# Characterizing transcriptomes using ngs data

T. Källman
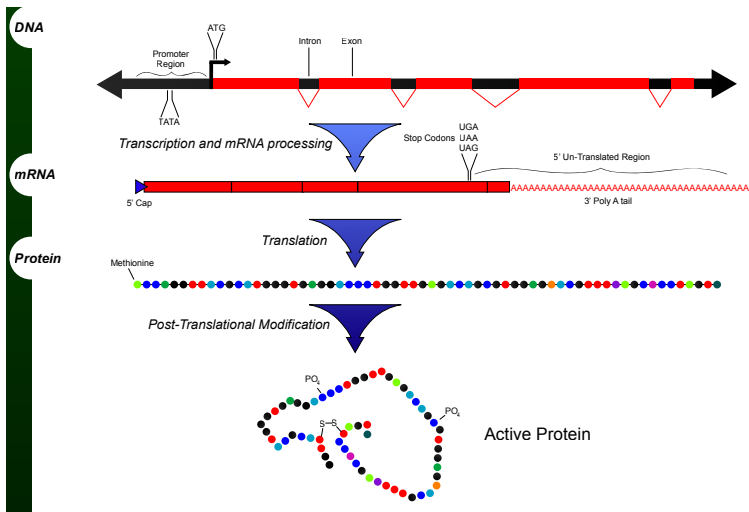
BILS/Scilife Lab/Uppsala University

Sep. 2015
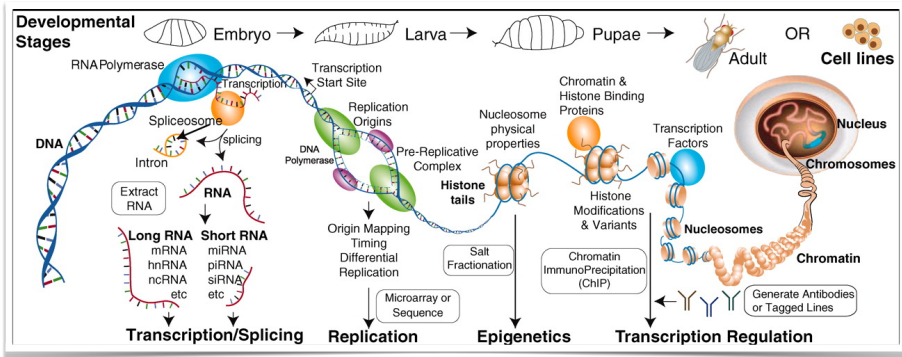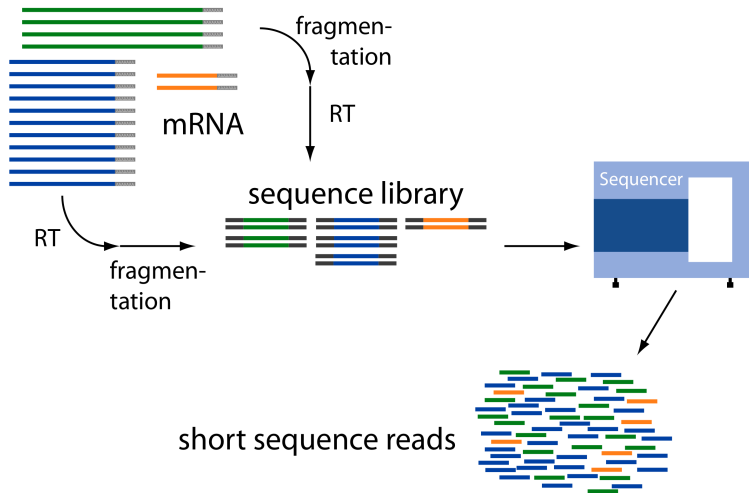
# Outline

