

Transcriptome Data Analysis

Federico Marini (marinif@uni-mainz.de)



IMBEI@



2022/07/20-21

 [@FedeBioinfo](https://twitter.com/FedeBioinfo)

Questions

What are the steps to process RNA-Seq data?

- How to convert RNA-seq reads into counts?
- How to perform quality control (QC) of RNA-seq reads?

How to identify differentially expressed genes across multiple experimental conditions?

- How to properly analyze RNA count data using DESeq2?
- How to perform quality control (QC) and exploratory data analysis (EDA) of RNA-seq count data?

What are the biological functions impacted by the differential expression of genes?

- How can I perform a gene ontology enrichment analysis?

How can I create neat visualizations of the data?

- How can I visualize the results for my enrichment analysis?

How can I generate interactive reports to summarise my analyses?

Questions

What are the steps to process RNA-Seq data?

- How to convert RNA-seq reads into counts?
- How to perform quality control (QC) of RNA-seq reads?

How to identify differentially expressed genes across multiple experimental conditions?

- How to properly analyze RNA count data using DESeq2?
- How to perform quality control (QC) and exploratory data analysis (EDA) of RNA-seq count data?

What are the biological functions impacted by the differential expression of genes?

- How can I perform a gene ontology enrichment analysis?

How can I create neat visualizations of the data?

- How can I visualize the results for my enrichment analysis?

How can I generate interactive reports to summarise my analyses?

What you will learn

- the basics of RNA-seq data
- the basics of RNA-seq data analysis
- to get familiar with the concepts of gene expression, high-dimensional data, expression quantification, differential expression analysis
- the importance to pose the right question, in order to get the right answer :)

This material has been developed together with Charlotte Soneson in the scope of the GTIPI Summer School
(<https://imbeimainz.github.io/GTIPI2022>)

Setup for practical sessions

Got R/RStudio?

Latest versions highly recommended!

See <https://imbeimainz.github.io/GTIPI2022/material.html> for details!

Decomposing the title

RNA

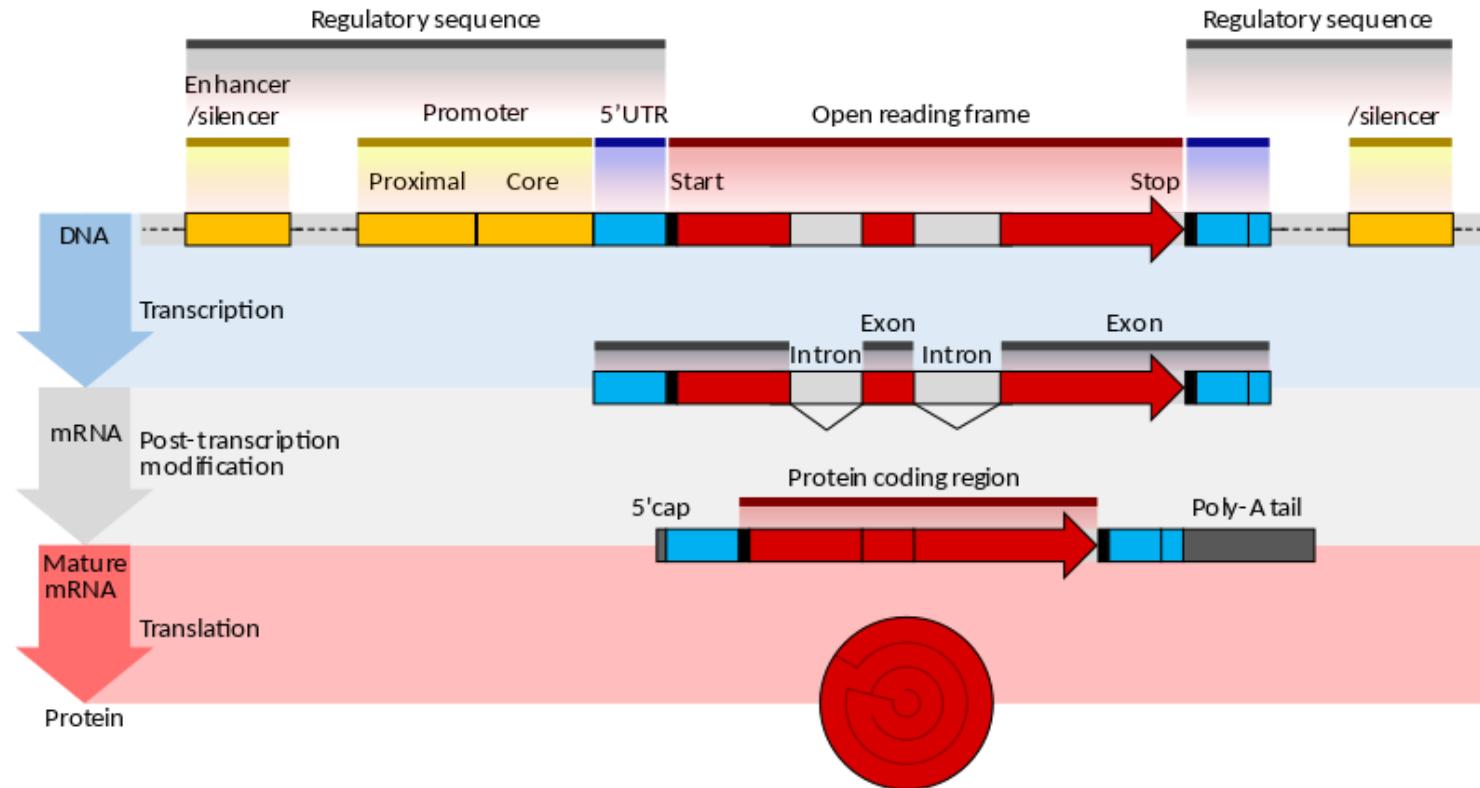
Sequencing

Bioinformatics

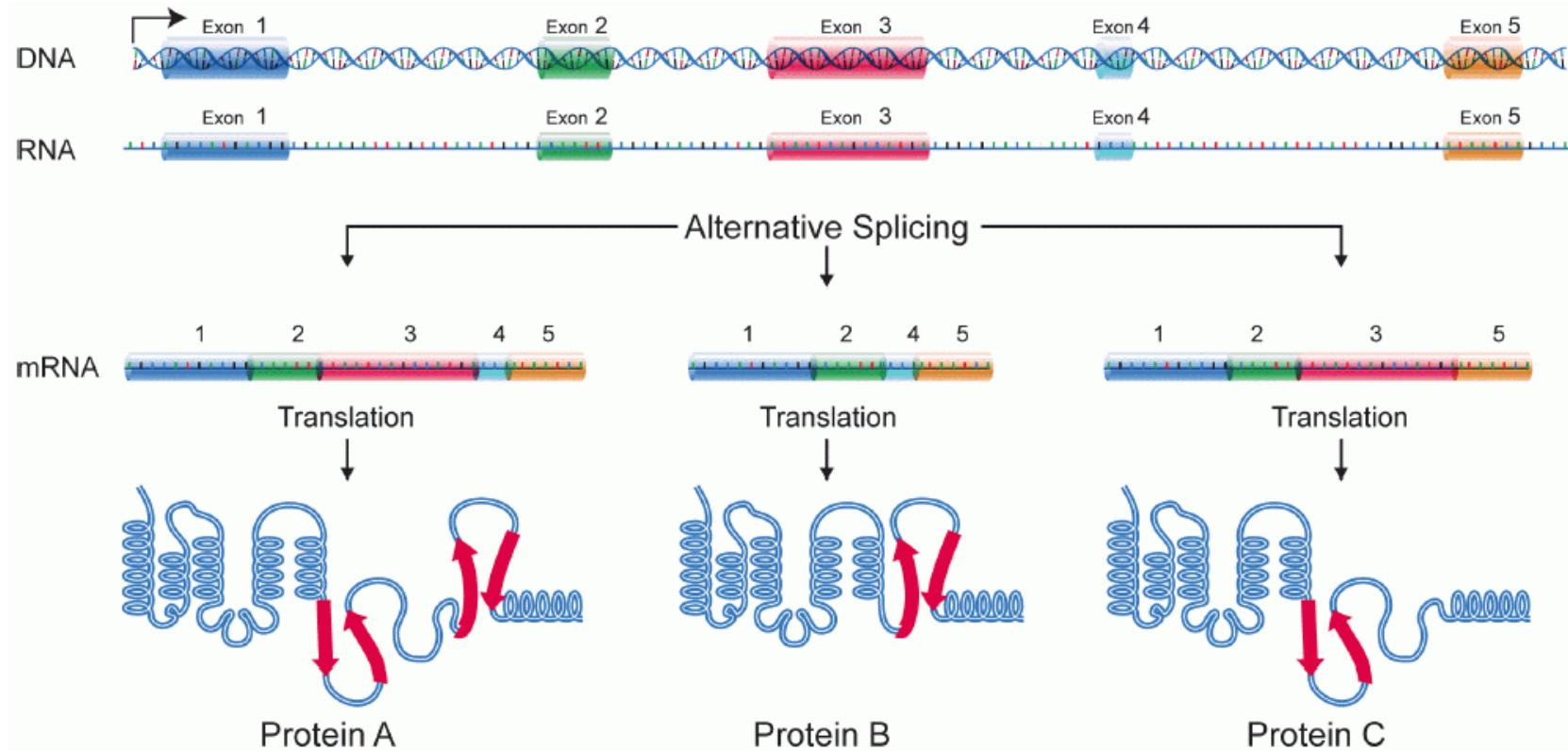
Transcriptome

Analysis

(messenger) RNA

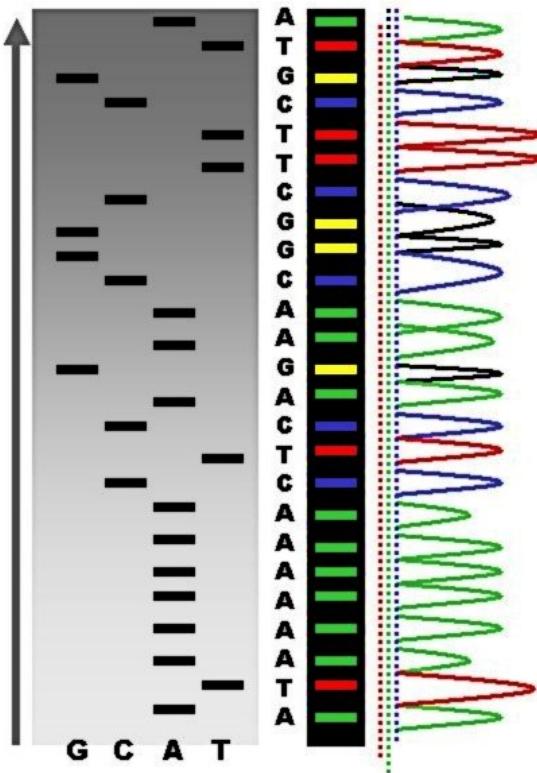


(messenger) RNA



Exons, introns, transcripts, isoforms

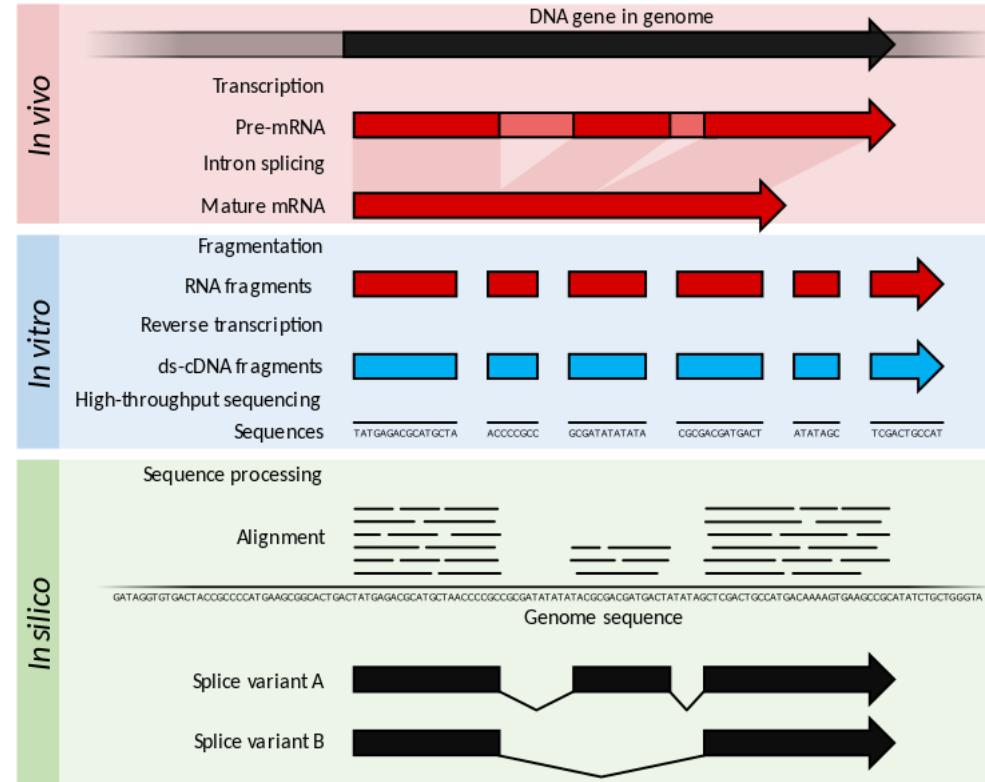
Sequencing



4 base pairs, 21 aminoacids

Excellent review on next-generation sequencing: <https://www.nature.com/articles/nrg.2016.49>

RNA-sequencing



- RNA quantification at single base resolution
- Cost efficient analysis of the whole transcriptome in a high-throughput manner

Challenges in RNA-seq

- Different origin for the sample RNA and the reference genome
- Presence of incompletely processed RNAs or transcriptional background noise
- Sequencing biases (e.g. PCR library preparation)

Benefits

- sensitive
- specific
- high-throughput
- cost-efficient
- basepair resolution

You can see: transcripts, splicing, lncRNA, circRNA, gene fusions

Bioinformatics

"an interdisciplinary field that develops methods and software tools for understanding biological data, in particular when the data sets are large and complex."

A combination of

- biology
- computer science
- information engineering
- mathematics
- statistics

... to analyze and interpret the biological data

Transcriptome

Gene expression is a fundamental level at which the results of various genetic and regulatory programs are observable.

RNA sequencing (RNA-seq) provides a quantitative and open system for profiling transcriptional outcomes on a large scale

Much has been learned about the characteristics of the RNA-seq data sets, as well as the performance of the myriad of methods developed

"RNA Sequencing Data: Hitchhiker's Guide to Expression Analysis" ->
<https://www.annualreviews.org/doi/abs/10.1146/annurev-biodatasci-072018-021255>

Analysis

There's data involved!

and these datasets have particular properties (how they are generated, ...)

How to make sense out of it?

There is a large diversity of applications to deal with

Among the most widely adopted workflows:

- Transcript discovery

Which RNA molecules are in my sample?

