

RNA-Seq: Differential Expression of Genes

Álvaro Cortés Calabuig

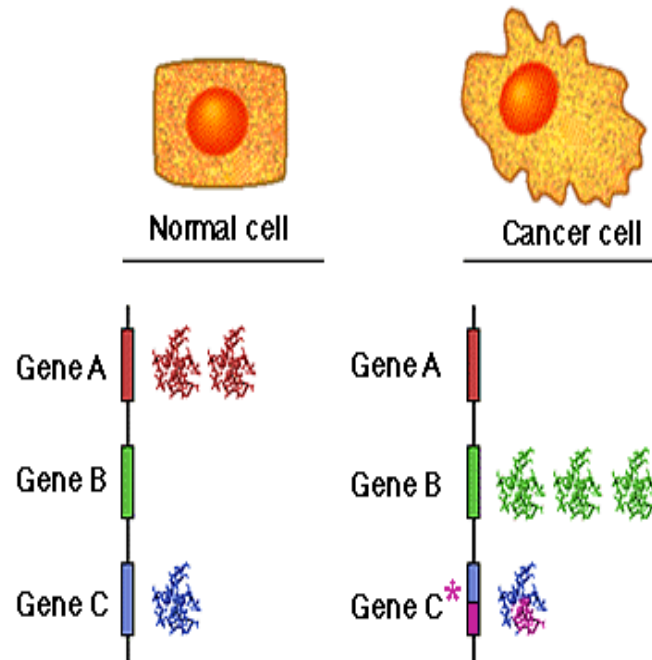
October 2019

Overview

- Differential Expression Principles
- Splice-aware Alignment and Counting
- Differential Expression Table

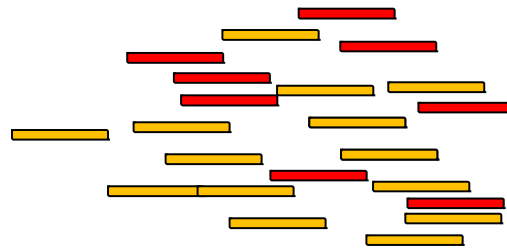
RNA-Seq: Differential Gene Expression

Discover changes in RNA abundance between conditions



RNA-Seq

Sequencer produces millions of **reads** and a qualification for each base call



What reads have to do with differential expression of genes?

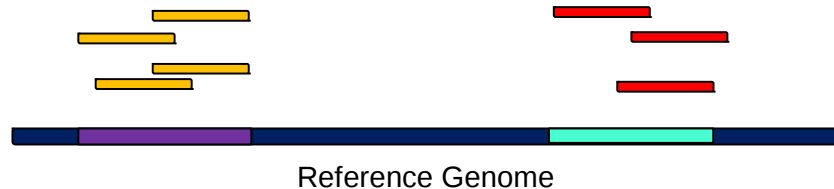
Read count is linearly related to the abundance of the target transcript

Differential Expression of Genes

Count the number of reads that fall into annotated genes

Align: Align reads w.r.t. a reference genome (transcriptome)

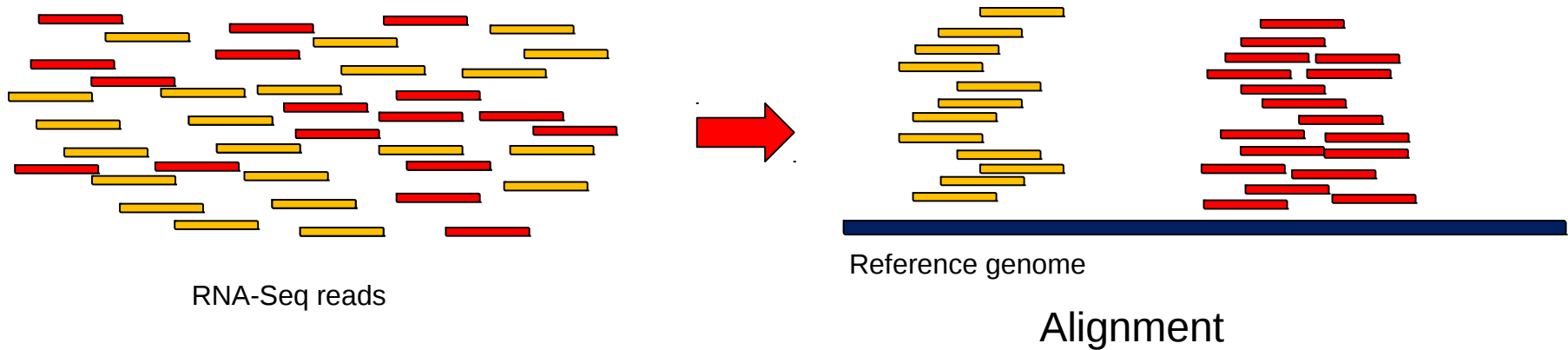
Count: number of aligned reads per feature (genes)



