

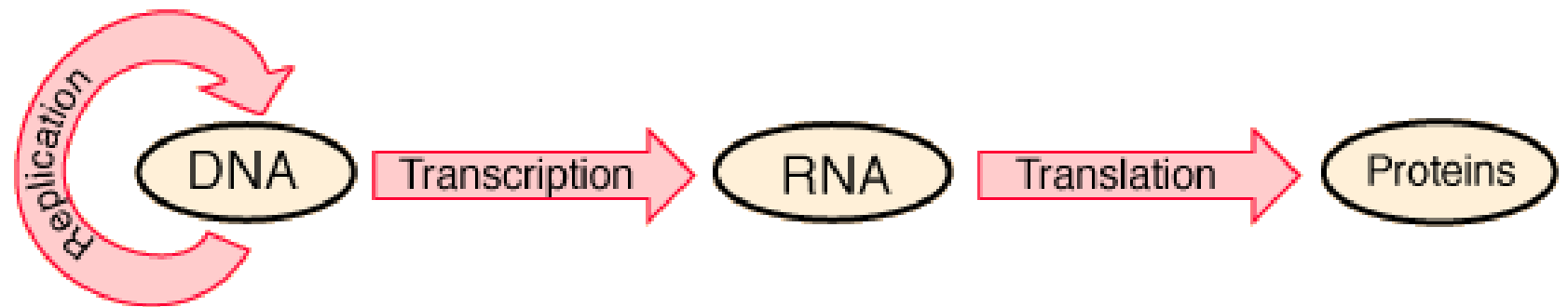
An introduction to gene expression analysis by RNA-seq: gene-level exploratory analysis and differential expression

Paolo Sonogo
Computational Biology Unit, Fondazione Edmund Mach
Via E. Mach 1, 38010, S. Michele a/A, Tn, Italy
office phone: +39 0461 615 645
mail: paolo.sonogo@fmach.it

