

Characterizing transcriptomes using ngs data

T. Källman

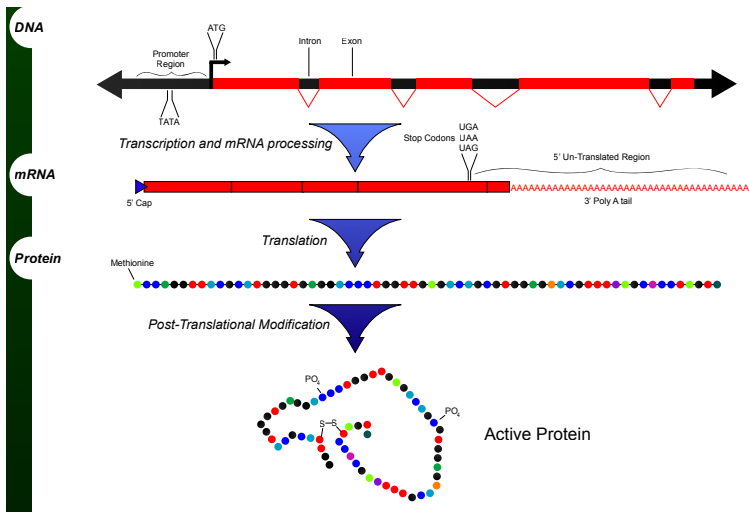
BILS/Scilife Lab/Uppsala University

May 2015

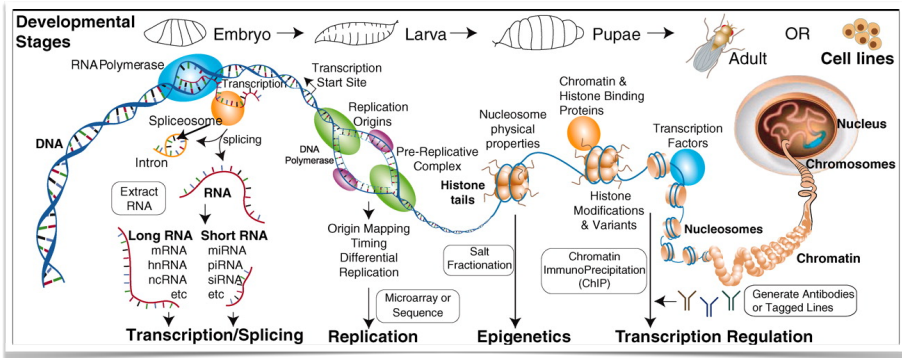
Outline

- 1 The transcriptome
- 2 RNA sequence technologies
- 3 RNA-seq analysis
 - Mapping based approach
 - Tools for working with ngs alignments
 - Gene expression from RNA-seq
 - de-novo assembly

