# Analysis of RNA-seq data
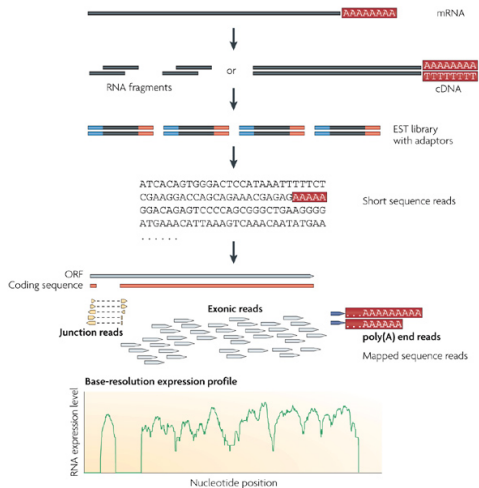## FOS 2017

Maarten van Iterson

Leiden University Medical Center
Department of Molecular Epidemiology

June 1, 2017

# A typical RNA-seq experiment



Wang *et al.* Nature Reviews Genetics 10, 57–63(2009)

# Experimental Design

- clear simple research question:
    - Good: Which genes are differentially expressed between disease and control, treated and untreated samples
    - Bad: Which genes are differentially expressed between treated and untreated samples with different dosages and different time points after treatment

- randomization
    - not all the controls on one day with one batch of chemicals and the cases on the other day with another batch of chemicals

- think ahead
    - sequencer (mostly Illumina nowadays)
    - aligner
    - annotation
    - statistical analysis

# Analysis of RNA-seq data

- Read mapping
    - Input: FASTQ file generated by the sequencer
    - Output: BAM file

- Summarization
    - Input: BAM file
    - Output: count table

- Quality control

- Normalization
    - Input: count table
    - Output: scale factors

- Differential Expression Analysis
    - Input: count table plus scale factors
    - Output: list of differentially expressed genes

# FASTQ file

file format for sequences plus quality scores
quality score indicates the probability that a given base is called incorrectly
by the sequencer

```
@SRR014849.1 EIXKN4201CFU84 length=93
GGGGGGGGGGGGGGGGGCTTTTTTTGTTTGGAACCGAAAGG
GTTTTGAATTTCAAACCCTTTTCGGTTTCCAACCTTCCAA
AGCAATGCCAATA
+SRR014849.1 EIXKN4201CFU84 length=93
3+&$#""""""""""""7F@71,";C?,B;?6B;:EA1EA
1EA59B:?:#9EA0D@2EA5:>5?:%A;A8A;?9B;D@
/=<?7=9<2A8==
```

@title and optional description
sequence line(s)
+optional repeat of title line

quality line(s)

---

Cock *et al.* Nucleic Acids Research,38(6), 1767–1771(2010)

# Read mapping



Align against genome or transcriptome

against transcriptome: easier, because no gapped alignment necessary

but: risk to miss possible alignments!

many tools available see e.g.,
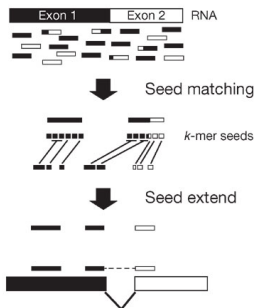http://massgenomics.org/short-read-aligners
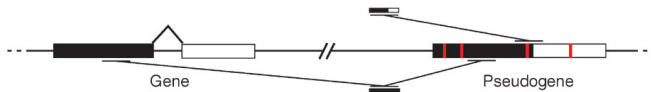
# Read mapping


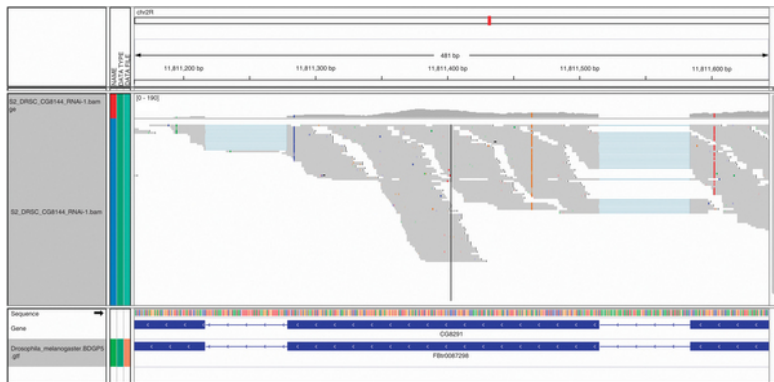
a Exon-first approach

b Seed-extend approach

c Potential limitations of exon-first approaches

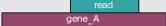Garber *et al.* Nature Methods 8, 469–477(2011)

# Summarizing mapped reads



Count each read at most once

Discard a read if

- it cannot be uniquely mapped
- its alignment overlaps with several genes
- the alignment quality score is bad