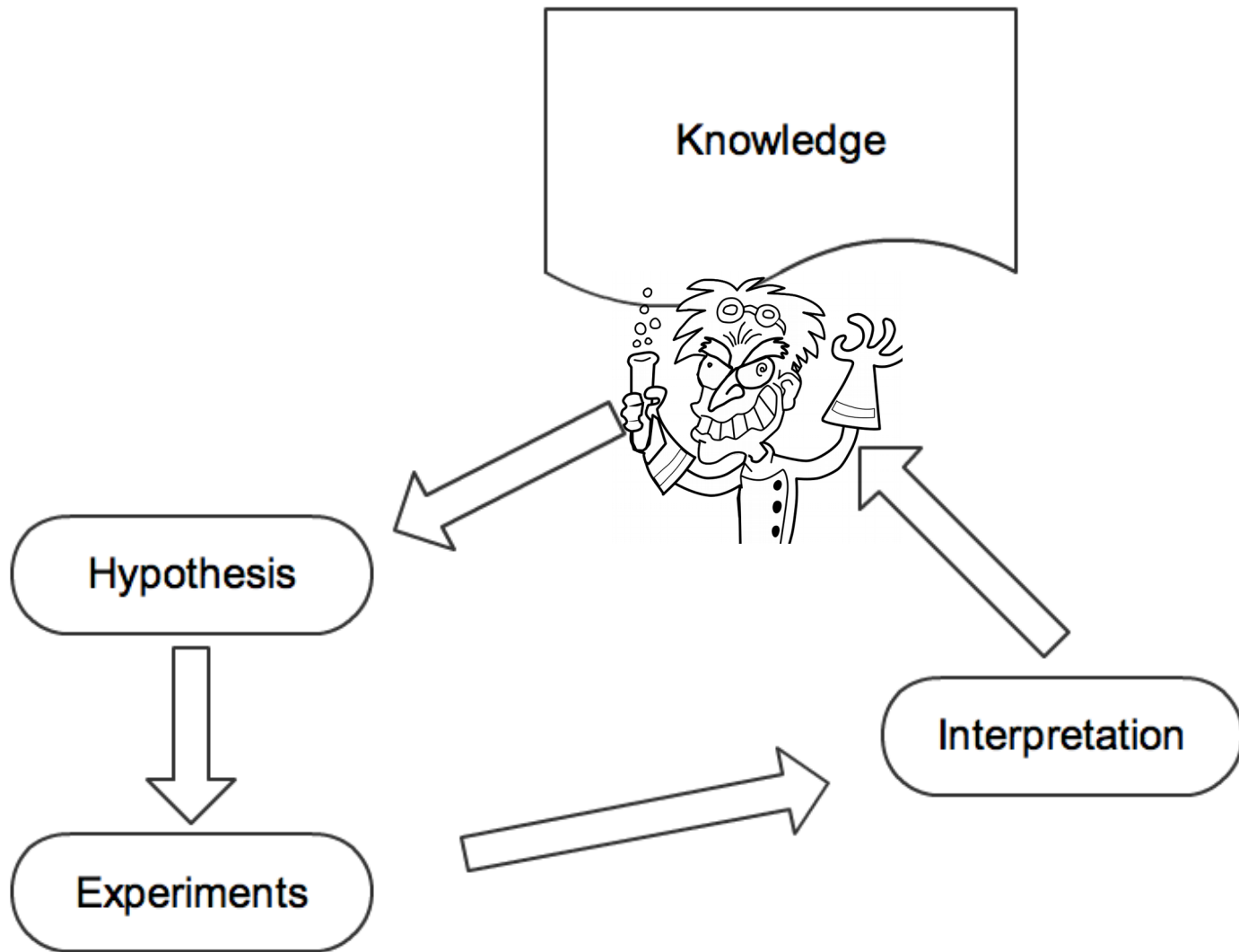


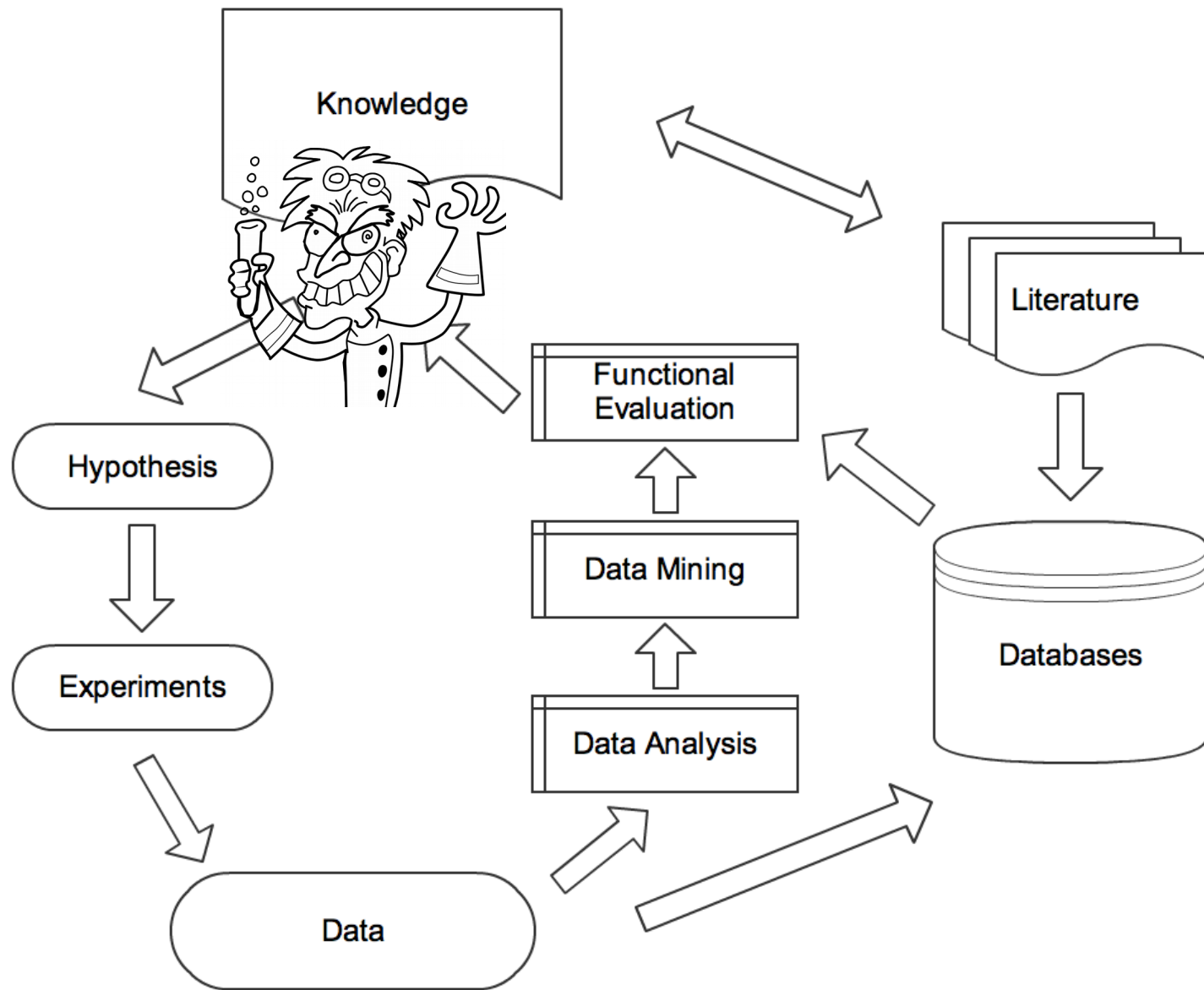
FONDAZIONE
EDMUND
MACH





Central dogma of molecular biology

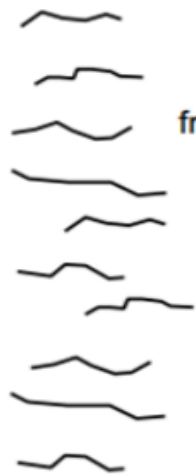




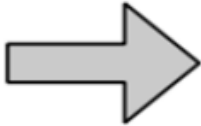
RNA-Seq

- RNA-seq works by sequencing every RNA molecule and profiling the expression of a particular gene by counting the number of time its transcripts have been sequenced.
- RNA or DNA is extracted from sample tissue/cells and fragmented. RNA is converted to cDNA by reverse transcription. DNA Fragments are converted into the library by ligation to sequencing adapters containing specific sequences designed to interact with the NGS platform. The next step involves clonal amplification of the library. The final step generates the actual sequence via the chemistries for each technology.

Sample
RNA



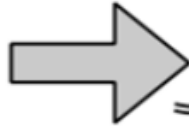
fragmentation



RNA
fragments



reverse
transcription +
amplification



cDNA
fragments



sequencing
machine



reads

CCTTCNCACTTCGTTTCCCAC
TTTTTNCAGAGTTTTTCTTG
GAACANTCCAACGCTTGGTGA
GGAAANAAGACCCTGTTGAGC
CCCGGNGATCCGCTGGGACAA
GCAGCATATTGATAGATAACT
CTAGCTACGCGTACGCGATCG
CATCTAGCATCGCGTTGCGTT
CCCGCGCGCTTAGGCTACTCG
TCACACATCTCTAGCTAGCAT
CATGCTAGCTATGCCTATCTA

Fastq files (raw data)

```
@SRR5227652.1 1 length=101
CCAGTGCCTTATTGACCTCAGATTTTCNNNNNNNNNNNNNNNNNNNNNNNNNNNNNTCCACGCANNNNNNNNNNNNAAATCATCAACAGACTCATCA
+SRR5227652.1 1 length=101
BBBF<FFFFFFFFFIIIBFIIFFFFF#####
@SRR5227652.2 2 length=101
CACTTGGATATGTTAATCTGTAGTCGATTANCNNNNNNNNNNNGAGAGGAATCANNNNNAAAAGGCCAANNNNNNNNNNAGAAAGAGGAAGCCAAAGATC
+SRR5227652.2 2 length=101
BBBFFFFFFFFFIIIIIIIIIIIFIIIIII#0#####00BFFIIFIIF####00<BFFFFBF#####00<BBFFFFFFFFFBFFFFFB
@SRR5227652.3 3 length=101
TACCCACTAGAGCCGGATACAGAGGTGNNNNNNNNNNNNNNNNNNNNNGNNNNNNNNNNNGGTAGCTGNNNNNNNNNNNNCTTAAGCCCCGCAATTTCCAT
+SRR5227652.3 3 length=101
BBBFFFFFFFFFIIIIIIIIIIIF#####0#####000<BFFF#####007BBFFFBFFFFFFFFFFFFF
@SRR5227652.4 4 length=101
ACCCAACACACAAGCACATTATAATTTNNNNNNNNNNNNNNNNNNNGNNNNNNNNNNNAAAGCTCTNNNNNNNNNNNNGCCAGGACATGCCATGGCCG
+SRR5227652.4 4 length=101
BBBFFFFFFFFFIIIIIIIIIIII#####0#####00<BFFFB#####00<BFFFFFFFFFFFFBFFFF
@SRR5227652.5 5 length=101
GCCCATATCCTTTTATGTCACTTAAAGANNNNNNNNNNNNNNNNNCTTTNANNNNNNNNNNAAAGATGANNNNNNNNNNNNCAGCCATCTGTCCTCAGCTTT
+SRR5227652.5 5 length=101
BBBFFFFFFFFFIIIIIIIIIIII#####00BF#0#####00<BFFFF#####007BBFFFFFFFFFFFFFBF
@SRR5227652.6 6 length=101
CACCTTTTGAATCCCATATACGTTGAGANNNNNNNNNNNNCACACACTGANNNNNNATTTTGTATTNNNNNNNNNNNCCCTGATATGTAATGCTTCAG
+SRR5227652.6 6 length=101
BBBFFFFFFFFFIIIIIFIIFIIFBFFF#####00BFFIIIII#####007BFFBFFF#####007<BBBFFFBFBFBFFFF
@SRR5227652.7 7 length=101
CTCAGTTTCTTCTCTTCAGATTTTGAGNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNATGCTGCTNNNNNNNNNNNNTTAACACCCAAATCGCGACA
+SRR5227652.7 7 length=101
BBBFFFFFFFFFIIIIIIIIIIIF#####00<BFFFF#####007BBFFFFFFFFFFFFF
@SRR5227652.8 8 length=101
AAGAAACACAGCTTCATTTGGGGTGGCAGTCANNNNNNNNTGGAATCATAGCNNNNNGAAGGGGCATTNNNNNNNNNNATGTTCCGGGATCAATTTAC
+SRR5227652.8 8 length=101
BBBFFFFFFFFFIIIIIIIIIBFIIIFI#####00<FFIIIIII#####00<BFFFFF#####00<BFFFFFFFFFBFFFFF
```

Sequence
Quality

Machine Learning for Gene Expression Analysis

- Class Discovery (Clustering):
 - Dividing samples into reproducible classes that have similar behavior or properties
 - Find genes whose expression profile is similar
 - Find groups of experiment whose expression profile is similar
- Class Prediction (Classification):
 - In classification, inputs are divided into two or more classes, and the learner must produce a model that assigns unseen inputs to one or more of these classes
- Class Comparison (Differential Expression):
 - Identify the genes which are up and down regulated among the experimental groups

Learning Objectives

- Understand the steps of the analysis from raw reads to expression counts, Class comparison and interpretation of gene lists using R/Bioconductor.
- Check both raw and processed data.
- Select the proper normalization/transformation according to the task performed (EDA/QC, Downstream analysis)
- Make sense of the results of the Class comparison

Experimental Design

- What is the purpose of this experiment?
- How many replicates?
- Technical or biological replicates?
- More reads or more replicates?