Perform **statistical analysis** on the counts to discover quantitative changes in expression levels between experimental groups

- Normalization of counts
- Probabilistic modeling of read counts
- Estimate differential expression

From Raw Reads to Alignment



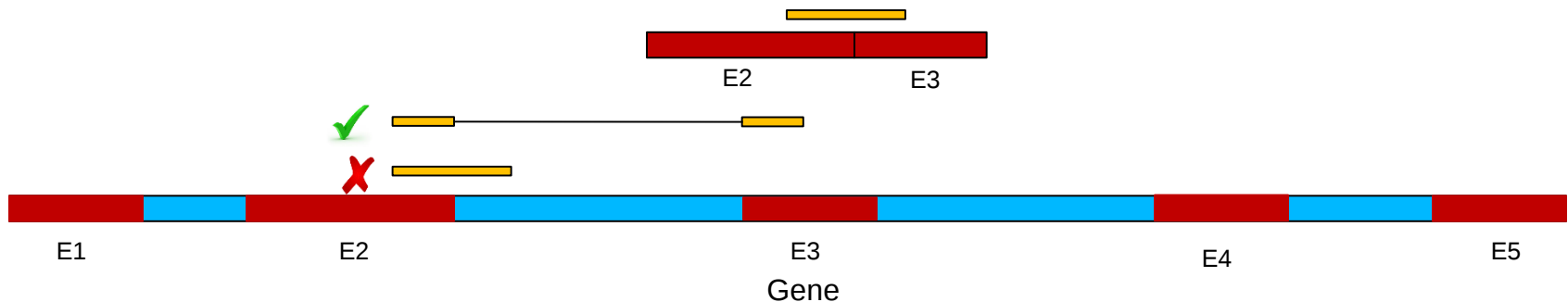
