
   - Select from Analysis tools / Visualisation
   - Save by right clicking on the file name and choosing "Export"

# Interactive visualizations



**Available actions:**

- Select genes and create a gene list
- Change titles, colors etc
- Zoom in/out
- Venn diagram: select genes with a list

# Static images produced by analysis tools

- ➢ **MA plot**
- ➢ **MDS plot**
- ➢ **Box plot**
- ➢ **Histogram**
- ➢ **Heatmap**
- ➢ **Idiogram**
- ➢ **Chromosomal position**
- ➢ **Correlogram**
- ➢ **Dendrogram**
- ➢ **K-means clustering**
- ➢ **SOM-clustering**
- ➢ **Dispersion plot**
- ➢ **etc**



**Dispersion plot**

# Acknowledgements to Chipster users and contibutors

# More info

➢ **chipster@csc.fi**

➢ **http://chipster.csc.fi**

➢ **Chipster tutorials in YouTube**

RNA-seq Data Analysis

Chapman & Hall/CRC
Mathematical and Computational Biology Series

**RNA-seq Data Analysis**

**A Practical Approach**

Korpelainen, Tuimala, Somervuo, Huss, and Wong.

**Eija Korpelainen, Jarno Tuimala,
Panu Somervuo, Mikael Huss, and Garry Wong**

CRC CRC Press
Taylor & Francis Group
A CHAPMAN & HALL BOOK

**GitHub** | This repository Search | Explore Features

chipster / **chipster**

Chipster is a user-friendly analysis software for high-throughput data.

⊕ 7,565 commits | ⅄ 18 branches | 🏷 123 releases | 👥 14 contributors

BMC Genomics

IMPACT FACTOR 4.21

home | journals A-Z | subject areas | advanced search | authors | reviewers | libraries | about | my BioMed Central

Software

**Highly accessed** **Open Access**

Chipster: user-friendly analysis software for microarray and other high-throughput data

M Aleksi Kallio ✉, Jarno T Tuimala ✉, Taavi Hupponen ✉, Petri Klemela ✉, Massimiliano Gentile ✉, Ilari Scheinin ✉, Mikko Koski ✉, Janne Kaki ✉ and Eija I Korpelainen ✉

BMC Genomics 2011, **12**:507    doi:10.1186/1471-2164-12-507

# Microarray data analysis

# Microarray data analysis workflow

- ➤ **Importing data to Chipster**
- ➤ **Normalization**
- ➤ **Describing samples with a phenodata file**
- ➤ **Quality control**
  - • Array level
  - • Experiment level
- ➤ **Filtering (optional)**
- ➤ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➤ **Annotation**
- ➤ **Pathway analysis**
- ➤ **Clustering**
- ➤ **Saving the workflow**

# Importing data

➢ **Affymetrix**
  • CEL-files are recognized by Chipster automatically

➢ **Illumina: two importing options**
  1. Import the GenomeStudio file as it is
     • All the samples need to be in one file.
     • Need columns AVG, BEAD_STDERR, Avg_NBEADS and DetectionPval
     • When imported this way, the data has to be normalized in Chipster using the lumi method
  2. Use Import tool to define the sample columns in the file(s)
     • Use the tool "Normalization / Illumina" to normalize the data
  → **The import option influences your normalization options later**

➢ **Agilent (and any other tab delimited files)**
  • Use Import tool to define the sample columns

# 1. Import tool: Select what to do

# 2. Import tool: Define rows (header, title, etc)

# 3. Import tool: Define columns (identifier,sample)

# Import tool - which columns should I mark?

➢ **http://chipster.csc.fi/manual/import-help.html**

➢ **Agilent**
- **Identifier (Probe<u>Name</u>, in case of miRNA arrays use GeneName)**
- **Annotation (Control type)**
- **Sample (rMeanSignal or rMedianSignal)**
- **Sample background (rBGMedianSignal)** ⎱ 1-color ⎱ 2-color
- **Control (gMeanSignal or gMedianSignal)**
- **Control background (gBGMedianSignal)**

➢ **Illumina BeadStudio version 3 file and GenomeStudio files**
- **Identifier (ProbeID)**
- **Sample (text "AVG")**

➢ **Illumina BeadStudio version 1-2 file**
- **Identifier (TargetID)**
- **Sample (text "AVG")**

# Importing <u>normalized</u> data

➤ **The data should be tab delimited and preferably log-transformed**
  • If your data is not log-transformed, you can transform it with the tool "Change interpretation"

➤ **Import the data file to Chipster using the Import tool. Mark the identifier column and all the sample columns.**

➤ **Run the tool <u>Normalize / Process prenormalized.</u> This**
  • Converts data to Chipster format by adding "chip." to expression column names
  • Creates the phenodata file. You need to indicate the chiptype using names given at http://chipster.csc.fi/manual/supported-chips.html

# Exercise 1. Start Chipster and open a session with Affymetrix .CEL-files

➢ **Log in to Chipster**

- Go to https://www.dkfz.de/gpcf/chipster0.html

- Log in with your normal DKFZ account

- If you don't have a DKFZ account, use user: gpcfproj and W110w110


➢ **Open session containing course data**

- Select **Open local session** and choose **Affymetrix_kidney_cancer**. The course data contains 17 samples from a kidney cancer study, measured using Affymetrix U133A chips. We want to find genes which are differentially expressed in cancer vs normal tissue.

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - • Array level
  - • Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Normalization

➤ **The goal is to make the arrays comparable to each other**
- Makes the expression value distributions similar
- Assumes that most genes don't change expression

➤ **After normalization the expression values are in <u>log2-scale</u>**
- Hence for example a fold change of 2 means 4-fold up

# Normalization of Affymetrix data

➢ **Normalization = background correction + expression estimation + summarization**

➢ **Methods**
- **RMA** (Robust Multichip Averaging) uses only PM probes, fits a model to them, and gives out expression values after quantile normalization and median polishing. Works nicely if you have more than a few chips
- **GCRMA** is similar to RMA, but takes also GC% content into account
- **MAS5** is the older Affymetrix method, **Plier** is a newer one
- **Li-Wong** is the method implemented in dChip

➢ **Custom chiptype parameter to use remapped probe information**
- Because some of the Affymetrix probe-to-transcript mappings can be outdated, probes have been remapped in the Bioconductor project.
- To use these remappings (alt CDF environments), select the matching chiptype from the Custom chiptype menu.

➢ **Variance stabilization option makes the variance similar over all the chips**
- Works only with MAS5 and Plier (the other methods log2-transform the data, which corrects for the same phenomenon)

# Quantile normalization procedure

| | Sample A | Sample B | Sample C |
|---|---|---|---|
| Gene 1 | 20 | 10 | 350 |
| Gene 2 | 100 | 500 | 200 |
| Gene 3 | 300 | 400 | 30 |

1. Raw data

| | Sample A | Sample B | Sample C | Median |
|---|---|---|---|---|
| Quantile 1 | 20 | 10 | 30 | 20 |
| Quantile 2 | 100 | 400 | 200 | 200 |
| Quantile 3 | 300 | 500 | 350 | 350 |

2. Rank data within sample and calculate median intensity for each row

| | Sample A | Sample B | Sample C | Median |
|---|---|---|---|---|
| Quantile 1 | 20 | 20 | 20 | 20 |
| Quantile 2 | 200 | 200 | 200 | 200 |
| Quantile 3 | 350 | 350 | 350 | 350 |

3. Replace the raw data of each row with its median (or mean) intensity

| | Sample A | Sample B | Sample C |
|---|---|---|---|
| Gene 1 | 20 | 20 | 350 |
| Gene 2 | 200 | 350 | 200 |
| Gene 3 | 350 | 200 | 20 |

4. Restore the original gene order

# Normalization of Agilent data

➢ **Background correction + averaging duplicate spots + normalization**

➢ **Background subtraction often generates negative values, which are coded as missing values after log2-transformation.**

- Using normexp + offset 50 will not generate negative values, and it gives good estimates

➢ **Loess removes curvature from the data (recommended)**



Before



After

# Agilent normalization parameters in Chipster

➢ **Background treatment**

- <u>Normexp</u>, Subtract, Edwards, None

➢ **Background offset**

- <u>50</u> or 0

➢ **Normalize chips**

- <u>Loess</u>, median, none

➢ **Chiptype**

- You must give this information in order to use annotation-based tools later

➢ **Normalize genes**

- None, scale (to median), quantile
- not needed for statistical analysis

# Illumina normalization: two analysis tools

1. **Illumina**
   - Normalization method
     - <u>Quantile</u>, vsn (variance stabilizing normalization), scale, none
   - Illumina software version
     - <u>GenomeStudio or BeadStudio3</u>, BeadStudio2, BeadStudio1
   - Chiptype
   - Identifier type
     - <u>Probe ID</u> (for BeadStudio version 3 data and newer), Target ID

2. **Lumi pipeline (data needs to be in one file, imported directly!)**
   - Normalization method
     - <u>Quantile</u>, vsn, rsn (robust spline normalization), loess, none
   - Transformation
     - <u>Log2</u>, vst (variance stabilizing transformation), none
   - Chiptype
     - human, mouse, rat
   - Background correction (usually done already in GenomeStudio)
     - <u>none</u>, bgAdjust.Affy

# Checking normalization

# Exercise 2: Normalize Affymetrix data

➢ Select all the CEL files by clicking on the box "**17**" in the Workflow view

➢ Select the tool **Normalisation / Affymetrix**, click **Show parameters**, set **Custom CDF annotation to be used** = **hgu133A**, and click **Run**.

➢ Repeat the process by setting **Custom CDF annotation to be used** = **Use original Affymetrix annotations**. When the result file **normalized.tsv** comes, rename it to **original_normalized.tsv**

➢ Open both normalized files and compare them. Do they have the same number of genes (rows)?

# Microarray data analysis workflow

➢ **Importing data to Chipster**

➢ **Normalization**

➢ **Describing samples with a phenodata file**

➢ **Quality control**

  • Array level

  • Experiment level

➢ **Filtering (optional)**

➢ **Statistical testing**

  • Parametric and non-parametric tests

  • Linear modeling

  • Multiple testing correction

➢ **Annotation**

➢ **Pathway analysis**

➢ **Clustering**

➢ **Saving the workflow**

# Phenodata file

➢ **Experimental setup is described with a phenodata file, which is created during normalization**

➢ **Fill in the group column with numbers describing your experimental groups**

- e.g. 1 = control sample, 2 = cancer sample
- necessary for the statistical tests to work
- note that you can sort a column by clicking on its title

➢ **Change sample names in Description column for visualizations**

# How to describe pairing, replicates, time, etc?

➢ **You can add new columns to the phenodata file**

➢ **How to describe different variables**

- **Time:** Use either real time values or recode with group codes
- **Replicates:** All the replicates are coded with the same number
- **Pairing:** Pairs are coded using the same number for each pair
- **Gender:** Use numbers
- **Anything else:** Use numbers

# Creating phenodata for <u>normalized</u> data

➤ **When you import data which has been already normalized, you need to create a phenodata file for it**

- Use Import tool to bring the data in
- Use the tool <u>Normalize / Process prenormalized</u> to create phenodata
  - Remember to give the chiptype
- Fill in the group column

➤ **Note: If you already have a phenodata file, you can import it too**

- Choose "Import directly" in the Import tool
- Right click on normalized data, choose "Link to phenodata"

# Exercise 3: Describe the experiment

➢ **Double click the phenodata file of the normalized.tsv**

➢ **In the phenodata editor, fill in the group column so that you enter**
- 1 for normal samples
- 2 for cancer samples

➢ **For the interest of visualizations later on, give shorter names for the samples in the Description column**
- Name the normal samples n1, n2,…
- Name the cancer samples c1, c2 ,…

# Microarray data analysis workflow

- ➤ **Importing data to Chipster**
- ➤ **Normalization**
- ➤ **Describing samples with a phenodata file**
- ➤ **Quality control**
  - • Array level
  - • Experiment level
- ➤ **Filtering (optional)**
- ➤ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➤ **Annotation**
- ➤ **Pathway analysis**
- ➤ **Clustering**
- ➤ **Saving the workflow**

# Array level quality control

➢ **Allows you to check if arrays are comparable to each other**
➢ **Tools in Chipster**

- Affymetrix basic: RNA degradation and Affy QC
- Affymetrix RLE and NUSE: fit a model to expression values
- Agilent 1-color: density plot and boxplot
- Agilent 2-color: MA-plot, density plot and boxplot
- Illumina: density plot and boxplot

# Affymetrix array level QC tools

➤ **Note that these tools use raw data (CEL files), not normalized data**

➤ **Affymetrix basic**
- Produces 3 plots:
    - QC stats plot
    - RNA degradation plot
    - Spike-in controls linearity plot
- Note that this tool uses the original probe set definitions from Affymetrix, not the alternative CDFs

➤ **Affymetrix RLE and NUSE**
- RLE (relative log expression)
- NUSE (normalized unscaled standard error plot)

➤ **Affymetrix RLE and NUSE for exon/gene arrays**

# Relative log expression, RLE



➤ RLE is the difference between log summarized expression of each chip to the log summarized expression on the median chip values.