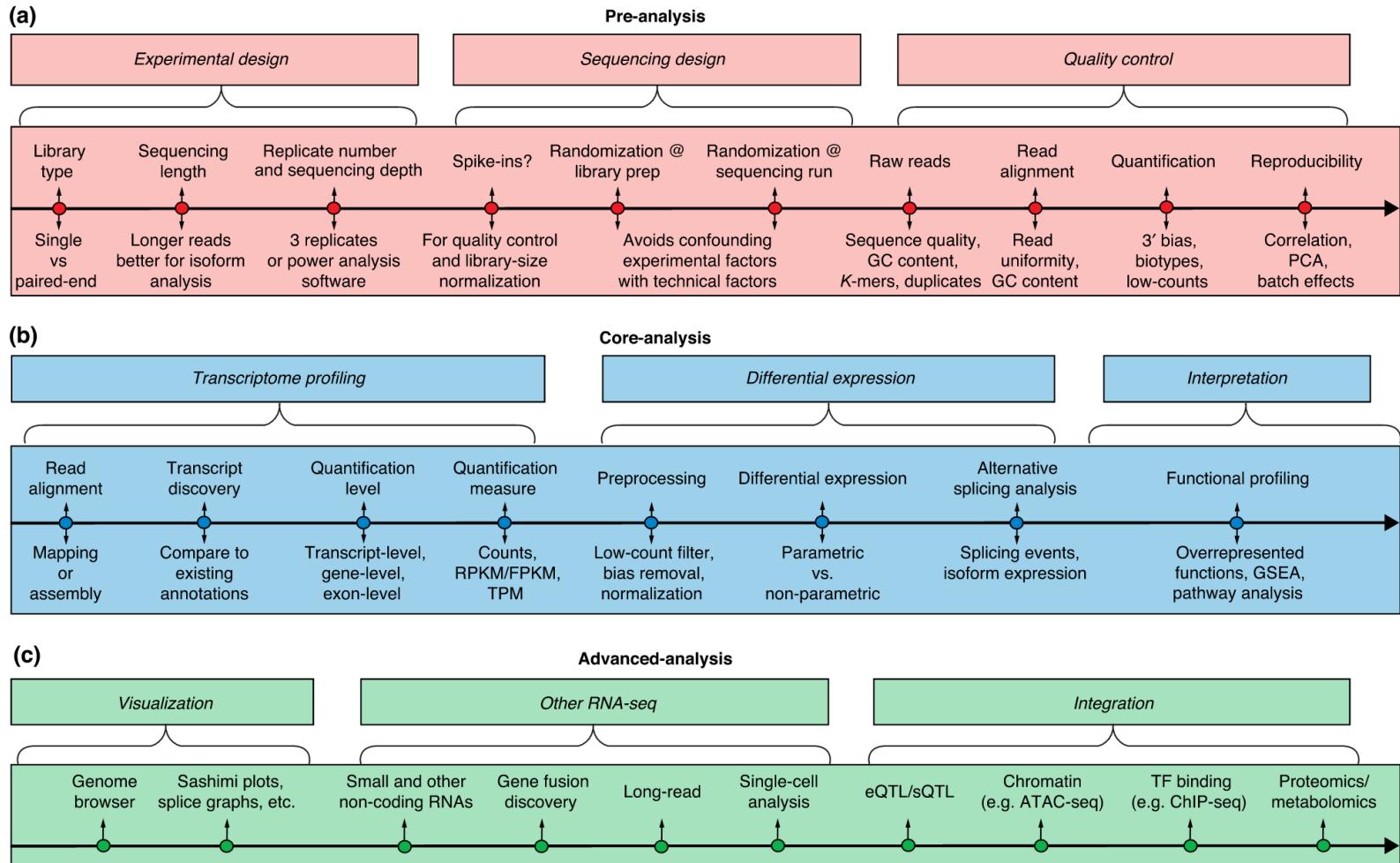
Novel isoforms and alternative splicing, Non-coding RNAs, Single nucleotide variations, Fusion genes

- RNA quantification

What is the concentration of RNAs?

Absolute gene expression (within sample), Differential expression (between biological samples)

Analysis - a bird's eye view



Differential analysis types for RNA-seq

No single available standardized workflow

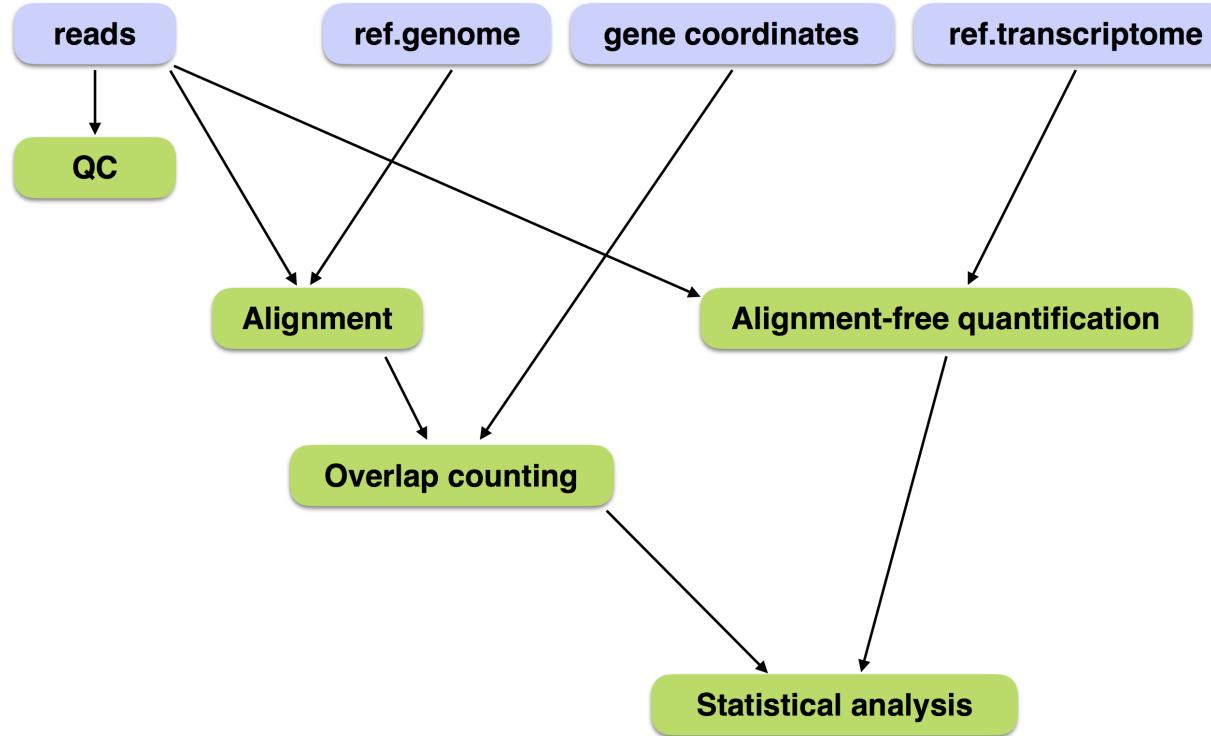
Multiple possible best practices for every dataset

To get the right answer, you have to pose the right question.

- Does the total output of a gene change between conditions? Differential Gene Expression (DGE)
- Does the expression of individual transcripts change? Differential Transcript Expression (DTE)
- Does any isoform of a given gene change? DTE+G
- Does the isoform composition for a given gene change? Differential Transcript Usage/Differential Exon Usage (DTU/DEU)

Each needs different computational approaches (quantifications + tests)

Overview of the processing workflow



Ingredients + operations

The raw data: sequencing reads

FASTQ files: sequence + base quality (phred score)

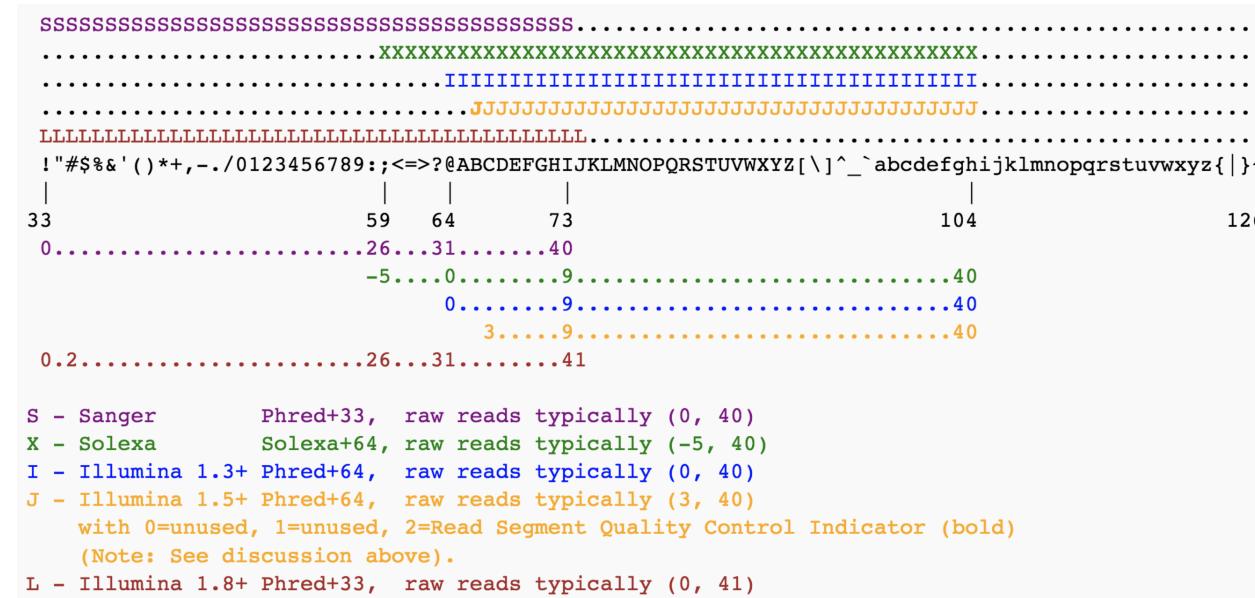
First lines of a FASTQ file

```
@SRR1055095.1 HWI-ST156:397:D09NJACXX:5:1101:1222:1915/1
NGCTGCTGGACTCCGAAGATGGCGGTATATCATCCACTGCTGACTCTN
+
#1=DDFFFHFHHHJJHIIGHIJJGIJAFFDHGGGIJJGJIJJJJIIH#
@SRR1055095.2 HWI-ST156:397:D09NJACXX:5:1101:1245:1920/1
NCTTTTCTTGTTCTCATCATCTTCAGGAGGGAGGGTCATCCTGTGN
+
#1=BDB:?:FFDFFDF?EF<FFF>B>?C@CF<1??CFB:09;09BFE9DB#
```

Repeat this tens of millions of times, and you'll have *one* sample

The raw data: sequencing reads

Different quality encodings exist



The raw data

Demo: quality control report, from FastQC

http://www.bioinformatics.babraham.ac.uk/projects/fastqc/good_sequence_short_fastqc.html

Your next best friend: MultiQC

Common sequence artefacts in NGS data:

- read errors
- base calling errors
- small insertions and deletions
- poor quality reads
- primer/adapter contamination

Solutions: Quality trimming & filtering (wide range of QC tools available)

Reference files

Reference genome sequences (in `fasta` format), required for genome alignment.

Think of the alignment as the address of each read in a 3-billion houses street, where the elements along the street can also end up repeating themselves.

What's out there?

Ensembl: <http://www.ensembl.org/info/data/ftp/index.html>

Gencode (human & mouse): <https://www.gencodegenes.org/>

UCSC: <http://hgdownload.cse.ucsc.edu/downloads.html>

iGenome: http://support.illumina.com/sequencing/sequencing_software/igenome.html

Reference files

Some critical points:

Be consistent!

Different chromosome identifiers!

Reference genomes and annotations are continuously refined, extended and improved

Keep track of version and be consistent!

Naming of genes can vary across versions, in some databases!

Reference files

An example:

http://ftp.ensembl.org/pub/release-106/fasta/homo_sapiens/dna/

... and one level up...

http://ftp.ensembl.org/pub/release-106/fasta/homo_sapiens/

```
>1 dna:chromosome chromosome:GRCh38:1:1:248956422:1 REF
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
...
TTGGTGCCAGTCCTCCAAGTCGATGGCACCTCCCTCCCTCAACCACTTGAGCAAAC
CCAAGACATCTTCTACCCCAACACCAGCAATTGTGCCAAGGGCATTAGGCTCTCAGCAT
GACTATTTTAGAGACCCGTGTGTCACTGAAACCTTTTGTTGGGAGACTATTCCCTC
CCATCTGCAACAGCTGCCCTGCTGACTGCCCTCTCCTCCCTCATCCCAGAGAAA
CAGGTCAGCTGGGAGCTTCTGCCCTGCCTAGGGACCAACAGGGCAGGAGGCAGTC
```

Reference files

The GTF format

```
chr1    unknown exon    11874    12227    .        +        .        gene_id "DDX11L1"; gene_name "DDX11L1";
transcript_id "NR_046018"; tss_id "TSS16107";
chr1    unknown CDS     3427347  3427466 .        -        2        gene_id "MEGF6"; gene_name "MEGF6";
p_id "P34437"; transcript_id "NM_001409"; tss_id "TSS31177";
```

- One line per "feature" (exon, transcript, gene, CDS, 3'UTR, 5'UTR, ...)
- One feature = 9 columns of data, plus optional track definition lines
- Essential for releasing annotation information

seqname – name of chromosome/scaffold

source – data/program source

feature – feature type name

start – positions of the feature

end

score – floating point value

strand – forward or reverse

frame – 0|1|2

attribute – semicolon-separated list of tag-value pairs, providing additional information

Bioconductor



Home Install Help Developers About

Search:

About Bioconductor

Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data. Bioconductor uses the R statistical programming language, and is open source and open development. It has two releases each year, and an active user community. Bioconductor is also available as an [AMI](#) (Amazon Machine Image) and a series of [Docker](#) images.

Bioc2020

Get the latest updates on the [Bioc2020 Conference!](#)

- Registration is Now Open! [Register Today!](#)
- Call for Abstracts! If you are interested in presenting a workshop, poster, or talk please [submit your proposal](#). Deadline March 3rd.
- Apply for [Travel Scholarships](#). Deadline March 3rd.

Bioconductor is hiring!

Bioconductor is hiring for a [full-time position](#) on the Bioconductor Core

Install »

- Discover [1823 software packages](#) available in *Bioconductor* release 3.10.

Get started with *Bioconductor*

- [Install Bioconductor](#)
- [Get support](#)
- [Latest newsletter](#)
- [Follow us on twitter](#)
- [Install R](#)

Learn »

Master *Bioconductor* tools

- [Courses](#)
- [Support site](#)
- [Package vignettes](#)
- [Literature citations](#)
- [Common work flows](#)
- [FAQ](#)
- [Community resources](#)
- [Videos](#)

Use »

Create bioinformatic solutions with *Bioconductor*

- [Software](#), [Annotation](#), and [Experiment](#) packages
- [Docker](#) and [Amazon](#) machine images
- Latest [release announcement](#)
- Use Bioconductor in the [AnVIL](#). Bioconductor [AnVIL Project Updates](#)
- [Community Slack](#) sign-up
- [Support site](#)

Develop »

Contribute to *Bioconductor*

- [Developer resources](#)
- [Use Bioc 'devel'](#)
- 'Devel' packages
- [Package guidelines](#)
- [New package submission](#)
- [Git source control](#)
- [Build reports](#)

Bioconductor - soon your best friend

- an open source project
- a repository of packages, focused on bioinformatics/computational biology
- a open development platform and community

Currently (May 2022)

- 2140 software packages
- 909 annotation packages
- 410 experiment packages
- 29 workflows
- 8 books

Aim: interdisciplinary research, collaboration and rapid development of scientific software

Documentation

- function manual pages, most of them with runnable examples
- package vignettes - mandatory here!
- workflows, documenting full analyses spanning multiple tools
- a very active support site

The real deal: Bioconductor's community

sign up / log in • about • faq • rss 

 Bioconductor
OPEN SOURCE SOFTWARE FOR BIOINFORMATICS

ASK QUESTION LATEST NEWS JOBS TUTORIALS TAGS USERS

Limit ▾ Sort ▾ Search

0 votes **2 answers** **26 views** **Use ddCT with gene specific efficiencies?**
ddct efficiency gene specific efficiencies
written 10 days ago by bettenbrock • 0 • updated 24 minutes ago by Zhang, Jitao David • 110

0 votes **0 answers** **18 views** **How to use panelcn.mops to detect CNVs from whole genome data by getting count windows from BED file?**
cn.mops cnv panelcn.mops cnv detection
written 3 days ago by metzgerlukas • 0

0 votes **0 answers** **23 views** **multifactor desing with interaction LRT or Wald test?**
deseq2 wald lrt
written 7 hours ago by jnaviapelaez • 0

0 votes **1 answer** **25 views** **Meaning of log2 fold change in 2 x 4 interaction**
deseq2
written 2 days ago by julia.chariker • 10 • updated 20 hours ago by Michael Love • 27k

0 votes **0 answers** **30 views** **DropletUtils::swappedDrops more cells than using Read10x with filtered_feature_bc_matrix files**
dropletutils single cell rna seq 10x genomics
written 1 day ago by gil.stelzer • 0

0 votes **1 answer** **25 views** **package to analyze generic protein microarray data**
microarray protein analysis proteomeprofiler
written 1 day ago by hcnbox • 0 • updated 1 day ago by Gordon Smyth • 40k

1 vote **1 answer** **40 views** **In DESeq2, how do you interpret results based on the order of variables in the "contrast" argument?**
deseq2
written 1 day ago by gpreising • 0 • updated 1 day ago by Kevin Blighe • 460

3 votes **1 answer** **51 views** **Conversion of LogFC to FC**
limma logfc fc
written 1 day ago by nia • 10 • updated 1 day ago by Gordon Smyth • 40k

Recent...