The Central Dogma



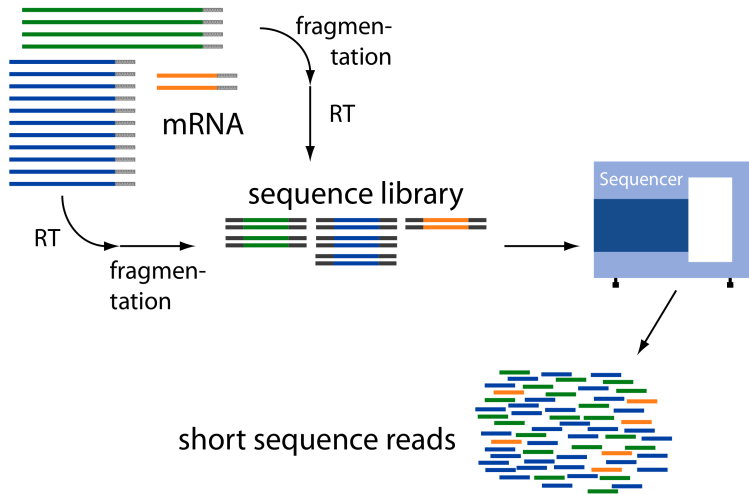
A more complex view



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



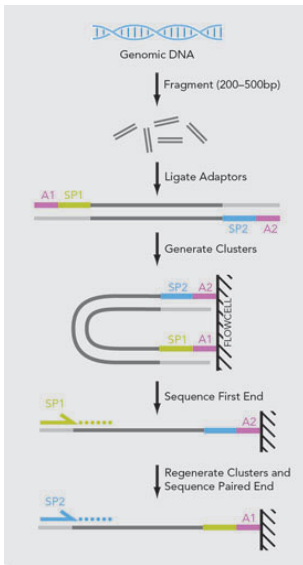
Machine output

```

@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@BBBBB@B9B@=BABA@A:@693:@B=
@SRR038845.53 HWI-EAS038:6:1:1:360 length=36
GTTCAAAAAGAACTAAATTGTGTCAATAGAAAACCTC
+SRR038845.53 HWI-EAS038:6:1:1:360 length=36
BBCBBBBBB@@BAB?BBBBCBC>BBBAA8>BBBAA@

```


Pair-end (PE) sequencing



Pair-end reads

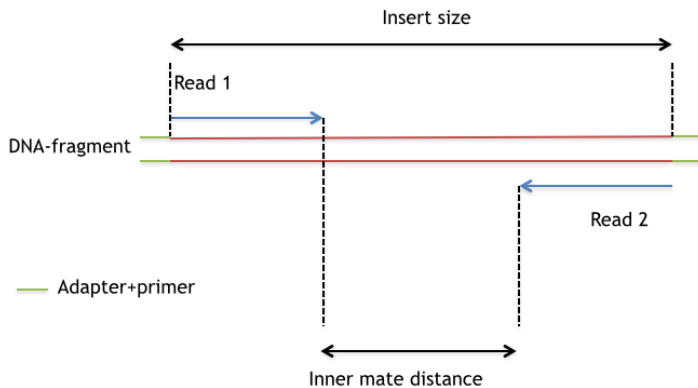
File format

- 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

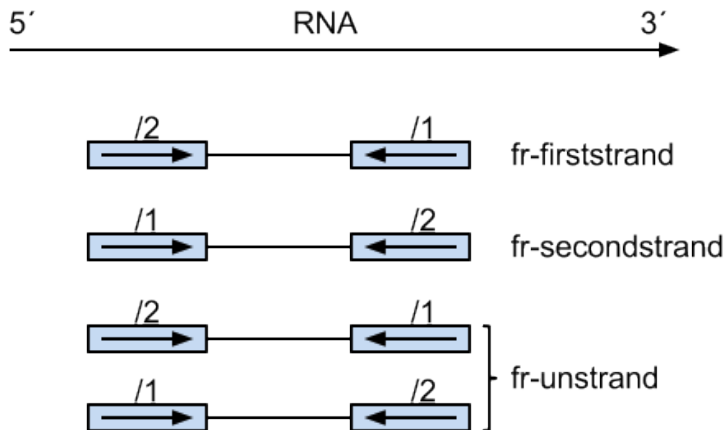
```
@61DFRAAXX100204:1:100:10494:3070/1
AAACAACAGGGCACATTGTCACTCTTGTATTTGAAAAACACTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCC@@CACCCCCA
```

```
@61DFRAAXX100204:1:100:10494:3070/2
ATCCAAGTTAAACAGAGGCCTGTGACAGACTCTTGGCCATCGTGTTGATA
+
_^_a^cccegcgghhgZc`ghhc^egggd^_[d]defcdfd^Z^OXWaq^ad
```