➤ Boxes should be centered near 0 and have similar spread.

# Normalized Unscaled Standard Error, NUSE



➢ NUSE is the individual probe error fitting the Probe-Level Model.

➢ Good chips have median values close to one, while bad ones have are above 1.1.

➢ Check also if some chips show higher spread of NUSE distri-bution than others.

# Affymetrix QC

probesets with present flag

average background on the chip

scaling factors for the chips
beta-actin 3':5' ratio
GADPH 3':5' ratio

Blue area shows where scaling
factors are less than 3-fold of the
mean.
•If the scaling factors or ratios fall
within this region (1.25-fold for
GADPH), they are colored **blue,**
otherwise **red**

△  actin3/actin5
ı  gapdh3/gapdh5

## QC Stats

microarray009.cel    47.03%    59.75

microarray008.cel    44.11%    72.05

microarray007.cel    41.26%    70.21

microarray006.cel    44.82%    70.54

microarray005.cel    43.84%    56.16

microarray004.cel    41.19%    55.39

microarray003.cel    44.4%    56.88

microarray002.cel    48.7%    55.2

microarray001.cel    42.24%    139.04

-3   -2   -1   0   1   2   3

# Affymetrix spike-ins and RNA degradation

**Spike-in linearity**

**RNA degradation plot**

# Density plot and box plot

# Agilent QC: MA-plot



> ➤ **Scatter plot of log intensity ratios M=log2(R/G) versus average log intensities A = log2 √(R*G), where R and G are the intensities for the sample and control, respectively**

> ➤ **M is a mnemonic for <u>m</u>inus, as M = log R − log G**

> ➤ **A is mnemonic for <u>a</u>dd, as A = (log R + log G) / 2**

# Exercise 4: Affymetrix array level quality control

➢ Select the **17 CEL files** and run the tool **Quality control / Affymetrix basic**. Please note that this tool uses the original probe set definitions from Affymetrix

- Inspect the three pdf image files. Are there outlier samples?

➢ Select the **17 CEL files** and run the tool **Quality control / Affymetrix – using RLE and NUSE** setting **Custom chiptype = hgu133ahsentrezg(hgu133a)**

- Inspect the RLE and NUSE images. Are there outlier samples?

➢ Select **normalized.tsv** and run the tool **Quality control / Illumina** which produces a boxplot and density plot

- Inspect the plots. Are there outlier samples?

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - Array level
  - Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - Parametric and non-parametric tests
  - Linear modeling
  - Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Experiment level quality control

➢ **Getting an overview of similarities and dissimilarities between samples allows you to check**

- Do the experimental groups separate from each other?
- Is there a confounding factor (e.g. batch effect) that should be taken into account in the statistical analysis?
- Are there sample outliers that should be removed?

➢ **Several methods available**

- NMDS (non-metric multidimentional scaling)
- PCA (principal component analysis)
- Clustering
- Dendrogram
- Correlogram

# Non-metric multidimensional scaling (NMDS)

➢ **Goal is to reduce dimensions from several thousands to two**
  • High dimensional space is projected into a 2-dimensional space
➢ **Check that the experimental groups separate on dimension 1**
  • Do the samples separate according to something else on dimension 2?

➢ **Method**

  • Computes a distance matrix for all genes

  • Constructs the dimensions so that the similarity of distances between the original and the 2-dimensional space is maximized

# Principal component analysis (PCA)

➤ **Goal is to reduce dimensions**
  - High dimensional space is projected into a lower dimensional space
➤ **Check the percentage of variance explained by each component**
  - If PC2 explains only a small percentage of variance, it can be ignored.

➤ **Method**
  - Computes a variance-covariance matrix for all genes

  - PC1, the first principal component, is the linear combination of variables that maximizes the variance

  - PC2 is a linear combination orthogonal to the previous one which maximizes variance.

  - etc

# PCA illustration

X

Y

Z

Z-Y

Z-X

X-Y

X is the first principal component of the pen

Explains most of the variability in the shape of the pen

# PCA illustration, continued

X

Y

Z

Y is the second principal component of the pen

Z-Y

Z-X

Explains most of the remaining variability in the shape of the pen

Y-X

# Dendrogram and correlogram

**Dendrogram**



**Correlogram**

# Exercise 5: Experiment level quality control

➢ **Run <u>Statistics / NMDS</u> for the normalized data (normalized.tsv)**
  - Do the groups separate along the first dimension?

➢ **Run <u>Statistics / PCA</u> on the normalized data.**
  - View **pca.tsv** as **3D scatter plot for PCA**. Can you see 2 groups?
  - Check in **variance.tsv** how much variance the first principal component explains? And the second one?

➢ **Run <u>Visualization / Dendrogram</u> for the normalized data**
  - Do the groups separate well?

➢ **Save the analysis session with name sessionKidneyCancer.zip**

# Microarray data analysis workflow

➢ **Importing data to Chipster**
➢ **Normalization**
➢ **Describing samples with a phenodata file**
➢ **Quality control**
  • Array level
  • Experiment level
➢ **Filtering (optional)**
➢ **Statistical testing**
  • Parametric and non-parametric tests
  • Linear modeling
  • Multiple testing correction
➢ **Annotation**
➢ **Pathway analysis**
➢ **Clustering**
➢ **Saving the workflow**

# Filtering

➢ **Why?**

- Reducing the number of genes tested for differential expression reduces the severity of multiple testing correction of p-values. As the p-values remain better, we detect more differentially expressed genes.

➢ **Why not?**

- Some statistical testing methods (inc. the empirical Bayes option in Chipster) need many genes, because they estimate variance by borrowing information from other genes which are expressed at similar level. Hence the more genes the better.

➢ **Filtering should**

- remove genes which don't have any chance of being differentially expressed: genes that are not expressed or don't change
- be <u>independent</u>: should not use the sample group information

# Filtering tools in Chipster

➢ **Filter by standard deviation (SD)**
  - Select the percentage of genes to be filtered out

➢ **Filter by coefficient of variation (CV = SD / mean)**
  - Select the percentage of genes to be filtered out

➢ **Filter by interquartile range (IQR)**
  - Select the IQR

➢ **Filter by expression**
  - Select the upper and lower cut-offs
  - Select the number of chips required to fulfil this rule

➢ **Filter by flag** (Affymetrix P, M and A flags)
  - Flag value and number of arrays

# Exercise 6: Filtering

➢ **Select the normalized data and play with the SD filter and CV filter.**

- Set the cutoffs so that you filter out 90% of genes (Percentage to filter out = 0.9).

- Preprocessing / Filter by SD

- Preprocessing / Filter by CV

➢ **Select the result files and compare them using the interactive Venn diagram visualization**

- Save the genes specific to SD filter to a new file. Rename it sd.tsv.

- Save the genes specific to CV filter to a new file. Rename it cv.tsv.

- View both as expression profiles. Is there a difference in expression levels of the two sets?

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - • Array level
  - • Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Statistical analysis: Why?

➢ **Distinguish the treatment effect from biological variability and measurement noise**

  • replicates
  • estimation of uncertainty (variability)

➢ **Generalisation of results**

  • representative sample
  • statistical inference

sampling

inference

# Parametric statistical methods

➢ **Comparing means of 1-2 groups**
- student's t-test

➢ **Comparing means of more than 2 groups**
- 1-way ANOVA

➢ **Comparing means in a multifactor experiment**
- 2-way ANOVA

# Parametric statistics



$$t = \frac{x_1 - x_2}{\sqrt{\dfrac{s_1^{\,2}}{n_1} + \dfrac{s_2^{\,2}}{n_2}}}$$

$H_0 : \mu_A = \mu_B \,,\; \mu_A - \mu_B = 0$

$H_1 : \mu_A \neq \mu_B$

Type 1 error, $\alpha$

Type 2 error, $\beta$

Power $= 1 - \beta$

# Non-parametric statistical methods

➢ **Comparing <u>ranks</u> of 2 groups**
  - Mann-Whitney

➢ **Comparing <u>ranks</u> of more than 2 groups**
  - Kruskal-Wallis

| Ranks | |
|:---:|:---:|
| group A | group B |
| 1 | 4 |
| 2 | 6 |
| 3 | 7 |
| 5 | 9 |
| 8 | 10 |

$$U_1 = n_1 * n_2 + \frac{n_1 * (n_1 + 1)}{2} - R_1$$

$$U_2 = n_1 * n_2 + \frac{n_2 * (n_2 + 1)}{2} - R_2$$

# Non-parametric tests compared to parametric

**Benefits**

- Do not make any assumptions on data distribution
  $\Rightarrow$ robust to outliers
  $\Rightarrow$ allow for cross-experiment comparisons


**Drawbacks**

- Lower power than parametric counterpart
- Granular distribution of calculated statistic
  $\Rightarrow$ many genes get the same rank
  $\Rightarrow$ requires at least 6 samples / group

# How to improve statistical power?

➢ **Need more accurate estimates of variability and effect size**

➢ **Improved analysis methods**
- Variance shrinking: Empirical Bayes method
- Partitioning variability: ANOVA, linear modeling

➢ **Improved experimental design**
- Increase number of biological replicates
- Use paired samples if possible
- Randomization
- Blocking

# Pairing = matched samples from the same individual

## Unpaired analysis

| | Before | After |
|---|---|---|
| | 2 | 3 |
| | 2 | 4 |
| | 3 | 2 |
| | 1 | 3 |
| Mean | 2 | 3 |
| Stdev | 0.8 | 0.8 |

## Paired analysis

| Before | After | Difference |
|---|---|---|
| 2 | 3 | 1 |
| 2 | 3 | 1 |
| 3 | 4 | 1 |
| 1 | 2 | 1 |

# Improving power with variance shrinking

➢ **Concept**
- Borrow information from other genes which are expressed at similar level, and form a pooled error estimate

➢ **How?**
- models the error - intensity dependence by comparing replicates
- uses a smoothing function to estimate the error for any given intensity
- calculates a weighted average between the observed gene specific variance and the model-derived variance (pooling)
- incorporates the pooled variance estimate in the statistical test (usually t- or F-test)

➢ **Available in Chipster**
- Two group test: Select empirical Bayes as the test
- Linear modeling tool

# Exercise 7: Statistical testing

➤ **Run different two group tests**

- Select the file **normalized.tsv** and **Statistics / Two group test.** What is the default value of the parameter "test"? How many differentially expressed genes do you get?

- Repeat the run but change **test = t-test**. Rename the result file to **t.tsv.** How many differentially expressed genes do you get now?

- Repeat the run but change **test = Mann-Whitney**. Rename the result file to **MW.tsv.** How many differentially expressed genes do you get now?

➤ **Compare the results with a Venn diagram**

- Do the gene lists overlap?