Replies

- C: Multiple testing across ... by fl • 0
- A: Use ddCT with gene speci... by Zhang, Jitao David • 110
- C: How to use panelcn.mops ... by Kevin Blighe • 460
- C: Conversion of LogFC to FC by Aaron Lun • 25k
- C: DESeq2: experiment with ... by Michael Love • 27k

Votes

- Error while uploading Bioc...
- A: Error while uploading Bi...
- C: Error while uploading Bi...
- A: Multiple testing across ...
- C: Conversion of LogFC to FC

Awards • All »

- Autobiographer  to jacknick1996 • 0
- Autobiographer  to smithcarter240191 • 0
- Autobiographer  to rk6415153 • 0
- Autobiographer  to lawyersforlandlords • 0
- Autobiographer  to casinolife24 • 0
- Scholar  to Mike Smith • 4.2k

Locations • All »

- Spain, 3 minutes ago
- Switzerland, 24 minutes ago
- EMBL Heidelberg / de.NBI, 27 minutes ago

Data processing



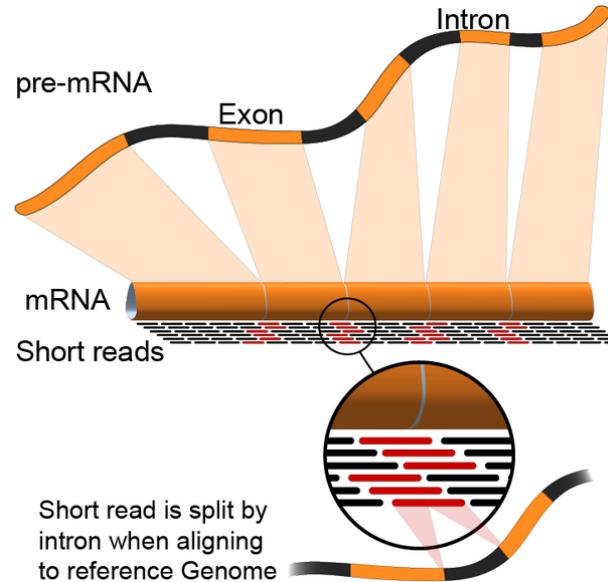
Data processing

Turning millions of text lines into properly structured abundance tables

Our aim: we often want to compare abundance (expression) of genes or other features between conditions

Splice-aware genome alignment vs "direct" transcript mapping and quantification

Alignment: not just simple mapping



For RNA-seq data, we need a splice-aware aligner

Common choices:

- STAR
- HISAT2

Alignment

File-format: [sam](#) (compressed into [bam](#))

```
SRR1055095.6079377 353 chr1    11167    0      50M    =    11751    634    CGCCCCTTGCTTGAGCCGGGCAC
TACAGGACCCGCTTGCTCACGGTGAA    CCCFFFFFFHHHHJJJJJJJJJJJJJJIGJJJJJJJJJJJJJJJDHJJCEHH
AS:i:-10    XN:i:0    XM:i:2    XO:i:0    XG:i:0    NM:i:2    MD:Z:48C0T0    YT:Z:UU    NH:i:20
CC:Z:chrY    CP:i:59361513    HI:i:0
```

Again: repeat this, one read at a time!

Entries:

QNAME - Query NAME of the read or the read pair

FLAG - Bitwise FLAG (pairing, strand, mate strand, etc.)

RNAME - Reference sequence NAME

POS - 1-Based leftmost POSition of clipped alignment

MAPQ - MAPping Quality (Phred-scaled)

CIGAR - Extended CIGAR string (operations: MIDNSHP)

MRNM - Mate Reference NaMe ('=' if same as RNAME)

MPOS - 1-Based leftmost Mate POSition

ISIZE - Inferred Insert SIZE

SEQ - Query SEQuence on the same strand as the reference

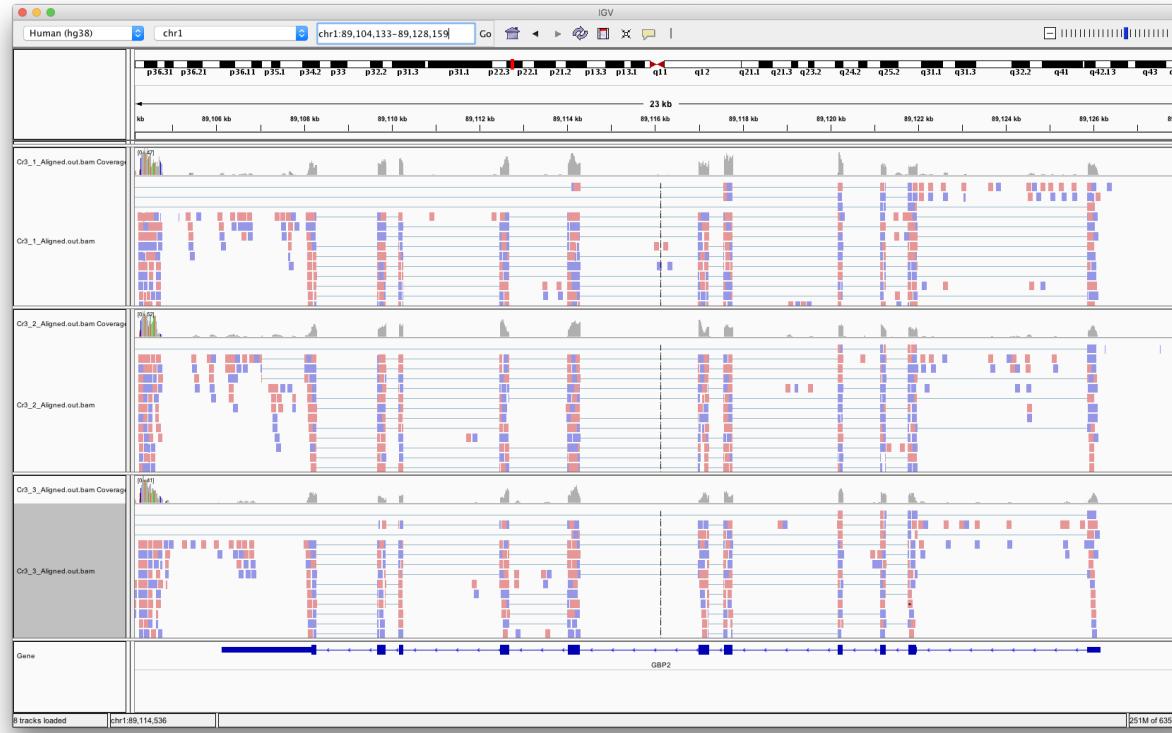
QUAL - Query QUALity (ASCII-33=Phred base quality)

Tags: used to store info about alignment

Before quantification

... and actually, always: Do visualize your data!

Options: UCSC Genome Browser, IGV, IGB - <http://software.broadinstitute.org/software/igv/>



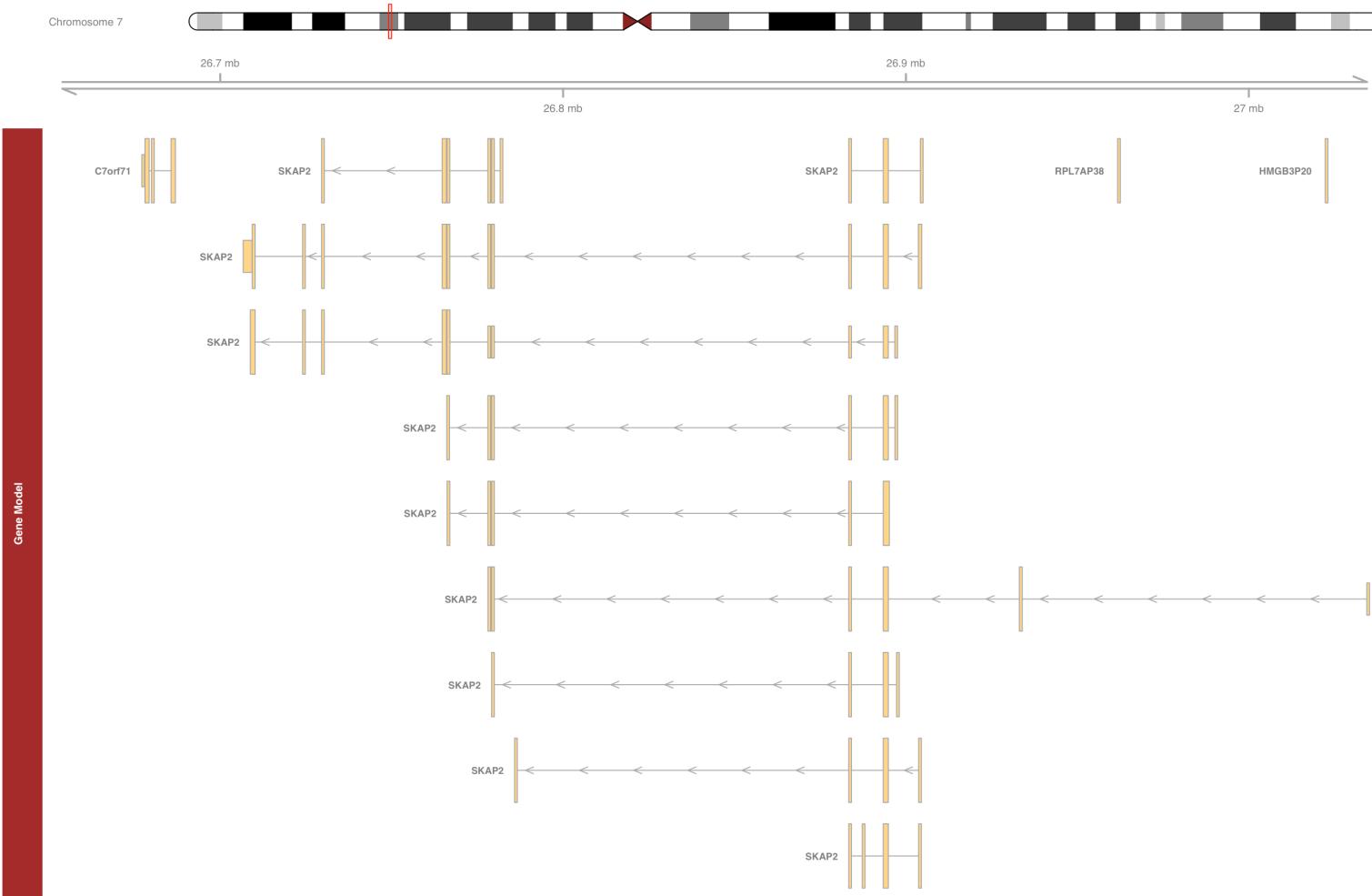
Quantification

Ingredients: BAM data + GTF annotation file

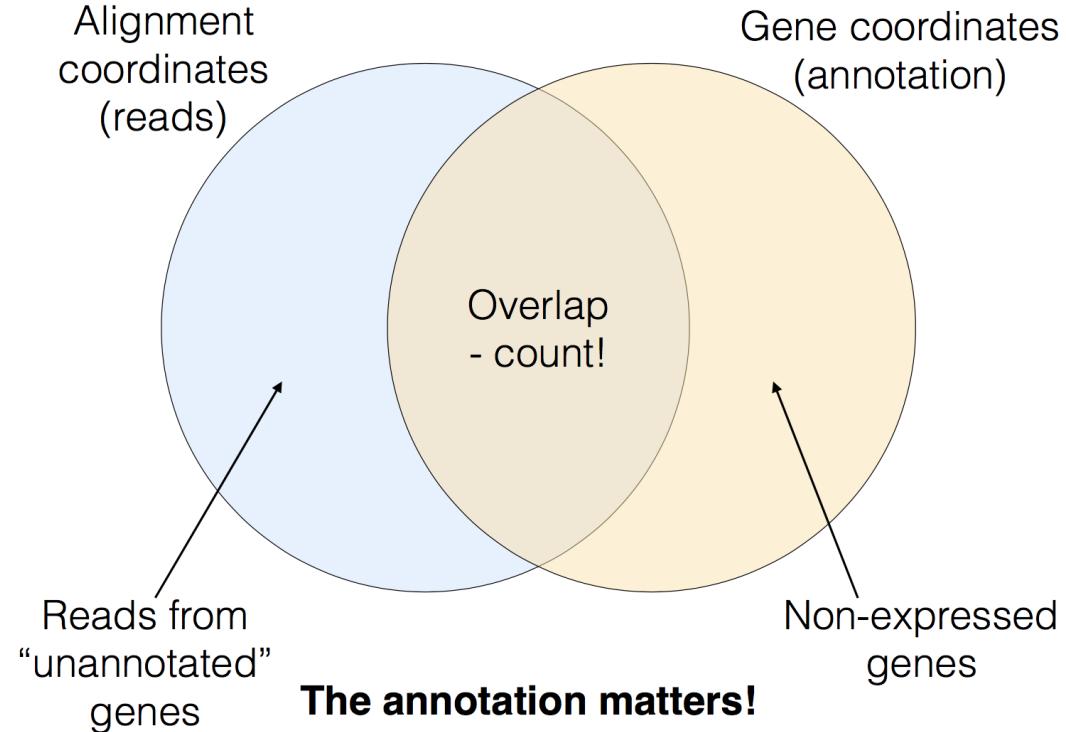
Output: number of reads overlapping known features (discrete, positive, skewed)

Gene-level counts, often obtained by genome alignment + overlap counting

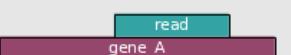
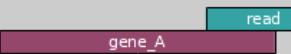
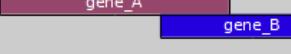
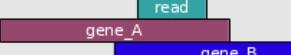
The role of annotation



The annotation matters!