RNA-Seq Workflow

Raw Sequence Read

Aligning/Mapping

Summarize

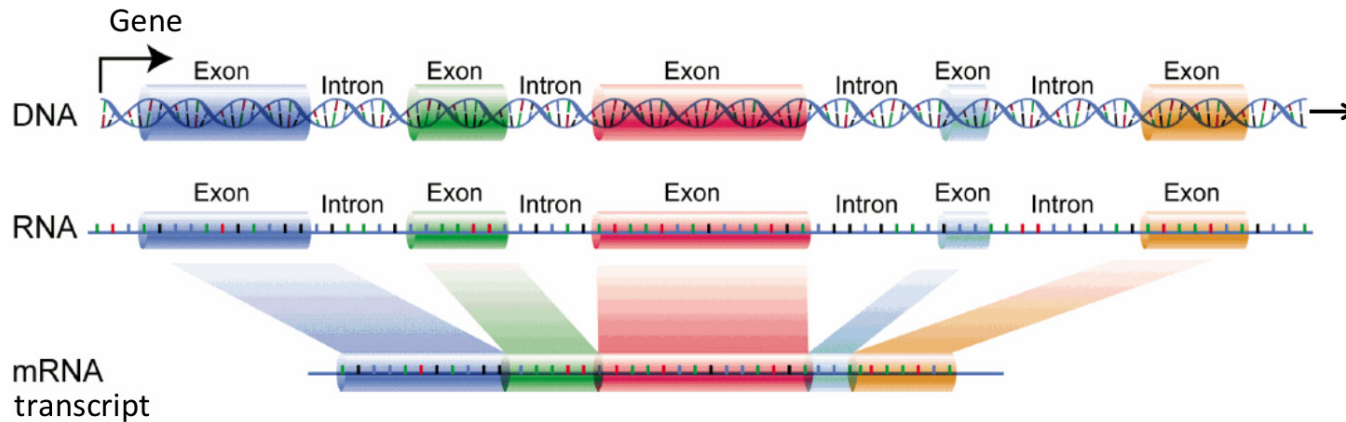
Differential Expression

System Biology

Alignment/Mapping

- There are several choices for this step, see, for example, Baruzzo et al. 2017:
- **alignement methods**
 - STAR
 - Subread/Rsubread
 - HISAT2
 - tophat
- **pseudo-alignement methods (transcript abundance quantification methods)**
 - Salmon
 - Kallisto

Summarization



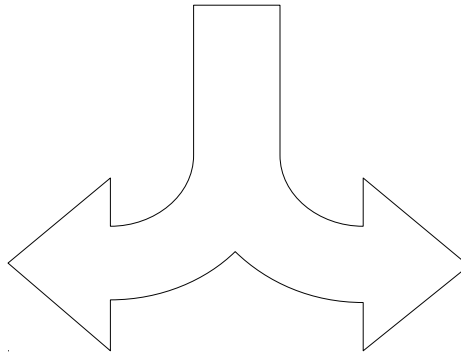
- Summarization is the process that assign mapped reads to genomic features such as genes, exons, promoters, etc. For gene-level differential expression the number of reads overlapping annotated exons of each gene can be used as a measure of the expression level of that gene.
- Typically, methods that perform this tasks take as input a set of files that contain read mapping results and an annotation file that includes genomic features and return a read count for each gene in each sample, producing a matrix of read counts (integer).

	SRR5227652_trimmed.bam	SRR5227653_trimmed.bam	SRR5227654_trimmed.bam	SRR5227655_trimmed.bam	SRR5227656_trimmed.bam
VIT_01s0010g00020	193	65	69	19	9
VIT_01s0010g00060	1016	675	356	369	430
VIT_01s0010g00240	2626	2175	770	260	255
VIT_01s0010g00330	142	126	69	196	186
VIT_01s0010g00340	32	32	23	79	93
VIT_01s0010g00360	12	10	4	7	7

Two Paths in RNA-Seq Analysis

	SRR5227652_trimmed.bam	SRR5227653_trimmed.bam	SRR5227654_trimmed.bam	SRR5227655_trimmed.bam	SRR5227656_trimmed.bam
VIT_01s0010g00020	193	65	69	19	9
VIT_01s0010g00060	1016	675	356	369	430
VIT_01s0010g00240	2626	2175	770	260	255
VIT_01s0010g00330	142	126	69	196	186
VIT_01s0010g00340	32	32	23	79	93
VIT_01s0010g00360	12	10	4	7	7

Transformations
and
Exploratory Data Analysis (EDA)



Differential Expression

Pre-processing

- Transformation for EDA
- Exploratory Data Analysis (EDA)
- Removing low expressed genes
- Normalization and transformation for DE analysis

Transformations for EDA

- The number of reads for a given gene is proportional to the expression level of the gene AND to its transcript length AND to the sequencing depth of the library.
- Counts Per Million (CPM): divide each read count per library size in millions.