# The Central Dogma

# A more complex view

# Transcriptomes vs genomes

- Dynamic, not the same over tissues and time points
- Smaller sequence space
- Less repetitive (but large gene families can be found)
- Fairly stable in size? (*eg.* 2-4 fold change among eukaryotes, whereas genome size can vary 1000-fold)
- Genes are often expressed in multiple different splice-variants
- RNA often from only one strand
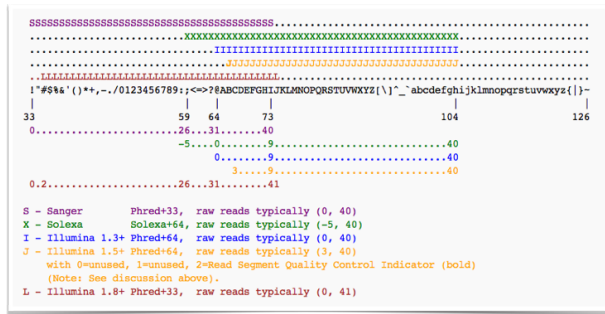
# NGS data

# Machine output

# Machine output

```
@SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
CAACGAGTTCACACCTTGGCCGACAGGCCCGGGTAA
+SRR038845.3 HWI-EAS038:6:1:0:1938 length=36
BA@7>B=>:>>7@7@>>9=BAA?;>52;>:9=8.=A
@SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
CCAATGATTTTTTTCCGTGTTTCAGAATACGGTTAA
+SRR038845.41 HWI-EAS038:6:1:0:1474 length=36
BCCBA@BB@BBBBBAB@B9B@=BABA@A:@693:@B=
@SRR038845.53 HWI-EAS038:6:1:1:360 length=36
GTTCAAAAAGAACTAAATTGTGTCAATAGAAAACTC
+SRR038845.53 HWI-EAS038:6:1:1:360 length=36
BBCBBBBBB@@BAB?BBBBCBC>BBBAA8>BBBAA@
```
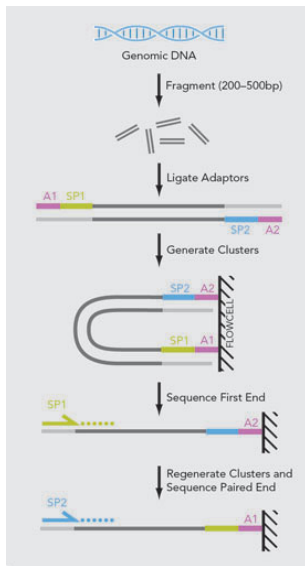
# Sequence quality

- Phred quality scores: Q = -10 x log P (High Q = high probability of the base being correct
- A Phred quality score of 20 to a base, means that the base is called incorrectly in 1 out of 100 times.