# Exercise 8: Visualize and filter results

➢ **Filter genes based on fold change**
  - Select **two-sample.tsv** and the tool **Utilities / Filter using a column value.** Keep genes whose expression changes more than 4-fold:
    - Column = FC
    - Cut-off = 2 (remember that the fold change values are in log2 scale)
    - Smaller or larger = outside (we want both up and down-regulated genes)

➢ **View results in interactive visualizations**
  - Select the **column-value-filter.tsv** and visualization method **Volcano plot**
  - Visualize the file also as **Expression profile**

# Exercise 9: Use paired samples in testing

➢ **Use pre-filled phenodata which contains more information about the samples**

- Select **normalized.tsv** and **phenodata.tsv**, right click, and select **Links between selected / Unlink**.

- Select **normalized.tsv** and right click to link it to **phenodata_moreSampleInfo**.

- Inspect the new phenodata for sample information. Note that sample pairing information is in the patient column.

➢ **Repeat statistical testing so that you include pairing information**

- Select the file **normalized.tsv** and **Statistics / Two group test** and set the parameter **Column with pairing information** = **patient**.

- Does the number of differentially expressed genes change?

- Rename the result file to **paired.tsv**

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - • Array level
  - • Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Linear modeling

➢ **Models the expression of a gene as a linear combination of explanatory factors (e.g. group, gender, time, patient,…)**

$$y = a + (b \cdot group) + (c \cdot gender) + (d \cdot group \cdot gender)$$

y = gene's expression

a, b, c and d = parameters estimated from the data

a = intercept (expression when factors are at "reference" level)

b and c = main effects

d = interaction effect

# Taking multiple factors into account

## 1 factor: treatment

| | Control | Treatment |
|---|---|---|
| | 2 | 5 |
| | 9 | 7 |
| | 1 | 3 |
| | 7 | 5 |
| | 8 | 4 |
| | 3 | 6 |
| **Mean** | **5** | **5** |

## 2 factors: treatment and <u>gender</u>

| | Control | Treatment |
|---|---|---|
| | 2 | 6 |
| **Males** | 3 | 7 |
| | 1 | 5 |
| **Mean** | **2** | **6** |
| | | |
| | 8 | 4 |
| **Females** | 9 | 5 |
| | 7 | 3 |
| **Mean** | **8** | **4** |

# Linear modeling: Interaction effect

# Linear modeling tool in Chipster

➢ **Linear modeling tool in Chipster can take into account**

- 3 main effects
- Their interactions
- Pairing
- Technical replication (one sample is hybridized to several arrays)

➢ **Main effects can be treated as**

- Linear = is there a trend towards higher numbers?
- Factor = are there differences between the groups?

If the main effect has only two levels (e.g. gender), selecting linear or factor gives the same result

➢ **Note that the result table contains all the genes, so in order to get the differentially expressed genes <u>you have to filter it</u>**

- Use the tool **Utilities / Filter using a column value**
- Select the column **p.adjusted** that corresponds to the comparison of your interest

# Exercise 10: Linear modeling

➢ **Perform linear modeling so that the analysis takes into account group and gender.**
- • Select **normalized.tsv** and **Statistics / Linear modelling**
- • Set **Main effect 2 = gender** and **treat both main effects as factors**.
- • Open **limma.tsv** and inspect the result columns.

➢ **Retrieve differentially expressed genes for the group comparison**
- • Select **limma.tsv** and the tool **Utilities / Filter using a column value.** Keep genes whose adjusted p-value < 0.05:
  - • Column = p.adjusted.main12
  - • Cut-off = 0.05
  - • Smaller or larger = smaller-than

➢ **Perform linear modeling so that the analysis takes into account group, gender and pairing.**
- • As above but include **pairing = patient**.
- • Open **limma.tsv** and inspect the result columns.
- • Retrieve differentially expressed genes as before.

# Microarray data analysis workflow

- ➤ **Importing data to Chipster**
- ➤ **Normalization**
- ➤ **Describing samples with a phenodata file**
- ➤ **Quality control**
  - • Array level
  - • Experiment level
- ➤ **Filtering (optional)**
- ➤ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
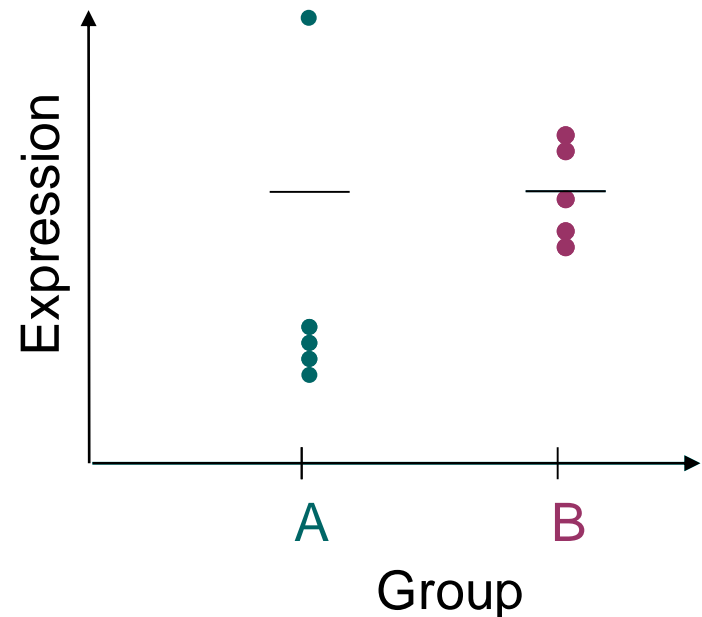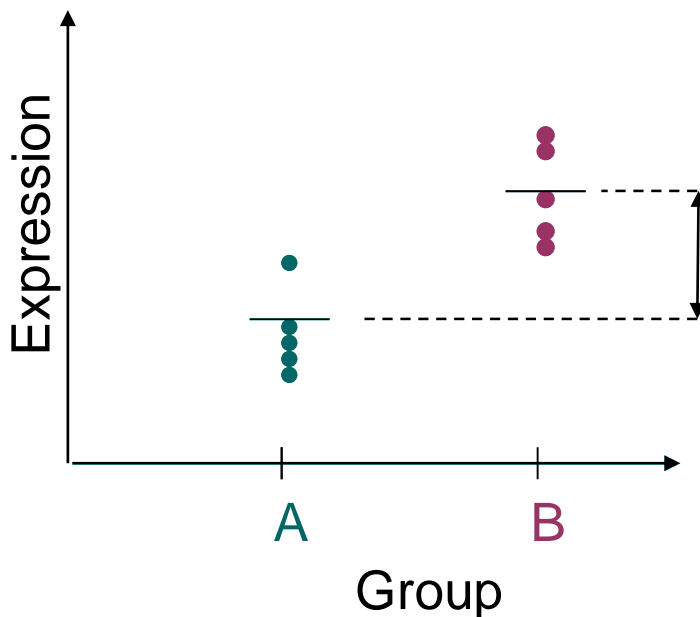- ➤ **Annotation**
- ➤ **Pathway analysis**
- ➤ **Clustering**
- ➤ **Saving the workflow**

# Multiple testing correction

➢ **Problem: When thousands of genes are tested for differential expression, a gene can get a good p-value just by chance.**

1 gene, $\alpha = 0.05$
$\Rightarrow$ false positive incidence = 1 / 20

30 000 genes, $\alpha = 0.05$
$\Rightarrow$ false positive incidence = 1500

➢ **Solution: Correct the p-values for multiple testing. Methods:**

- Bonferroni
- Holm (step down)
- Westfall & Young
- Benjamini & Hochberg

more false negatives

more false positives

# Benjamini & Hochberg method (BH)

➢ **How does it work?**
  • rank p-values from largest to smallest
  • largest p-value remains unaltered
  • second largest p-value = p * n / (n-1)
  • third largest p-value = p * n / (n-2)
  • …
  • smallest p-value = p * n / (n-n+1) = p * n

correction

n

1

raw p

➢ **Some adjusted p-values can become identical**
  • Adjusting should not change the order of p-values, so if $pa_{i+1} > pa_i$ then $pa_{i+1} = pa_i$

➢ **We can reduce the severity of multiple testing correction by reducing the number of genes tested (n)**
  • use independent filtering

➢ **The adjusted p-value is FDR (false discovery rate)**
  • Tells what proportion of <u>results</u> can be false positives

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - • Array level
  - • Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Annotation

➢ **Gene annotation = information about biological function, pathway involvement, chromosal location etc**

➢ **Annotation information is collected from different biological databases  to a single database by the Bioconductor project**

- Bioconductor provides annotation packages for many microarrays

➢ **Annotation package is required by many analysis tools**

-  Annotation, GO/KEGG enrichment, promoter analysis, chromosomal plots
- These tools don't work for those chiptypes which don't have Bioconductor annotation packages

# Annotations for the selected gene list

| Probe | Symbol | Description | Chromosome | Chromosome Location | GenBank | Gene | Cytoband | UniGene | PubMed | Gene Ontology | Pathway |
|-------|--------|-------------|------------|---------------------|---------|------|----------|---------|--------|---------------|---------|
| 205626_s_at | CALB1 | calbindin 1, 28kDa | 8 | -91140013 | NM_004929 | 793 | 8q21.3-q22.1 | Hs.65425 | 22 | locomotory behavior<br>cytoplasm<br>vitamin D binding<br>calcium ion binding<br>protein binding | |
| 220281_at | SLC12A1 | solute carrier family 12 (sodium/potassium/chloride transporters), member 1 | 15 | 46285789 | AI632015 | 6557 | 15q15-q21.1 | Hs.123116 | 13 | ion transport<br>potassium ion transport<br>sodium ion transport<br>chloride transport<br>membrane fraction<br>plasma membrane<br>membrane<br>integral to membrane<br>transporter activity<br>sodium:potassium:chloride symporter activity<br>symporter activity<br>potassium ion binding<br>sodium ion binding | |
| 206054_at | KNG1 | kininogen 1 | 3 | 187917813 | NM_000893 | 3827 | 3q27 | Hs.77741 | 86 | smooth muscle contraction<br>inflammatory response<br>negative regulation of cell adhesion<br>elevation of cytosolic calcium ion concentration<br>blood coagulation<br>diuresis<br>natriuresis<br>negative regulation of blood coagulation<br>vasodilation<br>positive regulation of apoptosis<br>extracellular region<br>cysteine protease inhibitor activity<br>receptor binding<br>heparin binding<br>zinc ion binding | Complement and coagulation cascades |
| | | | | | | | | | | behavior<br>gamma-aminobutyric acid catabolic process<br>neurotransmitter catabolic | Glutamate |

# Alternative CDF environments for Affymetrix

➢ **CDF is a file that links individual probes to gene transcripts**

➢ **Affymetrix default annotation uses older CDF files which may map many probes to wrong genes**

➢ **Alternative CDFs fix this problem**

➢ **In Chipster selecting "custom chiptype" in Affymetrix normalization takes altCDFs to use**

➢ **For more information see**

- Dai et al, (2005) Nuc Acids Res, 33(20):e175: *Evolving gene/transcript definitions significantly alter the interpretation of GeneChip data*

- http://brainarray.mbni.med.umich.edu/Brainarray/Database/CustomCDF/genomic_curated_CDF.asp

# Exercise 11: Annotation

➢ **Annotate genes**

- Select the file **two-sample.tsv**

- Run **Annotation / Agilent, Affymetrix or Illumina gene list** so that you include the FC and p-value information to the result file

- Run **Annotation / Add annotations to data**

# Microarray data analysis workflow

- ➤ Importing data to Chipster
- ➤ Normalization
- ➤ Describing samples with a phenodata file
- ➤ Quality control
  - Array level
  - Experiment level
- ➤ Filtering (optional)
- ➤ Statistical testing
  - Parametric and non-parametric tests
  - Linear modeling
  - Multiple testing correction
- ➤ Annotation
- ➤ **Pathway analysis**
- ➤ Clustering
- ➤ Saving the workflow

