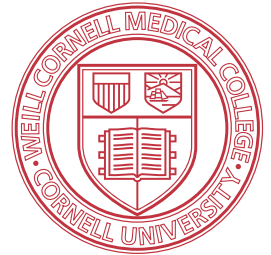


# Differential gene expression analysis using RNA-seq

---

Applied Bioinformatics Core, March 2018



Friederike Dündar with Luce Skrabanek & Paul Zumbo

# Day 2: Aligning reads

1. Experimental Design
2. Reference genome & transcript annotation
3. Alignment
  - STAR
  - BAM/SAM files
4. QC

# EXPERIMENTAL DESIGN

---

How to avoid spurious signals and drowning in noise

# Why do we need replicates?

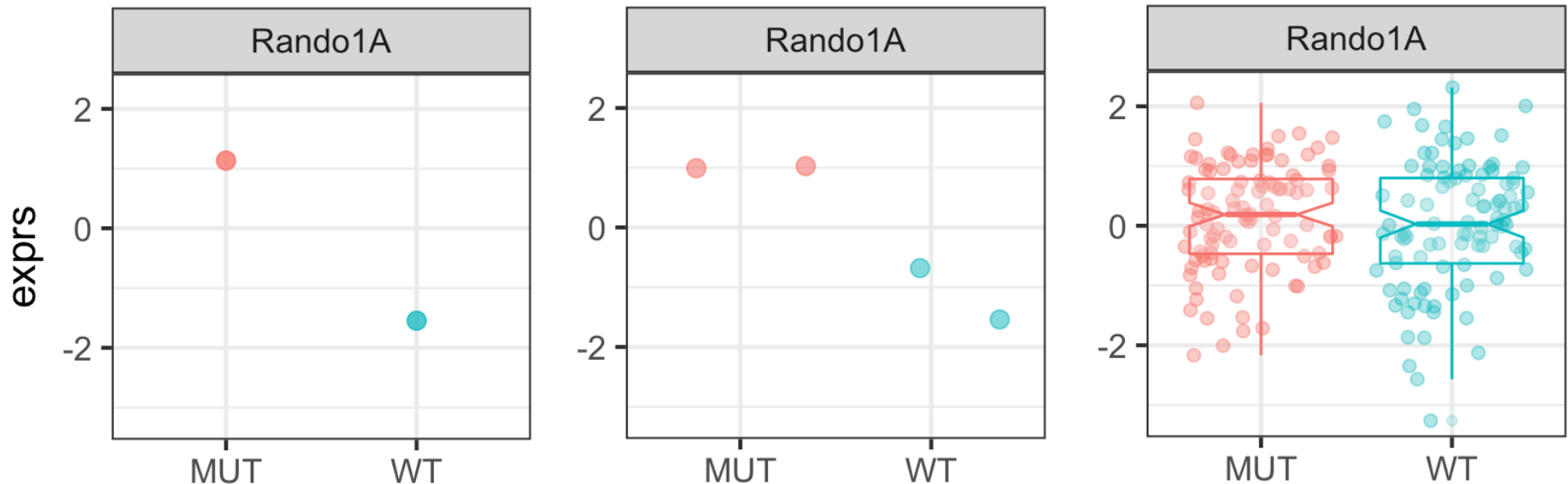
**Goal:** Identify differences in expression for every gene.

...and “differences” should preferably be due to our experiment, not noise!

*“Samples are our windows to the population, and their statistics are used to estimate those of the population.”*

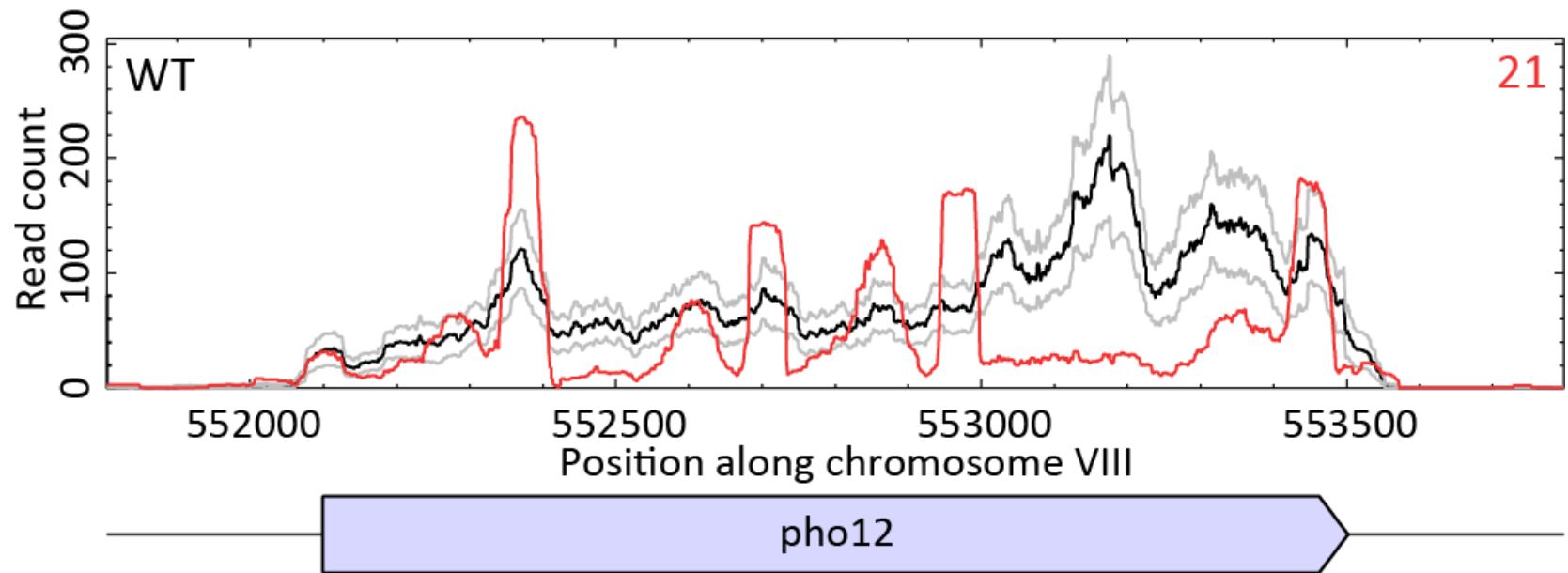
Martin Krzywinski & Naomi Altman

```
testdat <- data.frame(exprs = rnorm(200),  
                      condition = c("WT", "MUT"),  
                      gene_name = "Rando1A")
```



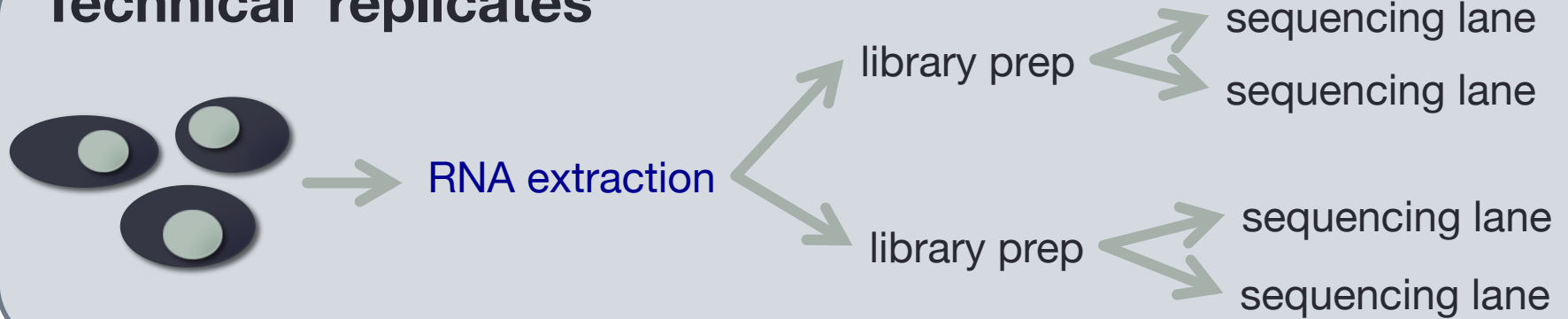
# Invest in replicates!