Splice-Aware Alignment

Individual reads are aligned to a reference genome

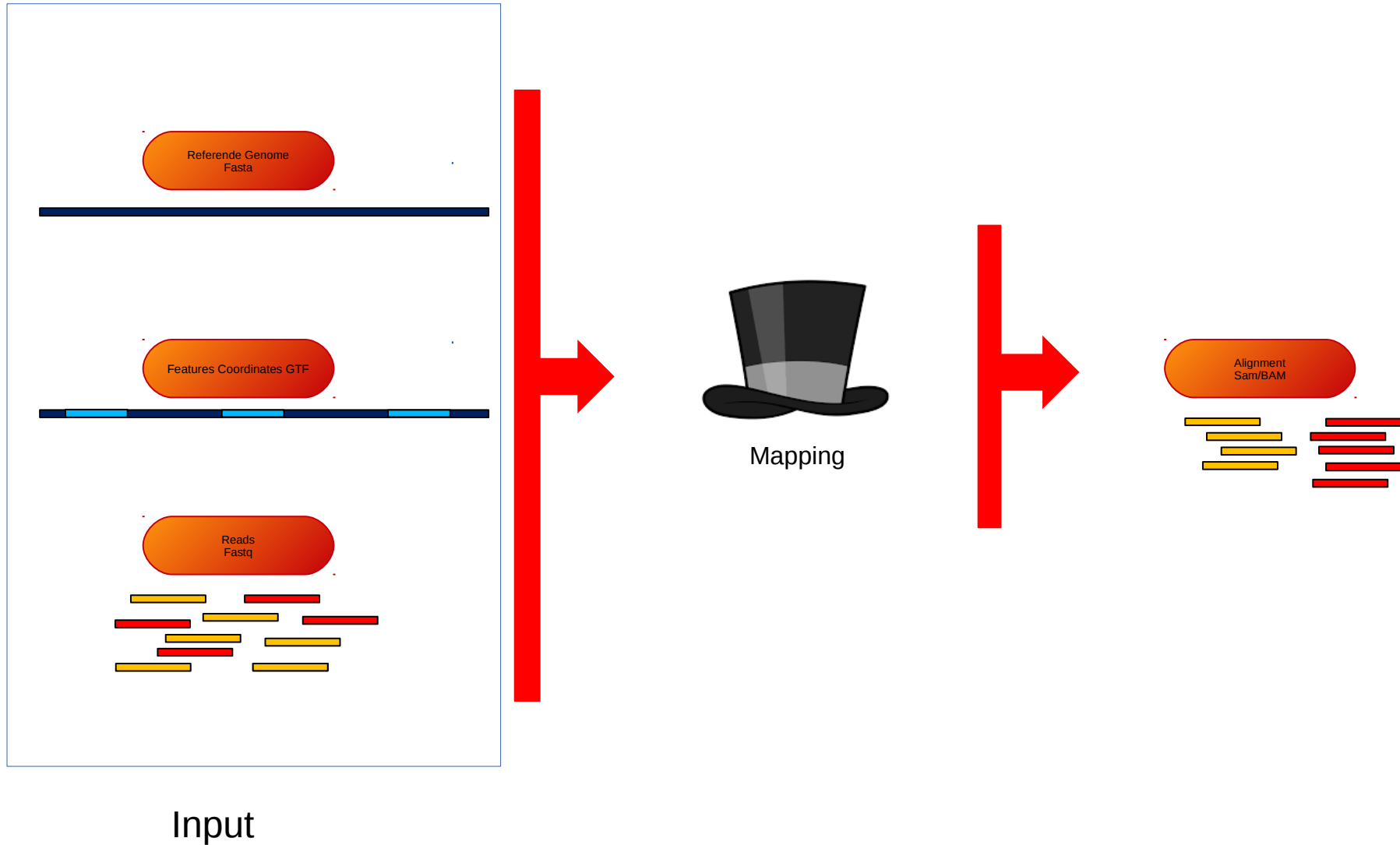
RNA-seq read alignment differs from standard alignment:

If an RNA-Seq read spans an exon boundary, part of the map **will not** map contiguously to the reference