# Pathway analysis – why?

➢ **Statistical tests can yield thousands of differentially expressed genes**

➢ **It is difficult to make "biological" sense out of the result list**

➢ **Looking at the bigger picture can be helpful, e.g. which pathways are differentially expressed between the experimental groups**

➢ **Databases such as KEGG, GO, Reactome and ConsensusPathDB provide grouping of genes to pathways, biological processes, molecular functions, etc**

➢ **Two approaches to pathway analysis**
- Gene set enrichment analysis
- Gene set test

# Approach I: Gene set enrichment analysis

1. **Perform a statistical test to find differentially expressed genes**
2. **Check if the list of differentially expressed genes is "enriched" for some pathways**

Our list and apoptosis (10 genes)

Our list (50 genes)

Apoptosis (200 genes)

$$H_0 : \frac{\dfrac{10}{50}}{\dfrac{200}{30000}} = 30 >> 1$$

Genome (30,000 genes)

# Approach II: Gene set test

1. **Do NOT perform differential <u>gene</u> expression analysis**

2. **Group genes to pathways and perform differential expression analysis <u>for the whole pathway</u>**

➢ **Advantages**

- More sensitive than single gene tests

- Reduced number of tests → less multiple testing correction

  → increased power



03050 –Proteasome

# ConsensusPathDB

➢ **One-stop shop: Integrates pathway information from 32 databases covering**

- biochemical pathways

- protein-protein, genetic, metabolic, signaling, gene regulatory and drug-target interactions

➢ **Developed by Ralf Herwig's group at the Max-Planck Institute in Berlin**

➢ **ConsensusPathDB over-representation analysis tool is integrated in Chipster**

- runs on the MPI server in Berlin

# GO (Gene Ontology)

➢ **Controlled vocabulary of terms for describing gene product characteristics**

➢ **3 ontologies**

- Biological process
- Molecular function
- Cellular component

➢ **Hierarchical structure**

```
⊡ all : all [841457 gene products]
  ⊞ ▯ GO:0008150 : biological_process [660879 gene products]
    ⊞ ▯ GO:0065007 : biological regulation [145630 gene products]
      ⊞ ▯ GO:0050789 : regulation of biological process [134091 gene products]
        ⊞ ▯ GO:0048518 : positive regulation of biological process [42078 gene products]
          ⊞ ▯ GO:0048522 : positive regulation of cellular process [34658 gene products]
            ⊞ ▯ GO:0031325 : positive regulation of cellular metabolic process [21272 gene products]
              ⊞ ▯ GO:0032270 : positive regulation of cellular protein metabolic process [6797 gene products]
                ⊞ ▯ GO:0031401 : positive regulation of protein modification process [5757 gene products]
                  ⊞ ▯ GO:0001934 : positive regulation of protein phosphorylation [4638 gene products]
                    ⊞ ▯ GO:0045860 : positive regulation of protein kinase activity [2860 gene products]
                      ⊞ ▯ GO:0032147 : activation of protein kinase activity [1745 gene products]
                        ⊞ ▯ GO:0000185 : activation of MAPKKK activity [82 gene products]
                      ⊞ ▯ GO:0071902 : positive regulation of protein serine/threonine kinase activity [1815 ge
                        ⊞ ▯ GO:0000185 : activation of MAPKKK activity [82 gene products]
                  ⊞ ▯ GO:0010562 : positive regulation of phosphorus metabolic process [6341 gene products]
```

# KEGG

➢ **Kyoto Encyclopedia for Genes and Genomes**

➢ **Collection of pathway maps representing molecular interaction and reaction networks for**

- metabolism
- cellular processes
- diseases, etc

RAS SIGNALING PATHWAY

# Exercise 12: Gene set enrichment analysis

➢ **Identify over-represented GO terms**

- Select the **two-sample.tsv** file and run **Pathways / Hypergeometric test for GO.** Open **hypergeo.html** and read about the first term. Check in **hypergeo.tsv** how many terms do you get.

➢ **Extract genes for a specific GO term**

- Copy the GO identifier for the top term (GO:0006082).
- Select **two-sample.tsv** and run tool **Utilities / Extract genes for GO term**, pasting the GO identifier in the parameter field.
- Open **extracted-from-GO.tsv.** How many genes do you get? Are they up- or down-regulated (use also Volcano plot and Expression profile)?

➢ **Identify over-represented ConsensusPathDB pathways**

- Select **two-sample.tsv** and run **Pathways / Hypergeometric test for ConsensusPathDB**.
- Click on the links in the **cpdb.html** file to read about the pathways.

# Exercise 13: Gene set test

➢ **Identify differentially expressed KEGG pathways**

- Select the **<u>normalized.tsv</u>** file and **Pathways / Gene set test**. Set the **Number of pathways to visualize = 4**

- Explore **global-test-result-table.tsv.** How many differentially expressed KEGG pathways do you get?

- Explore **multtest.pdf.** Which gene contributes most to the first pathway?

# Microarray data analysis workflow

- ➢ **Importing data to Chipster**
- ➢ **Normalization**
- ➢ **Describing samples with a phenodata file**
- ➢ **Quality control**
  - • Array level
  - • Experiment level
- ➢ **Filtering (optional)**
- ➢ **Statistical testing**
  - • Parametric and non-parametric tests
  - • Linear modeling
  - • Multiple testing correction
- ➢ **Annotation**
- ➢ **Pathway analysis**
- ➢ **Clustering**
- ➢ **Saving the workflow**

# Clustering in Chipster

➢ **Hierarchical**

- Includes reliability checking of the resulting tree with bootstrapping

➢ **K-means**

- Additional tool to estimate K

➢ **Quality threshold**

➢ **Self-organizing maps**

➢ **K-nearest neighbor (KNN)**

- Classification aka class prediction

# Hierarchical clustering

➢ **Provides stable clusters**
➢ **Assumes pairwise correlations**
➢ **Early mistakes cannot be corrected**
➢ **Computationally intensive**
➢ **Drawing methods**
  • Single / average / complete linkage
➢ **Distance methods**
  • Euclidean distance
  • Pearson / Spearman correlation

# Hierarchical clustering: distance methods

One can either calculate the <u>distance</u> between two pairs of data sets (e.g. samples) or the <u>similarity</u> between them

Euclidean distance                    Pearson correlation

# Distance methods can yield very different results

## Distances

- the Correlation distance
  - red-blue is 0.006
  - red-gray is 0.768
  - blue-gray is 0.7101
- Euclidean distance:
  - red-blue is 9.45
  - red-gray is 10.26
  - blue-gray is 3.29

# Correlations are sensitive to outliers  (use Spearman)!

# Hierarchical clustering: drawing methods

single linkage     average linkage     complete linkage

# Hierarchical clustering (euclidean distance)

calculate distance matrix

|        | gene 1 | gene 2 | gene 3 | gene 4 |
|--------|--------|--------|--------|--------|
| gene 1 | 0      |        |        |        |
| gene 2 | 2      | 0      |        |        |
| gene 3 | 8      | 7      | 0      |        |
| gene 4 | 10     | 12     | 4      | 0      |

calculate averages of most similar

|          | gene 1,2 | gene 3 | gene 4 |
|----------|----------|--------|--------|
| gene 1,2 | 0        |        |        |
| gene 3   | 7.5      | 0      |        |
| gene 4   | 11       | 4      | 0      |

calculate averages of most similar

|          | gene 1,2 | gene 3,4 |
|----------|----------|----------|
| gene 1,2 | 0        |          |
| gene 3,4 | 9.25     | 0        |

# Hierarchical clustering (avg. linkage)

calculate
distance
matrix

|         | gene 1 | gene 2 | gene 3 | gene 4 |
|---------|--------|--------|--------|--------|
| gene 1  | 0      |        |        |        |
| gene 2  | 2      | 0      |        |        |
| gene 3  | 8      | 7      | 0      |        |
| gene 4  | 10     | 12     | 4      | 0      |

calculate averages of
most similar

|          | gene 1,2 | gene 3 | gene 4 |
|----------|----------|--------|--------|
| gene 1,2 | 0        |        |        |
| gene 3   | 7.5      | 0      |        |
| gene 4   | 11       | 4      | 0      |

calculate averages of
most similar

|          | gene 1,2 | gene 3,4 |
|----------|----------|----------|
| gene 1,2 | 0        |          |
| gene 3,4 | 9.25     | 0        |

**Dendrogram**

| 1 | 2 | 3 | 4 |

When assessing similarity, look at the branching pattern instead of sample order

# Bootstrap resampling



Cluster dendrogram with AU/BP values (%)

Distance: correlation
Cluster method: average

- – checks uncertainty in hierarchical cluster analysis
- – AU = approximately unbiased p-value, computed by multiscale bootstrap resampling. Clusters with AU larger than 95% are strongly supported by data.
- – BP = bootstrap probability p-value, computed by normal bootstrap resampling

# K-means clustering

# Quality threshold clustering

# K nearest neighbour clustering

# Exercise 14: Hierarchical clustering

➢ **Cluster genes**

- Select the **column-value-filter.tsv** and run **Clustering / Hierarchical**.

- View the resulting file **hc.tre** as **Hierarchical clustering**.

➢ **Cluster genes and samples**

- Select the **column-value-filter.tsv** and run the tool **Visualization / Heatmap.**

- Select the **column-value-filter.tsv** and run the tool **Visualization / Annotated heatmap,** using parameters
  - Coloring scheme = Blue - white – red
  - Cluster samples only = no

# Microarray data analysis workflow

➢ **Importing data to Chipster**
➢ **Normalization**
➢ **Describing samples with a phenodata file**
➢ **Quality control**
  • Array level
  • Experiment level
➢ **Filtering (optional)**
➢ **Statistical testing**
  • Parametric and non-parametric tests
  • Linear modeling
  • Multiple testing correction
➢ **Annotation**
➢ **Pathway analysis**
➢ **Clustering**
➢ **Saving the workflow**

# Saving and using workflows



➢ **Select the <u>starting point</u> for your workflow**

➢ **Select "Workflow/ Save starting from selected"**

➢ **Save the workflow file on your computer with a meaningful name**
  - Don't change the ending (.bsh)

➢ **To run a workflow on another dataset, select**
  - Workflow → Open and run
  - Workflow → Run recent (if you saved the workflow recently).

# Exercise 15: Saving a workflow

➢ **Prune your workflow if necessary by removing**

- cyclic structures
- files produced by visual selection (gray boxes)

➢ **Save the workflow**

- Select **normalized.tsv** and click on **Workflow / Save starting from selected**. Give your workflow a meaningful name and save it.

# Tools for analysis of miRNA data

➢ **Normalize Agilent miRNA arrays**

- Averages the signal for probes targeting same miRNA
- Ability to exclude control probes (mark as "annotation")
- Support for Human, Mouse and Rat

➢ **Correlate miRNA with target expression**

- Requires corresponding gene expression data set
- Target genes for a list of miRNAs fetched from TargetScan and PicTar
- The correlation between expression of each miRNA and its target genes is calculated with parametric or non-parametric tests
- Finds both **positive** and **negative** correlations

➢ **Up/Down analysis of miRNA targets**

- Requires corresponding gene expression data set
- Target genes for a list of miRNAs fetched from TargetScan and PicTar
- miRNA fold-changes are calculated and **oppositely** behaving target genes are identified

# Exercise 16. Import Agilent miRNA data

➢ **Import folder using Import tool**

- Select **File / Import folder** and select the folder **Agilent_miRNA**
- Set the **Action** to **Use Import tool** for each file and click **OK**
- Click **Mark header** and paint the first **9** rows
- Click **Mark title row** and click on row **10**. Click **Next**.
- Make the Import tool full screen so that you have more space.
- Mark the following columns by clicking first on the button and then <u>in</u> the column so that the column gets colored:
  - Identifier: GeneName
  - Sample: gMeanSignal
  - SampleBG: gBGMedianSignal
  - Annotation: ControlType
- Click **Finish**

# Exercise 17. Normalize Agilent miRNA data

➢ **Normalize Agilent miRNA data and fill in the phenodata**
- Select the **6 files** and the tool **Normalize / Agilent miRNA**, and set
  - **Remove control probes = yes**
  - **Chiptype = Human**.
- The experiment compares miRNA expression in intestinal and diffuse gastric cancer samples. Fill in the phenodata accordingly.
- How would you perform quality control on these samples?