- recommended: **6 biological replicates per condition** for DGE of strongly changing genes ( $\log FC \geq 2$ ) [based on insights from the fairly simple yeast transcriptome]



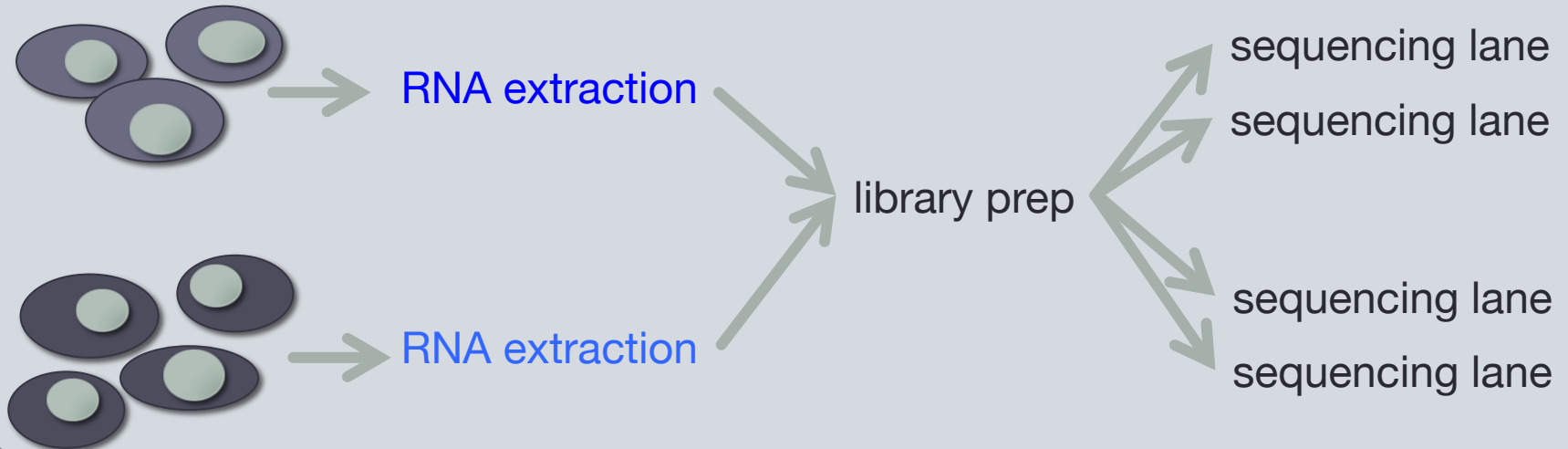
# Replicate types

## Technical replicates

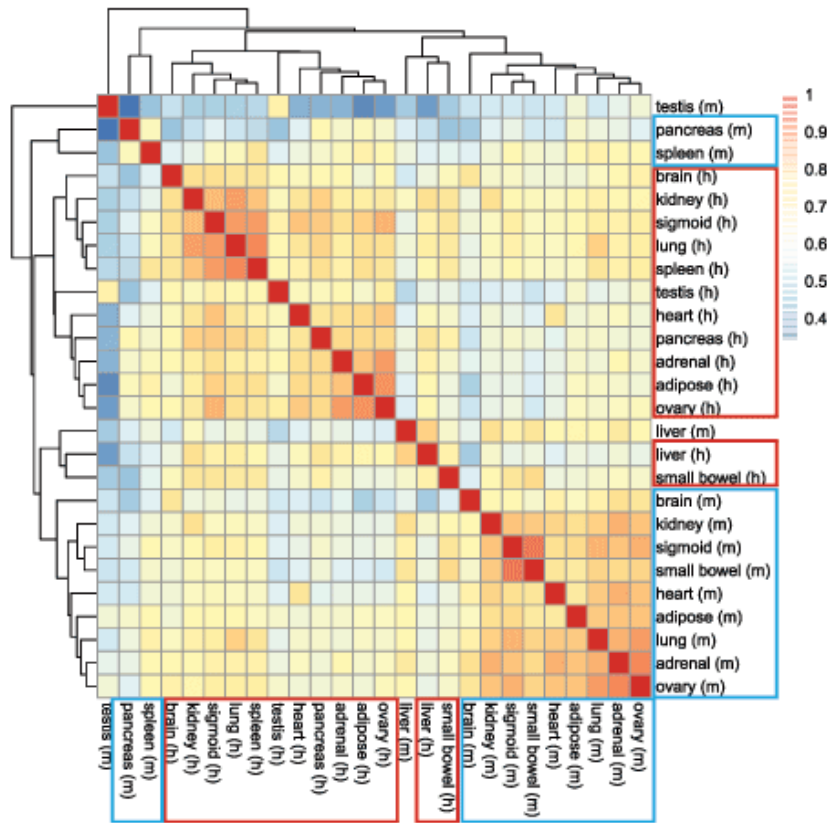


## Biological replicates

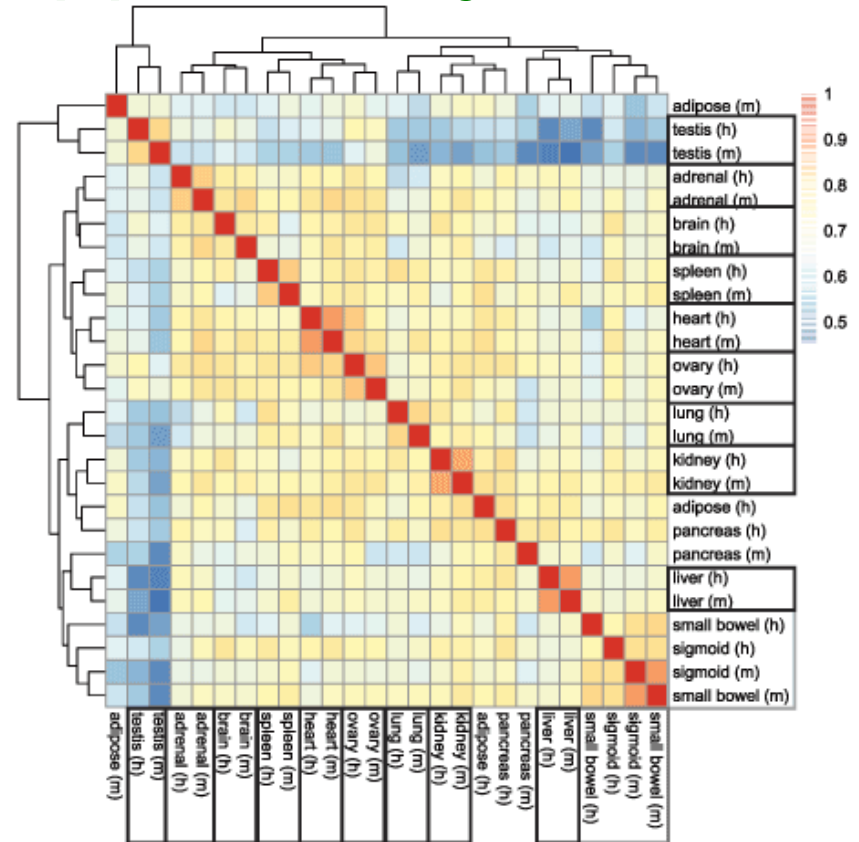
RNA from an independent growth of cells/tissue