$$\text{CPM}_i = \frac{X_i}{\frac{N}{10^6}} = \frac{X_i}{N} \cdot 10^6$$

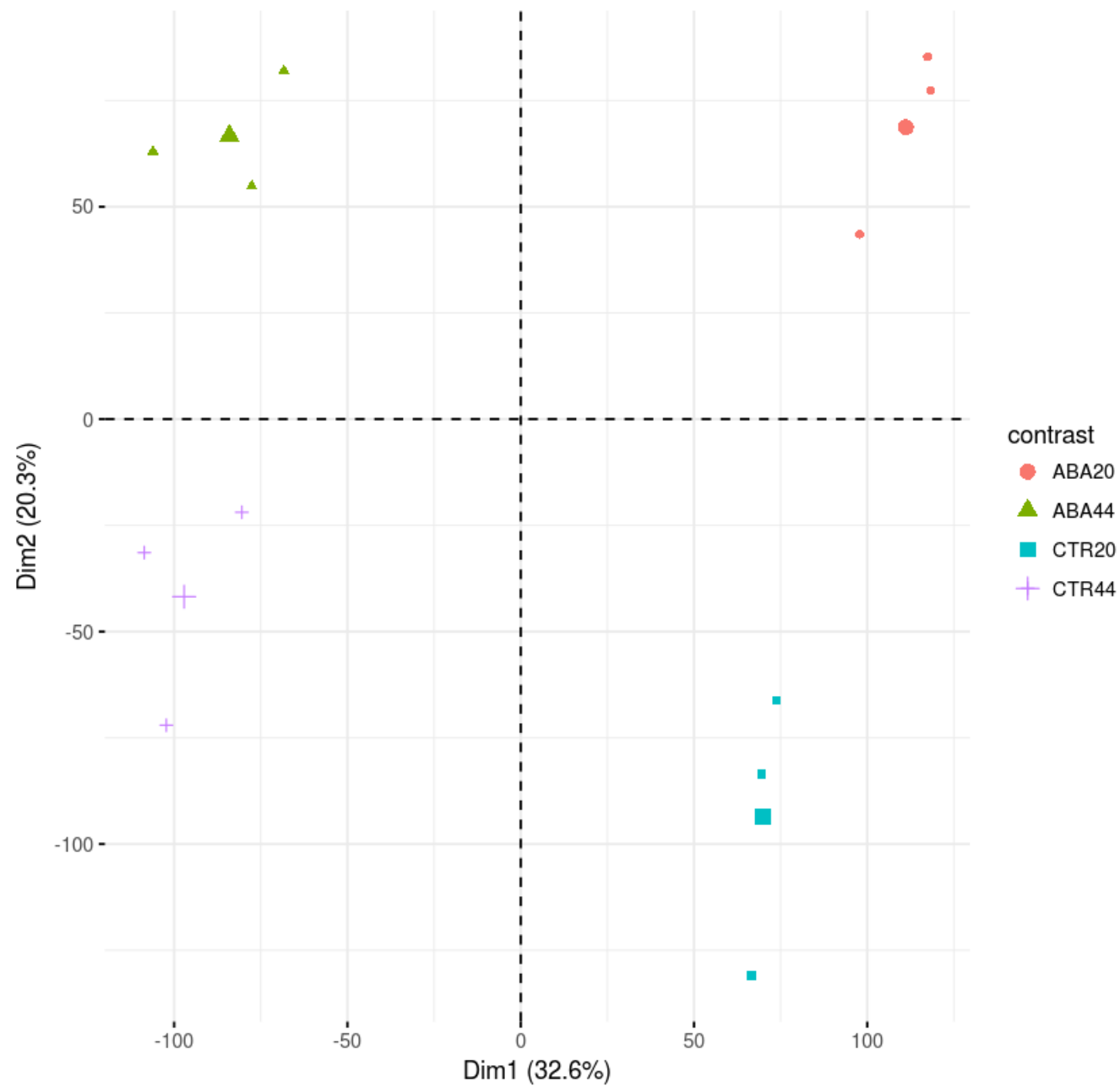
- Reads/Fragments Per Kilobase per Million (RPKM/FPKM): counts are divided by the transcript length (kb) times the total number of millions of mapped reads.

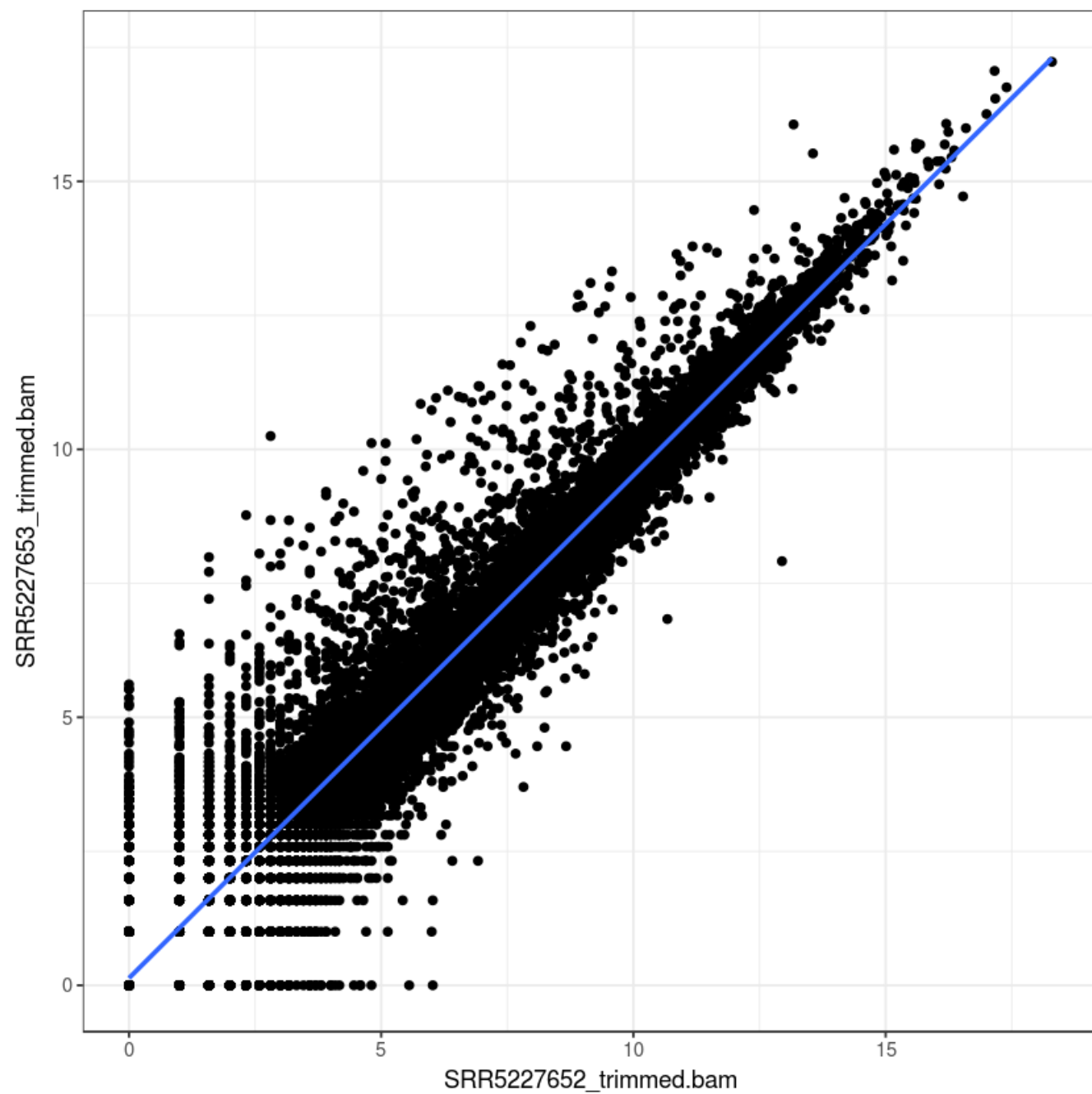
$$\text{FPKM}_i = \frac{X_i}{\left(\frac{\tilde{l}_i}{10^3}\right) \left(\frac{N}{10^6}\right)} = \frac{X_i}{\tilde{l}_i N} \cdot 10^9$$