This causes the standard mapping procedure to fail



Reference-Based Alignment



Reference Sequence in Fasta Format

- genome.fa **human-readable** nucleotide sequence
- Mouse genome: 2.6GB
- And it looks like this:

```
AATAAGTCAATGGCCTTTCTCTACACAAAGAATAAACAGGCTGAGAAAGAAATTAGGGAA
ACAACACCCTTCTCAATAGTCACAAATAATATAACATATCTCGGCGTGACTCTAACTAAG
GAAGTGAAAGATCTGTATGATAAAAACTTCAAGTCTCTGAAGAAAGAAATTAAAGAAGAT
CTCAGAAGATGGAAAGATCTCCCATGCTCATGGATTGGCAGGATCAATATTGTAAAAATG
GCTATCTTGCCAAAAGCAATCTACAGATTCAATGCAATCCCCATCAAAATTCCAACTCAA
TTCTTCAACGAATTAGAAGGAGCAATTTGCAAAATTCATCTGTAAATAACAAAAACCTAGG
ATAGCAAAAAGTCTTCTCAAGGATAAAAGAACCTCTGGTGGAATCACCATGCCTGACCTA
AAGCTTTACTACAGAGCAATTGTGGTAAAACTGCATGGTACTGGTATAGAGACAGACAA
GTAGACCAATGGAATAGAATTGAAGACCCAGAAATGAACCCACACACCTATGGTCACTTG
ATCTTCGACAAGGGAGCTAAAACCATCCAGTGGAAGAAAGACAGCATTTTCAACAAATGG
TGCTGGCACAACCTGGTTGTTATCATGTAGAAGAATGCGAATCGATCCATACCTTATCTCCT
TGTAATAAGGTCAAATCTAAATGGATCAAAGAAGTTCACATAAAACCAGAGACACTGAAA
CTTATAGAGGAGAAAGTGGGGAAAAGCCTTGAAGATATGGGCACAGGGGAAAAATTCCTG
AACAGAACAGCAATGGCTTGCTGTAAGATTGAGAATTGACAAATGGGACCTAATGAAA
CTCCAAAGTTTCTGCAAGGCAAAAGACACCGTCAATAAGAGAAAGAGACCACCAACAGAT
TGGGAAAGGATCTTTACCTATCCTAAATCAGATAGGGGACTAATATCCAACATATATAAA
GAACTCAAGAAGGTGGACTTCAGAAAATCAAACAACCCATTAAAAAATGGGGCTCAGAA
CTGAACAAAGAATTCTCACCTGAGTTATACCGAATGGCAGAGAAGCACCTGAAAAATGC
TCAACATCCTTAATCATCAGGGAATGCAAAATCAAAACAACCCCTGAGATTCCACCTCACA
CCAGTCAGAATGTCTAAGATCAAAATTCAGGTGACAGCAGATGCTGGCGAGGATGTGGA
GAAAGAAGAACACTCCTCCATTGTTGGTGGGATTGCAGGCTTGTACAACCACTCTGGAAA
TCCGTCTGGCGGTTCTCAGAAAATTTGGACATAGTACTACCGGAGGATCCAGCAATACCT
CTCCTGGGCATATATCCAGAAGATGCCCCAACTGGTAAGAAGGACACATGCTCCTATG
TTCATAGCAGCCTTATTTATAATAGCCAGAAGCTGGAAGAAGCCAGATGCCCTCAACA
GAGGAATGGATACAGAAAATGTGGTACATCTACACAATGGAGTACTACTCAGCTATTAAA
AAGAATGAATTTATGAAATTCAGCCAAATGGATGGACCTGGAGGCATCATCCTGAGT
```

GTF Files: Gene Transfer Format

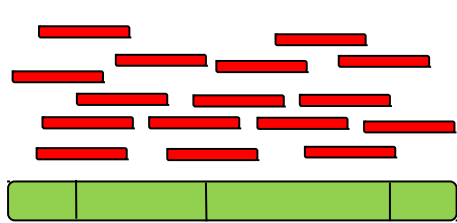


```

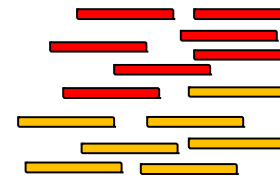
9   havana gene 44107301 44108960 . - . gene_id "ENSMUSG000000092258"; gene_version "1"; gene_name "Gm20444"; gene_source "havana"; gene_biotype "processed_transcript";
9   ensembl gene 44112923 44121947 . + . gene_id "ENSMUSG000000097467"; gene_version "1"; gene_name "Gm26737"; gene_source "ensembl"; gene_biotype "lincRNA";
9   ensembl gene 44123768 44134485 . - . gene_id "ENSMUSG000000097617"; gene_version "1"; gene_name "Gm10687"; gene_source "ensembl"; gene_biotype "lincRNA";
  
```

column-number	content	values/format
1	chromosome name	chr{1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17,18,19,20,21,22,X,Y,M}
2	annotation source	{ENSEMBL,HAVANA}
3	feature-type	{gene,transcript,exon,CDS,UTR,start_codon,stop_codon,Selenocysteine}
4	genomic start location	integer-value (1-based)
5	genomic end location	integer-value
6	score (not used)	.
7	genomic strand	{+,-}
8	genomic phase (for CDS features)	{0,1,2,.}