# Batch effects can happen everywhere



*“Overall, our results indicate that there is **considerable RNA expression diversity between humans and mice**, well beyond what was described previously, likely reflecting the fundamental physiological differences between these two organisms.”*



*“Once we accounted for the batch effect (...), the comparative gene expression data no longer clustered by species, and instead, we observed **a clear tendency for clustering by tissue.**”*

# ENCODE's\* study design was not optimal

D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7)	D87PMJN1 (run 253, flow cell D2GUAACXX , lane 8)	D4LHBFN1 (run 276, flow cell C2HKJACXX , lane 4)	MONK (run 312, flow cell C2GR3ACXX , lane 6)	HWI-ST373 (run 375, flow cell C3172ACXX , lane 7)
heart	adipose	adipose	heart	brain
kidney	adrenal	adrenal	kidney	pancreas
liver	sigmoid colon	sigmoid colon	liver	brain
small bowel	lung	lung	small bowel	spleen
spleen	ovary	ovary	testis	● Human
testis		pancreas		● Mouse

Tissue	Human	Mouse
adipose	FEMALE	MALE
adrenal	MALE	FEMALE
brain	FEMALE	MALE
heart	FEMALE	FEMALE
kidney	MALE	FEMALE
liver	MALE	FEMALE
lung	FEMALE	FEMALE
ovary	FEMALE	FEMALE
pancreas	FEMALE	FEMALE
sigmoid colo	MALE	FEMALE
small bowel	FEMALE	FEMALE
spleen	FEMALE	MALE
testis	MALE	MALE

Tissue was confounded with (at least):

- sequencer
- sex
- age
- tissue handling

human data: deceased organ donors  
mouse data: 10-week-old littermates

A very good read (including the reviews and comments) that discusses many scientific as well as ethical issues: <https://f1000research.com/articles/4-121/v1>



# Avoiding bias

## Completely randomized design

STRESS	A	B	A	A	B	A	B	A	A	B	B	B
DIET	1	2	1	2	2	1	1	2	2	1	2	1

## Restricted randomized design

GENOTYPE	A	A	A	A	A	A	B	B	B	B	B	B
DIET	1	2	1	2	2	1	1	2	1	1	2	2

## Blocked & randomized design

GENOTYPE	A	A	B	B	A	A	B	B	A	A	B	B
DIET	1	2	1	2	1	2	1	2	1	2	1	2
WEIGHT	•	•	•	•	•	•	•	•	•	•	•	•



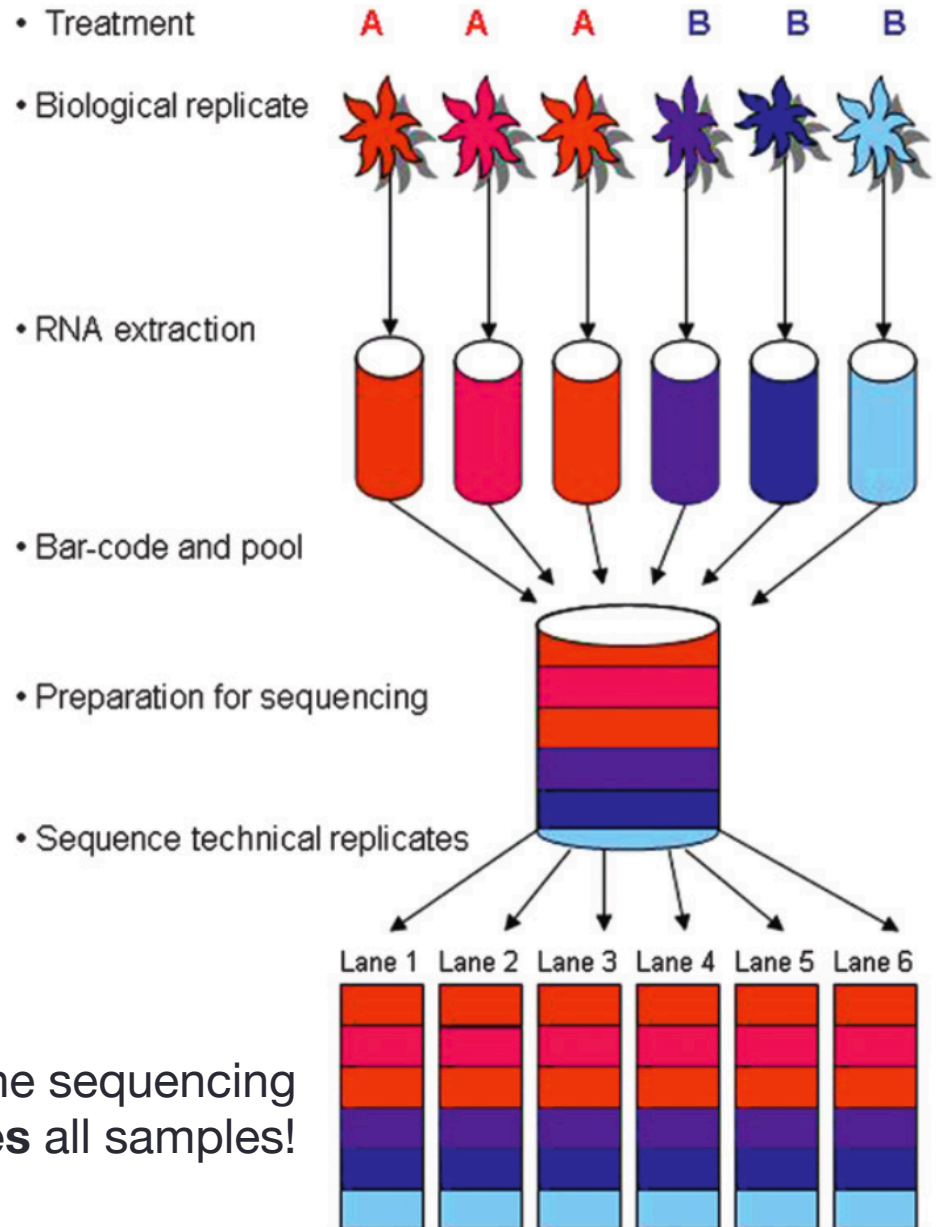
**Block** what you can,  
**randomize** what you cannot.

*What factors are of interest? Which ones might introduce noise?  
Which nuisance factors do you absolutely need to account for?*

# Typical RNA-seq set-up

- keep the **technical nuisance** factors (harvest date, RNA extraction kit, sequencing date...) to a **minimum**
- cover only as much of the **biological variation as needed** (just keep possible restrictions about your conclusions in mind for later)

Make sure the sequencing core **multiplexes** all samples!