EDA Visualizations

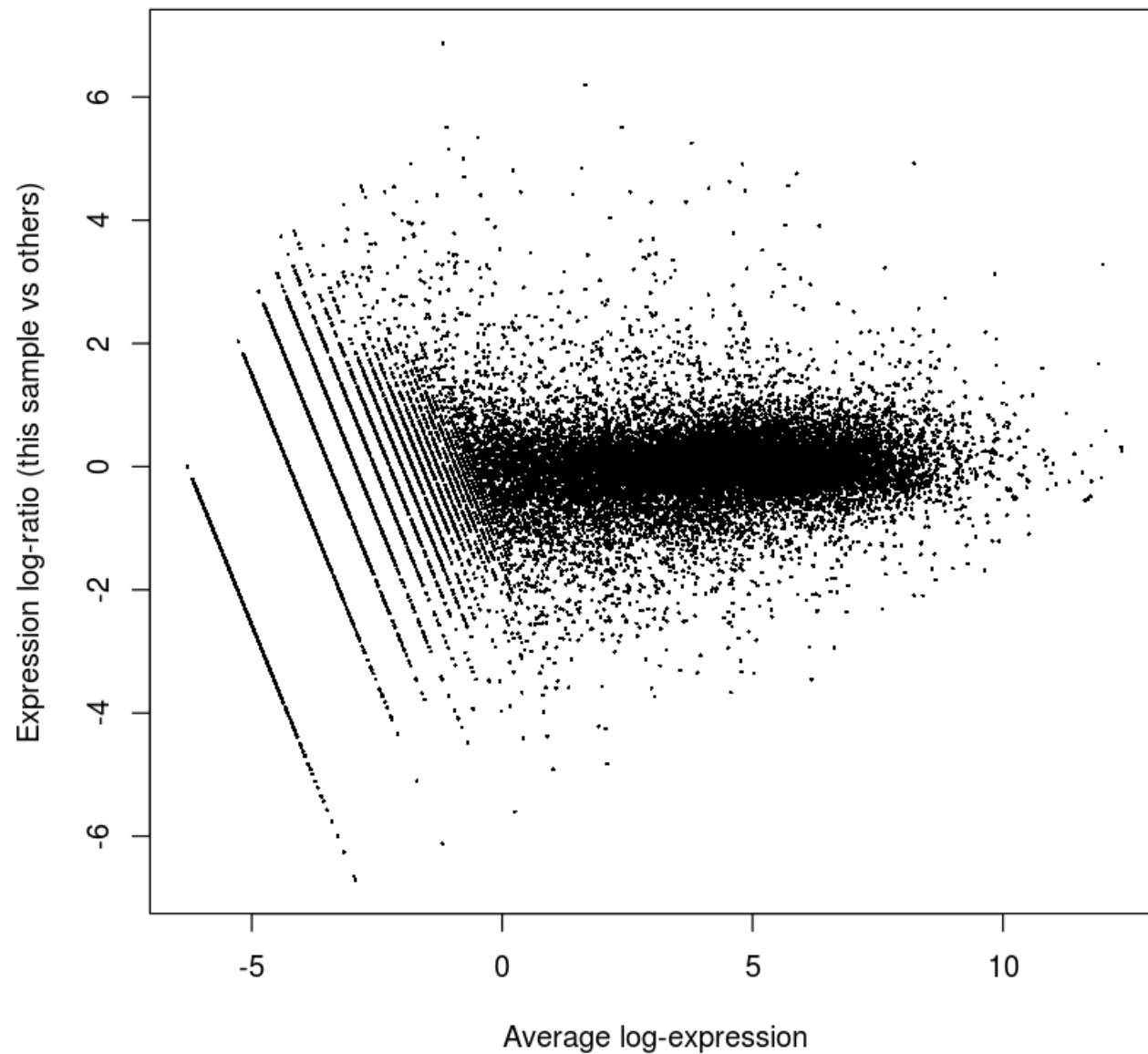
- PCA/MDS
- Clustering
- Scatterplot
- Box-plot
- MA-plot (Mean-Difference MD-plot)

Individuals - PCA





SRR5227652_trimmed.bam



Remove low expressed genes

- Genes that are not expressed at a biologically meaningful level in any condition should be discarded to reduce the subset of genes to those that are of interest, and to reduce the number of tests carried out downstream when looking at differential expression.
- Although any sensible value can be used as the expression cutoff, typically a CPM value of 1 is used in analyses as it separates expressed genes from unexpressed genes well for most datasets. Here, a CPM value of 1 means that a gene is “expressed” if it has at least 20 counts in the sample with library size ≈ 20 million .

Differential Expression Analysis

- Design matrix
- Contrast matrix
- Trimmed mean of M-values (TMM) Normalization
- Voom transformation
- Limma

Design

Design matrix

CTRL20	ABA20	CTRL44	ABA44
1	0	0	0
1	0	0	0
1	0	0	0
0	1	0	0
0	1	0	0
0	1	0	0
0	0	1	0
0	0	1	0
0	0	1	0
0	0	0	1
0	0	0	1
0	0	0	1

Contrast matrix

	ABA20vsCTRL20	ABA44vsCTRL44
CTRL20	-1	0
ABA20	1	0
CTRL44	0	-1
ABA44	0	1

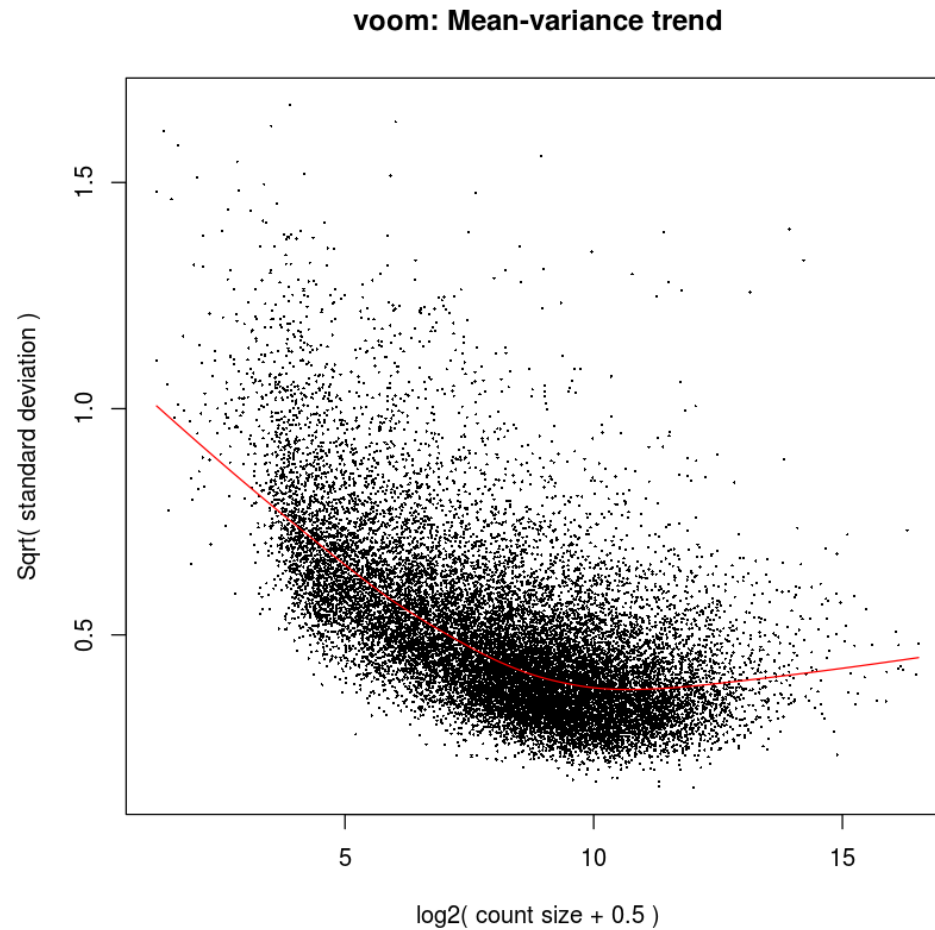
Normalization for DE Analysis

- Trimmed mean of M-Values (TMM) normalization is a normalization method commonly performed to eliminate composition biases between libraries. This generates a set of normalization factors, where the product of these factors and the library sizes defines the effective library size.
- A normalization factor below one indicates that the library size will be scaled down, as there is more suppression (i.e., composition bias) in that library relative to the other libraries. This is also equivalent to scaling the counts upwards in that sample. Conversely, a factor above one scales up the library size and is equivalent to downscaling the counts.