The annotation matters!

	union	intersection _strict	intersection _nonempty
 A single read overlaps with gene A.	gene_A	gene_A	gene_A
 A single read overlaps with gene A, but only a portion of the read is within the gene's boundaries.	gene_A	no_feature	gene_A
 A single read is entirely contained within gene A.	gene_A	no_feature	gene_A
 A single read overlaps with two genes, A and B.	gene_A	gene_A	gene_A
 A single read overlaps with two genes, A and B, and both genes are annotated.	gene_A	gene_A	gene_A
 A single read overlaps with two genes, A and B, and both genes are annotated.	ambiguous (both genes with --nonunique all)	gene_A	gene_A
 A single read overlaps with two genes, A and B, and only gene B is annotated.	ambiguous (both genes with --nonunique all)		
 A single read overlaps with two genes, A and B, and neither gene is annotated.	alignment_not_unique (both genes with --nonunique all)		

Alignment-free quantifications

Some recently developed methods:

- salmon (Patro et al, Nat Methods 2017)
- kallisto (Bray et al, Nat Biotechnol 2016)

return...

- Transcript-level counts and TPM (transcripts-per-million) estimates, which can be summed up to get
- Gene-level counts and TPM estimates

Pros & cons

- considerably faster than traditional alignment+counting -> allow bootstrapping
- more highly resolved estimates (transcripts rather than gene) + can be aggregated
- can use a slightly larger fraction of the reads since multi-mapping reads are not excluded
- don't return precise alignments (bam files, for e.g. visualization in genome browser)

Which way to go?

Based on genome alignment - mainly gene-level quantification: combine exons, "ignoring" splice variants

- Simple, powerful, yet in some cases inaccurate
- Tools:
 - `htseq-count`, `featureCounts` for estimating expression levels (counts)
 - `edgeR`, `DESeq2`, `voom+limma` for statistical modeling

Based on transcriptome mapping - transcript- and gene-level quantification: 'assign' reads (or rather, estimate most likely expression level) based on probabilistic modeling

- Potentially cleaner, but high degree of uncertainty on the transcript level!
- Tools:
 - `bitSeq`, `RSEM`, `salmon`, `kallisto` for (pseudo)alignment/quantification
 - `DESeq2`, `edgeR`, `voom+limma`, `swish`, `DRIMseq`, `DEXSeq`, `sleuth` for modeling (depending on the question of interest)

What it would look like - STAR + featureCounts:

... due to time constraints

Index the genome

```
$ STAR --runThreadN 24 \
    --runMode genomeGenerate \
    --genomeDir my_genome \
    --genomeFastaFiles my_genome.fa \
    --sjdbGTFfile my_genes.gtf \
    --sjdbOverhang 99
```

What it would look like - STAR + featureCounts:

... due to time constraints

Map each file

```
$ STAR --runThreadN 24 \
    --runMode alignReads \
    --genomeDir my_genome \
    --readFilesIn my_sample_read1.fastq.gz \
                  my_sample_read2.fastq.gz \
    --readFilesCommand zcat \
    --outFileNamePrefix output/S1/ \
    --outSAMtype BAM SortedByCoordinate \
    --quantMode GeneCounts
```

What it would look like - STAR + featureCounts:

 SRR1039508	▶	 SRR1039508_Aligned.sortedByCoord.out.bam
 SRR1039509	▶	 SRR1039508_Log.final.out
 SRR1039512	▶	 SRR1039508_Log.out
 SRR1039513	▶	 SRR1039508_Log.progress.out
 SRR1039516	▶	 SRR1039508_ReadsPerGene.out.tab
 SRR1039517	▶	 SRR1039508_SJ.out.tab
 SRR1039520	▶	
 SRR1039521	▶	

What it would look like - STAR + featureCounts:

... due to time constraints

Quantify

```
featureCounts(files = bamfiles,  
              annot.ext = "my_genes.gtf",  
              isGTFAnnotationFile = TRUE,  
              GTF.featureType = "exon",  
              GTF.attrType = "gene_id",  
              useMetaFeatures = TRUE,  
              isPairedEnd = TRUE,  
              strandSpecific = 0)
```

Directly generates a count matrix in your R session.

What it would look like - salmon

... due to time constraints

Create an index of the transcriptome

```
$ salmon index -i my_transcripts.idx \
    -t <(cat my_transcripts.fasta my_genome.fasta) \
    -d chromosome_names.txt
```

The genome acts as a 'decoy' sequence, to collect reads truly arising from intronic or intergenic locations.

What it would look like - salmon

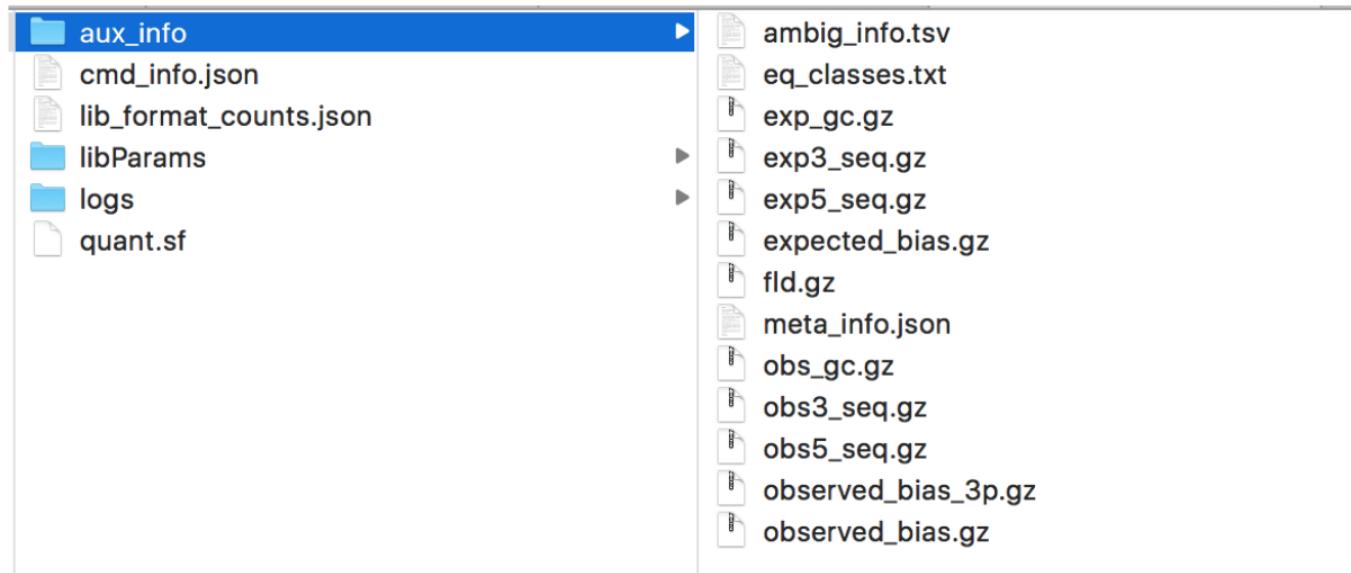
... due to time constraints

Quantify a sample at the transcript level

```
$ salmon quant -i my_transcripts.idx -l A \  
    -1 my_sample_read1.fastq.gz -2 my_sample_read2.fastq.gz \  
    -p 10 -o results/sample1 --validateMappings \  
    --numBootstraps 30 --seqBias --gcBias
```

What it would look like - salmon

... due to time constraints



What it would look like - salmon

... due to time constraints

Salmon
[quant.sf]

Name	Length	EffectiveLength	TPM	NumReads
ENST00000406070	2025	1869.81	0	0
ENST00000446844	2227	2071.81	0.137334	3.71695
ENST00000599620	686	530.936	0	0
ENST00000471557	505	350.256	0.731211	3.3457
ENST00000338761	1456	1300.81	0	0
ENST00000417509	1444	1288.81	7.58582e-08	1.27717e-06
ENST00000484946	610	455.039	2.87905	17.1142
ENST00000490656	660	504.969	1.46703	9.67744
ENST00000439537	1161	1005.81	1.47611	19.3952
ENST00000493251	641	485.994	0.597774	3.79512
ENST00000460127	408	253.708	0	0

Importing salmon quantifications into R

You can follow (offline) the instructions of the `tximport` package -
<https://bioconductor.org/packages/tximport/>

`tximeta`: another precious assistant on the way to be consistent and to keep track of provenance identification (we'll see it in action during the exercises)

What does our data look like now?

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516	SRR1039517	SRR1039520	SRR1039521
ENSG00000000003	679	448	873	408	1138	1047	770	572
ENSG00000000005	0	0	0	0	0	0	0	0
ENSG00000000419	467	515	621	365	587	799	417	508
ENSG00000000457	260	211	263	164	245	331	233	229
ENSG00000000460	60	55	40	35	78	63	76	60
ENSG00000000938	0	0	2	0	1	0	0	0
ENSG00000000971	3251	3679	6177	4252	6721	11027	5176	7995
ENSG00000001036	1433	1062	1733	881	1424	1439	1359	1109
ENSG00000001084	519	380	595	493	820	714	696	704
ENSG00000001167	394	236	464	175	658	584	360	269
ENSG00000001460	172	168	264	118	241	210	155	177
ENSG00000001461	2112	1867	5137	2657	2735	2751	2467	2905
ENSG00000001497	524	488	638	357	676	806	493	475
ENSG00000001561	71	51	211	156	23	38	134	172
ENSG00000001617	555	394	905	415	727	697	618	599
ENSG00000001626	10	2	9	2	10	6	5	5
ENSG00000001629	1660	1251	2259	1079	2462	2514	1888	1660
ENSG00000001630	59	54	66	23	84	87	31	59
ENSG00000001631	729	692	943	475	1034	1163	731	744
ENSG00000002016	201	161	256	99	268	257	160	137
ENSG00000002079	3	0	3	1	4	0	0	1
ENSG00000002330	206	174	184	111	194	260	156	177

Some challenges in RNA-seq data analysis

- 1 - Choosing an appropriate statistical distribution
- 2 - Normalization between samples
- 3 - Few samples available make it difficult to estimate parameters (e.g., variance)
- 4 - Many genes, many tests - high dimensionality

Some challenges in RNA-seq data analysis - 1

Choosing an appropriate statistical distribution

Variance depends on the mean count

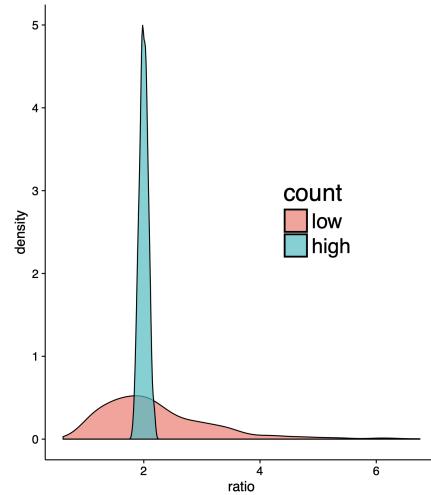
Counts are non-negative and often highly skewed

This means you can't just use t-tests, ANOVA - no prob, `glm`s to the rescue!

Poisson -> negative binomial, better captures variability across biological replicates

"Why do we not just take the ratios?"

Fold changes, relative abundances



Ex: ratio between two Poisson distributed variables

Low: mean = 20 vs mean = 10

High: mean = 2000 vs mean = 1000

Which one would you trust more? Why?

This goes back to having appropriate statistical frameworks that nicely model your datasets (and how these get generated)

Some challenges in RNA-seq data analysis - 2

Normalization between samples

Observed counts depend on:

- abundance level
- gene/transcript length
- sequencing depth
- sequencing biases

"As-is" estimates not directly comparable across samples

Normalization aims to ensure our expression estimates are

- comparable across features (genes, isoforms, etc)
- comparable across libraries (different samples)
- on a human-friendly scale (interpretable magnitude)

Most RNA-seq methods (e.g., edgeR, DESeq2, voom) need raw counts (or equivalent) as input

Don't provide these methods with (e.g.) RPKMs, FPKMs, TPMs, CPMs, log-transformed counts, normalized counts, ...

Read documentation carefully!

Digression: Normalization expression units

- RPKM/FPKM (Reads/Fragments per kilobase of transcript per million reads of library)
 - Corrects for total library coverage
 - Corrects for gene length
 - Comparable between different genes within the same dataset
- TPM (transcripts per million)
 - normalizes to transcript copies instead of reads - gives an idea of the proportion of transcripts
 - Corrects for cases where the total RNA output differs between samples
 - More appropriate for between sample comparisons (the sum of all TPMs in each sample are the same)

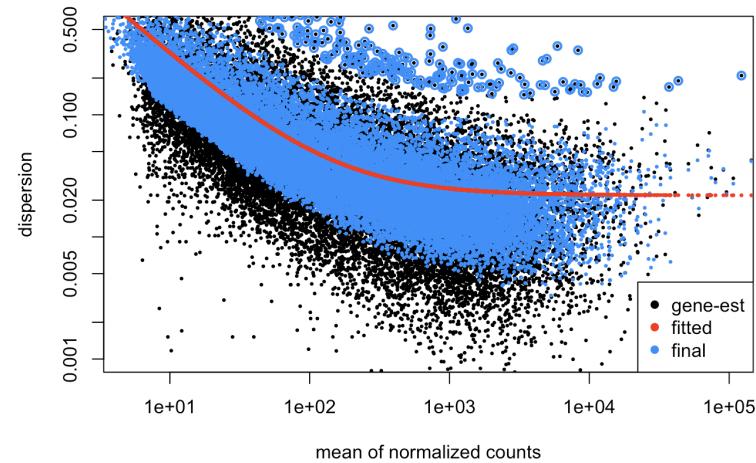
For DE analysis you have to work with discrete counts...

... and for comparisons you can use normalized counts (median ratio/TMM methods are robust across all genes!)

Some challenges in RNA-seq data analysis - 3

Few samples available make it difficult to estimate parameters (e.g., variance)

You can take advantage of the large number of genes



-> Shrink the gene-wise estimates towards a center value defined by the observed distribution of dispersions across...

- all genes (common dispersion estimate)
- genes with similar expression (trended dispersion estimate)

Some challenges in RNA-seq data analysis - 4

Many genes, many tests - high dimensionality

FDR for multiple test correction

Some more ideas:

- filter out genes that have little chance of showing significance (without looking at the test results)
- independent hypothesis weighting

... all nicely implemented for the DESeq framework

Do not test AND filter on logFC post-hoc!

Think of the null you're testing against - $\beta = 0$ by default, useful to adapt if you want to focus on larger effect sizes

Exploratory analysis and visualization

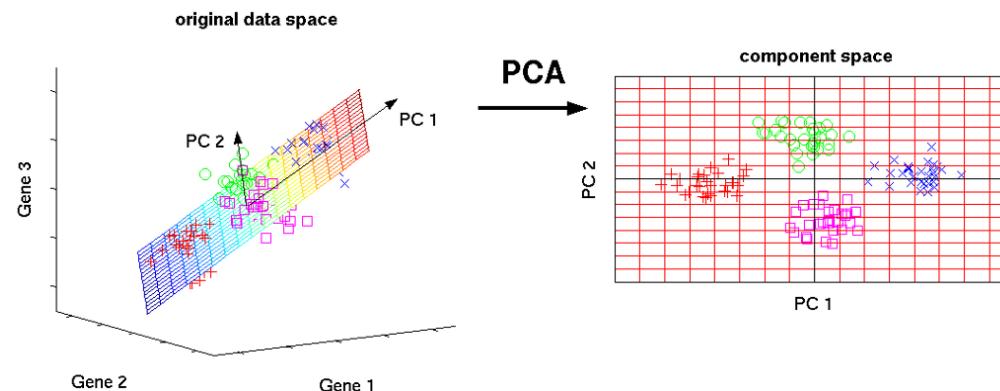
Is the data of good quality?

Quality as "fitness for purpose": DE

Very important: Transforming the data is often required for better further explorations (here: variance stabilization, regularized logarithm...)

One notable example, commonly used: Principal Components Analysis.