# How deep is deep enough?

for DGE ( $\log FC \sim 2$ ) in mammals:  
20 – 50 mio SR, 75 bp

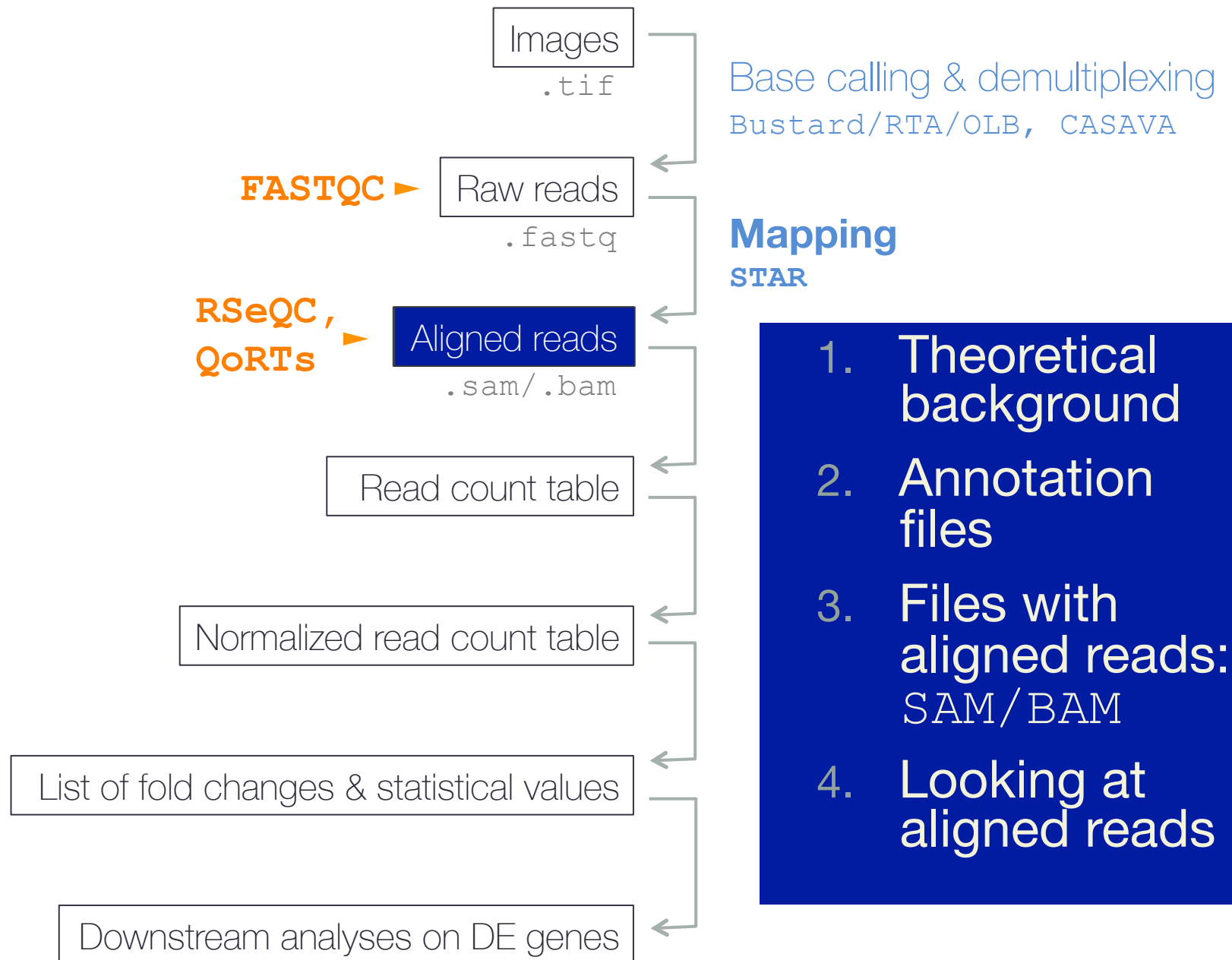
Goals that require **more**, **longer**, and possibly **paired-end** reads:

- quantification of **lowly expressed** genes
- identification of genes with **small changes** between conditions
- investigation of **alternative splicing**/isoform quantification
- identification of **novel transcripts**, chimeric transcripts
- *de novo* **transcriptome assembly**