Limma: voom function

- Unlike methods such as edgeR or DESeq(2) that model counts using a Negative Binomial distribution, limma performs linear modelling on the log-CPM values assumed to be normally distributed and the relationship between mean and variance is taken care using precision weights calculated by the voom function.
- The read counts of gene g in sample i do not follow a Normal distribution however the transformed (logCPM) response variable converges quickly to normality. voom() models the log counts per million and fits a loess trend line to the scatterplot of variance vs. mean to create weight that are then fed into a standard limma analysis. The variance model is at the observation level and the loess trend is robust against highly variable genes.

limma/voom: mean-variance trend



Limma: Fitting linear models for comparisons of interest

- Fitting a separate linear model to the expression values for each gene.
- eBayes moderated t-test which borrows information across all genes to obtain precise estimate of gene-wise variability.
- The method aims to remove the dependency of the variance from the mean expression level.

	ID	logFC	AveExpr	t	P.Value	adj.P.Val	B
VIT_02s0025g04330	VIT_02s0025g04330	-5.967939	3.903416	-26.01021	1.311512e-11	2.419346e-07	16.32686
VIT_06s0004g05460	VIT_06s0004g05460	-3.438538	6.257815	-21.98922	8.794390e-11	3.955772e-07	15.26717
VIT_19s0093g00550	VIT_19s0093g00550	-6.354018	3.811761	-24.29892	2.839847e-11	2.619333e-07	15.26463
VIT_13s0019g02200	VIT_13s0019g02200	-2.411133	7.400294	-21.56152	1.097765e-10	3.955772e-07	15.09002
VIT_04s0044g00130	VIT_04s0044g00130	-3.106525	5.974479	-20.83956	1.611715e-10	3.955772e-07	14.71222
VIT_01s0026g00220	VIT_01s0026g00220	-2.121529	7.314954	-20.77744	1.666848e-10	3.955772e-07	14.68831
VIT_17s0000g00430	VIT_17s0000g00430	-3.373297	7.214778	-20.72442	1.715519e-10	3.955772e-07	14.64948
VIT_03s0038g01380	VIT_03s0038g01380	-3.059084	6.422147	-20.45446	1.988447e-10	3.992256e-07	14.51462
VIT_14s0171g00360	VIT_14s0171g00360	-4.849779	4.405474	-21.29683	1.261908e-10	3.955772e-07	14.45467
VIT_02s0236g00130	VIT_02s0236g00130	-1.799988	5.688626	-20.14180	2.364765e-10	3.992256e-07	14.34860

logFC log2 fold change between compared groups

AveExpr average log2 expression (across complete dataset)

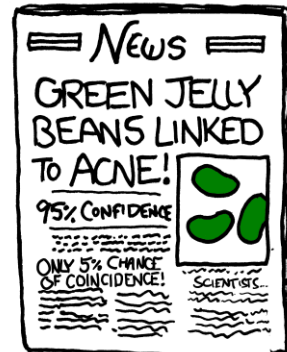
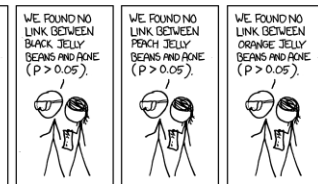
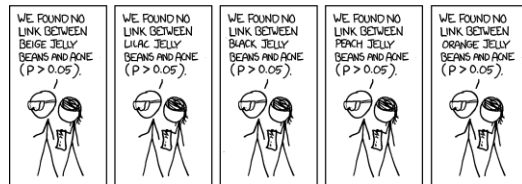
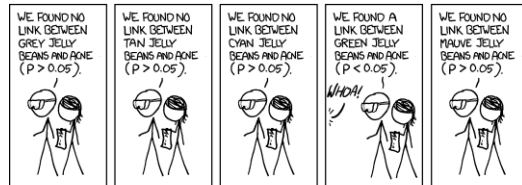
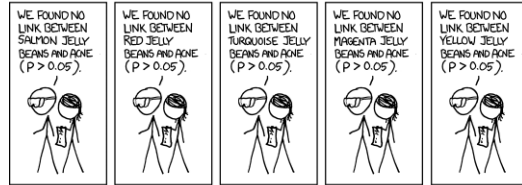
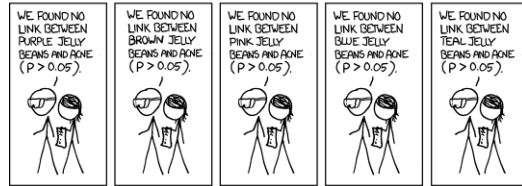
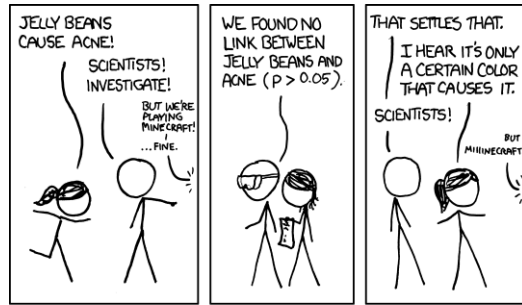
t moderated t-statistics

P.Value raw p-value

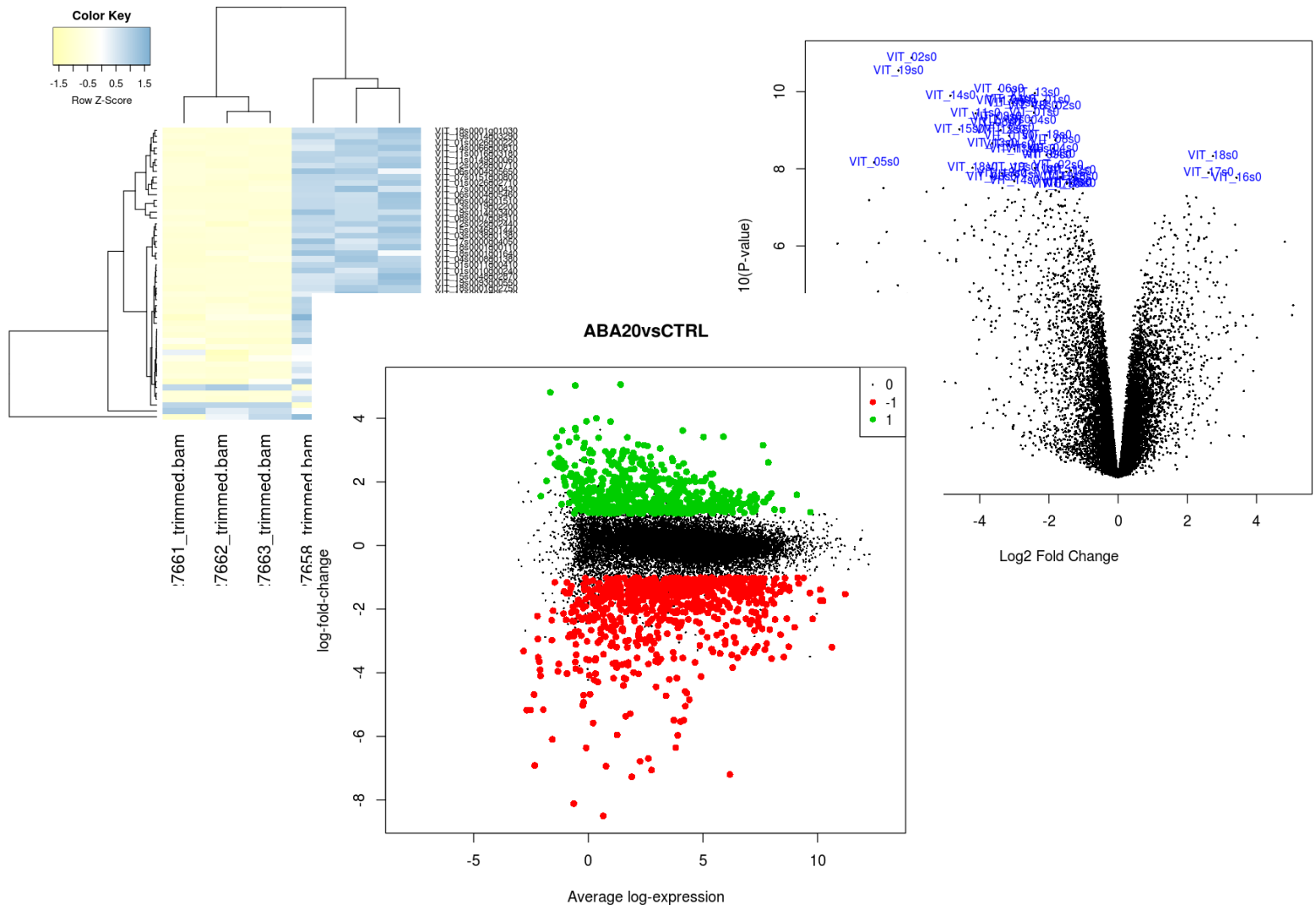
adj.P.Val adjusted p-value for multiple testing issue (FDR correction by default)