# Microarray data analysis summary

- **Normalization**
  - RMA for Affy
- **Quality control at array level: are there outlier arrays?**
  - RLE, NUSE
- **Quality control at experiment level: do the sample groups separate? Are there batch effects or outliers?**
  - PCA, NMDS, dendrogram
- **(Independent filtering of genes)**
  - e.g. 50% based on coefficient of variation
  - Depends on the statistical test to be used later
- **Statistical testing**
  - Empirical Bayes method (two group test / linear modeling)
- **Annotation, pathway analysis, promoter analysis, clustering, classification…**

# Introduction to RNA-seq

# What can I investigate with RNA-seq?

- ➢ **Differential expression**
- ➢ **Isoform switching**
- ➢ **New genes and isoforms**
- ➢ **New transcripts and transcriptomes**
- ➢ **Variants**
- ➢ **Allele-specific expression**
- ➢ **Etc etc**

# Is RNA-seq better than microarrays?

**+** **Wider detection range**

**+** **Can detect new genes and isoforms**

**-** **Data analysis is not as established as for microarrays**

**-** **Data is voluminous**

A comprehensive assessment of RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control Consortium

SEQC/MAQC-III Consortium

Affiliations | Contributions | Corresponding authors

# How was your data produced?



extraction of poly-A RNAs

PolyA purification

conversion into ds-cDNA and shearing

cDNA generation & fragmentation

amplification and adapter ligation

Library construction

Size selection

sequencing

single end (SET)

paired-end (PET)

# RNA-seq data analysis

# RNA-seq data analysis: typical steps

Raw data (reads)

Align reads to reference genome

**Gene A**          **Gene B**

Match alignment positions with known gene positions

A = 6          B = 11

Count how many reads each gene has

| | Control 1 | Control 2 | Control 3 | Sample 1 | Sample 2 | Sample 3 |
|---|---|---|---|---|---|---|
| **Gene A** | **6** | 5 | 7 | 170 | 100 | 110 |
| **Gene B** | **11** | 11 | 10 | 3 | 4 | 2 |
| **Gene C** | 200 | 150 | 355 | 50 | 1 | 3 |
| **Gene D** | 0 | 1 | 0 | 2 | 0 | 1 |

Compare sample groups: differential expression analysis

# RNA-seq data analysis: steps, tools and files



| STEP | TOOL | FILE |
|------|------|------|
| Quality control | FastQC | FASTQ |
| Pre-processing | Trimmo-matic | FASTQ |
| Alignment | HISAT2 | BAM |
| Quality control | RSeQC | |
| Quantitation | HTSeq | Read count file (TSV) |
| Combine count files to table | Define NGS experiment | Read count table (TSV) |
| Quality control | PCA, clustering | |
| Differential expression analysis | DESeq2, edgeR | Gene lists (TSV) |

|  | Control 1 | Control 2 | Control 3 | Sample 1 | Sample 2 | Sample 3 |
|--------|-----------|-----------|-----------|----------|----------|----------|
| Gene A | 6 | 5 | 7 | 170 | 100 | 110 |
| Gene B | 11 | 11 | 10 | 3 | 4 | 2 |
| Gene C | 200 | 150 | 355 | 50 | 1 | 3 |
| Gene D | 0 | 1 | 0 | 2 | 0 | 1 |

# RNA-seq data analysis workflow

# RNA-seq data analysis workflow

➢ **Quality control of raw reads**
➢ **Preprocessing if needed**
➢ **Alignment (=mapping) to reference genome**
➢ **Alignment level quality control**
➢ **Quantitation**
➢ **Experiment level quality control**
➢ **Differential expression analysis**
➢ **Visualization of reads and results in genomic context**

# What and why?

➤ **Potential problems**

- low confidence bases, Ns

- sequence specific bias, GC bias

- adapters

- sequence contamination

- …

**Knowing about potential problems in your data allows you to**

➤ **correct for them before you spend a lot of time on analysis**

➤ **take them into account when interpreting results**

# Raw reads: FASTQ file format

➢ **Four lines per read:**

@read name

GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT

+ read name

!"*((((***+))%%%++)(%%%%).1***-+*"))**55CCF>>>>>>CCCCCCC65

➢ **http://en.wikipedia.org/wiki/FASTQ_format**

➢ **Attention: Do not unzip FASTQ files**

• Chipster's analysis tools can cope with zipped files (.gz)

# Base qualities

➢ **If the quality of a base is 20, the probability that it is wrong is 0.01.**

- Phred quality score Q = -10 * $\log_{10}$ (probability that the base is wrong)

```
T   C   A   G   T   A   C   T   C   G
40  40  40  40  40  40  40  40  37  35
```

➢ **"Sanger" encoding: numbers are shown as ASCII characters so that 33 is added to the Phred score**

- E.g. 39 is encoded as "H", the 72nd ASCII character (39+33 = 72)

- Note that older Illumina data uses different encoding

  - Illumina1.3: add 64 to Phred
  - Illumina 1.5-1.7: add 64 to Phred, ASCII 66 "B" means that the whole read segment has low quality

# Base quality encoding systems

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...............................

..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL........................................
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmn
 |                                   |       |           |                              |
33                                  59      64          73                            104
 0.................................26...31.......40



 0.2...............................26...31........41

S - Sanger          Phred+33,  raw reads typically (0, 40)



L - Illumina 1.8+ Phred+33,  raw reads typically (0, 41)
```

# Per position base quality (FastQC)



Quality scores across all bases (Illumina 1.5 encoding)

good

ok

bad

Position in read (bp)

# Per position base quality (FastQC)



Quality scores across all bases (Illumina 1.5 encoding)

# Per position sequence content (FastQC)

# Per position sequence content (FastQC)



Sequence content across all bases

➢ **Enrichment of k-mers at the 5' end due to use of random hexamers or transposases in the library preparation**

➢ **Typical for RNA-seq data**

➢ **Can't be corrected, doesn't usually effect the analysis**

# RNA-seq data analysis workflow

➢ **Quality control of raw reads**

➢ **Preprocessing (trimming / filtering) if needed**

➢ **Alignment (=mapping) to reference genome**

➢ **Manipulation of alignment files**

➢ **Alignment level quality control**

➢ **Quantitation**

➢ **Experiment level quality control**

➢ **Visualization of reads and results in genomic context**

➢ **Differential expression analysis**

# Filtering and trimming

➢ **Filtering removes the entire read, trimming removes only the bad quality bases**

- It can remove the entire read, if all bases are bad

➢ **Trimming makes reads shorter**

- This might not be optimal for some applications

➢ **Base quality threshold for trimming is a trade-off between having good quality reads and having enough sequence**

➢ **Paired end data: the matching order of the reads in the two files has to be preserved**

- If a read is removed, its pair has to removed as well

# Was your data made with stranded protocol?

➤ **Several protocols available**

- TruSeq <u>stranded</u>, NEB Ultra <u>Directional</u>, Agilent SureSelect <u>Strand-Specific</u>…

➤ **You need to select the right parameters in tools later on!**

- HISAT2, HTSeq

➤ **See http://chipster.csc.fi/manual/library-type-summary.html**

➤ **The tool <u>Quality control / RNA-seq strandedness inference and inner distance estimation using RseQC</u> allows you to check the type of strandedness in your fastq files**

# RNA-seq data analysis workflow

➤     **Quality control of raw reads**

➤     **Preprocessing (trimming / filtering) if needed**

➤     **Alignment (=mapping) to reference genome**

➤     **Manipulation of alignment files**

➤     **Alignment level quality control**

➤     **Quantitation**

➤     **Experiment level quality control**

➤     **Visualization of reads and results in genomic context**

➤     **Differential expression analysis**

# Alignment to reference genome

➢ **Goal is to find out where a read originated from**

- Challenge: variants, sequencing errors, repetitive sequence

➢ **Many organisms have introns, so RNA-seq reads map to genome non-contiguously → spliced alignments needed**

- But sequence signals at splice sites are limited and introns can be thousands of bases long

➢ **Splice-aware aligners**

- HISAT2, TopHat
- STAR

# HISAT2

- **HISAT = <u>H</u>ierarchical <u>I</u>ndexing for <u>S</u>pliced <u>A</u>lignment of <u>T</u>ranscripts**
- **Fast spliced aligner with low memory requirement**
- **Reference genome is indexed for fast searching**
- **Uses two types of indexes**
  - One global index: used to anchor each alignment (28 bp is enough)
  - Thousands of small local indexes, each covering a genomic region of 56 Kbp: used for rapid extension of alignments (good for reads with short anchors over splice sites)
- **Uses splice site information found during the alignment of earlier reads in the same run**
- **You can use the reference genomes available in Chipster, or provide your own in fasta format**

148

# HISAT2 parameters



| Analysis tools - Alignment - HISAT2 for single end reads | |
|---|---|
| Genome | Homo_sapiens.GR... ▼ |
| Base quality encoding used | Sanger - Phred+33 ▼ |
| Minimum intron length | 20 |
| Maximum intron length | 500000 |
| Library type | fr-unstranded ▼ |
| How many hits is a read allowed to have | 5 |
| Disallow soft-clipping | Use soft-clipping ▼ |
| Require long anchor lengths for subsequent assembly | Don't require ▼ |

➢ **Remember to set the strandedness (library type) correctly!**

➢ **Require long anchors (> 16 bp) if you are going to do transcript assembly**

➢ **Note that there can an alignment that is better than the 5 reported ones**

➢ **Soft-clipping = read ends don't need to align to the genome, if this maximizes the alignment score**

149

# Do you have several FASTQ files <u>per sample</u>?

➢ **Align all of them in one HISAT2 run**

➢ **Single end data: Select all the (zipped) FASTQ files for the sample**

➢ **Paired end data: Make filename lists first**
- Select all read1 files and run the tool "Utilities / Make a list of file names"
- Select all read2 files and and do as above
- Select the FASTQ files and the filename lists and run HISAT2 (check that the files have been assigned correctly)

# STAR

- ➢ **STAR (Spliced Transcripts Alignment to a Reference) uses a 2-pass mapping process**
  - splice junctions found during the 1st pass are inserted into the genome index, and all reads are re-mapped in the 2nd mapping pass
  - this doesn't increase the number of detected novel junctions, but it allows more spliced reads mapping to novel junctions.
- ➢ **Maximum alignments per read -parameter sets the maximum number of loci the read is allowed to map to**
  - Alignments (all of them) will be output only if the read maps to no more loci than this. Otherwise no alignments will be output.
- ➢ **Chipster offers an Ensembl GTF file to detect annotated splice junctions**
  - you can also give your own. For example the GENCODE GTF
- ➢ **Two log files**
  - Log_final.txt lists the percentage of uniquely mapped reads etc.
  - Log_progress.txt contains process summary

# Mapping quality

➢ **Confidence in read's point of origin**

➢ **Depends on many things, including**

- uniqueness of the aligned region in the genome

- length of alignment

- number of mismatches and gaps

➢ **Expressed in Phred scores, like base qualities**

- $Q = -10 * \log_{10}$ (probability that mapping location is wrong)

➢ **Values differ in different aligners. E. g. unique mapping is**

- 60 in HISAT2

- 255 in STAR

- 50 in TopHat

- https://sequencing.qcfail.com/articles/mapq-values-are-really-useful-but-their-implementation-is-a-mess/

# File format for mapped reads: BAM/SAM

➢ SAM (Sequence Alignment/Map) is a tab-delimited text file containing aligned reads. BAM is a binary (and hence more compact) form of SAM.



Visualisation

Method: BAM viewer ▼    Redraw    Restore