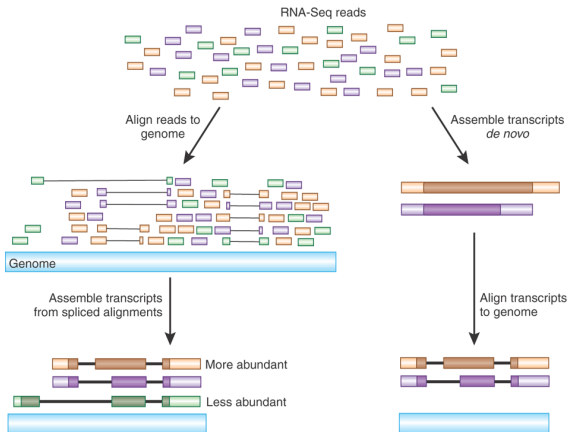
Pair-end data



Stranded or not



Two main routes for analysis



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

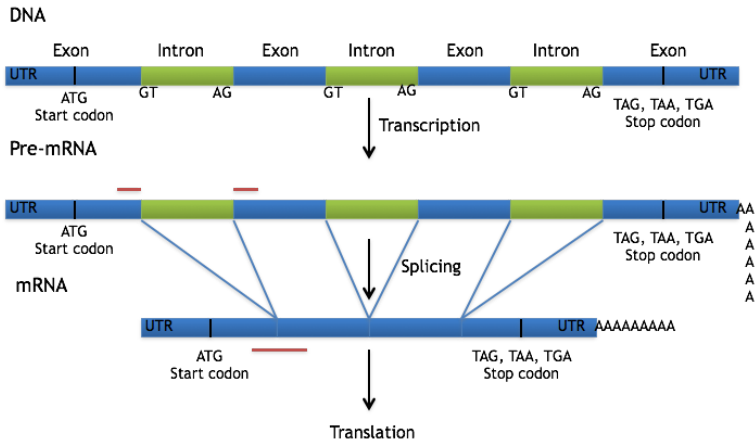
Aligning short reads from RNA to genomes

- 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)

Galaxy		Analyze Data	Workflow	Shared Data	Visualization	Help	User	Using 0 bytes
Seqname	Source	Feature	Start	End	Score	Strand	Frame	Attributes
chr12	unknown	exon	87984	88017	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	88257	88392	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	88570	88771	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	88860	89018	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	89675	89927	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	90587	90655	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	90796	91263	-	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS58200";
chr12	unknown	exon	147946	148509	-	-	.	gene_id "FAM1380"; gene_name "FAM1380"; transcript_id "NR_026823"; tss_id "TSS11802";
chr12	unknown	exon	148612	148814	-	-	.	gene_id "FAM1380"; gene_name "FAM1380"; transcript_id "NR_026823"; tss_id "TSS11802";
chr12	unknown	exon	149052	149412	-	-	.	gene_id "FAM1380"; gene_name "FAM1380"; transcript_id "NR_026823"; tss_id "TSS11802";
chr12	unknown	CDS	176049	176602	-	0	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	176049	176602	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	start_codon	176049	176051	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	186542	186878	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P13619"; transcript_id "NM_015232"; tss_id "TSS12565";
chr12	unknown	CDS	208312	208380	-	1	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	208312	208380	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P13619"; transcript_id "NM_015232"; tss_id "TSS12565";
chr12	unknown	exon	208312	208380	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	CDS	234790	235078	-	1	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	234790	235078	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	246577	246793	-	+	.	gene_id "LOC574538"; gene_name "LOC574538"; transcript_id "NR_033859"; tss_id "TSS17153";
chr12	unknown	CDS	247433	248520	-	0	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	exon	247433	248520	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P13619"; transcript_id "NM_015232"; tss_id "TSS12565";
chr12	unknown	exon	247433	248520	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";
chr12	unknown	CDS	247439	248520	-	0	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P13619"; transcript_id "NM_015232"; tss_id "TSS12565";
chr12	unknown	start_codon	247439	247441	-	+	.	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P13619"; transcript_id "NM_015232"; tss_id "TSS12565";