The data points (here, the samples) are projected onto the 2D plane such that they spread out in the two directions that explain most of the differences (the variability)



Exploratory analysis and visualization

"Perspective matters"



Gene identifiers

Take for example **GBP2**, guanylate binding protein 2 (4 transcripts available)

Ensembl ID: ENSG00000162645 (human), ENSMUSG00000028270 (mouse)

<http://www.ensembl.org/id/ENSG00000162645>

Entrez ID: 2634

HGNC ID: 4183

RefSeq ID: NM_004120

UCSC ID: uc001dmz.3

Official symbol: GBP2

Synonyms: - (sometimes makes the ambiguity even bigger!!!)

Gene symbols can change over time!

Typically, no 1:1 mapping between different ID types

Working with Excel

Don't do it - and not just because I do like R

Correspondence | [Open Access](#) | Published: 23 June 2004

Mistaken Identifiers: Gene name errors can be introduced inadvertently when using Excel in bioinformatics

[Barry R Zeeberg](#), [Joseph Riss](#), [David W Kane](#), [Kimberly J Bussey](#), [Edward Uchio](#), [W Marston Linehan](#), [J Carl Barrett](#) & [John N Weinstein](#) 

[BMC Bioinformatics](#) **5**, Article number: 80 (2004) | [Cite this article](#)

113k Accesses | **43** Citations | **515** Altmetric | [Metrics](#)

Working with Excel

There are excellent reasons *not* to do it

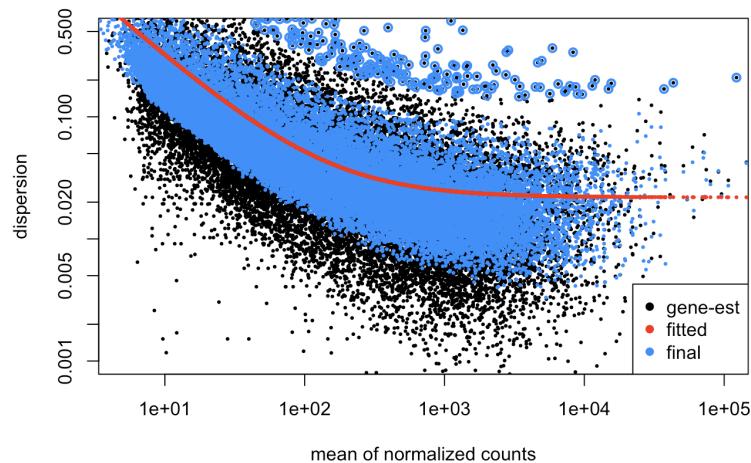
The screenshot shows a Microsoft Excel spreadsheet titled "excel.gene2date.xls". The data is organized into several columns:

	A	B	C	D	E	F	G	H	I	J	K
1	APR-1	35885	1-Apr	OCT-1	36068	1-Oct	SEP2	36039	2-Sep		
2	APR-2	35886	2-Apr	OCT-2	36069	2-Oct	SEP3	36040	3-Sep		
3	APR-3	35887	3-Apr	OCT-3	36070	3-Oct	SEP4	36041	4-Sep		
4	APR-4	35888	4-Apr	OCT-4	36071	4-Oct	SEP5	36042	5-Sep		
5	APR-5	35889	5-Apr	OCT-6	36073	6-Oct	SEP6	36043	6-Sep		
6	DEC-1	36129	1-Dec	OCT1	36068	1-Oct	SEPT1	36038	1-Sep		
7	DEC-2	36130	2-Dec	OCT11	36078	11-Oct	SEPT2	36039	2-Sep		
8	DEC1	36129	1-Dec	OCT2	36069	2-Oct	SEPT3	36040	3-Sep		
9	DEC2	36130	2-Dec	OCT3	36070	3-Oct	SEPT4	36041	4-Sep		
10	MAR1	35854	1-Mar	OCT4	36071	4-Oct	SEPT5	36042	5-Sep		
11	MAR2	35855	2-Mar	OCT6	36073	6-Oct	SEPT6	36043	6-Sep		
12	MAR3	35856	3-Mar	OCT7	36074	7-Oct	SEPT7	36044	7-Sep		
13	NOV1	36099	1-Nov	SEP-1	36038	1-Sep	SEPT8	36045	8-Sep		
14	NOV2	36100	2-Nov	SEP-2	36039	2-Sep	SEPT9	36046	9-Sep		
15				SEP1	36038	1-Sep					

Differential expression analysis, with DESeq2

There is one main function, `DESeq`, that does the following

- estimation of size factors
- estimation of the dispersion values for each gene
- fitting of the generalized linear model and performing statistical testing



What do our results look like?

	id	baseMean	log2FoldChange	IfcSE	stat	pvalue	padj
1	ENSG00000152583	997.43977	-4.5749187	0.18405609	-24.856111	2.220933e-136	4.003898e-132
2	ENSG00000165995	495.09291	-3.2910618	0.13317370	-24.712551	7.839410e-135	7.066444e-131
3	ENSG00000120129	3409.02938	-2.9478099	0.12143769	-24.274258	3.666925e-130	2.203577e-126
4	ENSG00000101347	12703.38706	-3.7669954	0.15543799	-24.234715	9.583815e-130	4.319425e-126
5	ENSG00000189221	2341.76725	-3.3535799	0.14178235	-23.653014	1.098955e-123	3.962392e-120
6	ENSG00000211445	12285.61515	-3.7304027	0.16583058	-22.495263	4.618318e-112	1.387651e-108
7	ENSG00000157214	3009.26322	-1.9767725	0.08998232	-21.968454	5.769975e-107	1.486016e-103
8	ENSG00000162614	5393.10168	-2.0356649	0.09418231	-21.614091	1.323849e-103	2.983294e-100
9	ENSG00000125148	3656.25278	-2.2109786	0.10564078	-20.929214	2.902397e-97	5.813824e-94
10	ENSG00000154734	30315.13547	-2.3456037	0.11580806	-20.254235	3.259713e-91	5.876611e-88
11	ENSG00000139132	1223.46614	-2.2289030	0.11283944	-19.752872	7.578284e-87	1.242012e-83
12	ENSG00000162493	1099.62587	-1.8912170	0.09577959	-19.745511	8.767336e-87	1.317146e-83
13	ENSG00000134243	5510.95826	-2.1957116	0.11200768	-19.603224	1.451321e-85	2.012647e-82
14	ENSG00000179094	776.59667	-3.1917499	0.16396292	-19.466291	2.120834e-84	2.731028e-81
15	ENSG00000162692	508.17023	3.6926606	0.19018590	19.416059	5.646022e-84	6.785765e-81
16	ENSG00000163884	561.10717	-4.4591282	0.23486369	-18.986027	2.225442e-80	2.507517e-77
17	ENSG00000178695	2649.85015	2.5281746	0.13443012	18.806607	6.667027e-79	7.070186e-76
18	ENSG00000198624	2057.19173	-2.9184357	0.16129634	-18.093626	3.577533e-73	3.583098e-70
19	ENSG00000107562	25136.30757	1.9116698	0.10574810	18.077582	4.786167e-73	4.541317e-70
20	ENSG00000148848	1365.17399	1.8145431	0.10048365	18.058094	6.813522e-73	6.141709e-70

I would like to understand more what is going on here

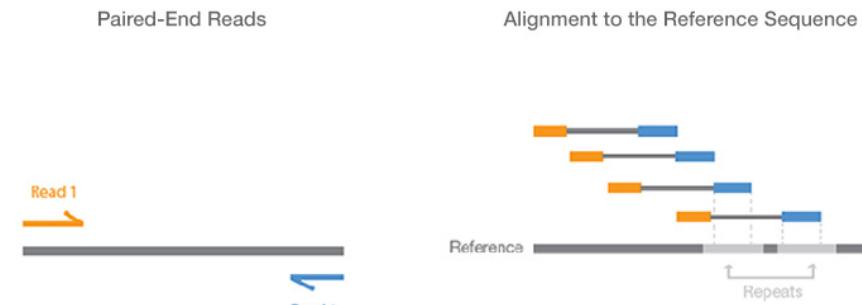
Practical tips for your (next) RNA-seq dataset

Single- vs paired-end sequencing

Each fragment can be sequenced from one end only, or from both ends

Single-end cheaper and faster

Paired-end provide improved ability to localize the fragment in the genome and resolve mapping close to repeat regions - less multimapping reads



Paired-end sequencing enables both ends of the DNA fragment to be sequenced. Because the distance between each paired read is known, alignment algorithms can use this information to map the reads over repetitive regions more precisely. This results in much better alignment of the reads, especially across difficult-to-sequence, repetitive regions of the genome.

Practical tips for your (next) RNA-seq dataset

Strand-specificity

In “standard” protocols, we don’t know from which strand a read stems

Various “strand-specific” protocols allow us to keep this information

Strand-specificity leads to lower number of ambiguous reads (overlapping multiple genes)

Practical tips for your (next) RNA-seq dataset

Sequencing depth and number of replicates

Recommendation: At least 3 biological replicates - Conesa et al, Genome Biol, 2016

Recommendation: At least 6 biological replicates - Schurch

Diminishing returns in increasing the read depth - go for replicates given a fixed budget!

I am not kidding you



zack chiang
@z_chiang · [Follow](#)



doing bulk sequencing analysis in 2022



7:28 PM · May 31, 2022



547

Reply

Copy link

[Read 3 replies](#)

Practical tips for your (next) RNA-seq dataset

Bulk or single cell?

That depends on your question!

Did someone say heterogeneity?

Questions

What are the steps to process RNA-Seq data?

- How to convert RNA-seq reads into counts?
- How to perform quality control (QC) of RNA-seq reads?

How to identify differentially expressed genes across multiple experimental conditions?

- How to properly analyze RNA count data using DESeq2?
- How to perform quality control (QC) and exploratory data analysis (EDA) of RNA-seq count data?

What are the biological functions impacted by the differential expression of genes?

- How can I perform a gene ontology enrichment analysis?

How can I create neat visualizations of the data?

- How can I visualize the results for my enrichment analysis?

How can I generate interactive reports to summarise my analyses?

Questions

What are the steps to process RNA-Seq data?

- How to convert RNA-seq reads into counts?
- How to perform quality control (QC) of RNA-seq reads?

How to identify differentially expressed genes across multiple experimental conditions?

- How to properly analyze RNA count data using DESeq2?
- How to perform quality control (QC) and exploratory data analysis (EDA) of RNA-seq count data?

What are the biological functions impacted by the differential expression of genes?

- How can I perform a gene ontology enrichment analysis?

How can I create neat visualizations of the data?

- How can I visualize the results for my enrichment analysis?

How can I generate interactive reports to summarise my analyses?

Goals

after constructing and running a differential gene expression analysis...

- Perform a functional enrichment analysis
- Perform and visualize the enrichment analysis for Gene Ontology terms
- Get familiar with representations of the relationship between genes and gene sets

Prerequisites

- a basic understanding of R
- familiarity with gene expression data (RNA-seq)
- familiarity with the concept of differential expression analysis

What do our results look like?

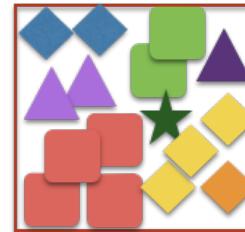
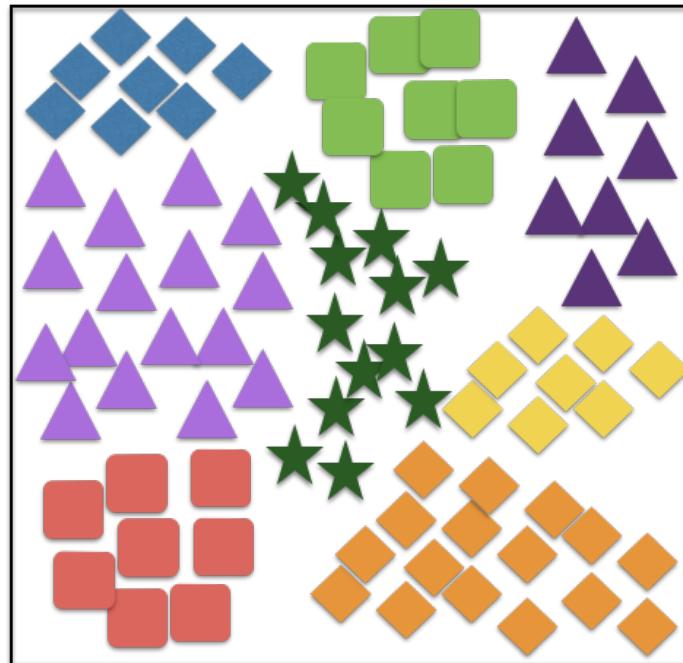
	id	baseMean	log2FoldChange	IfcSE	stat	pvalue	padj
1	ENSG00000152583	997.43977	-4.5749187	0.18405609	-24.856111	2.220933e-136	4.003898e-132
2	ENSG00000165995	495.09291	-3.2910618	0.13317370	-24.712551	7.839410e-135	7.066444e-131
3	ENSG00000120129	3409.02938	-2.9478099	0.12143769	-24.274258	3.666925e-130	2.203577e-126
4	ENSG00000101347	12703.38706	-3.7669954	0.15543799	-24.234715	9.583815e-130	4.319425e-126
5	ENSG00000189221	2341.76725	-3.3535799	0.14178235	-23.653014	1.098955e-123	3.962392e-120
6	ENSG00000211445	12285.61515	-3.7304027	0.16583058	-22.495263	4.618318e-112	1.387651e-108
7	ENSG00000157214	3009.26322	-1.9767725	0.08998232	-21.968454	5.769975e-107	1.486016e-103
8	ENSG00000162614	5393.10168	-2.0356649	0.09418231	-21.614091	1.323849e-103	2.983294e-100
9	ENSG00000125148	3656.25278	-2.2109786	0.10564078	-20.929214	2.902397e-97	5.813824e-94
10	ENSG00000154734	30315.13547	-2.3456037	0.11580806	-20.254235	3.259713e-91	5.876611e-88
11	ENSG00000139132	1223.46614	-2.2289030	0.11283944	-19.752872	7.578284e-87	1.242012e-83
12	ENSG00000162493	1099.62587	-1.8912170	0.09577959	-19.745511	8.767336e-87	1.317146e-83
13	ENSG00000134243	5510.95826	-2.1957116	0.11200768	-19.603224	1.451321e-85	2.012647e-82
14	ENSG00000179094	776.59667	-3.1917499	0.16396292	-19.466291	2.120834e-84	2.731028e-81
15	ENSG00000162692	508.17023	3.6926606	0.19018590	19.416059	5.646022e-84	6.785765e-81
16	ENSG00000163884	561.10717	-4.4591282	0.23486369	-18.986027	2.225442e-80	2.507517e-77
17	ENSG00000178695	2649.85015	2.5281746	0.13443012	18.806607	6.667027e-79	7.070186e-76
18	ENSG00000198624	2057.19173	-2.9184357	0.16129634	-18.093626	3.577533e-73	3.583098e-70
19	ENSG00000107562	25136.30757	1.9116698	0.10574810	18.077582	4.786167e-73	4.541317e-70
20	ENSG00000148848	1365.17399	1.8145431	0.10048365	18.058094	6.813522e-73	6.141709e-70

I need to understand more what is going on here

Functional enrichment analysis

Our problem, at a glance

All known genes in a species
(categorized into groups)



DEGs

Functional enrichment analysis

Test whether known biological functions or processes are over-represented (= enriched) in an experimentally-derived gene list, e.g. a list of differentially expressed (DE) genes.

