

Characterizing transcriptomes using ngs data

T. Källman

BILS/Scilife Lab/Uppsala University

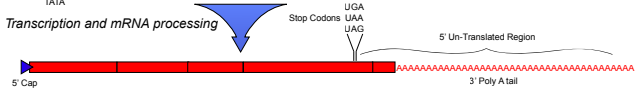
Sep. 2015

Outline

DNA



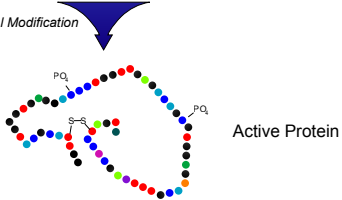
mRNA



Protein



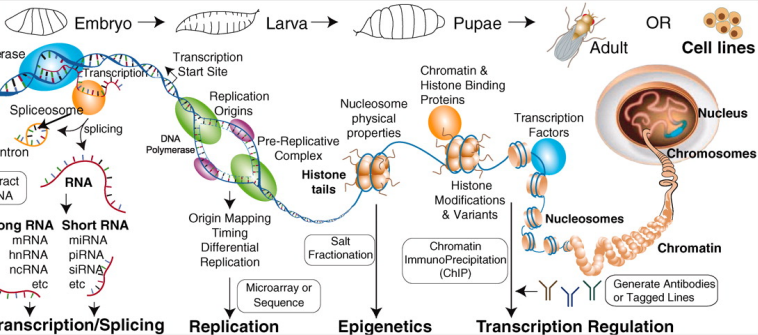
Post-Translational Modification



The transcriptome

al Dogma

Developmental Stages

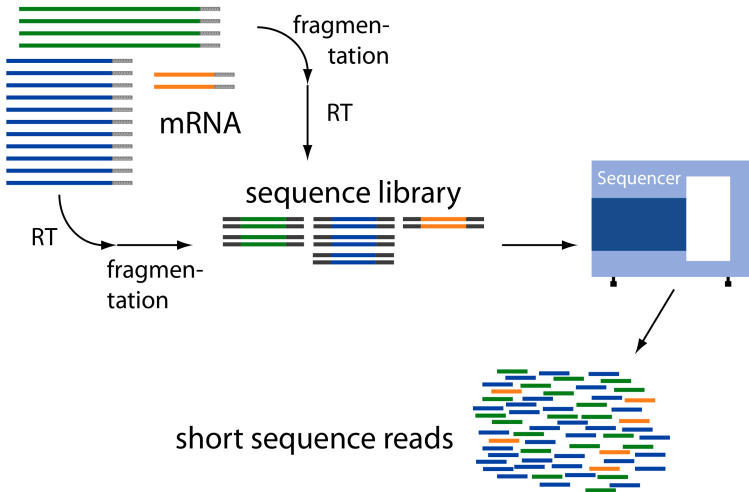


The transcriptome

Complex view

The transcriptome

omes vs genomes



RNA sequence technologies

RNA sequence technologies

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@BBBBAB@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
BBCBBBBBBB@@BAB?BBBBBCBC>BBBAA8>BBBAA@

RNA sequence technologies

output

RNA sequence technologies

quality



Genomic DNA



Fragment (200–500bp)



Ligate Adaptors



Generate Clusters



Sequence First End



Regenerate Clusters and
Sequence Paired End



RNA sequence technologies

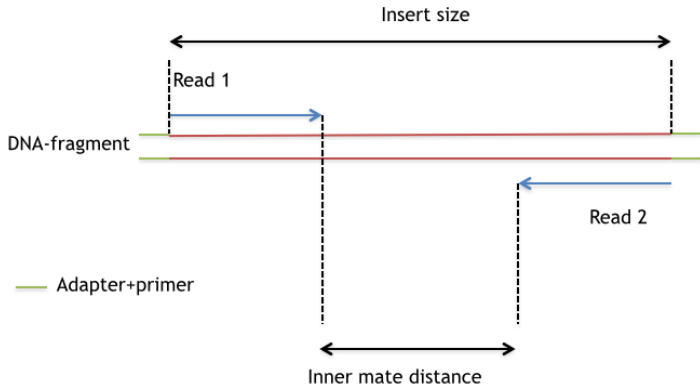
PE) sequencing

@61DFRAAXX100204:1:100:10494:3070/1
AAACAACAGGGCACATTGTCACCTCTTGTATTTGAAAAACACTTTCCGGCCAT
+
ACCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCBC?CCCCCCCCCCCC@@CACCCCCA

@61DFRAAXX100204:1:100:10494:3070/2
ATCCAAGTTAAACAGAGGCCTGTGACAGACTCTTGGCCCATCGTGTTGATA
+
_ ^ _ a ^ c c c e g c g g h h g Z c ` g h h c ^ e g g g d ^ _ [d] d e f c d f d ^ Z ^ 0 X W a Q ^ a d