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

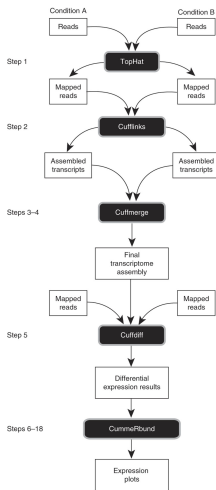
```
read_distribution.py -i Paired_StrandSpecific_51mer_Human_hg19.bam -r hg19.refseq.bed12
```

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

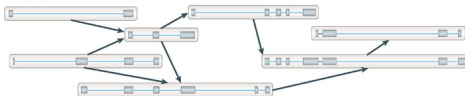
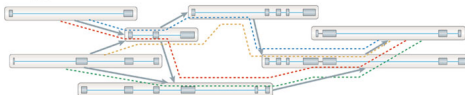


Trapnell *et al.* (2012), Nature Protocols 7, 562–578

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 run using the splice information

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**

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

Samtools

- Program to work with ngs alignment files (SAM, BAM, CRAM)
- Can be used to view data, calculate basic info, extract subsets of alignments and convert between file formats
- <http://www.htslib.org>

Picard

- A set of Java command line tools with the same (or similar functionality as samtools)
- Note that even though they largely aim at doing similar functions Picard and Samtools is not always generating compatible file formats
- <http://broadinstitute.github.io/picard/>

IGV: Integrative Genomics Viewer

Viewer

- Home
- Downloads
- Documents
 - Hosted Genomes
 - FAQ
 - IGV User Guide
 - File Formats
 - Release Notes
 - IGV for iPad
 - Credits
- Contact

Search website

search

[Broad Home](#)

[Cancer Program](#)

BROAD
INSTITUTE

© 2013 Broad Institute



What's New

September 2014. The IGV iPad app can now be installed from the Apple App Store. **IGV for iPad** is a lightweight genomic data viewer that provides some of the functionality available in our regular desktop IGV. See the [IGV for iPad documentation](#) for details.

June 2014. We're hiring! See the [job description](#) on the Broad Institute careers website.

Overview

The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated genomic datasets. It supports a wide variety of data types, including array-based and next-generation sequence data, and genomic annotations.

Downloads

Please [register](#) to download IGV. After registering, you can log in at any time using your email address. Permission to use IGV is granted under

Citing IGV

To cite your use of IGV in your publication:

Heiga Thorvaldsdóttir, James T. Robinson, Jill P. Mesirov. **Integrative Genomics Viewer (IGV): high-performance genomics data visualization and exploration.** *Briefings in Bioinformatics* 2012.

James T. Robinson, Heiga Thorvaldsdóttir, Wendy Winckler, Mitchell Gutman, Eric S. Lander, Gad Getz, Jill P. Mesirov. **Integrative Genomics Viewer.** *Nature Biotechnology* 29, 24–26 (2011)

Funding

Development of IGV is made possible by funding from the [National Cancer Institute](#), the [National Institute of General Medical Sciences](#) of the [National Institutes of Health](#), and the [Starr Cancer Consortium](#).

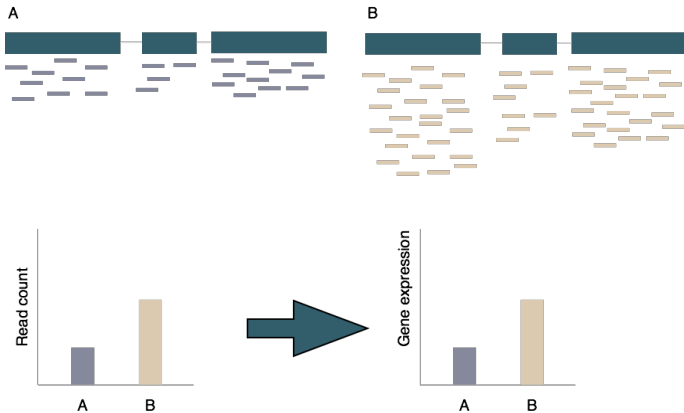
IGV participates in the [GenomeSpace Initiative](#), which is funded by the [National Human Genome Research Institute](#).