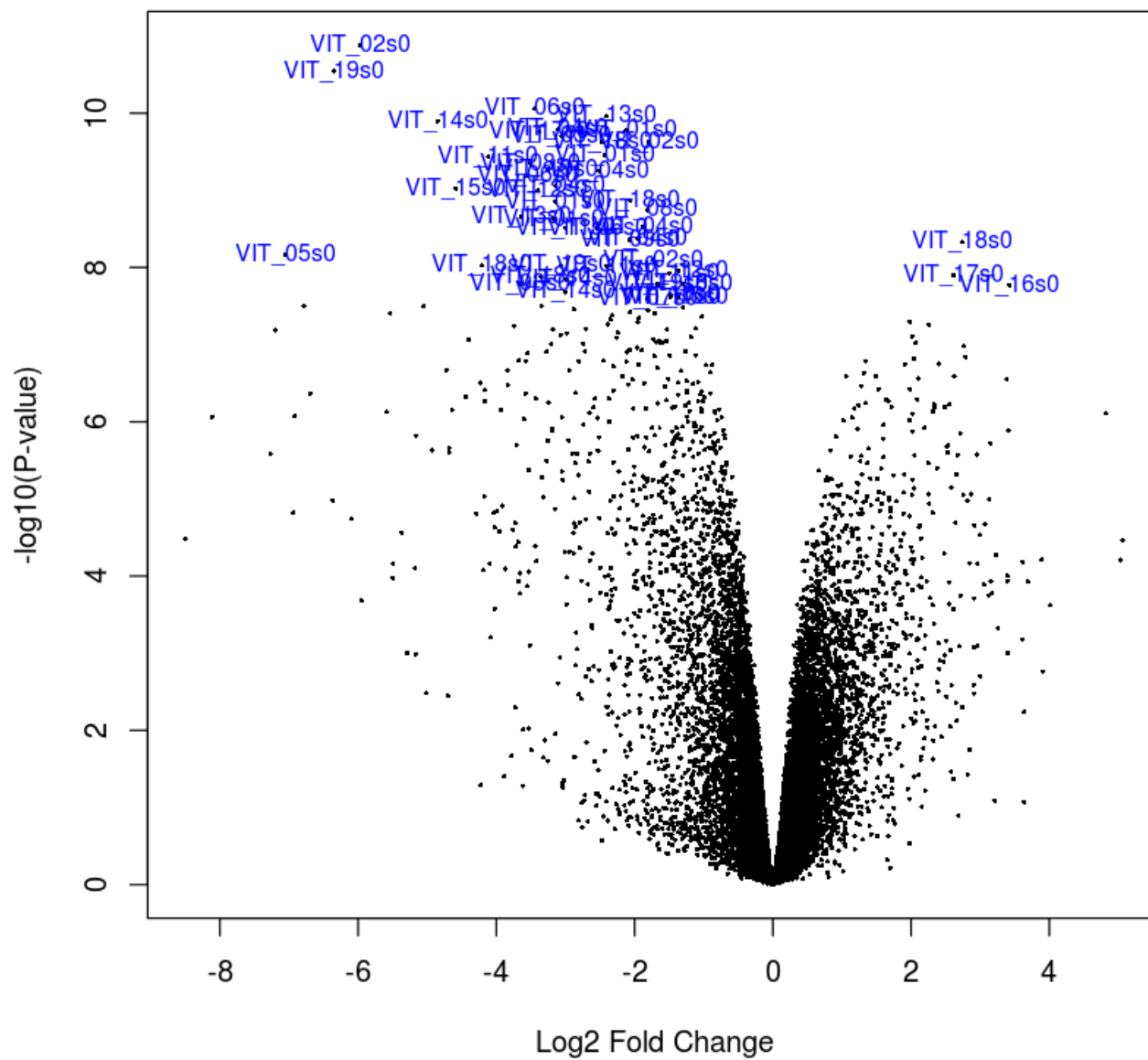
B B-statistics, log-odds that the gene is differentially expressed

<u>P-VALUE</u>	<u>INTERPRETATION</u>
0.001	HIGHLY SIGNIFICANT
0.01	
0.02	
0.03	
0.04	SIGNIFICANT
0.049	
0.050	OH CRAP. REDO CALCULATIONS.
0.051	ON THE EDGE OF SIGNIFICANCE
0.06	
0.07	HIGHLY SUGGESTIVE, SIGNIFICANT AT THE $P < 0.10$ LEVEL
0.08	
0.09	
0.099	HEY, LOOK AT THIS INTERESTING SUBGROUP ANALYSIS
≥ 0.1	