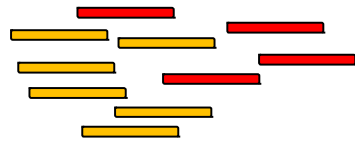
Alignment Strategy – GTF



Transcriptome



Unmapped reads



Unmapped reads

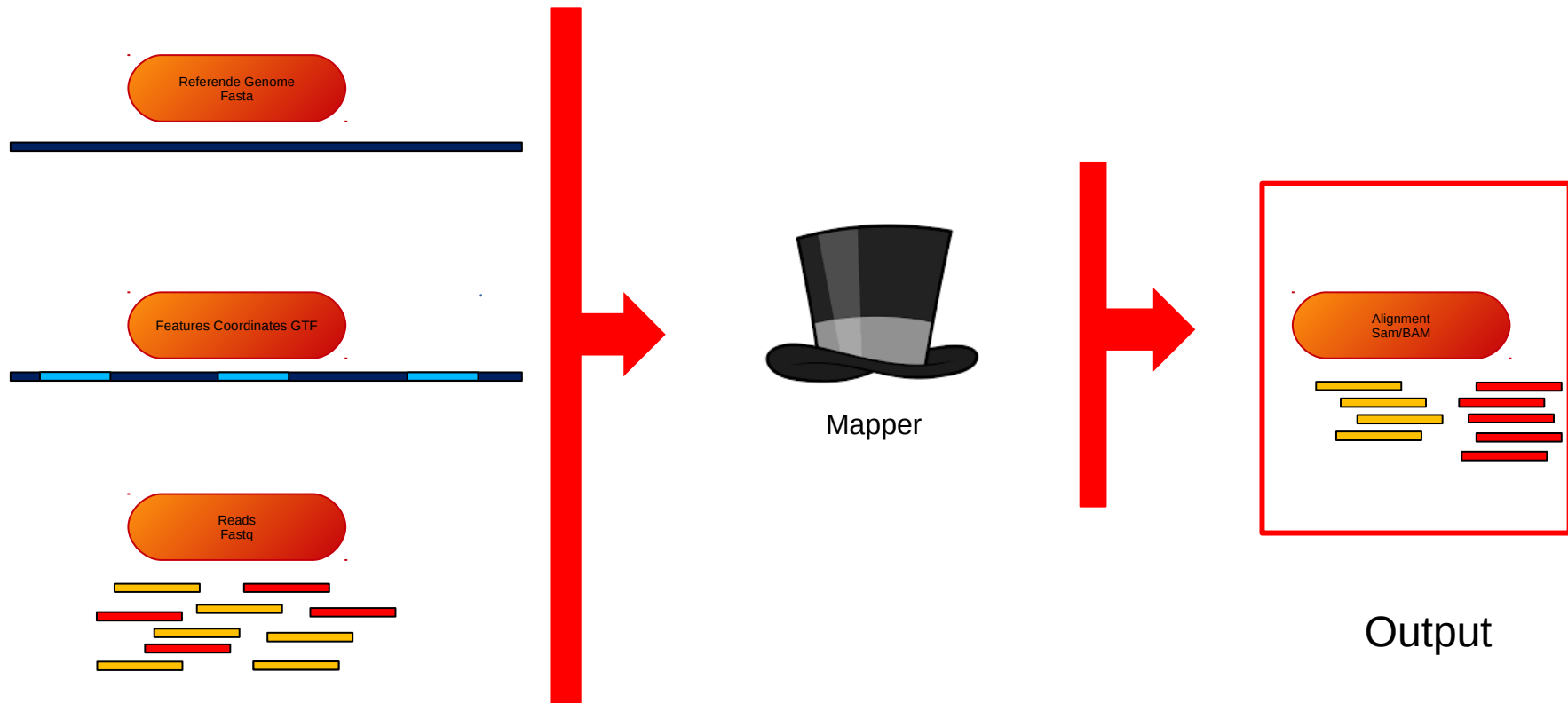


Reference Genome

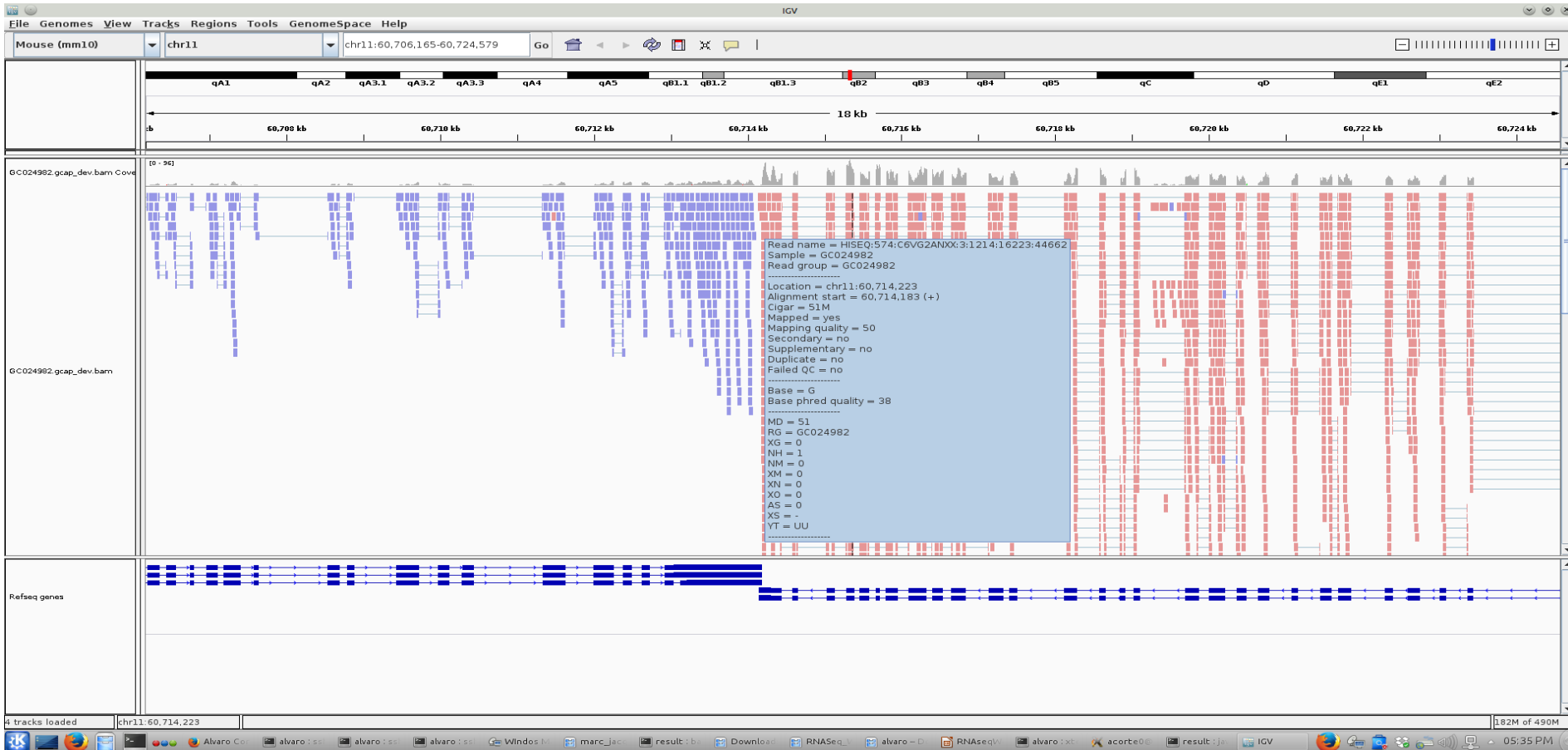


Read are split-aligned

Reference-Based Alignment



RNA-Seq BAM/SAM



Counting Reads in HTSeq



Method:

- Count each read only **once**
- Multi-mapping reads and reads overlapping multiple features **discarded**