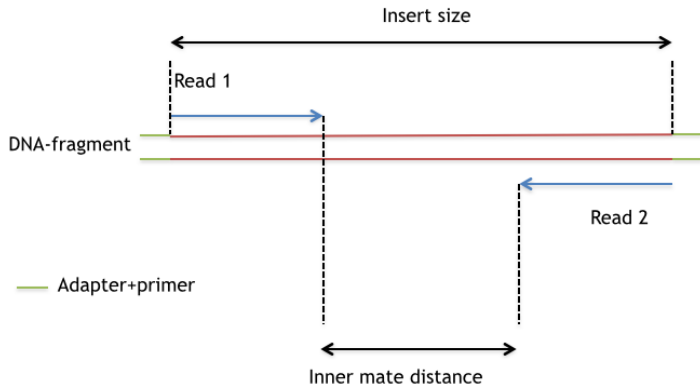# Pair-end (PE) sequencing

# Pair-end reads

## File format

- Two files are created
- The order in files identical and naming of reads are the same with the exception of the end
- The way of naming reads are changing over time so the read names depend on software version
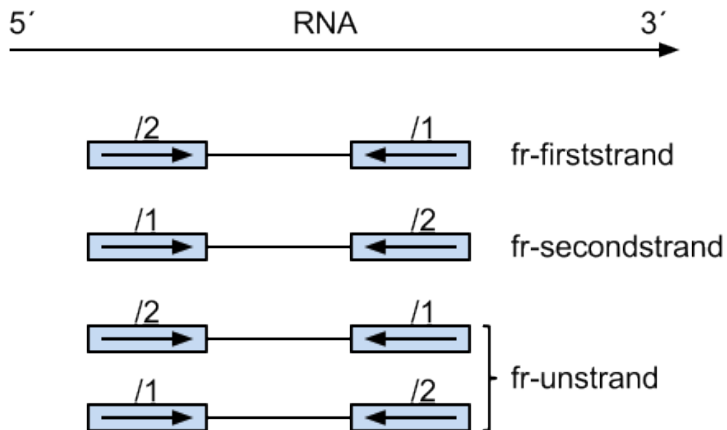
```
@61DFRAAXX100204:1:100:10494:3070/1
AAACAACAGGGCACATTGTCACTCTTGTATTTGAAAAACACTTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCCC@@CACCCCCA
```

```
@61DFRAAXX100204:1:100:10494:3070/2
ATCCAAGTTAAAACAGAGGCCTGTGACAGACTCTTGGCCCATCGTGTTGATA
+
_^_a^cccegcgghhgZc`ghhc^egggd^_[d]defcdfd^Z^OXWaQ^ad
```
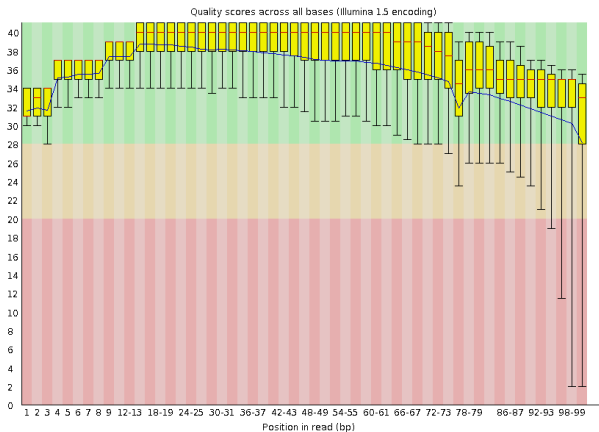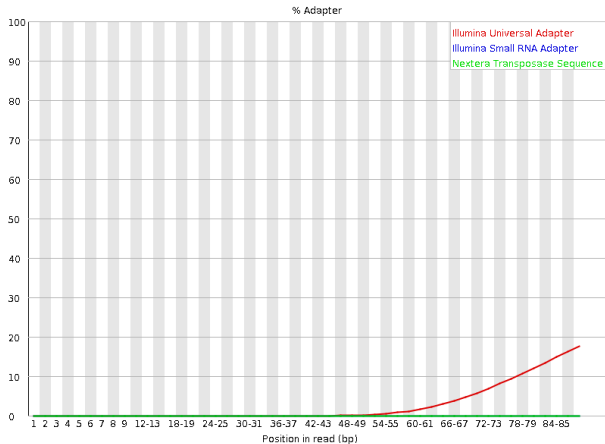
# Pair-end data

# Stranded or not

# Basic quality control of raw reads

- FastQC



Quality scores across all bases (Illumina 1.5 encoding)

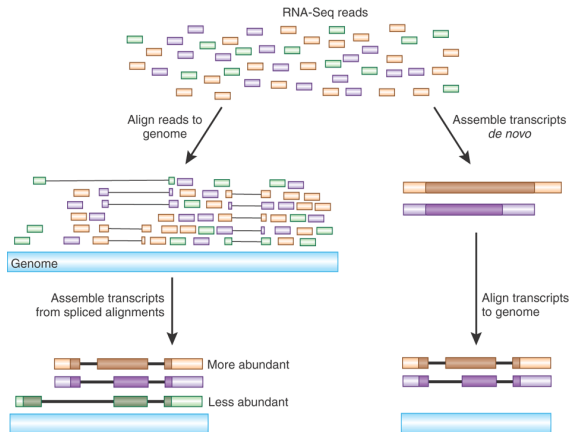# Basic quality control of raw reads

- FastQC

# Basic quality control of raw reads

- RNA-seq is not random sample from the genome eg. GC content might be different
- Highly expressed genes can be frequent and create warnings in quality controls that assumes whole genome data
- Random hexamer in cDNA synthesis might create 'biases' in base frequencies in the beginning of reads