# Summarizing mapped reads: Counting rules



| | 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 |

http://www-huber.embl.de/users/anders/HTSeq/doc/count.html

# Quality Control: on 'raw' reads

- Basic information (total reads, sequence length, etc.)
- Per base sequence quality
- Overrepresented sequences (e.g., ribosomal RNAs)
- GC content
- Duplication level
- Etc

Tools: samtools, Fastqc, Fastx, Galaxy fastq tools, $\cdots$

# Quality Control: on aligned reads

- Percentage of reads properly mapped or uniquely mapped
- Among the mapped reads, the percentage of reads in exon, intron, and intergenic regions.
- 5' or 3' bias
- The percentage of expressed genes

Tools: RSeQC[1], RNA-SeQC, $\cdots$

---

[1]Wang *et al.* Bioinformatics, 28(16), 2184–2185(2012)

# Normalization

Each sample (library) will have different number of total reads
Total count, Counts per million and Reads Per Kilobase per Million
mapped reads (RPKM)
For differential expression these are not appropriate!!!

Toy example:

|  | sample A | sample B |
|---|---|---|
| gene 1 | 100 | 80 |
| gene 2 | 100 | 80 |
| . . . | . . . | . . . |
| gene 100 | 100 | 80 |
| gene 101 | 0 | 2.000 |
| Total Counts: | 10.000 | 10.000 |

# Normalization

Using Counts per million $\frac{X_{ij}}{X_{.j}}10^6$ with $X_{.j} = \sum_{i=1}^{101} X_{ij}$

|          | sample A | sample B |
|---------:|---------:|---------:|
| gene 1   | 10.000   | 8.000    |
| gene 2   | 10.000   | 8.000    |
| . . .    | . . .    | . . .    |
| gene 100 | 10.000   | 8.000    |
| gene 101 | 0        | 20.000   |

All genes are differentially expressed!
Is this really true?

# Normalization

Using TMM (trimmed mean of M-values)

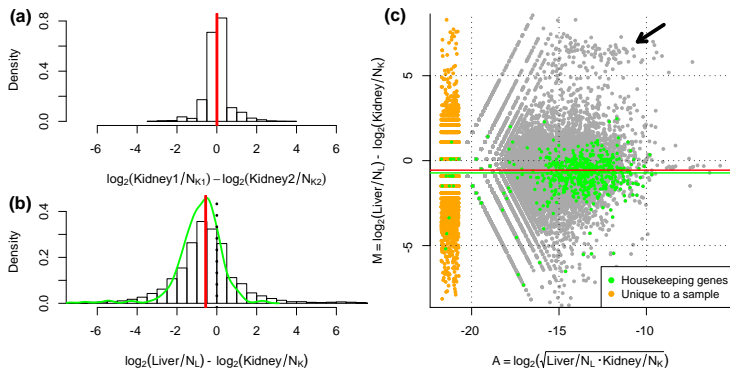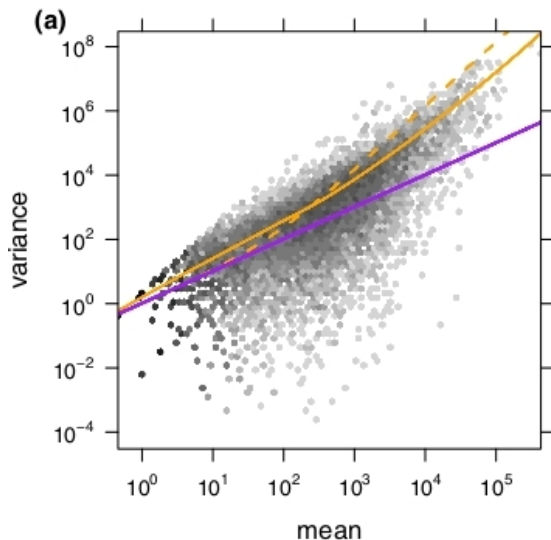|          | sample A | sample B |
|---------:|---------:|---------:|
| gene 1   | 10.000   | 10.000   |
| gene 2   | 10.000   | 10.000   |
| ...      | ...      | ...      |
| gene 100 | 10.000   | 10.000   |
| gene 101 | 0        | 250.000  |

One gene differentially expressed!
Seems more realistic

# TMM



Figure : (a) technical replicates and (b) liver versus kidney (c) An M versus A plot comparing liver and kidney.

Robinson *et al.* Genome Biology, 11(3), (2010)

# Testing for differential expression

- counts are discrete: $0, 1, 2, \cdots$
- large dynamic range $[0, > 100000]$
- Poisson or Negative Binomial distributed
- mean is approximately equal to the variance
- generalized linear model
- likelihood ratio test

Tools: edgeR, DESeq2, Cuffdiff, Myrna, $\cdots$

# Mean variance relationship



Anders *et al.* Genome Biology, 11(10), (2010)

From differential expression to Biology

- gene set enrichment (GO and KEGG)
- network construction (co-regulated genes)
- data integration (eQTL, meQTL, $\cdots$)
- $\cdots$

Other things you can do with RNAseq

- allele specific expression
- isoform (transcript) expression
- variant detection
- $\cdots$