	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

Simon Anders

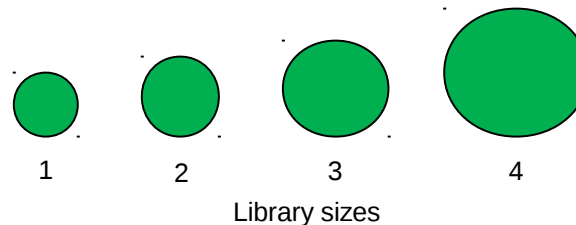
Counting Reads in HTSeq(*)

(*) Simon Anders, Wolfgang Huber (EMBL)

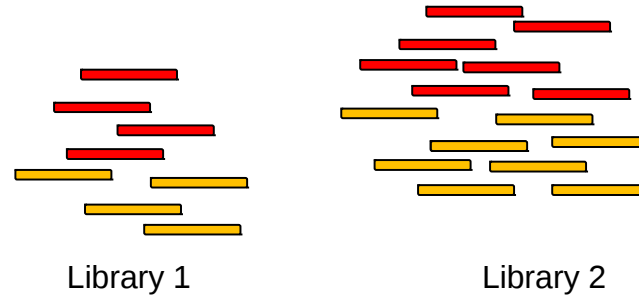
Resulting table of **absolute counts** for each sample and feature:

GeneID	Sample 1	Sample 2	Sample 3	Sample 4
FBgn0000003	0	0	0	1
FBgn0000008	76	70	88	70
FBgn0000014	0	0	0	0
FBgn0000015	1	2	0	0
FBgn0000017	3564	3150	3072	3334
FBgn0000018	245	310	299	308

But!