# Two main routes for analysis



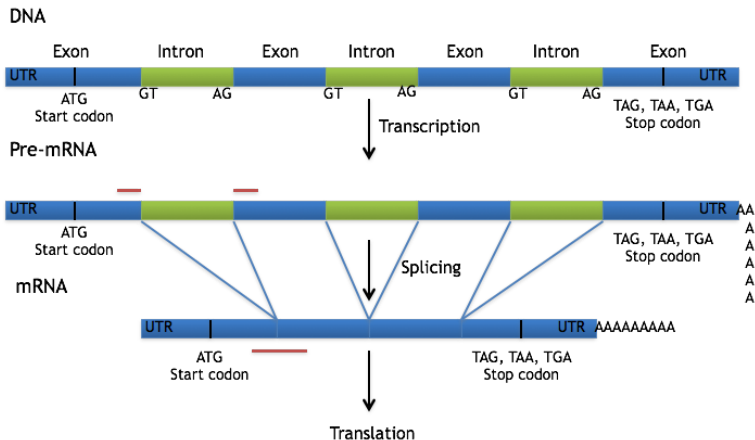Haas & Zody (2010), Nature Biotechnology 28, 421–423

# Aligning short reads from RNA to genomes

- If available map to the genome sequence
- If no genome sequence one can also map to transcriptome reference
- Make use of available genome annotation (GTF, GFF, BED files)

# Aligning short reads from RNA to genomes

- Large number of programs available: Star, Tophat, Subread etc
- Important feature: Allow for spliced mapping

# Aligning short reads from RNA to genomes

- After mapping perform QC of the output

```
read_distribution.py  -i Pairend_StrandSpecific_51mer_Human_hg19.bam -r hg19.refseq.bed12
```

Output:

| Group | Total_bases | Tag_count | Tags/Kb |
|---|---|---|---|
| CDS_Exons | 33302033 | 20002271 | 600.63 |
| 5'UTR_Exons | 21717577 | 4408991 | 203.01 |
| 3'UTR_Exons | 15347845 | 3643326 | 237.38 |
| Introns | 1132597354 | 6325392 | 5.58 |
| TSS_up_1kb | 17957047 | 215331 | 11.99 |
| TSS_up_5kb | 81621382 | 392296 | 4.81 |
| TSS_up_10kb | 149730983 | 769231 | 5.14 |
| TES_down_1kb | 18298543 | 266161 | 14.55 |
| TES_down_5kb | 78900674 | 729997 | 9.25 |
| TES_down_10kb | 140361190 | 896882 | 6.39 |

# Example workflow

- Tophat: Aligns reads to genome (allows for spliced read mapping)
- Cufflinks: Extract transcripts from spliced read alignments
- Cuffmerge: Merge results from multiple Cufflinks results
- Cuffdiff: Detect differential gene expression



Trapnell *et al.* (2012), Nature Protocols 7, 562–578

# Tophat

1. Efficient and fast alignment to the genome using bowtie2
2. Create a data base of putative splice junctions from the reads mapping in step 1
3. Map reads that did not map in step 1 using the splice information

# QC of mapped reads
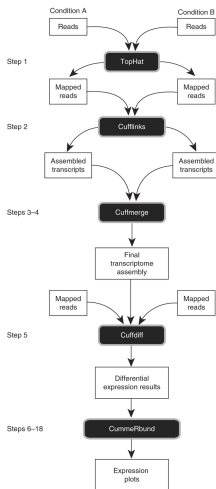
Reads should mostly map to known (annotated genes)

```
read_distribution.py  -i Pairend_StrandSpecific_51mer_Human_hg19.bam -r hg19.refseq.bed12
```

Output:

| Group | Total_bases | Tag_count | Tags/Kb |
|---|---|---|---|
| CDS_Exons | 33302033 | 20002271 | 600.63 |
| 5'UTR_Exons | 21717577 | 4408991 | 203.01 |
| 3'UTR_Exons | 15347845 | 3643326 | 237.38 |
| Introns | 1132597354 | 6325392 | 5.58 |
| TSS_up_1kb | 17957047 | 215331 | 11.99 |
| TSS_up_5kb | 81621382 | 392296 | 4.81 |
| TSS_up_10kb | 149730983 | 769231 | 5.14 |
| TES_down_1kb | 18298543 | 266161 | 14.55 |
| TES_down_5kb | 78900674 | 729997 | 9.25 |
| TES_down_10kb | 140361190 | 896882 | 6.39 |

# QC of mapped reads

Most splice event should be known and canonical (GU-AG)