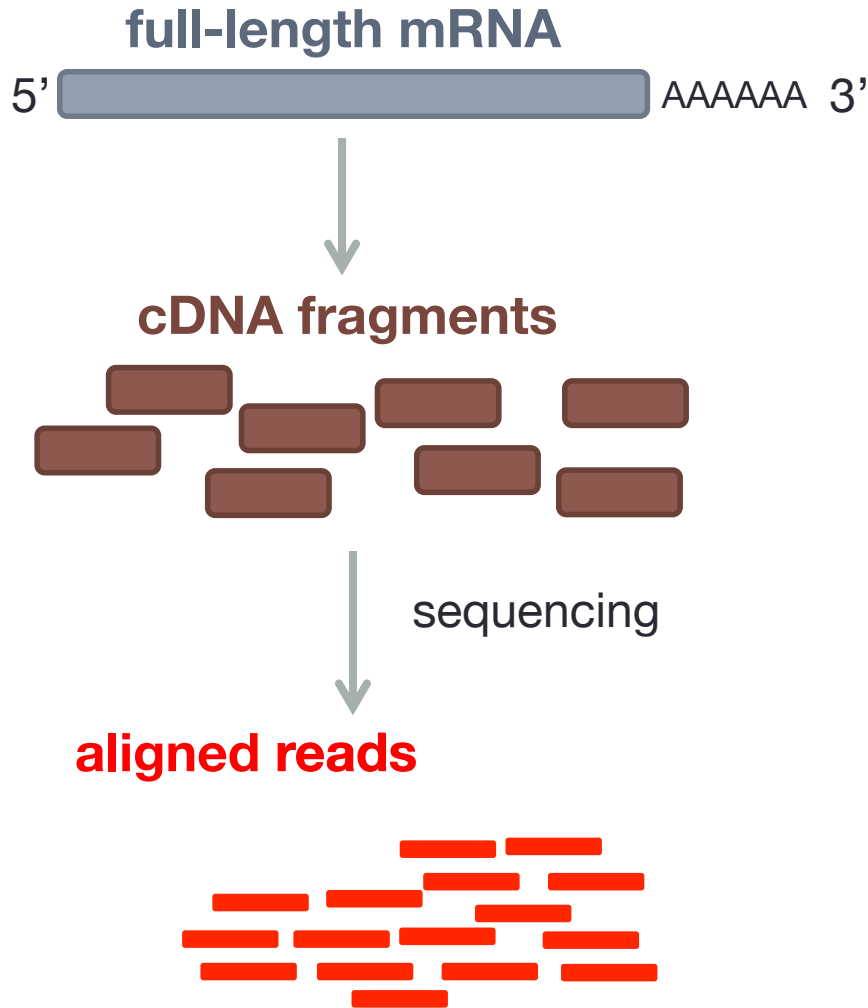
# ALIGNMENT

---

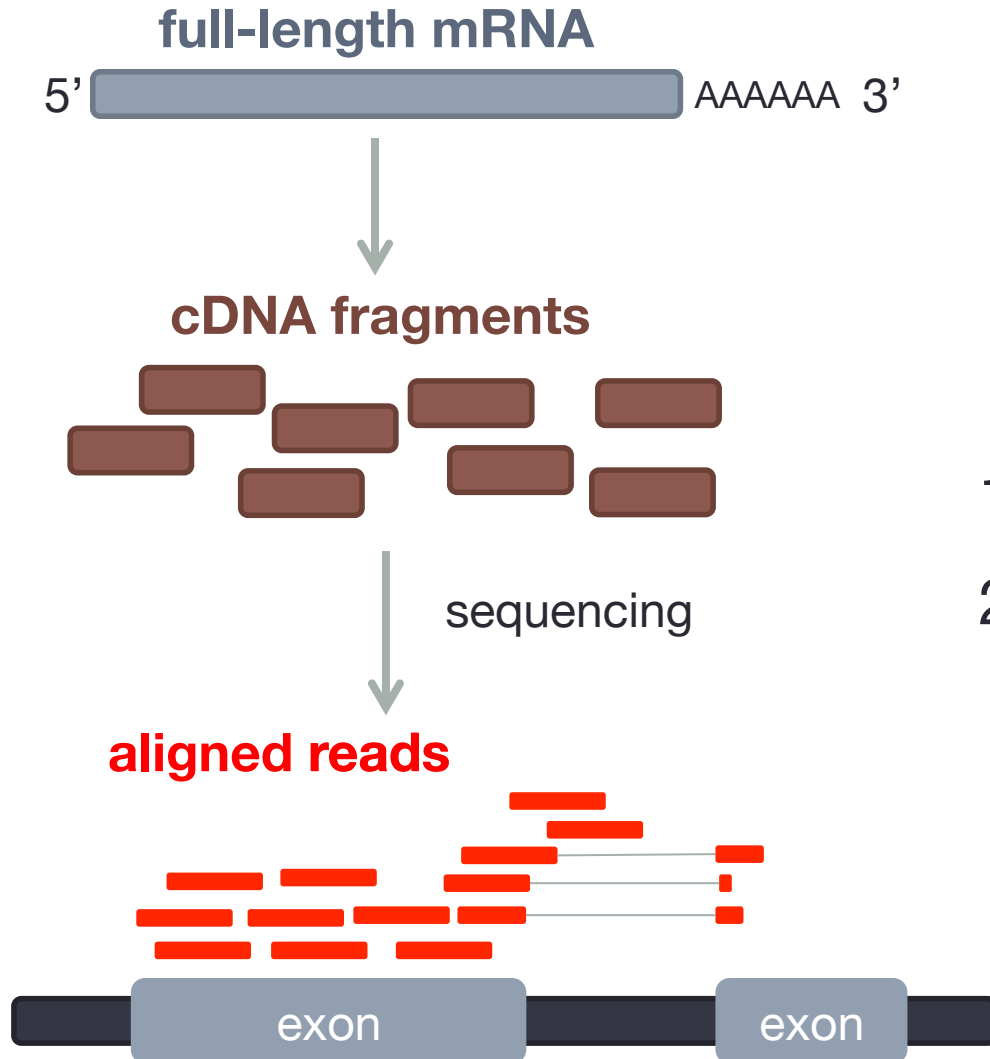
Finding out where the reads came from



# Aligning short RNA-seq reads



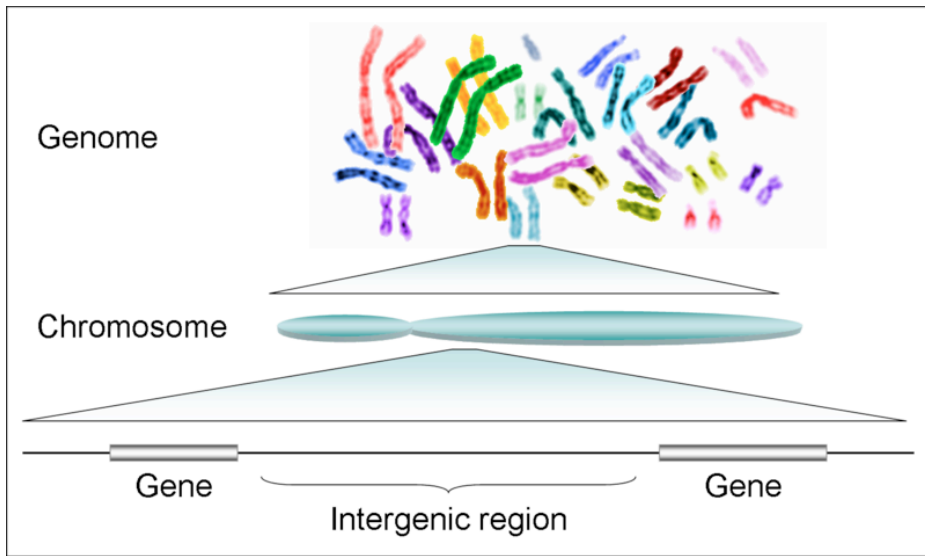
# Aligning short RNA-seq reads



Spliced alignment tools  
usually need:

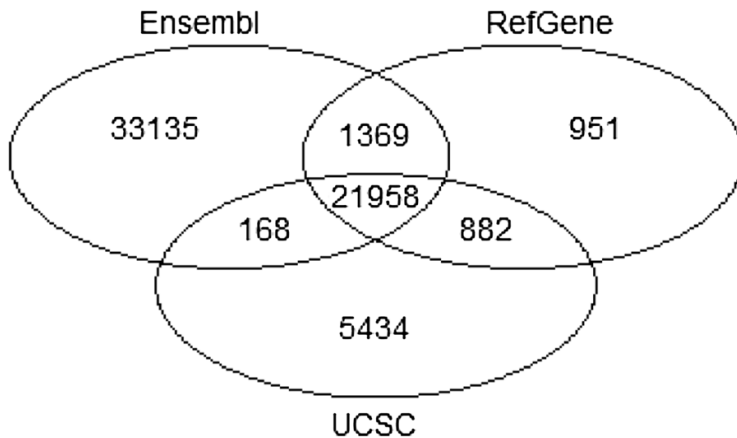
- 1) reference genome
- 2) annotation

# Annotation & their data bases



[https://commons.wikimedia.org/wiki/File:Human\\_genome\\_to\\_genes.png](https://commons.wikimedia.org/wiki/File:Human_genome_to_genes.png)

- Annotation is **dynamic!** (sequence, coordinates, types of elements)
- Automated vs. manual curation (“evidence-based”)



RefSeq [ncbi.nlm.nih.gov/refseq](http://ncbi.nlm.nih.gov/refseq)

UCSC Known Genes [genome.ucsc.edu](http://genome.ucsc.edu)

Ensembl/Gencode [gencodegenes.org](http://gencodegenes.org)

1/3 protein-coding genes  
> 17,000 non-coding RNAs  
> 15,000 pseudogenes



# Which annotation should one use?

*“More sensitive annotations, such as **Ensembl** (...) **should be preferred** over more specific annotations, such as RefSeq (...) if the aim is to obtain accurate expression estimates.”*

Janes et al. (Briefings in Bioinformatics, 2015). doi:  
10.1093/bib/bbv007

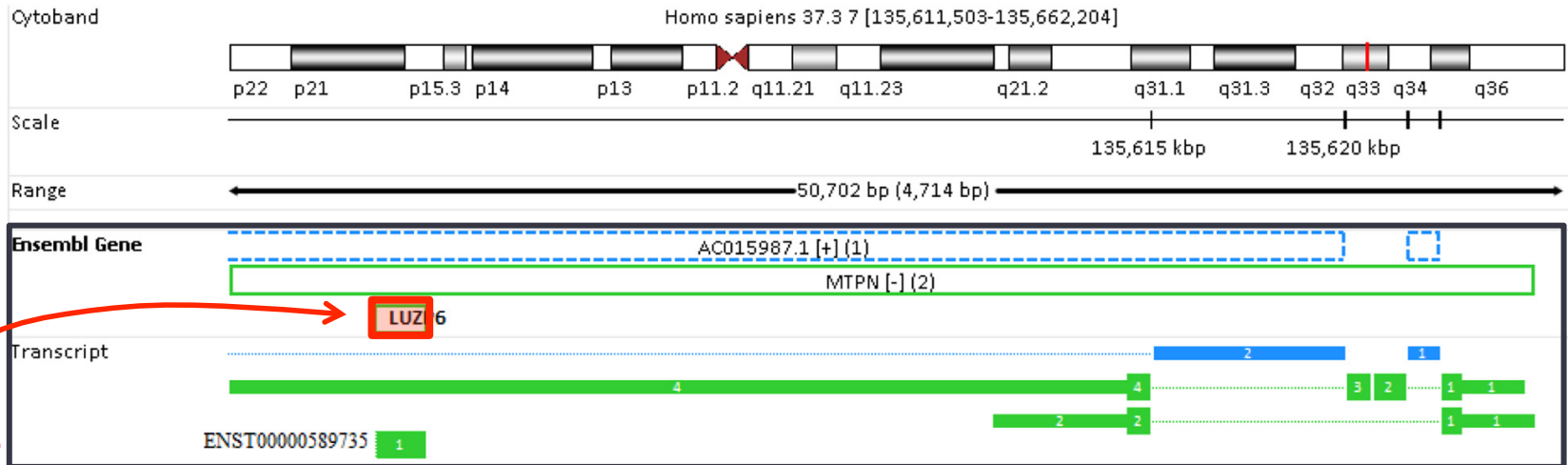
*“We observe that **RefSeq Genes** produces the **most accurate fold-change measures** with respect to a ground truth of RT-qPCR gene expression estimates. “*