Normalization of Counts: Impact Library Size



A bigger library produces more reads...But in slide 3:

“Read Counts is linearly related to the abundance of the target transcript”

Normalization of Counts: Size Factors

Adjust for library sizes to produce count values in a **common scale**:

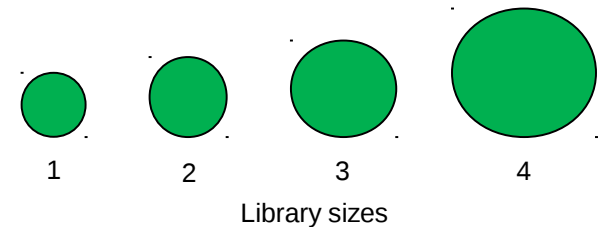
Sample 1	Sample 2	Sample 3	Sample 4
0.873	1.011	1.022	1.115

Size factors

	Sample 1	Sample 2	Sample 3	Sample 4
FBgn0000003	0	0	0	1
FBgn0000008	76	70	88	70
FBgn0000014	0	0	0	0
FBgn0000015	1	2	0	0
FBgn0000017	3564	3150	3072	3334
FBgn0000018	245	310	299	308

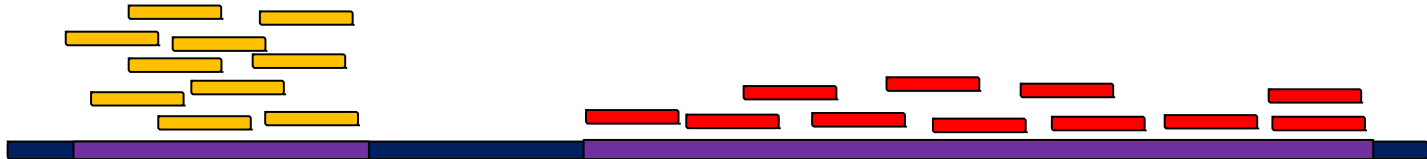
Normalized counts

	Sample 1	Sample 2	Sample 3	Sample 4
FBgn0000003	0.00	0.00	0.0	0.897
FBgn0000008	87.05	69.27	86.1	62.803
FBgn0000014	0.00	0.00	0.0	0.000
FBgn0000015	1.15	1.98	0.0	0.000
FBgn0000017	4082.02	3116.93	3004.5	2991.238
FBgn0000018	280.61	306.75	292.4	276.335



$$76/0.873 = 87.05$$

More on Counting and Normalization



- Possible extra normalization:
 - Longer transcripts are more likely to have sequences mapped to their genes
 - **Higher counts**, biasing comparisons between transcripts of different lengths.
- **RPKM** - Reads per kilo base per million mapped reads

“based on three real mRNA and one miRNA-seq datasets, we confirm previous observations that RPKM and TC, both of which are still widely in use [40,41], are ineffective and should be definitively abandoned in the context of differential analysis”

Briefings in Bioinformatics Advance Access published September 17, 2012
BRIEFINGS IN BIOINFORMATICS, page 1 of 11 doi:10.1093/bib/bbs046