Example: Transcriptomic study, in which 12671 genes have been tested for differential expression between two sample conditions and 529 genes were found DE

Among the DE genes, 28 are annotated to a specific functional gene set, which contains in total 170 genes. This setup corresponds to a 2x2 contingency table:

```
deTable <- matrix(  
  c(28, 142, 501, 12000),  
  nrow = 2,  
  dimnames = list(c("DE", "Not_DE"),  
                 c("In_gene_set", "Not_in_gene_set")))
```

```
deTable  
#           In_gene_set Not_in_gene_set  
# DE             28          501  
# Not_DE        142         12000
```

```
fisher.test(deTable, alternative = "greater")  
#  
#      Fisher's Exact Test for Count Data
```

Some terminology: Gene sets, pathways, networks

Gene sets are simple lists of usually functionally related genes without further specification of relationships between genes. **Any *a priori*** classification of genes into biologically relevant groups. Sets do not need to be exhaustive or disjoint

Pathways can be interpreted as specific gene sets, typically representing a group of genes that work together in a biological process.

Gene regulatory networks describe the interplay and effects of regulatory factors (such as transcription factors and microRNAs) on the expression of their target genes.

Is the expression of genes in a gene set associated with the experimental conditions?

Genes categories	Organism-specific Background	DE results	Over-represented?
Functional category 1	35/13000	25/1000	Likely
Functional category 2	56/13000	4/1000	Unlikely
Functional category 3	90/13000	8/1000	Unlikely
Functional category 4	15/13000	10/1000	Likely
...			
...			

Functional enrichment analysis

Primarily developed and applied on transcriptomic data, now extended and applied also in other fields of genomic and biomedical research (proteomic and metabolomic data, genomic regions, ...)

Many methods available!

- Over-representation analysis (ORA) - are differentially expressed (DE) genes in the set more common than expected? Based on (variations) of the 2x2 contingency table method
- Functional class scoring (FCS) - summarize gene set (= functional class) scores by summarizing statistic of DE of genes in a set, and compare to null
- Pathway topology (PT) - explicitly taking into account interactions between genes as defined in signaling pathways and gene regulatory networks

Topology-based methods appear to be most realistic, but: features that are not-detectable on the transcriptional level + insufficient network knowledge

Cautious interpretation of results is required to derive valid conclusions!

Collections of gene sets

Gene Ontology ([GO](#)) Annotation (GOA)

- CC Cellular Components
- BP Biological Processes
- MF Molecular Function

Pathways

- [MSigDb](#)
- [KEGG](#)
- [reactome](#)
- [PantherDB](#)
- ...

GO and KEGG annotations are most frequently used for the enrichment analysis of functional gene sets - likely due to their long-standing curation and availability for a wide range of species

Software for functional enrichment

- DAVID
- GSEA and its variations
- inside Bioconductor: [reactome](#), [topGO](#), [pathview](#), [GOSeq](#), ...
- Ingenuity Pathway Analysis

Not straightforward, but amazing way of helping generating hypothesis for the bench scientist

Friendly tip: wrappers for topGO and goseq are in [pcaExplorer](#) and in [ideal!](#)

General recommendations

Never forget to visualize your data!

Pick candidates, sets of candidates, plot instead of guessing!

The choice of a background matters!

Background matters

... with a slight modification of the example above

```
deTable_detected <- matrix(  
  c(8, 142, 501, 12000),  
  nrow = 2,  
  dimnames = list(c("DE", "Not_DE"),  
                 c("In_gene_set", "Not_in_gene_set")))  
fisher.test(deTable_detected, alternative = "greater")
```

versus...

```
deTable_allgenes <- matrix(  
  c(8, 142, 501, 32000),  
  nrow = 2,  
  dimnames = list(c("DE", "Not_DE"),  
                 c("In_gene_set", "Not_in_gene_set")))  
fisher.test(deTable_allgenes, alternative = "greater")
```

Data analysis: where's the real bottleneck?

- ✓ Efficient methods to process your data (entire workflows available, see e.g. [the Bioconductor workflow packages](#))
- ✓ Compelling ways of visualizing your data (better QC, better hypotheses, better answers → see <http://bioconductor.org/packages/iSEE/> & <http://bioconductor.org/packages/iSEEu/>)
- ✓ Powerful framework to communicate, interactively and reproducibly

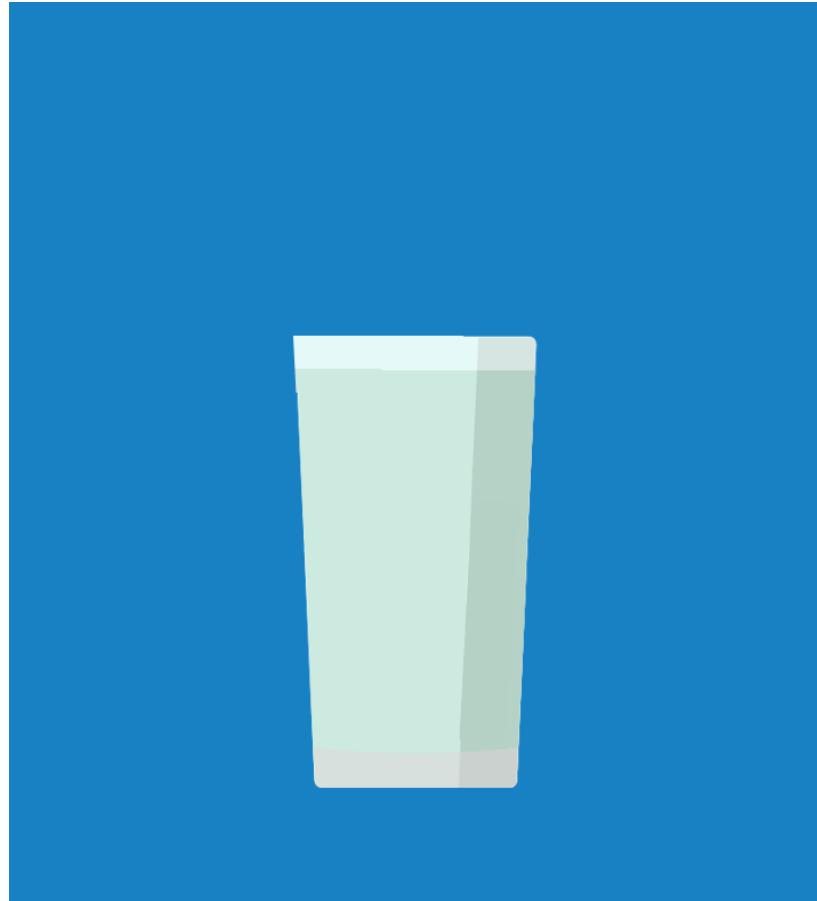
🔥 Data interpretation



... just not in plain sight

Show me your data...

A cocktail recipe



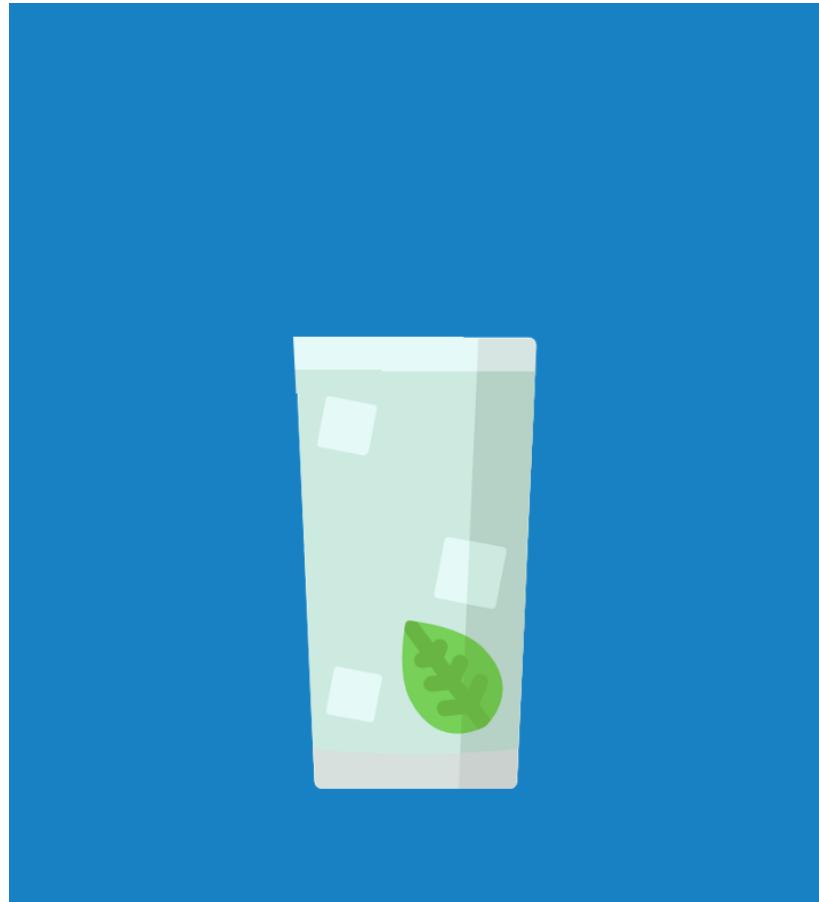
*You might (and should) have a set of standardized objects for your datasets and results (Excel does **not** count)*

A cocktail recipe



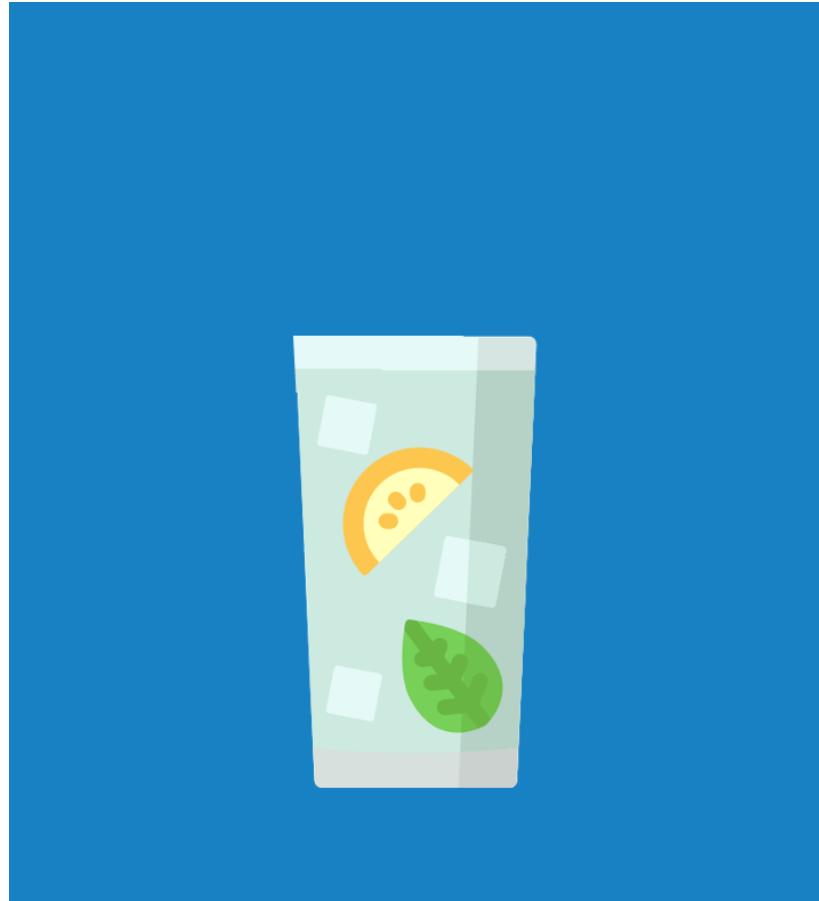
A container for the expression matrix, `dds`

A cocktail recipe



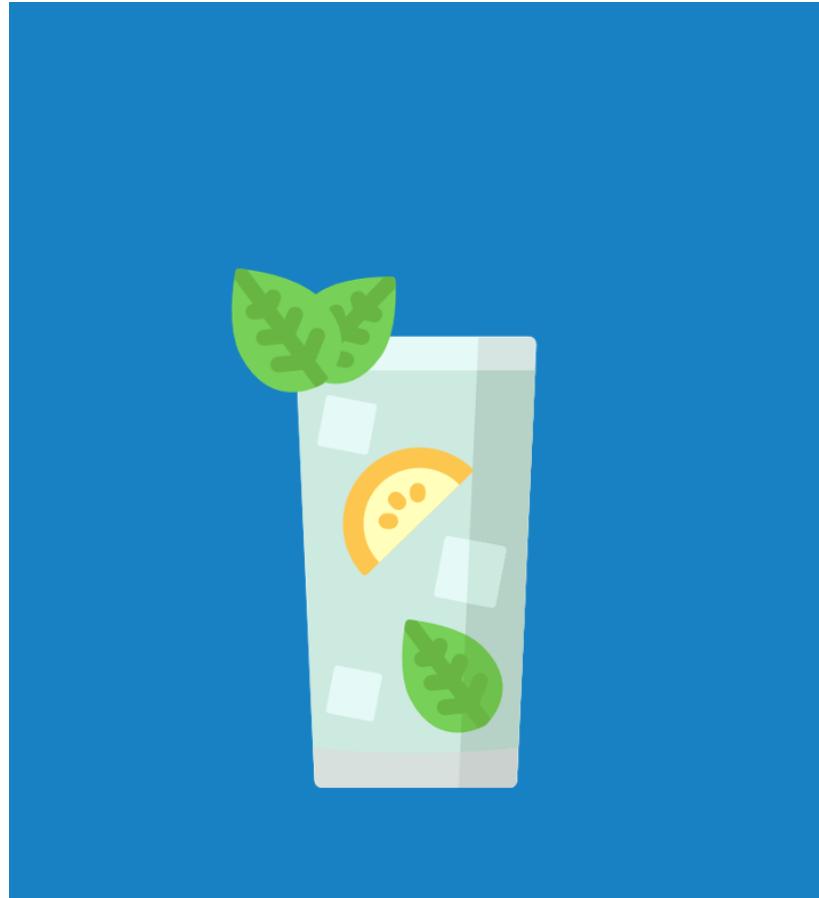
A container for the results of differential expression, `res_de`

A cocktail recipe



A table for the result of functional enrichment, `res_enrich`

A cocktail recipe



Decorate with an annotation table, [anno_df](#)

A cocktail recipe



Shaken, not stirred...

A cocktail recipe



GeneTonic is now on Bioconductor!

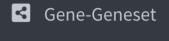
http://127.0.0.1:5184 | Open in Browser | [Bookmark](#)

~/Development/GeneTonic - Shiny

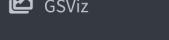
 GeneTonic 

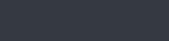
≡ GeneTonic 

Welcome! 

Gene-Geneset 

Enrichment Map 

Overview 

GSViz 

Bookmarks

First steps with a full glass



Overview on the provided input

Expression Matrix + DE results +

Functional analysis results + Annotation info +

58294 genes x 24 samples
dds object 

928 DE genes
res object 

500 functional categories
func enrich object 

2 feature identifiers for 58294 features
annotation object 



GeneTonic is a project developed by [Federico Marini](#) in the Bioinformatics division of the [IMBEI](#) - Institute for Medical Biostatistics, Epidemiology and Informatics
License: [MIT](#) - The GeneTonic package is developed and available on [GitHub](#)

Welcome!