RNA sequence technologies

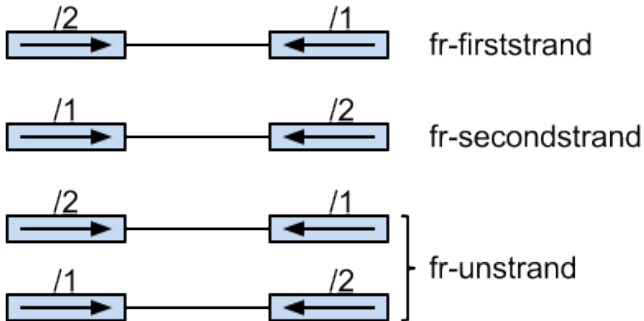
reads



RNA sequence technologies

data

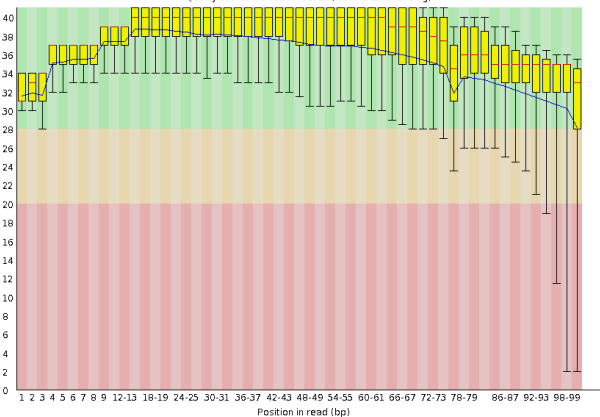
5' RNA 3'



RNA sequence technologies

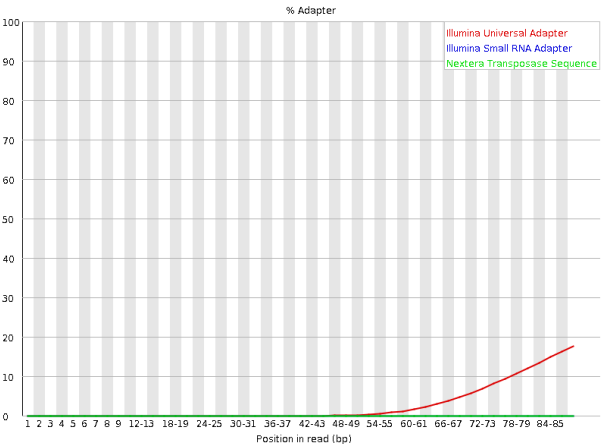
or not

Quality scores across all bases (Illumina 1.5 encoding)



RNA sequence technologies

Quality control of raw reads



RNA sequence technologies

Quality control of raw reads

RNA sequence technologies

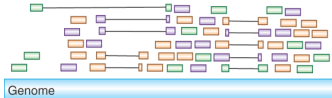
Quality control of raw reads

RNA-Seq reads



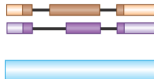
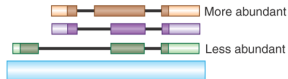
Align reads to genome

Assemble transcripts *de novo*



Assemble transcripts from spliced alignments

Align transcripts to genome



RNA-seq analysis

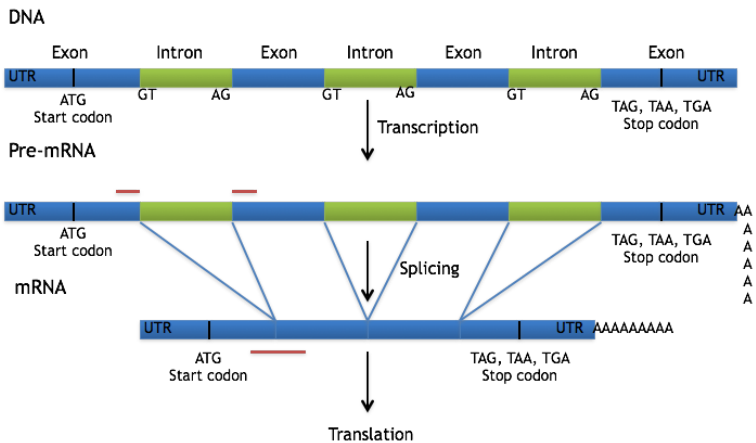
routes for analysis

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 "TSS8200";	
chr12	unknown	exon	88257	88392	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	88570	88771	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	88860	89018	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	89675	89827	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	90587	90655	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	90796	91263	.	+	.	gene_id "LOC100288778"; gene_name "LOC100288778"; transcript_id "NR_028269"; tss_id "TSS8200";	
chr12	unknown	exon	147946	148509	.	-	.	gene_id "FAM138D"; gene_name "FAM138D"; transcript_id "NR_026823"; tss_id "TSS11862";	
chr12	unknown	exon	148612	148814	.	-	.	gene_id "FAM138D"; gene_name "FAM138D"; transcript_id "NR_026823"; tss_id "TSS11862";	
chr12	unknown	exon	149052	149412	.	-	.	gene_id "FAM138D"; gene_name "FAM138D"; transcript_id "NR_026823"; tss_id "TSS11862";	
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	234799	235078	.	+	1	gene_id "IQSEC3"; gene_name "IQSEC3"; p_id "P5442"; transcript_id "NM_001170738"; tss_id "TSS17433";	
chr12	unknown	exon	234799	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";	