Wu et al. (BMC Bioinfo, 2013). doi:  
10.1186/1471-2105-14-S11-S8

*“In practice, there is **no simple answer to this question**, and it depends on the purpose of the analysis. (...) When choosing an annotation database, researchers should keep in mind that **no database is perfect** and **some gene annotations might be inaccurate or entirely wrong.**”*

Zhao & Zhang (BMC Genomics, 2015). doi:10.1186/s12864-015-1308-8

# Gene models can differ dramatically



1

**ENSEMBL**

0 unambiguous reads for LUZP6;  
x number of reads for MTPN

**VS.**

2

**RefSeq**

0 reads for either  
or  
the same number of reads  
for both genes

pay attention to the **source** as well as to the version of the genome/annotation **build!**

# Storing annotation information

*see the course notes for details*

- representing genome coordinates + description/name
- various formats (all are plain text files): GFF2, GFF3, GTF, BED, SAF...

## GTF (“GFF2.5”)

1. reference coordinate
2. source
3. annotation type
4. start position
5. end position
6. score
7. strand
8. frame/phase
9. attributes: <TYPE VALUE>

```
1 # GFF-version 2
2 IV      curated exon      5506900 5506996 . + .      Transcript B0273.1
3 IV      curated exon      5506026 5506382 . + .      Transcript B0273.1
4 IV      curated exon      5506558 5506660 . + .      Transcript B0273.1
5 IV      curated exon      5506738 5506852 . + .      Transcript B0273.1
6
7 # GFF-version 3
8 ctg123  .  exon  1300  1500  .  +  .  ID=exon00001
9 ctg123  .  exon  1050  1500  .  +  .  ID=exon00002
10 ctg123  .  exon  3000  3902  .  +  .  ID=exon00003
11 ctg123  .  exon  5000  5500  .  +  .  ID=exon00004
12 ctg123  .  exon  7000  9000  .  +  .  ID=exon00005
```

```
# example for the 9th field of a GTF file
gene_id "Em:U62.C22.6"; transcript_id "Em:U62.C22.6.mRNA"; exon_number 1
```

# 0 vs. 1 based conventions

## one-based, fully-closed



ATG location: 7 - 9 or [7,9]  
Cut site: 11^12 or (11,12)  
Interval length = stop - start + 1

GFF format

## zero-based, half open

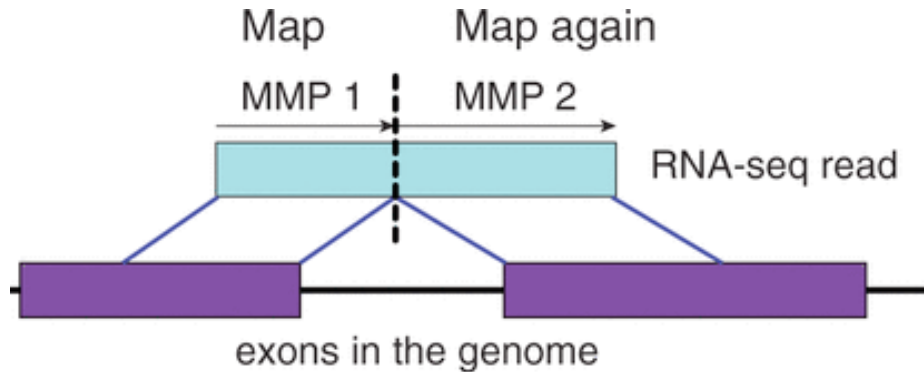


ATG location: 6 - 9 or [6,9)  
Cut site: 11-11 or [11,11)  
Interval length = stop - start

BED format

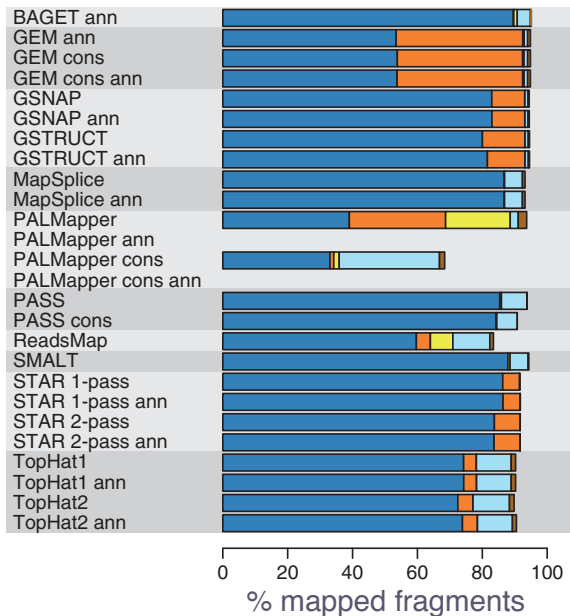
[http://  
alternateallele.blogspot.com/  
2012/03/genome-coordinate-cheat-  
sheet.html](http://alternateallele.blogspot.com/2012/03/genome-coordinate-cheat-sheet.html)

# Spliced Transcriptome Alignment to Reference (STAR)



- accurate & sensitive
- very fast
- memory intensive!  
(use it on the server!)

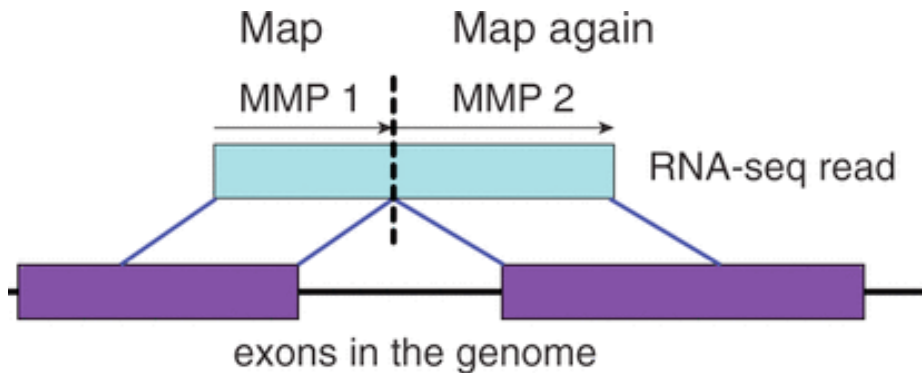
## Spliced alignment programs



Engström et al. (2013) Nature Methods, 10(12), 1185–1191. doi:10.1038/nmeth.2722

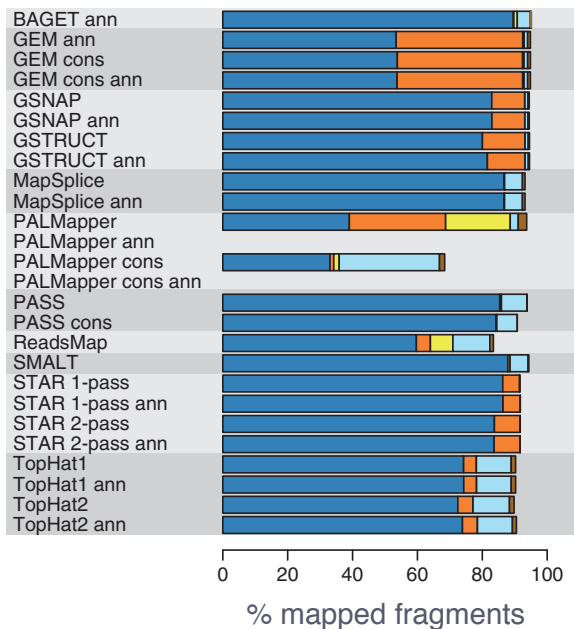
- MMP = maximal mappable prefix  
(aka maximum matching portion)
- reads are split when a continuous alignment is not possible
- the remaining unmappable portion is then aligned again
- finally, aligned portions of the original full-length reads are stitched together

# STAR spliced alignment



- accurate & sensitive
- very fast
- memory intensive!