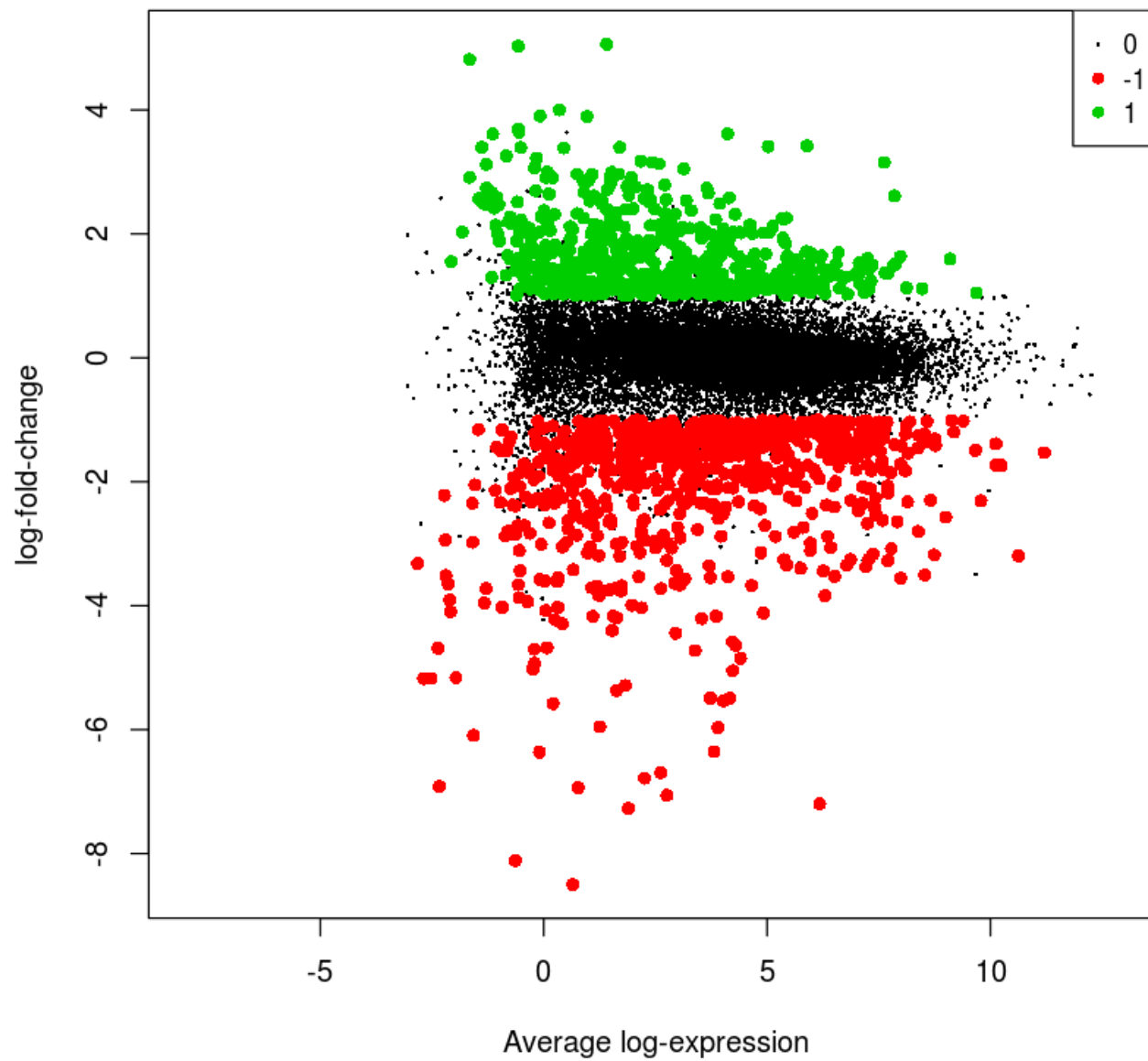


DE genes visualizations





ABA20vsCTRL



Beyond the gene list

- GO enrichment(overrepresentation) analysis:
 - Gene Ontology (BP,CC,MF) categorizes genes into functional groups.
 - Investigate whether gene sets associated with particular biological functions as represented by Gene Ontology annotations are statistically overrepresented in the identified gene groups: statistical test are used to determine which GO terms appear more frequently than would be expected by chance when examining the set of terms annotated to the input genes.
- Integration of transcriptomics data (microarray and RNA-Seq) with other -omics.
- One more thing ...

Gene Expression Compendia

The image displays two web interfaces for gene expression analysis. The top interface is Vespucci, featuring a navigation bar with links to home, download, tutorials, help, about, and feedback. It includes a workspace area with a data dropdown and a quick search bar. The main content area is divided into two sections: 'Click below to choose a gene selection method' and 'Click below to choose a contrast selection method'. The bottom interface is Colombos, also with a navigation bar and workspace. It shows a workspace with a data dropdown and a quick search bar. The main content area is divided into two sections: 'Click below to choose a gene selection method' and 'Click below to choose a contrast selection method'. The bottom interface is Colombos, also with a navigation bar and workspace. It shows a workspace with a data dropdown and a quick search bar. The main content area is divided into two sections: 'Click below to choose a gene selection method' and 'Click below to choose a contrast selection method'.

vespucci home download tutorials help about feedback

Workspace
Data Quick search
Organism: *Vitis vinifera* Reset

Click below to choose a **gene** selection method

Click below to choose a **contrast** selection method

By gene name/locus tag
Click here to add a custom list of genes to the module, by entering either common gene names or locus tags.

By contrast id
Click here to add contrasts based on the original sample/hybridization identifiers.

By GO term

By experiment
Click here to add contrasts whose contrasts you want to add to the module.

Contrast annotation
Click here to add contrasts that measure a change in certain condition properties.

login

colombos home download tutorials help about feedback

Workspace
Data Quick search
Escherichia coli
Quick module m1

Overview Heatmap Network Edit Split Genes Contrasts Enrichment

Quick module m1
Escherichia coli
Click on plot to unfreeze info panel content

b3164
pnp
Description: polynucleotide phosphorylase
polyadenylase
Transcription unit: metY-yhbC-nusA-infB-rbfA-truB-rpsO-pnp,rpsO-pnp
Regulon: CRP,Sigma70,Fis,ArgR
Pathway:
tRNA processing
Gene Ontology:
GO:0006396|RNA processing
GO:0006402|mRNA catabolic process
GO:0006950|response to stress

test_E-MEXP-1370_reiA_spoT_ILE_40m (E-MEXP-1370.RELASPOT_ILE_6_REP.ch1)
Reference: ref_E-MEXP-1370_wt_40m (E-MEXP-1370.E-COLI_2.ch1)
Experiment: E-MEXP-1370
Platform: enterocha520147F_custom,
Publication: 0
Test annotations:
STRAIN MG1655

Sort contrasts by: expression

Back to overview

colombos v3.0 beta | contact us

<http://colombos.net/>
<http://vespucci.colombos.fmach.it/>

Reproducibility

- Bioconductor (16 years of development)
- (R)Markdown
- Python Notebook
- github

References

- Papers:
 - Andrews, S. 2010. FastQC: A quality control tool for high throughput sequence data (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>)
 - Liao Y, Smyth GK and Shi W (2013). The Subread aligner: fast, accurate and scalable read mapping by seed-and-vote. *Nucleic Acids Research*, 41(10):e108
 - Liao Y, Smyth GK and Shi W (2014). featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. *Bioinformatics*, 30(7):923-30.
 - MD Robinson, A Oshlack. 2010. A scaling normalization method for differential expression analysis of RNA-seq data, *Genome biology*,
 - Ritchie, M. E., B. Phipson, D. Wu, Y. Hu, C. W. Law, W. Shi, and G. K. Smyth. 2015. "limma powers differential expression analyses for RNA-sequencing and microarray studies." *Nucleic Acids Research*, January, gkv007. doi:10.1093/nar/gkv007.
 - Law, Charity W, Yunshun Chen, Wei Shi, and Gordon K Smyth. 2014. "Voom: precision weights unlock linear model analysis tools for RNA-seq read counts." *Genome Biology* 15 (2): R29. doi:10.1186/gb-2014-15-2-r29.
- Web sites:
 - <http://bioconductor.org>
 - <https://f1000research.com/articles/5-1408/>
 - <https://dockflow.org/workflow/rnaseq-gene/>
 - <http://robpatro.com/blog/?p=235>
 - <https://haroldpimentel.wordpress.com/2014/05/08/what-the-fpkm-a-review-rna-seq-expression-units/>