```
@HD     VN:1.4    SO:coordinate
@SQ     SN:chr1   LN:249250621
@SQ     SN:chr10  LN:135534747
@SQ     SN:chr11  LN:135006516
@SQ     SN:chr12  LN:133851895
@SQ     SN:chr13  LN:115169878
@SQ     SN:chr14  LN:107349540
@SQ     SN:chr15  LN:102531392
@SQ     SN:chr16  LN:90354753
@SQ     SN:chr17  LN:81195210
@SQ     SN:chr18  LN:78077248
@SQ     SN:chr19  LN:59128983
@SQ     SN:chr2   LN:243199373
@SQ     SN:chr20  LN:63025520
@SQ     SN:chr21  LN:48129895
@SQ     SN:chr22  LN:51304566
@SQ     SN:chr3   LN:198022430
@SQ     SN:chr4   LN:191154276
@SQ     SN:chr5   LN:180915260
@SQ     SN:chr6   LN:171115067
@SQ     SN:chr7   LN:159138663
@SQ     SN:chr8   LN:146364022
@SQ     SN:chr9   LN:141213431
@SQ     SN:chrM   LN:16571
@SQ     SN:chrX   LN:155270560
@SQ     SN:chrY   LN:59373566
@PG     ID:TopHat VN:2.0.9  CL:/opt/chipster/tools/tophat2/tophat -p 2 --read-mismatches 2 -a 8 -m 0 -i 70 -I 500000 -g 20 --library-type fr-unstranded
--transcriptome-index=/opt/chipster/tools/bowtie2/indexes/hg19.ti --no-novel-juncs /opt/chipster/tools/bowtie2/indexes/hg19 reads1.fq
HWI-EAS229_1:4:82:1371:1147    272    chr1    18378    1    2M6358N73M*    0    0
TCCTGCTGAAGATGTCTCCAGAGACCTTCTGCAGGTACTGAAGGGCATCCGCCATCTGCTGGACGGCCTCCTCTC    5661525416816488666(6(6(6261?8==(B=513);(/BB=141=>6>?=<=?B>9B?>BA<66>BA>BBB
CC:Z:chr15 MD:Z:40C34 XG:i:0    NH:i:3    HI:i:0    NM:i:1    XM:i:1    XN:i:0    XO:i:0    CP:i:102506354    AS:i:0    XS:A:-    YT:Z:UU
```

optional header section

alignment section: one line per read, containing 11 mandatory fields, followed by optional tags

# Fields in BAM/SAM files

- **read name**       HWI-EAS229_1:2:40:1280:283
- **flag**            272
- **reference name**  1
- **position**        18506
- **mapping quality** 0
- **CIGAR**           49M6183N26M
- **mate name**       *
- **mate position**   0
- **insert size**     0
- **sequence**
  AGGGCCGATCTTGGTGCCATCCAGGGGGCCTCTACAAGGAT
  AATCTGACCTGCTGAAGATGTCTCCAGAGACCTT
- **base qualities**
  ECC@EEF@EB:EECFEECCCBEEEE;>5;2FBB@FBFEEFCF@F
  FFFCEFFFFEE>FFEFC=@A;@>1@6.+5/5
- **tags**            MD:Z:75  NH:i:7  AS:i:-8  XS:A:-

# BAM index file (.bai)

➢ **BAM files can be sorted by chromosomal coordinates and indexed for efficient retrieval of reads for a given region.**

➢ **The index file must have a matching name (e.g. reads.bam and reads.bam.bai)**

➢ **Genome browser requires both BAM and the index file.**

➢ **The alignment tools in Chipster automatically produce sorted and indexed BAMs.**

➢ **When you import BAM files, Chipster asks if you would like to preproces them (convert SAM to BAM, sort and index BAM).**

# RNA-seq data analysis workflow

- ➢ **Quality control of raw reads**
- ➢ **Preprocessing (trimming / filtering) if needed**
- ➢ **Alignment (=mapping) to reference genome**
- ➢ **Alignment level quality control**
- ➢ **Quantitation**
- ➢ **Experiment level quality control**
- ➢ **Differential expression analysis**
- ➢ **Visualization of reads and results in genomic context**

# Annotation-based quality metrics

- ➤ **Saturation of sequencing depth**
  - Would more sequencing detect more genes and splice junctions?
- ➤ **Read distribution between different genomic features**
  - Exonic, intronic, intergenic regions
  - Coding, 3' and 5' UTR exons
  - Protein coding genes, pseudogenes, rRNA, miRNA, etc
- ➤ **Is read coverage uniform along transcripts?**
  - Biases introduced in library construction and sequencing
    - polyA capture and polyT priming can cause 3' bias
    - random primers can cause sequence-specific bias
    - GC-rich and GC-poor regions can be under-sampled
  - Genomic regions have different mappabilities (uniqueness)

# Quality assessment with RseQC

➤ **Checks coverage uniformity, saturation of sequencing depth, novelty of splice junctions, read distribution between different genomic regions, etc.**

➤ **Takes a BAM file and a BED file**

- Chipster has BED files available for several organisms
- You can also use your own BED if you prefer

# BED file format

➢ **BED (Browser extensible data) file format is used for reporting location of features (e.g. genes and exons) in a genome**

➢ **5 obligatory columns: chr, start, end, name, score**

➢ **0-based, like BAM**

| column0 | column1 | column2 | column3 | column4 |
|---------|---------|---------|--------------|----|
| chr22 | 21022480 | 21024796 | JUNC00000001 | 1 |
| chr19 | 201609 | 201783 | JUNC00000002 | 5 |
| chr19 | 281478 | 282180 | JUNC00000003 | 3 |
| chr19 | 282242 | 282811 | JUNC00000004 | 21 |
| chr19 | 282751 | 287541 | JUNC00000005 | 37 |
| chr19 | 287705 | 288084 | JUNC00000006 | 6 |
| chr19 | 288105 | 291354 | JUNC00000007 | 18 |
| chr19 | 307484 | 308600 | JUNC00000008 | 1 |
| chr19 | 308603 | 308858 | JUNC00000009 | 2 |
| chr19 | 308868 | 311907 | JUNC00000010 | 13 |
| chr19 | 311872 | 312256 | JUNC00000011 | 26 |
| chr19 | 312205 | 313558 | JUNC00000012 | 22 |
| chr19 | 313575 | 325706 | JUNC00000013 | 68 |

# Own BED? Check chromosome names

➢ **RseQC needs the same chromosome naming in BAM and BED**

➢ **Chromosome names in BED files can have the prefix "chr"**

- e.g. chr1

➢ **Chipster BAM files are Ensembl-based and don't have the prefix**

- If you use your own BED (e.g. from UCSC Table browser) you need to remove the prefix (chr1 → 1)

➢ **Use the tool Utilities / Modify text with the following parameters:**

- Operation = Replace text
- Search string = chr
- Input file format = BED

# QC tables by RseQC

```
#================================================
#All numbers are READ count  (alignment, actually…)
#================================================

Total records:                          103284

QC failed:                              0
Optical/PCR duplicate:                  0
Non primary hits:                       18476
Unmapped reads:                         0
mapq < mapq_cut (non-unique):           4208
     Default=30
mapq >= mapq_cut (unique):              80600
Read-1:                                 0
Read-2:                                 0
Reads map to '+':                       48292
Reads map to '-':                       32308
Non-splice reads:                       50919
Splice reads:                           29681
Reads mapped in proper pairs:           0
Proper-paired reads map to different chrom:0
```

```
read_distribution:

Total Reads              84808
Total Tags               116738
Total Assigned Tags      111352
================================================================
Group          Total_bases        Tag_count        Tags/Kb
CDS_Exons      2211343            90961            41.13
5'UTR_Exons    529860             1662             3.14
3'UTR_Exons    1415234            12423            8.78
Introns        25801210           5349             0.21
TSS_up_1kb     1295771            31               0.02
TSS_up_5kb     5332522            321              0.06
TSS_up_10kb    8804879            584              0.07
TES_down_1kb   1292506            217              0.17
TES_down_5kb   5108821            344              0.07
TES_down_10kb  8282641            373              0.05
================================================================
```

| | |
|---|---|
| Total records: | 7 |
| Non primary hits: | 4 |
| Total reads: | 3 |
| Total tags: | 8 |

**Read B**    Read B

Read A    Read A    Read A    **Read A**    **Read C**

Reference

# Did I accidentally sequence ribosomal RNA?

➢ **The majority of RNA in cells is rRNA**

➢ **Typically we want to sequence protein coding genes, so we try to avoid rRNA**

- polyA capture
- Ribominus kit (may not work consistently between samples)

➢ **How to check if we managed to avoid rRNA?**

- RseQC might not be able to tell, if the rRNA genes are not in the BED file (e.g. in human the rRNA gene repeating unit has not been assigned to any chromosome yet)

- You can map the reads to human ribosomal DNA repeating unit sequence (instead of the genome) with the Bowtie aligner, and check the alignment percentage

# RNA-seq data analysis workflow

➢ **Quality control of raw reads**
➢ **Preprocessing (trimming / filtering) if needed**
➢ **Alignment (=mapping) to reference genome**
➢ **Alignment level quality control**
➢ **Quantitation**
➢ **Experiment level quality control**
➢ **Differential expression analysis**
➢ **Visualization of reads and results in genomic context**

# Counting reads per genes with HTSeq

➢ **Given a BAM file and a list of gene locations, counts how many reads map to each gene.**

- A gene is considered as the union of all its exons.
- Reads can be counted also per exons.

➢ **Locations need to be supplied in GTF file**

- Note that GTF and BAM must use the same chromosome naming

➢ **Multimapping reads and ambiguous reads are not counted**

➢ **3 modes to handle reads which overlap several genes**

- Union (default), Intersection-strict, Intersection-nonempty

➢ **Attention: was your data made with stranded protocol?**

- You need to select the right counting mode!

# Not unique or ambiguous?



Ambiguous

Stranded data
→ Not ambiguous

Multimapping
(not unique)

# HTSeq count modes

| | union | intersection_strict | intersection_nonempty |
|---|---|---|---|
| read / gene_A | gene_A | gene_A | gene_A |
| read / gene_A | gene_A | no_feature | gene_A |
| read / gene_A — gene_A | gene_A | no_feature | gene_A |
| read — read / gene_A — gene_A | gene_A | gene_A | gene_A |
| read / gene_A / gene_B | gene_A | gene_A | gene_A |
| read / gene_A / gene_B | ambiguous | gene_A | gene_A |
| read / gene_A / gene_B | ambiguous | ambiguous | ambiguous |

# GTF file format

➢ **9 obligatory columns: chr, source, name, start, end, score, strand, frame, attribute**

➢ **1-based**

➢ **For HTSeq to work, all exons of a gene must have the same gene_id**
  - Use GTFs from Ensembl, avoid UCSC

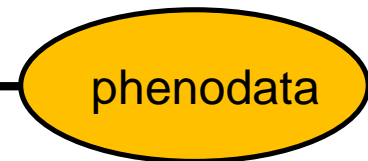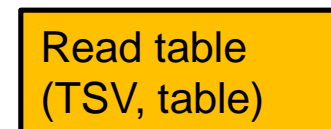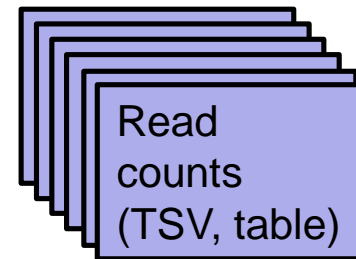| chr1 | unknown | exon | 14362 | 14829 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
|------|---------|------|-------|-------|---|---|---|---|
| chr1 | unknown | exon | 14970 | 15038 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 15796 | 15947 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 16607 | 16765 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 16858 | 17055 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 17233 | 17368 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 17606 | 17742 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 17915 | 18061 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 18268 | 18366 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 24738 | 24891 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |
| chr1 | unknown | exon | 29321 | 29370 | . | – | . | gene_id "WASH7P"; gene_name "WASH7P"; transcript_id "NR_024540"; tss_id "TSS7245"; |

# Combine individual count files into a count table

➢ **Select all the count files and run "Utilities / Define NGS experiment"**

➢ **This creates a table of counts and a phenodata file, where you can describe experimental groups**

# Phenodata file: describe the experiment

➢ **Describe experimental groups, time, pairing etc <u>with numbers</u>**
- e.g. 1 = control, 2 = cancer

➢ **Define sample names for visualizations in the Description column**

| sample | original_name | description | patient | group | treatment | time | hours |
|--------|---------------|-------------|---------|-------|-----------|------|-------|
| ngs001.tsv | SRR479052 | 1_C_24 | 1 | 1 | Control | 1 | 24h |
| ngs002.tsv | SRR479053 | 1_C_48 | 1 | 1 | Control | 2 | 48h |
| ngs003.tsv | SRR479054 | 1_DP_24 | 1 | 2 | DPN | 1 | 24h |
| ngs004.tsv | SRR479055 | 1_DP_48 | 1 | 2 | DPN | 2 | 48h |
| ngs007.tsv | SRR479058 | 2_C_24 | 2 | 1 | Control | 1 | 24h |
| ngs008.tsv | SRR479059 | 2_C_48 | 2 | 1 | Control | 2 | 48h |
| ngs009.tsv | SRR479060 | 2_DP_24 | 2 | 2 | DPN | 1 | 24h |
| ngs011.tsv | SRR479062 | 2_DP_48 | 2 | 2 | DPN | 2 | 48h |
| ngs015.tsv | SRR479066 | 3_C_24 | 3 | 1 | Control | 1 | 24h |
| ngs016.tsv | SRR479067 | 3_C_48 | 3 | 1 | Control | 2 | 48h |
| ngs017.tsv | SRR479068 | 3_DP_24 | 3 | 2 | DPN | 1 | 24h |
| ngs018.tsv | SRR479069 | 3_DP_48 | 3 | 2 | DPN | 2 | 48h |

# What if somebody gives you a count table?

➤ **Make sure that the filename ending is tsv**
➤ **When importing the file to Chipster select "Use Import tool"**
➤ **In Import tool**
- Mark the title row
- Mark the identifier column and the count columns

➤ **Select the imported files and run the tool "Utilities / Preprocess count table"**
- This creates a count table and a phenodata file for it

# RNA-seq data analysis workflow

➢ **Quality control of raw reads**

➢ **Preprocessing (trimming / filtering) if needed**

➢ **Alignment (=mapping) to reference genome**

➢ **Alignment level quality control**

➢ **Quantitation**

➢ **Experiment level quality control**

➢ **Differential expression analysis**

➢ **Visualization of reads and results in genomic context**

# Experiment level quality control

➢ **Getting an overview of similarities and dissimilarities between samples allows you to check**

- Do the experimental groups separate from each other?
- Is there a confounding factor (e.g. batch effect) that should be taken into account in the statistical analysis?
- Are there sample outliers that should be removed?

➢ **Several methods available**

- MDS (multidimensional scaling)
- PCA (principal component analysis)
- Clustering

# PCA plot by DESeq2



➢ **The first two principal components, calculated after variance stabilizing transformation**

➢ **Indicates the proportion of variance explained by each component**

- If PC2 explains only a small percentage of variance, it can be ignored

# MDS plot by edgeR

➢ **Distances correspond to the logFC or biological coefficient of variation (BCV) between each pair of samples**

➢ **Calculated using 500 most heterogenous genes (= have largest dispersion when treating all samples as one group)**

# Sample heatmap by DESeq2



➢ **Euclidean distances between the samples, calculated after variance stabilizing transformation**

# Exercise 18. Experiment-level QC of RNA-seq

➢ **Description of exercise data:**

- RNAseq data from a Drosophila experiment where the gene for splicing factor Pasilla was knocked down with RNAi. There are 4 control and 3 knock-down samples. We want to find differentially expressed genes. Note that some samples were sequenced single end and some paired end.

- Reads were aligned to the Drosophila genome and counted per genes. The count table contains all the 7 samples, described in the phenodata file.

➢ **Open session**

- Select **Open local session** and **RNAseq_drosophila.zip**.

- Read the experiment description in the **phenodata.tsv**.

- Look at the contents of the file **counts.tsv**.

➢ **Perform experiment-level quality control**

- Select **counts.tsv** and **Quality control / PCA and heatmap of samples with DESeq2.** Do the groups separate along PC1? How much variance does PC2 explain?

- Repeat the run so that you set **Phenodata column for the shape of samples in PCA plot = readtype**. What does the PCA plot tell you?
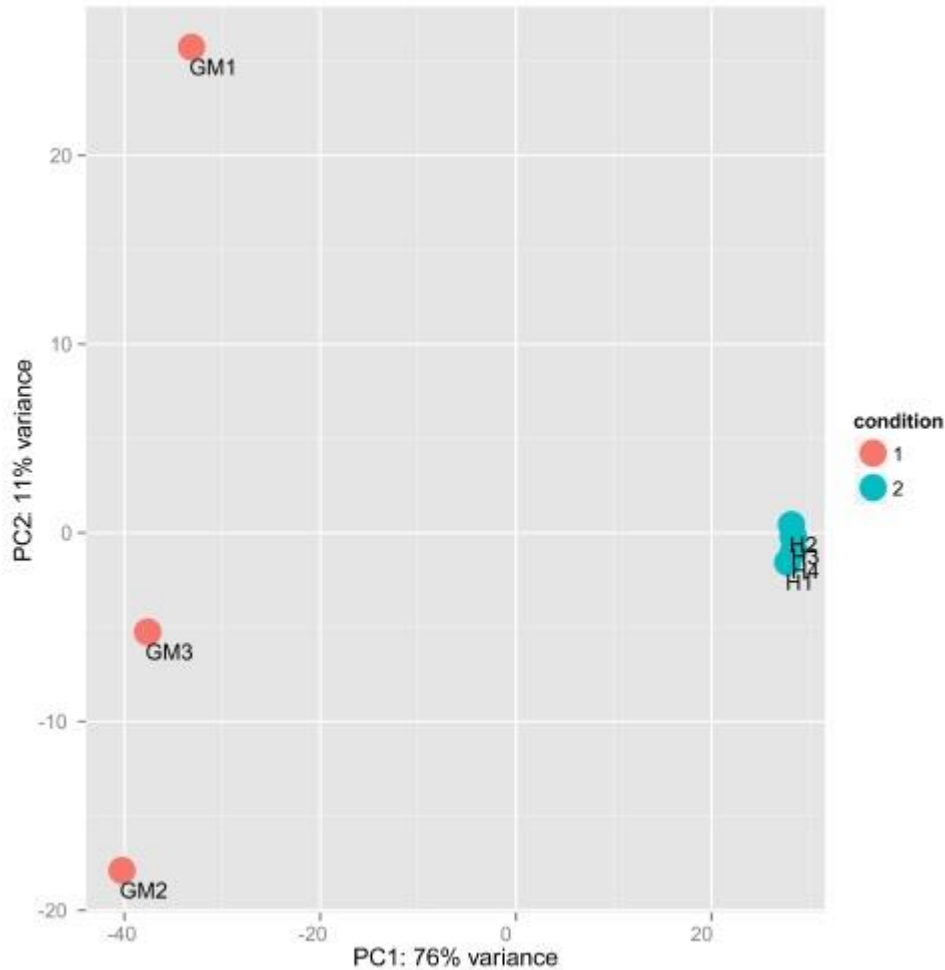
# RNA-seq data analysis workflow

- ➢ **Quality control of raw reads**
- ➢ **Preprocessing (trimming / filtering) if needed**
- ➢ **Alignment (=mapping) to reference genome**
- ➢ **Manipulation of alignment files**
- ➢ **Alignment level quality control**
- ➢ **Quantitation**
- ➢ **Experiment level quality control**