## Spliced alignment programs



Engström et al. (2013) Nature Methods, 10(12), 1185–1191. doi:10.1038/nmeth.2722

STAR has numerous options! tune them to meet your needs

Current Protocols in Bioinformatics  
(Sept 2015)

DOI: 10.1002/0471250953.bi1114s51

and  
STARmanual.pdf

# 2 main STAR modules

## 1. generate **genome index**

```
--runMode genomeGenerate  
--genomeFastaFiles sacCer3.fa  
--sjdbGTFfile sacCer3.gtf
```

needs to be done just  
1x per transcriptome!

## 2. align

2.1. align to *reference* & identify  
novel splice junctions

```
$runSTAR --genomeDir STARindex/ \  
--readFilesIn $FASTQ_FILES \  
--readFilesCommand zcat \  
--twopassMode
```

2.2 *re-run* alignment including  
the novel splice junctions

--twopassMode

must be done for  
every sample

*Let's align the reads for WT\_1!*

# Storing aligned reads: SAM/BAM

@HD VN:

@SQ SN: LN:

@RG ID: SM:

@PG ID:

@CO

(theoretically) optional  
HEADER SECTION  
general information about the file

1	2	3	4	5	6	7	8	9	10	11	>11
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT

Paired read?  
Unmapped?  
Mapped to rev.  
strand?  
1<sup>st</sup> in pair?  
2<sup>nd</sup> in pair?  
Failed QC?  
...

M (mis)match  
I insertion  
D deletion  
N skipped  
S soft clipped  
H hard clipped  
P padding

<TAG>:<TYPE>:<VALUE>  
AS A  
BC i  
NH f  
NM z  
... H

ALIGNMENT  
SECTION  
1 line per locus

QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT



# Storing aligned reads: SAM/BAM

1	2	3	4	5	6	7	8	9	10	11	>11
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT

2<sup>nd</sup> field: binary FLAG

Binary (Decimal)	Hex	Description
000000000001 (1)	0x1	Is the read paired?
000000000010 (2)	0x2	Are both reads in a pair mapped “properly” (i.e., in the correct orientation with respect to one another)?
000000000100 (4)	0x4	Is the read itself unmapped?
000000001000 (8)	0x8	Is the mate read unmapped?
000000010000 (16)	0x10	Has the read been mapped to the reverse strand?
000000100000 (32)	0x20	Has the mate read been mapped to the reverse strand?
000001000000 (64)	0x40	Is the read the first read in a pair?
000010000000 (128)	0x80	Is the read the second read in a pair?
000100000000 (256)	0x100	Is the alignment not primary? (A read with split matches may have multiple primary alignment records.)
010000000000 (512)	0x200	Does the read fail platform/vendor quality checks?
100000000000 (1024)	0x400	Is the read a PCR or optical duplicate?

most common FLAGS for SR: 0; 4; 16

<https://broadinstitute.github.io/picard/explain-flags.html>

# Storing aligned reads: SAM/BAM

1	2	3	4	5	6	7	8	9	10	11	>11
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT

6<sup>th</sup> field: CIGAR string – which hoops did the aligner have to jump through to align the read to the genome locus that it thought was the best fit?

**M** alignment (match or **mismatch!!**)  
**I (N)** insertion (large insertions)  
**D** deletion  
**S/H** clipping

spliced out introns = sequences are missing in the read, i.e., they need to be inserted in order to align the read to the genome

Reference sequence with aligned reads	CIGAR string	Explanation
C T G C A T G T T A G A T A A * * G A T A G C T G T G C T A		
A A G G A T A * C T G	<b>1M2I4M1D3M</b>	Insertion & Deletion
G A T A A * G G A T A	<b>5M1P1I4M</b>	Padding & Insertion
T G T T A [blue bar] T G C T A	<b>5M15N5M</b>	Spliced read
a a a C A T G T T A G	<b>3S8M</b>	Soft clipping
A A A C A T G T T A G	<b>3H8M</b>	Hard clipping

# Storing aligned reads: SAM/BAM

1	2	3	4	5	6	7	8	9	10	11	>11
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT

after 11<sup>th</sup> field: OPTIONAL information

AS:i Alignment score  
BC:Z Barcode sequence  
HI:i Query is *i*-th hit stored in the file  
NH:i Number of reported alignments for the query sequence  
NM:i Edit distance of the query to the reference  
MD:Z String that contains the exact positions of mismatches (should complement the CIGAR string)  
RG:Z Read group (should match the entry after ID if @RG is present in the header)

**<TAG>:<TYPE>:<VALUE>**  
tags are not standardized!

NH HI NM MD have standard meaning as defined in the SAM format specifications.

AS is the local alignment score (paired for paired-end reads).

nM is the number of mismatches per (paired) alignment, not to be confused with NM, which is the number of mismatches in each mate.

jM:B:c,M1,M2,... intron motifs for all junctions (i.e. N in CIGAR): 0: non-canonical; 1: GT/AG, 2: CT/AC, 3: GC/AG, 4: CT/GC, 5: AT/AC, 6: GT/AT. If splice junctions database is used, and a junction is annotated, 20 is added to its motif value.

jI:B:I,Start1,End1,Start2,End2,... Start and End of introns for all junctions (1-based).

jM jI attributes require samtools 0.1.18 or later, and were reported to be incompatible with some downstream tools such as Cufflinks.

# Basic QC of aligned reads

How many reads were aligned? What were reasons for lack of alignment?

Do you have enough paired mates (for PE sequencing)?

- aligner output (e.g., Log.final.out)
- `samtools flagstat`
- RSeQC's `bam_stat`
- QoRTs

visual  
inspection!

(almost) all of these can be summarized  
using MultiQC!

→ Section 3.4.1 of the course notes

# Integrative Genomics Viewer

`http://software.broadinstitute.org/software/igv/download`

## Integrative Genomics Viewer (IGV) (Version 2.3)

### Install IGV

Options for installing and running IGV:

1. (Mac only) Download and run the Mac application; or
2. (Windows) Download and run the self-extracting archive; or
3. (All systems) Use the Java Web Start buttons (Mac users: see below for limitations); or
4. (All systems) Download the binary distribution and run IGV from the command line.

**Note:** IGV 2.3.x requires Java 7. Users with Java 6 (JRE 1.6) should first try to upgrade Java to the latest version. If this is not possible you will need to run a 2.2.x version available in the [archive](#).