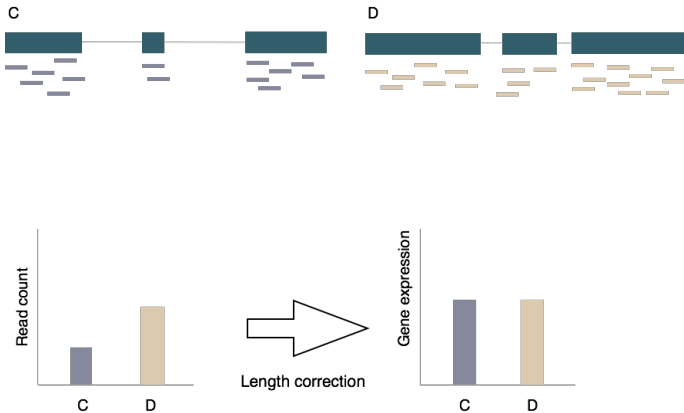
IGV: Integrative Genomics Viewer



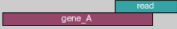
From counts to gene expression



From counts to gene expression



Not all reads are the same

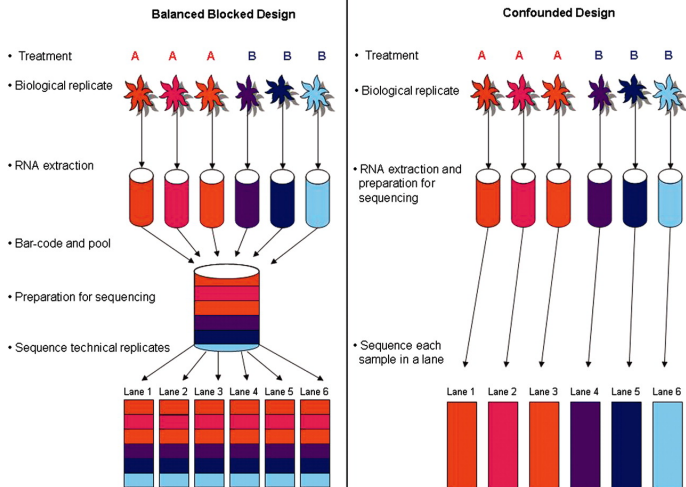
	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

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

Normalized expression Values

- Transcript-mapped read counts are normalized for both length of the transcript 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

Major challenges in relation to genome assembly

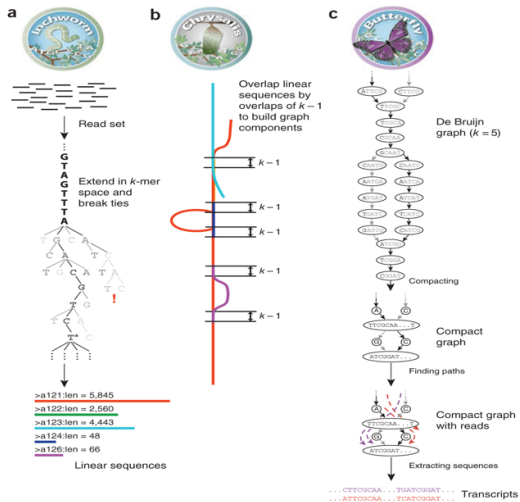
- Genes show different levels of gene expression, hence uneven coverage among genes
- Many genes are expressed in different isoforms
- As sequence depth increase detected number of loci increase. (What is actually expressed?)
- Sequence error from highly expressed genes might be seen more often than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



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 Bayesemblem)
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
- Replicate, replicate....!

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)