Gene-Geneset

Enrichment Map

DEview

GeneSets

Bookmarks

About

Overview on the provided input

Expression Matrix

+

DE results

+

Functional analysis results

+

Annotation info

+

58294 genes x 24 samples



dds object

928 DE genes



res object

500 functional categories



func enrich object

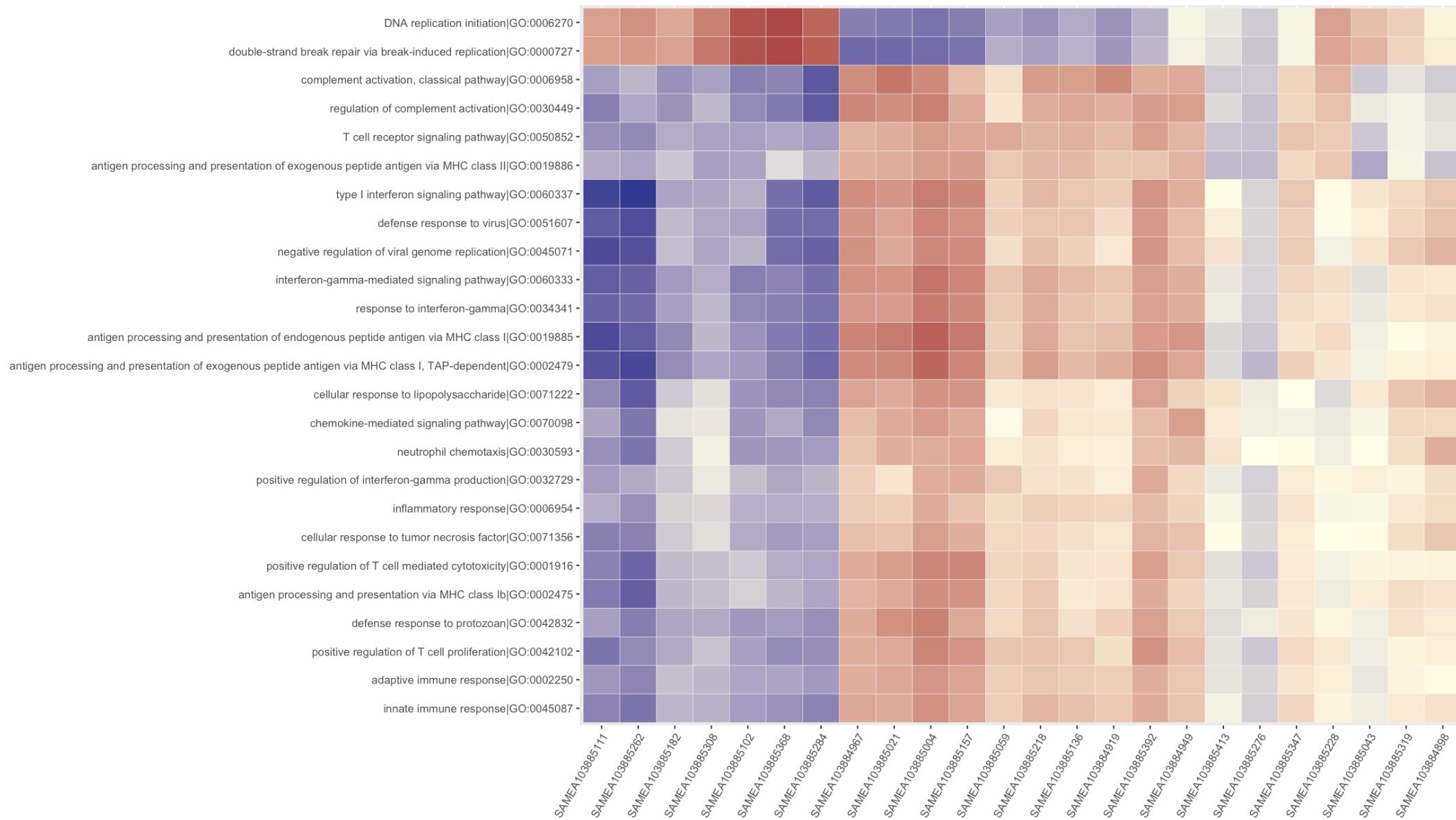
2 feature identifiers for 58294 features