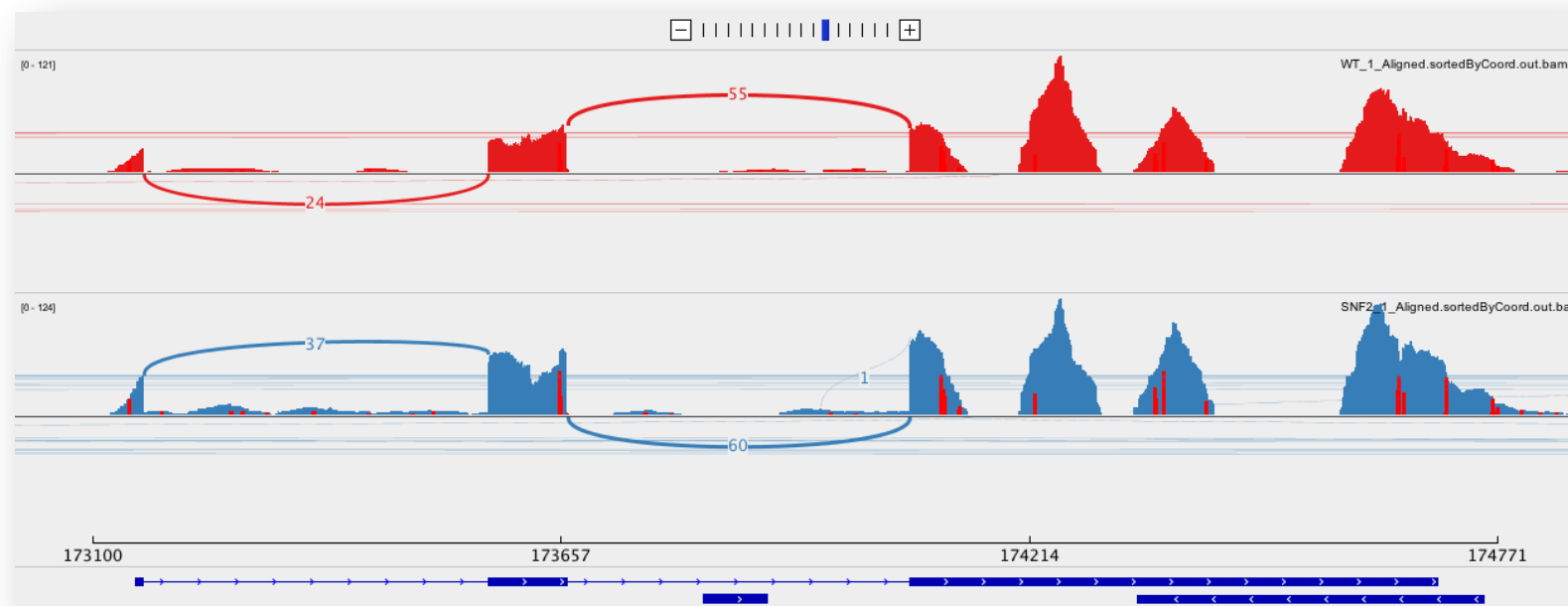
### Mac

Download and unzip the Mac App archive, then double-click the IGV application to run it. The application can be moved to the "Applications" folder, or anywhere else.

Download  
Mac App

# Integrative Genomics Viewer

- load BAM file(s) from URL (“File” -> “Open URL...”):  
<http://chagall.med.cornell.edu/RNASEQcourse/>  
<http://www.trii.org/courses/rnaseq.html>
- take a snapshot of the reads around gene *YPL198W*

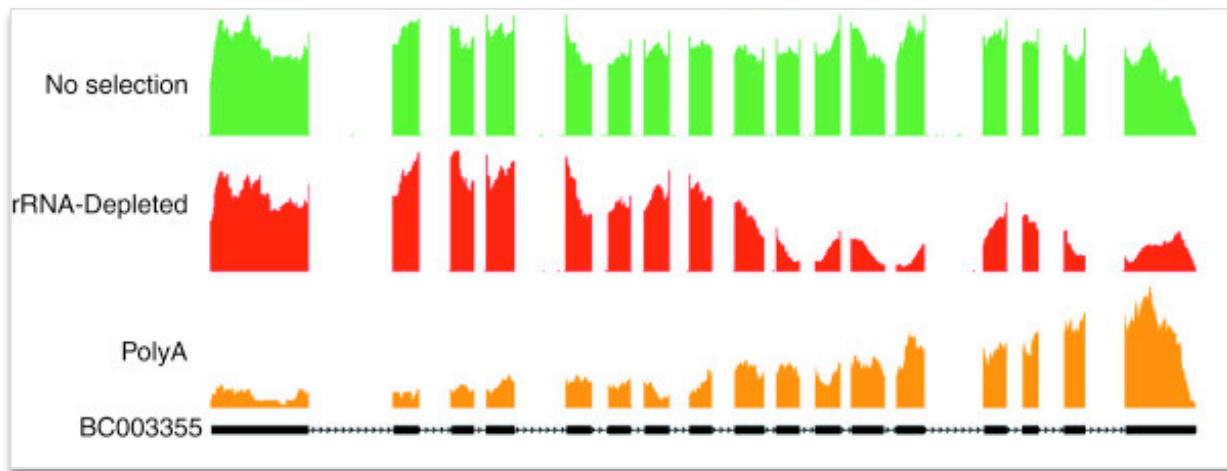
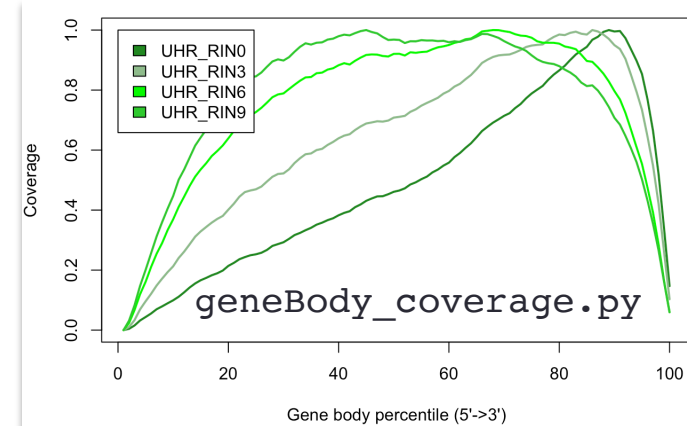


starting with  
IGV 2.3,  
**Sashimi** plots  
can easily be  
created

<http://software.broadinstitute.org/software/igv/Sashimi>

# Identifying RNA-seq-specific biases

- lack of **gene diversity**: dominance of rRNAs, tRNAs (and other highly abundant transcripts)
- **read distribution** (`read_distribution.py`)
  - high **intron** coverage ~ incomplete poly(A) enrichment step or many immature transcripts
  - many **intergenic** reads: possibly gDNA contamination
- **gene coverage: 3' bias** → RNA degradation, poly(A) enrichment protocol



popular QC tools:  
RSeQC, QoRTs

(see course notes)

# removing rRNAs

Can be done at virtually every step of the analysis. Choose the version that makes most sense to you.

- **sortMeRNA:** <http://bioinfo.lifl.fr/RNA/sortmerna/>
  - input: reads in fastq file + rRNA sequences
  - will extract those reads that do not match to the rRNA sequences
  - <https://www.ncbi.nlm.nih.gov/nuccore/U13369> (human rRNA),  
<https://www.ncbi.nlm.nih.gov/nuccore/BK000964> (mouse)
- make a “**genome**” **index for rRNAs only** (and perhaps tRNAs), then align your reads and only use those that do not map for the next round of alignment
- do your alignment and counting as is, simply **ignore the rRNA genes** in your subsequent downstream analysis



# Summary Day 2

- aligning unspliced reads is not too hard, but it takes a long time
- spliced reads are quite tricky, and identifying novel splice junctions is error-prone and far from being solved
- the file format for storing aligned reads (SAM/BAM) is fairly standardized, but the optional fields (and how alignment tools interpret some of the mandatory entries) leave lots of room for variability
- the file format(s) for storing genome annotation (e.g. genes, transcripts) tend to be even stricter defined and even less well followed (aka it's a mess!)
- historically, samtools are the most widely used tools when it comes to exploring and manipulating SAM/BAM files (although there are alternatives, e.g. bamtools)
- **QC of aligned read files is at least as important as QC of the raw reads!**