- ➢ **Differential expression analysis**
- ➢ **Visualization of reads and results in genomic context**

# Statistical testing for differential expression: things to take into account

➢ **Biological replicates are important!**

➢ **Normalization is required in order to compare expression between samples**

- Different library sizes

- RNA composition bias caused by sampling approach

➢ **Raw counts are needed to assess measurement precision**

- Counts are the "the units of evidence" for expression

- No FPKMs thanks!

➢ **Multiple testing problem**

# Differential gene expression analysis

➢ **Normalization**

➢ **Dispersion estimation**

➢ **Log fold change estimation**

➢ **Statistical testing**

➢ **Filtering**

➢ **Multiple testing correction**

# Normalization by edgeR and DESeq

➢ **Aim to make normalized counts for non-differentially expressed genes similar between samples**

   • Do not aim to adjust count distributions between samples

➢ **Assume that**

   • Most genes are not differentially expressed

   • Differentially expressed genes are divided equally between up- and down-regulation

➢ **Do not transform data, but use normalization factors within statistical testing**

# Normalization by edgeR and DESeq – how?

- ➢ **DESeq(2)**
  - Take geometric mean of gene's counts across all samples
  - Divide gene's counts in a sample by the geometric mean
  - Take median of these ratios → sample's normalization factor (applied to read counts)

- ➢ **edgeR**
  - Select as reference the sample whose upper quartile is closest to the mean upper quartile
  - Log ratio of gene's counts in sample vs reference → M value
  - Take weighted trimmed mean of M-values (TMM) → normalization factor (applied to library sizes)
    - Trim: Exclude genes with high counts or large differences in expression
    - Weights are from the delta method on binomial data

Differential expression analysis:
Dispersion estimation

# Dispersion

➢ **When comparing gene's expression levels between groups, it is important to know also its within-group variability**

➢ **Dispersion = (BCV)$^2$**
  - BCV = gene's biological coefficient of variation
  - E.g. if gene's expression typically differs from replicate to replicate by 20% (so BCV = 0.2), then this gene's dispersion is $0.2^2$ = 0.04

➢ **Note that the variability seen in counts is a sum of 2 things:**
  - Sample-to-sample variation (dispersion)
  - Uncertainty in measuring expression by counting reads
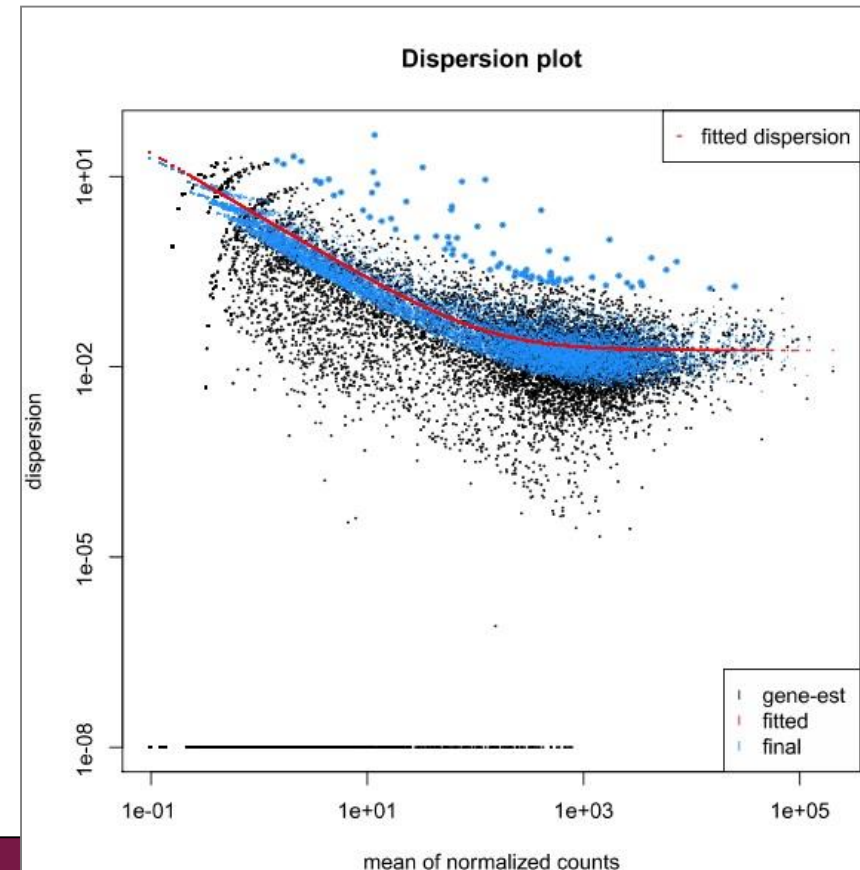
# How to estimate dispersion reliably?

➢ **RNA-seq experiments typically have only few replicates**
   **→ it is difficult to estimate within-group variability**

➢ **Solution: pool information across genes which are expressed at similar level**
   • assumes that genes of similar average expression strength have similar dispersion

➢ **Different approaches**
   • edgeR
   • DESeq2

# Dispersion estimation by DESeq2

➢ **Estimates genewise dispersions using maximum likelihood**

➢ **Fits a <span style="color:red">curve</span> to capture the dependence of these estimates on the average expression strength**

➢ **Shrinks <span style="color:blue">genewise values towards the curve</span> using an empirical Bayes approach**

- The amount of shrinkage depends on several things including sample size

- Genes with high gene-wise dispersion estimates are dispersion outliers (blue circles above the cloud) and they are not shrunk



**Dispersion plot**

# Differential expression analysis:
# Statistical testing

# Generalized linear models

➢ **Model the expression of each gene as a linear combination of explanatory factors (eg. group, time, patient)**

- y = a + (b · group) + (c · time) + (d · patient) + e

  y = gene's expression

  a, b, c and d = parameters estimated from the data

  a = intercept (expression when factors are at reference level)

  e = error term

➢ **<u>Generalized</u> linear model (GLM) allows the expression value distribution to be different from normal distribution**

- Negative binomial distribution used for count data
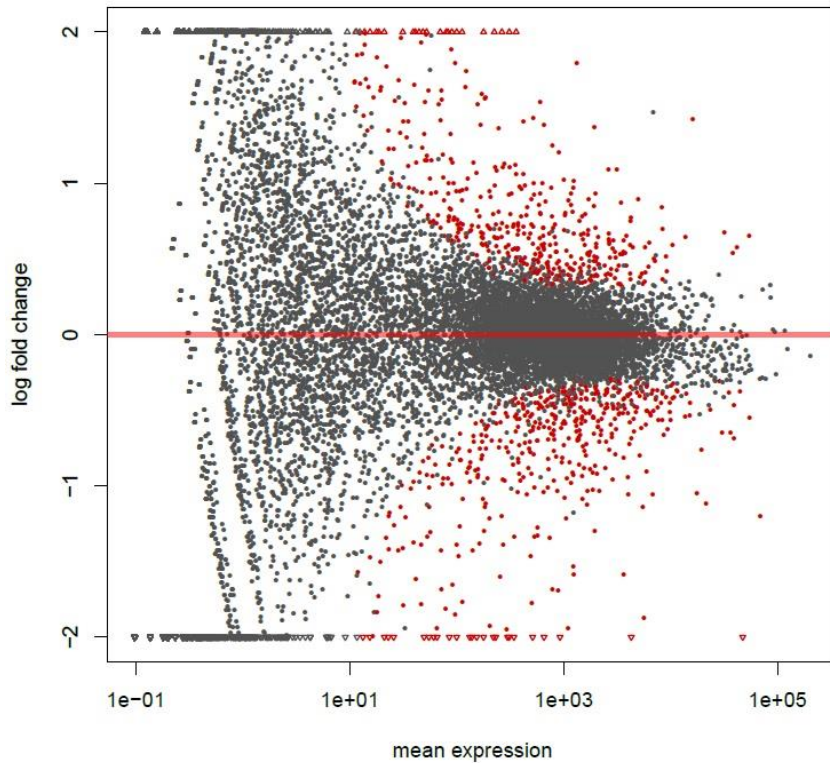
# Statistical testing

- ➢ **edgeR**
  - Two group comparisons
    - Exact test for negative binomial distribution.
  - Multifactor experiments
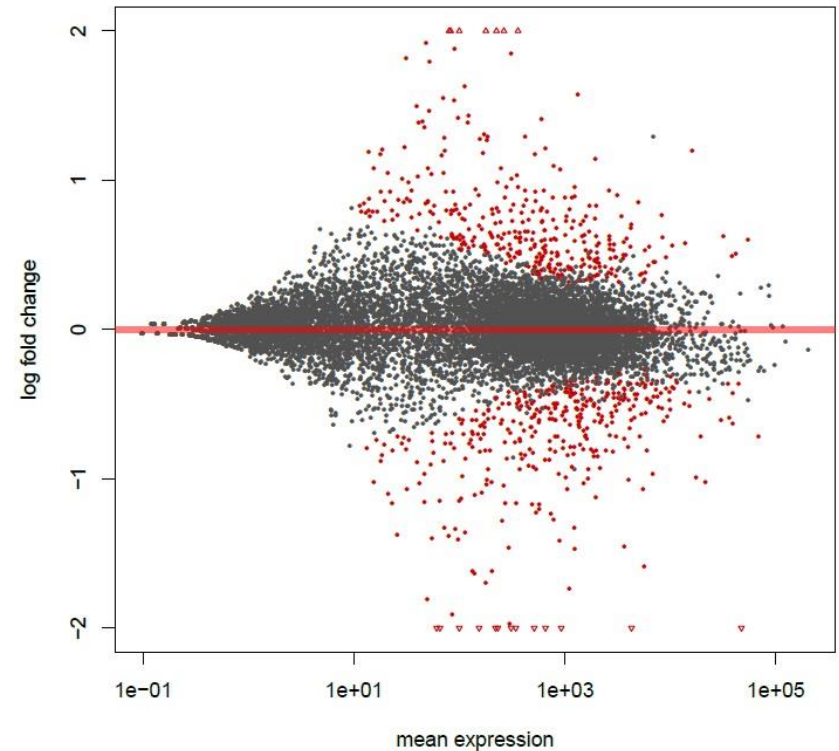    - Generalized linear model, likelyhood ratio test.

- ➢ **DESeq2**
  - Shrinks log fold change estimates toward zero using an empirical Bayes method
    - Shrinkage is stronger when counts are low, dispersion is high, or there are only a few samples
  - Generalized linear model, Wald test for significance
    - Shrunken estimate of log fold change is divided by its standard error and the resulting z statistic is compared to a standard normal distribution

# Fold change shrinkage by DESeq2

# Multiple testing correction

➤ **We tests thousands of genes, so it is possible that some genes get good p-values just by chance**

➤ **To control this problem of false positives, p-values need to be corrected for multiple testing**

➤ **Several methods are available, the most popular one is the Benjamini-Hochberg correction (BH)**
  - largest p-value is not corrected
  - second largest p = (p *n)/ (n-1)
  - third largest p = (p * n)/(n-2)
  - …
  - smallest p = (p * n)/(n- n+1) = p * n

➤ **The adjusted p-value is FDR (false discovery rate)**

# Filtering

- ➢ **Reduces the severity of multiple testing correction by removing some genes (makes n smaller)**

- ➢ **Filter out genes which have little chance of showing evidence for significant differential expression**
  - genes which are not expressed
  - genes which are expressed at very low level (low counts are unreliable)

- ➢ **Should be independent**
  - do not use information on what group the sample belongs to

- ➢ **DESeq2 selects filtering threshold automatically**

# DSeq2 result table

➤ baseMean = mean of counts (divided by size factors) taken over all samples

➤ **log2FoldChange = log2 of the ratio meanB/meanA**

➤ lfcSE = standard error of log2 fold change

➤ stat = Wald statistic

➤ pvalue = raw p-value

➤ **<u>padj = Benjamini-Hochberg adjusted p-value</u>**

| | baseMean | log2FoldChange | lfcSE | stat | pvalue | padj |
|---|---|---|---|---|---|---|
| FBgn0026562 | 47282.42 | -2.4 | 0.08 | -30.26 | 4.159e-201 | 3.309e-197 |
| FBgn0039155 | 924.27 | -4.46 | 0.16 | -27.04 | 4.476e-161 | 1.781e-157 |
| FBgn0029167 | 4287.44 | -2.21 | 0.08 | -26.75 | 1.107e-157 | 2.937e-154 |
| FBgn0035085 | 654.94 | -2.5 | 0.11 | -22.08 | 5.278e-108 | 1.050e-104 |
| FBgn0034736 | 231.7 | -3.29 | 0.18 | -18.28 | 1.261e-74 | 2.006e-71 |
| FBgn0000071 | 359.53 | 2.6 | 0.14 | 17.98 | 2.741e-72 | 3.635e-69 |
| FBgn0034434 | 153.84 | -3.69 | 0.21 | -17.26 | 9.008e-67 | 1.024e-63 |
| FBgn0039827 | 342.77 | -3.83 | 0.23 | -16.54 | 1.742e-61 | 1.733e-58 |
| FBgn0029896 | 513.08 | -2.34 | 0.14 | -16.29 | 1.168e-59 | 1.033e-56 |
| FBgn0052407 | 220.26 | -2.2 | 0.15 | -14.99 | 8.597e-51 | 6.841e-48 |
| FBgn0037754 | 299.03 | -2.23 | 0.15 | -14.94 | 1.916e-50 | 1.386e-47 |

# Exercise 19. Find differentially expressed genes

- ➢ Select **counts.tsv** and **RNA-seq / Differential expression with DESeq2.**
  - • Open **de-list-deseq2.tsv** and check how many differentially expressed genes do you get.
  - • Check in **summary.txt** how many genes were filtered out because they had too low counts and what was the low count threshold used?
- ➢ **Correct for the confounding factor (different sequencing type)**
  - • Run as before but set **Column describing additional experimental factor = readtype**.
  - • How many differentially expressed genes do you get now?
  - • Did the low count filtering threshold change?

# Summary of differential expression analysis steps and files

➢ **Quality control / Read quality with FastQC** → **html report**

➢ (Preprocessing / Trim reads with Trimmomatic → FASTQ)

➢ (Utilities / Make a list of file names → txt)

➢ **Alignment / HISAT2 for paired end reads** → **BAM**

➢ **Quality control / RNA-seq quality metrics with RseQC** → **pdf**

➢ **RNA-seq / Count aligned reads per genes with HTSeq** → **tsv**

➢ **Utilities / Define NGS experiment** → **tsv**

➢ **Quality control / PCA and heatmap of samples with DESeq2** → **pdf**

➢ **RNA-seq / Differential expression using DESeq2** → **tsv**

➢ **Utilities / Annotate Ensembl identifiers** → **tsv**

# RNA-seq data analysis workflow

➢ **Quality control of raw reads**

➢ **Preprocessing (trimming / filtering) if needed**

➢ **Alignment (=mapping) to reference genome**

➢ **Alignment level quality control**

➢ **Quantitation**

➢ **Experiment level quality control**

➢ **Differential expression analysis**

➢ **Visualization of reads and results in genomic context**

# Chipster Genome Browser

- ➢ **Integrated with Chipster analysis environment**
- ➢ **Automatic sorting and indexing of BAM, BED and GTF files**
- ➢ **Automatic coverage calculation (total and strand-specific)**
- ➢ **Zoom in to nucleotide level**
- ➢ **Highlight variants**
- ➢ **Jump to locations using BED, GTF, VCF and tsv files**
- ➢ **View details of selected BED, GTF and VCF features**
- ➢ **Several  views (reads, coverage profile, density graph)**