# Cufflinks



**a** Splice-align reads to the genome

**b** Build a graph representing alternative splicing events

**c** Traverse the graph to assemble variants

**d** Assembled isoforms

Nature Reviews | **Genetics**

# Cuffdiff

- Program that estimate expression levels and identify differentially expressed genes from ngs alignments
- Basically uses the read data to estimate dispersion parameters (the amount of deviation from a Poisson distr.)
- Genes that show patterns deviating from the above expectations are differentially expressed between treatments
- Will work also for detection of isoform differential expression

# From counts to gene expression

# From counts to gene expression

# Not all reads are the same



| | 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 |

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

# Normalized expression Values

- Mapped read counts are normalized for both length of the transcript they map to and total depth of sequencing.
- Count data is hence converted to: Reads/Fragments per kb of transcript length and million mapped reads (RPKM or FPKM)

# Experimental design

# Experimental design

- Count reads (convert to RPKM/FPKM?)
- Small number of reads (= low RPKM/FPKM values) often non-significant
- Remember that Fold change is not the same as significance

|  | Condition 1 | Condition 2 | Fold_Change | Significant? |
|---|---|---|---|---|
| **Gene A** | 1 | 2 | 2-fold | No |
| **Gene B** | 100 | 200 | 2-fold | Yes |

# Two main routes for analysis



RNA-Seq reads

Align reads to genome

Assemble transcripts *de novo*

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant

Less abundant

Haas & Zody (2010), Nature Biotechnology 28, 421–423

# Major challenges in relation to genome assembly

- Genes show different levels of gene expression, hence uneven coverage among genes
- Many genes are expressed in different isoforms
- As sequence depth increase detected number of loci increase. (What is actually expressed?)
- Sequence error from highly expressed genes might be seen more often than "true" sequences from lowly expressed genes

# Several programs available

- SOAP-denovo TRANS
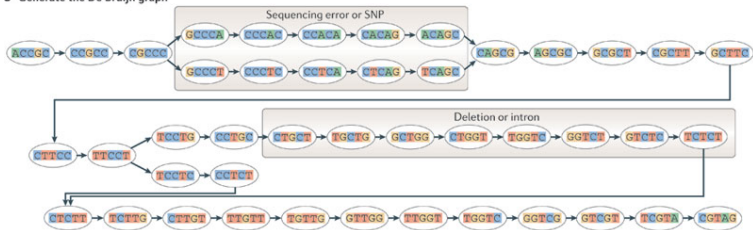- Oases
- Trans-ABYSS
- Trinity

All of them uses de Bruijn graphs to cope with the data and many of them have been developed from a genome assembly program
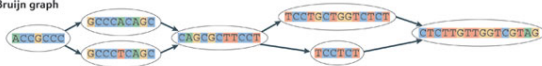
# Trinity



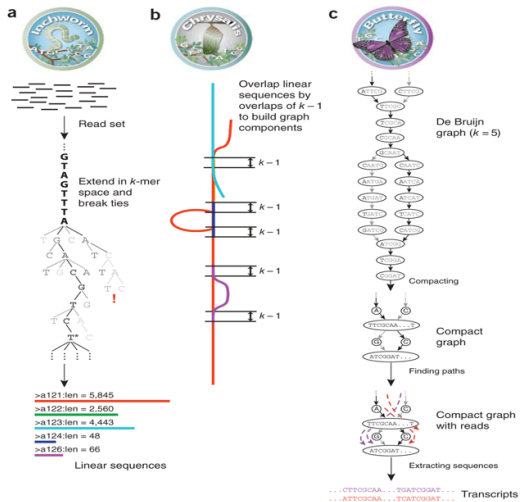**a** Generate all substrings of length k from the reads

**b** Generate the De Bruijn graph

**c** Collapse the De Bruijn graph

# Trinity

# Summary - with ref.

- Map to genome allow for spliced alignment
- If novel transcripts of interest: use method that can re-create transcripts from mapped reads (Cufflinks, Scripture or Bayesembler)
  NB! In well annotated genomes most reads should map to known genes
- If interest is expression of known genes/exons: Use available annotation for analysis
- Spend time on experimental design and more replicates gives more power in gene expression analysis

# Summary - without ref.

- Assemble using your favourite assembler
- Spend lots of time in assessing the results (compare to related species, look for ORFs etc)
- Often large number of partial transcripts (hence often large number of contigs).
- Merge with other data from transcripts?