A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis

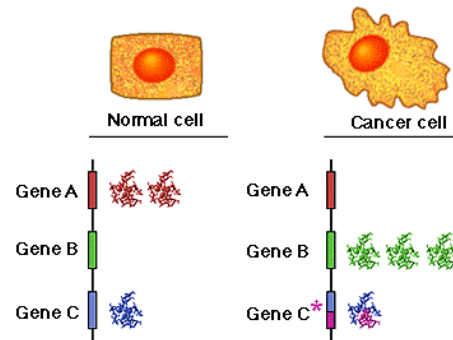
Mario-Agnès Dillies*, Andrea Rau*, Julie Aubert*, Christelle Hannequet-Antier*, Marine Jeanmougin*, Nicolas Servant*, Céline Keime*, Guillemette Marot, David Castel, Jordi Estelle, Gregory Guernec, Bernd Jagla, Luc Jouneau, Denis Labé, Caroline Le Gall, Brigitte Schaeffer, Stéphane Le Crom*, Mickaël Guedj*, Florence Jaffrézic* and on behalf of The French StatOmique Consortium

Submitted: 15th April 2012; Received (in revised form): 29th June 2012

Differential Expression

Statistical test: Decide whether for a given gene, an observed difference in reads counts is significant

Comparison between different **biological conditions**



- Read counts per gene are modeled by a **probability distribution**

Read Count Distribution Assumption

Method		Read Count Distribution Assumption
EdgeR		Negative binomial distribution
DESeq		Negative binomial distribution
Cuffdiff2 (CuffLinks)		Beta negative binomial distribution

Package Selection

DESeq or edgeR?

Box 2 | Differences between DESeq and edgeR

The two packages described in this protocol, DESeq and edgeR, have similar strategies to perform differential analysis for count data. However, they differ in a few important areas. First, their look and feel differs. For users of the widely used limma package⁶⁰ (for analysis of microarray data), the data structures and steps in edgeR follow analogously. The packages differ in their default normalization: edgeR uses the trimmed mean of M values⁵⁶, whereas DESeq uses a relative log expression approach by creating a virtual library that every sample is compared against; in practice, the normalization factors are often similar. Perhaps most crucially, the tools differ in the choices made to estimate the dispersion. edgeR moderates feature-level dispersion estimates toward a trended mean according to the dispersion-mean relationship. In contrast, DESeq takes the maximum of the individual dispersion estimates and the dispersion-mean trend. In practice, this means DESeq is less powerful, whereas edgeR is more sensitive to outliers. Recent comparison studies have highlighted that no single method dominates another across all settings^{27,61,62}.

PROTOCOL

Count-based differential expression analysis of RNA sequencing data using R and Bioconductor

Simon Anders¹, Davis J McCarthy^{2,3}, Yunshun Chen^{4,5}, Michal Okoniewski⁶, Gordon K Smyth^{4,7}, Wolfgang Huber¹ & Mark D Robinson^{8,9}

Differential Expression Table (DeSeq)

id	baseMean	baseMeanA	baseMeanB	foldChange	log2FoldChange	pval	padj
TTp9	61.2142079613	114.5039786047	7.9244373179	0.0692066548	-3.8529454185	3.18278695454205E-038	7.8910836963961E-034
Roums4	111.7253852962	3.3863424303	220.064428162	64.9858756734	6.0220542852	3.14004006066747E-016	3.89255066120643E-012
Serinc3	5049.0292624521	2783.4853979737	7314.5731269304	2.6278467752	1.3938811573	8.88017567331908E-011	7.33887318228667E-007
Apoba	687.9674505131	478.6234549012	897.3114461251	1.8747753311	0.9067177166	3.02197214770456E-005	0.1833853013
Psiga	318.5920273479	219.6643392529	417.519715443	1.9007168704	0.926543645	3.69832818348108E-005	0.1833853013
Igdkv12-98	3.4254621585	0.697231647	6.1536926701	8.8258940866	3.1417424346	7.3351076233012E-005	0.2589065108
Ssdlc12a5	17.8775194124	28.617735224	7.1373036008	0.249401413	-2.0034584562	8.04863286246797E-005	0.2589065108
Treb3l3	31.0676067843	13.9815013501	48.1537122185	3.444101675	1.7841277338	8.35418096347497E-005	0.2589065108
Tm16316	1.205992866	2.4119857321	0	0	-Inf	0.0001326348	0.3653795179
Ces1d	553.0344051691	424.629437072	681.4393732662	1.604785994	0.6823809198	0.0002487961	0.6168402322
Vsig1	134.705417607	168.1043611548	101.3064740591	0.6026403679	-0.7306307796	0.0005224394	1

Thanks!