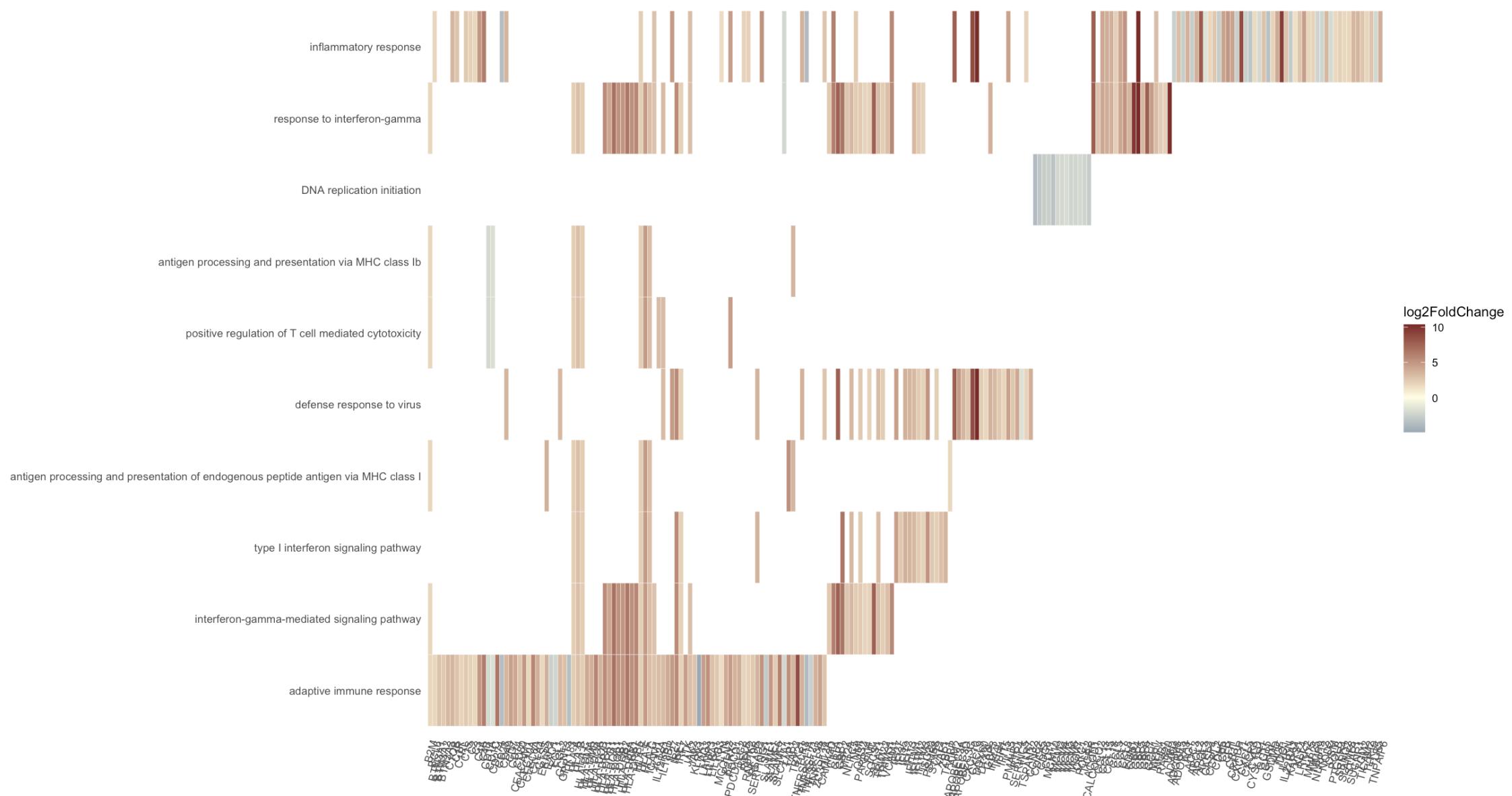


annotation object



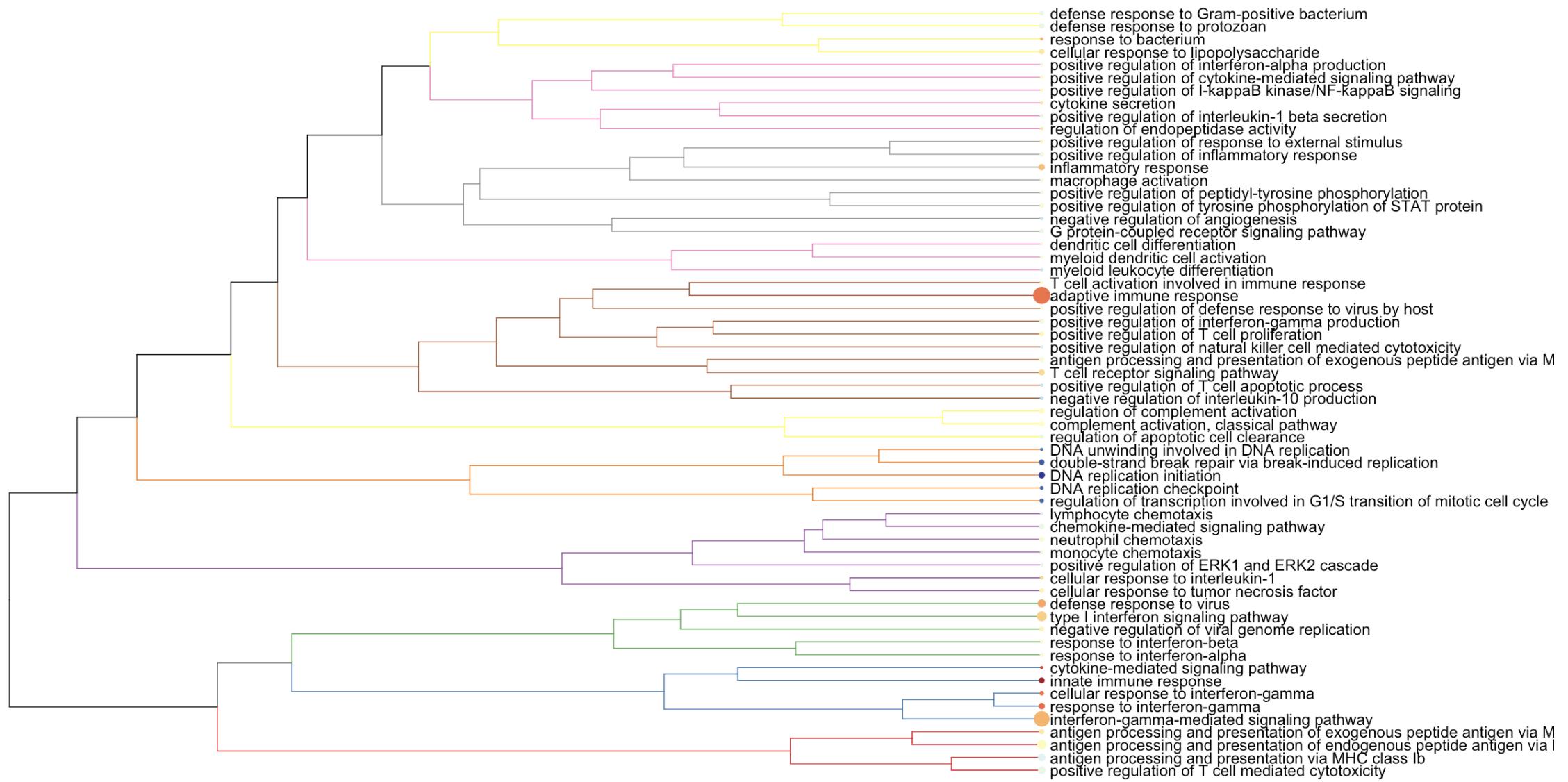




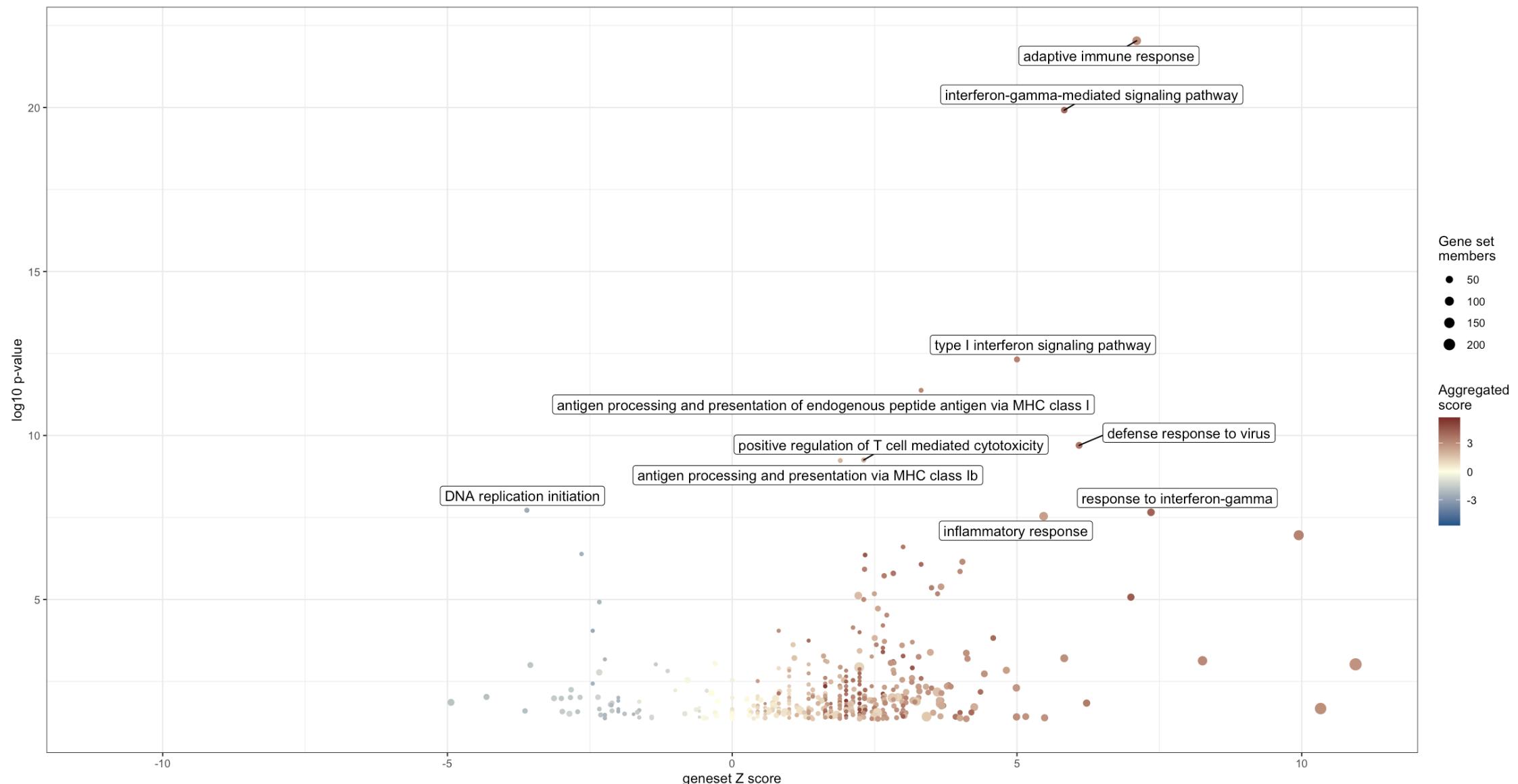


Geneset MDS plot - condition IFNg vs naive

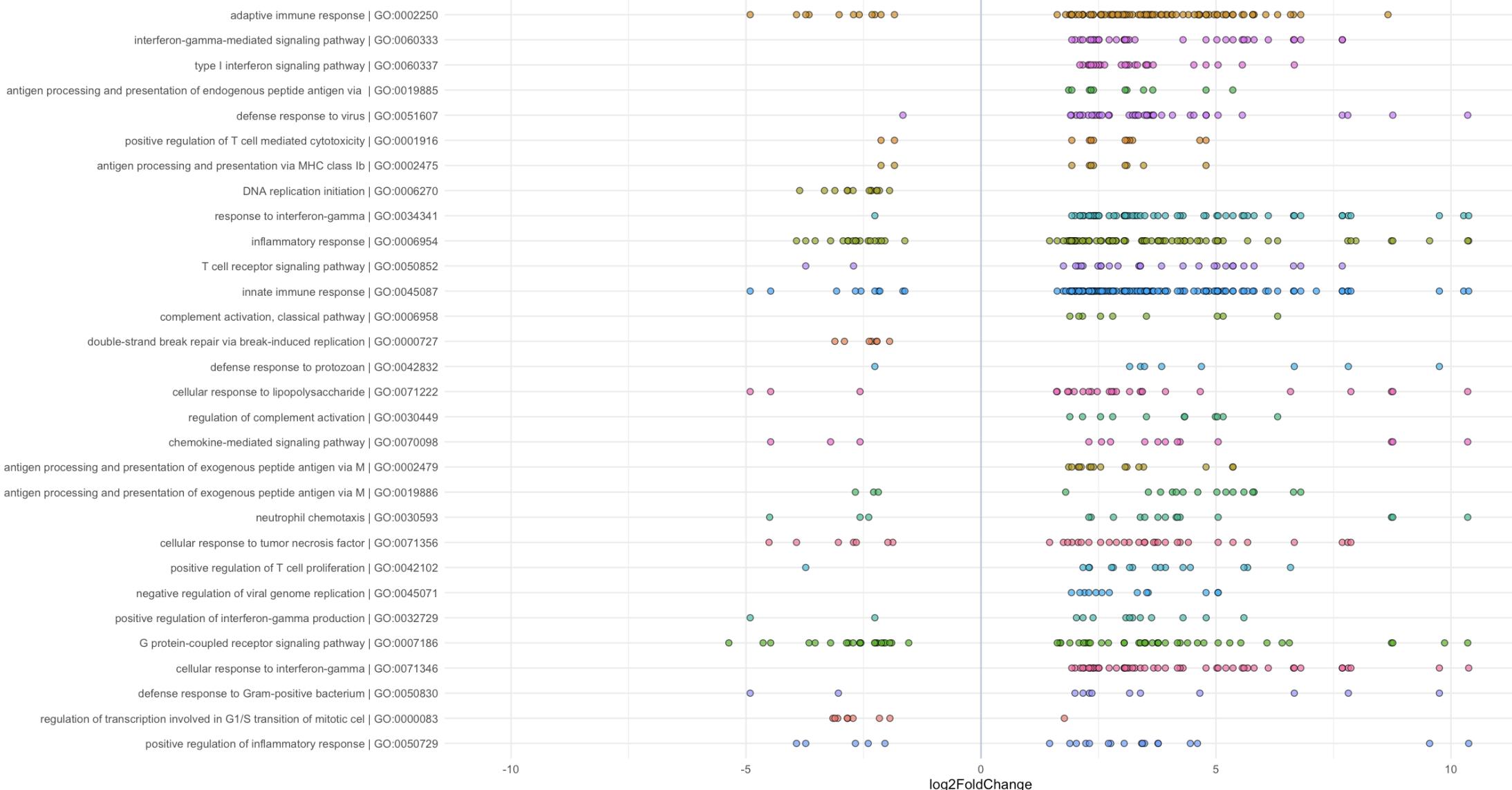




Geneset volcano



Enrichment overview - condition IFNg vs naive



Welcome!

Gene-Geneset

Enrichment Map

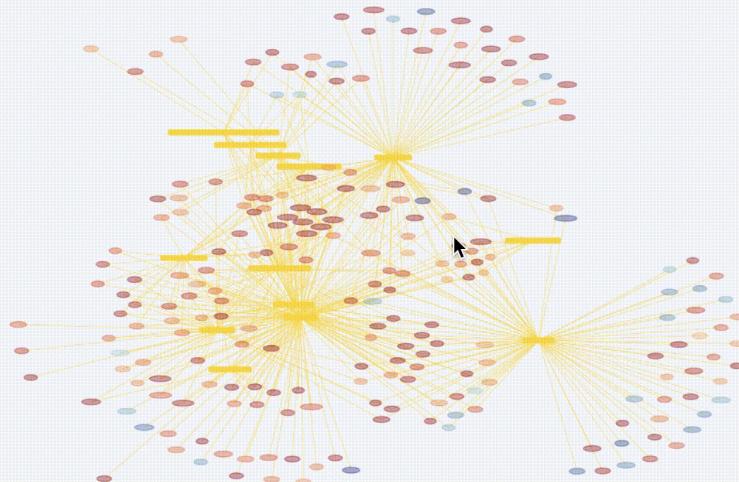
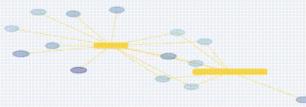
DView

GeneSets

Bookmarks

About

Select by id



Genesetbox

Please select a gene set.



Genebox

Please select a gene/feature.



The interplay of genes and genesets

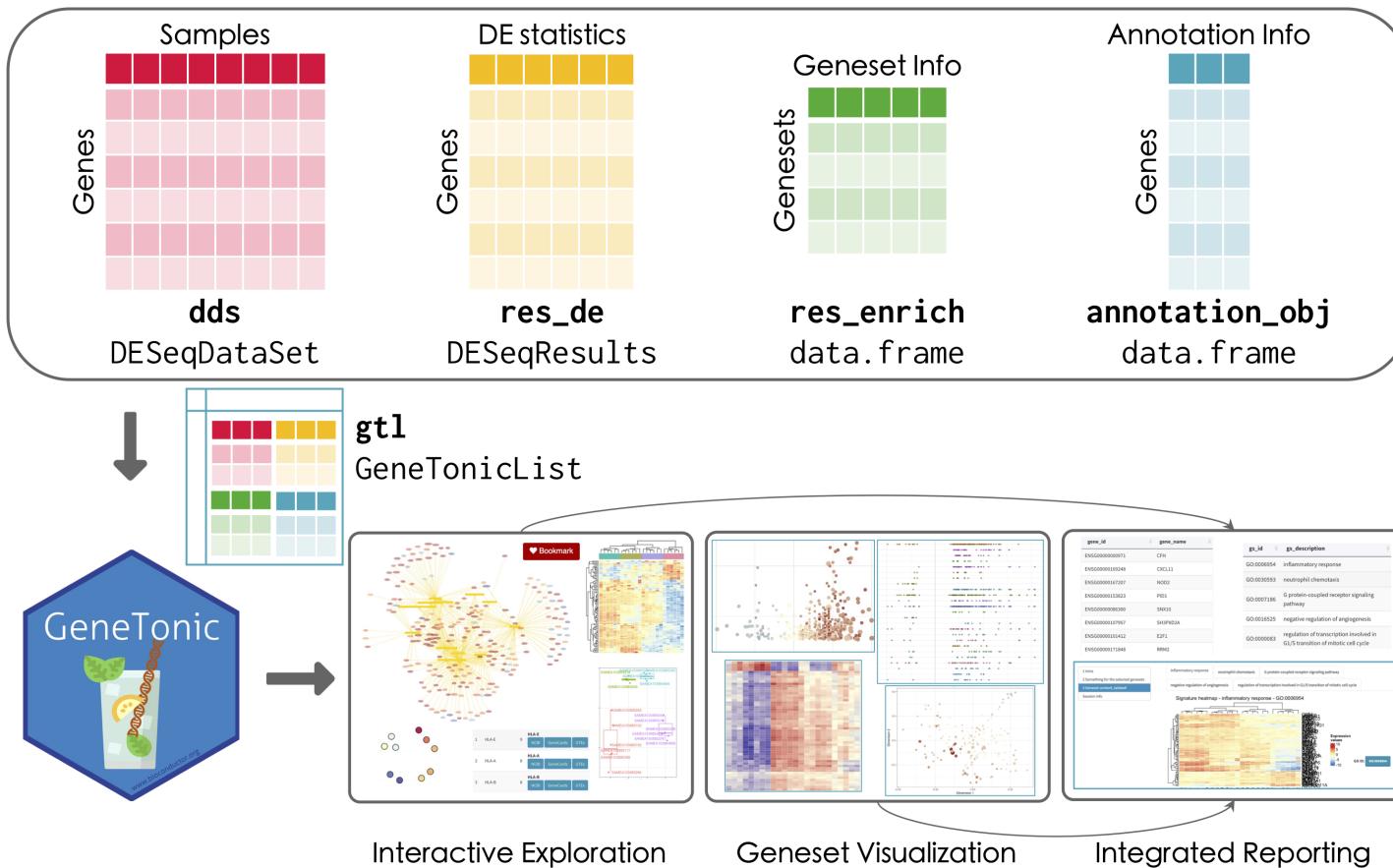
A quick demo

```
BiocManager::install("GeneTonic")
library("GeneTonic")
example(GeneTonic, ask = FALSE)
```


GeneTonic



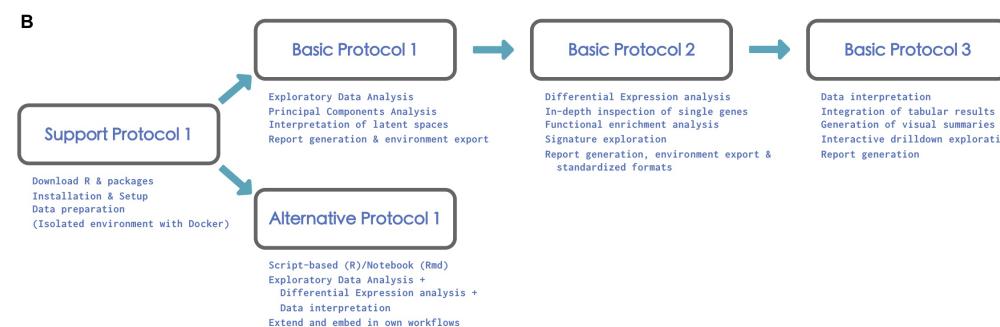
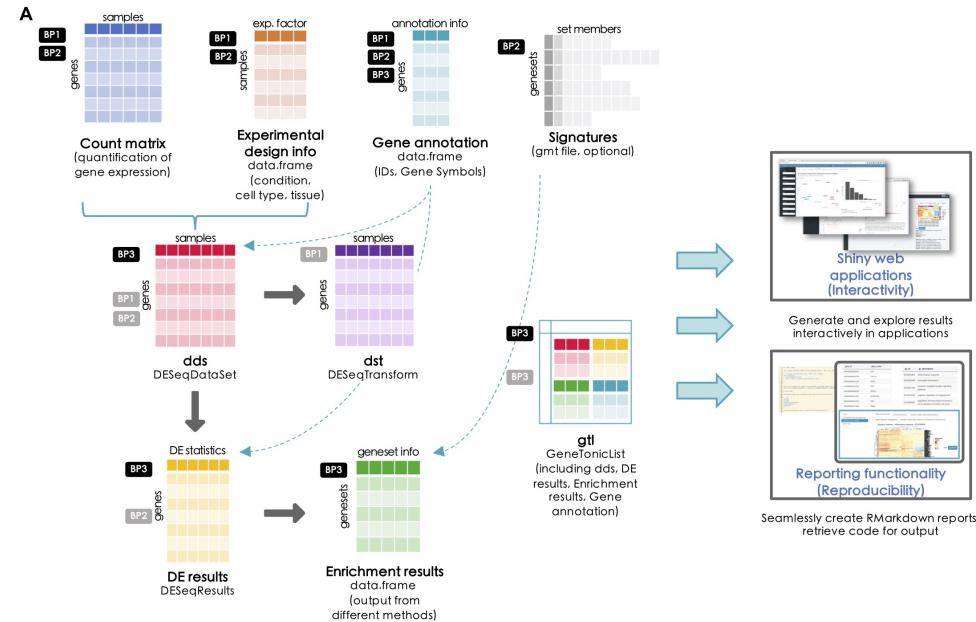
A way to enjoy the data interpretation...



Suddenly i SEE



Putting it all together



Interfaces for RNA-seq

Systematically evaluating interfaces for RNA-seq analysis from a life scientist perspective FREE

Alicia Poplawski, Federico Marini, Moritz Hess, Tanja Zeller, Johanna Mazur,
Harald Binder

Briefings in Bioinformatics, Volume 17, Issue 2, March 2016, Pages 213–223,
<https://doi.org/10.1093/bib/bbv036>

Published: 23 June 2015 Article history ▾

Interfaces have been developed for making such analysis steps accessible to life scientists without extensive knowledge of command line tools.

Systematic search and evaluation of such interfaces

Definition of criteria for evaluation (ease of configuration, documentation, usability, computational demand, reporting)

There's no **one tool fits all** winner!

A woman with short black hair and bangs, wearing a white zip-up top and patterned pants, sits cross-legged on a grey sofa. She is smiling and looking towards the camera. The background features a light-colored wooden wall with a shelf holding small framed pictures and a stack of books. There are several potted plants in the room, including a large one on the left and a tall one on the right. A pink blanket is draped over the back of the sofa.

Does the analysis/exploration/interpretation of your data
spark joy?

Lessons learnt for me (so far)

It takes two to tango

Results & Interpretation

Bioinformatician & Experimental scientist

Reproducibility & Interactivity

Stop using Excel for performing Bioinformatics data analysis :)

Our contributions to the Bioconductor project:

`pcaExplorer` - Exploratory Data Analysis

`ideal` - Differential Expression analysis

`GeneTonic` - Interpretation of DE results

`iSEE` - Exploration of single cell datasets

...

`countsimQC` - comparing count data sets

`iCOBRA` - calculation and visualization of performance metrics

`ExploreModelMatrix`

`tximeta` - tx importing with metadata, *for free*

`alevinQC`

Practical session

References:

- Marini and Binder (2016) - Development of Applications for Interactive and Reproducible Research: a Case Study (Genomics Computational Biology) [10.18547/gcb.2017.vol3.iss1.e39](https://doi.org/10.18547/gcb.2017.vol3.iss1.e39)
- Marini and Binder (2019) - pcaExplorer: an R/Bioconductor package for interacting with RNA-seq principal components (BMC Bioinformatics) [10.1186/s12859-019-2879-1](https://doi.org/10.1186/s12859-019-2879-1)
- Marini, Linke, Binder (2020) - ideal: an R/Bioconductor package for Interactive Differential Expression Analysis (BMC Bioinformatics) [10.1186/s12859-020-03819-5](https://doi.org/10.1186/s12859-020-03819-5)
- Rue-Albrecht, Marini, Soneson, Lun (2018) - iSEE: Interactive SummarizedExperiment Explorer (F1000 Research) <https://doi.org/10.12688/f1000research.14966.1>
- Marini, Ludt, Linke, Strauch (2021) - GeneTonic: an R/Bioconductor package for streamlining the interpretation of RNA-seq data (BMC Bioinformatics) <https://doi.org/10.1186/s12859-021-04461-5>
- Ludt, Ustjanzew, Binder, Strauch, Marini (2022) - Interactive and Reproducible Workflows for Exploring and Modeling RNA-seq Data with pcaExplorer, Ideal, and GeneTonic (Current Protocols) <https://doi.org/10.1002/cpz1.411>

References: (cont'd)

- Srivastava, Malik, Sarkar, Zakeri, Almodaresi, Soneson, Love, Kingsford, Patro (2020) - Alignment and mapping methodology influence transcript abundance estimation (Genome Biology)
<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-020-02151-8>
- Conesa, Madrigal, Tarazona, Gomez-Cabrero, Cervera, McPherson, Szcześniak, Gaffney, Elo, Zhang, Mortazavi (2016) - A survey of best practices for RNA-seq data analysis (Genome Biology)
<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-016-0881-8>

Resources

The whole `rnaseqGene` workflow is detailed here: <https://www.bioconductor.org/packages/rnaseqGene/>

... thank you for your attention!

marinif@uni-mainz.de -  @FedeBioinfo charlotte.soneson@fmi.ch -  @CSoneson