RNA-seq analysis

Mapping based approach

Short reads from RNA to genomes



RNA-seq analysis

Mapping based approach

Short reads from RNA to genomes

```
read_distribution.py -i Pairend_StrandSpecific_5lmer_Human_hg19.bam -r hg19.refseq.bed12
```

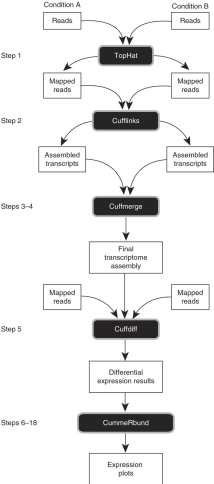
Output:

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

RNA-seq analysis

Mapping based approach

Short reads from RNA to genomes



RNA-seq

workflow

- Aligns reads to

RNA-seq

```
read_distribution.py -i Paired_StrandSpecific_5lmer_Human_hg19.bam -r hg19.refseq.bed12
```

Output:

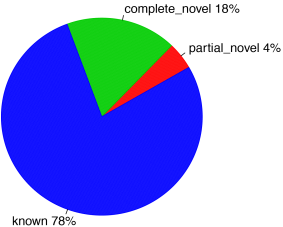
Group	Total_bases	Tag_count	Tags/Kb
CDS_Exons	33302033	20002271	600.63
5'UTR_Exons	21717577	4408991	203.01
3'UTR_Exons	15347845	3643326	237.38
Introns	1132597354	6325392	5.58
TSS_up_1kb	17957047	215331	11.99
TSS_up_5kb	81621382	392296	4.81
TSS_up_10kb	149730983	769231	5.14
TES_down_1kb	18298543	266161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

RNA-seq analysis

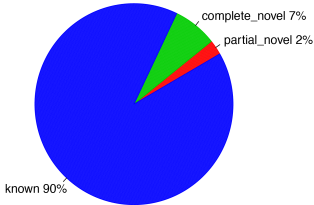
Mapping based approach

mapped reads

splicing junctions



splicing events

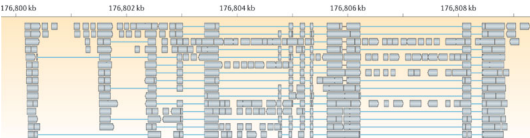


RNA-seq analysis

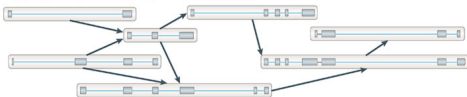
Mapping based approach

mapped reads

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



RNA-seq analysis

Mapping based approach

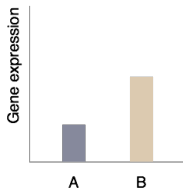
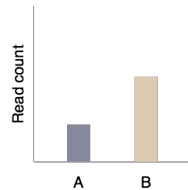
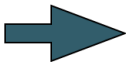
RNA-seq analysis

Mapping based approach

A



B



RNA-seq analysis

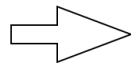
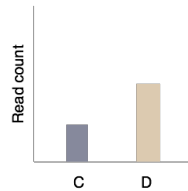
Gene expression from RNA-seq

nts to gene expression

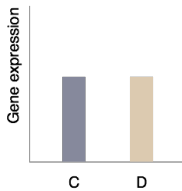
C



D



Length correction



RNA-seq analysis

Gene expression from RNA-seq

nts to gene expression

	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
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

RNA-seq analysis

Gene expression from RNA-seq

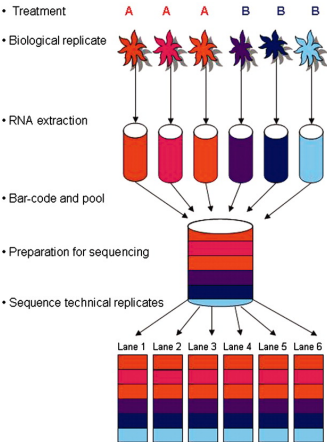
reads are the same

RNA-seq analysis

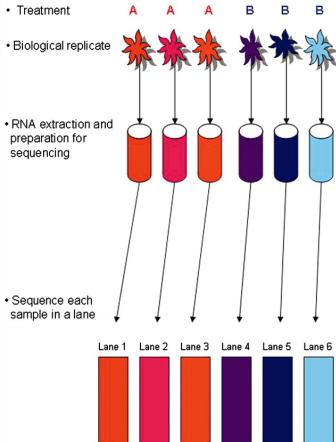
Gene expression from RNA-seq

ed expression Values

Balanced Blocked Design



Confounded Design



RNA-seq analysis

Gene expression from RNA-seq

Experimental design

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

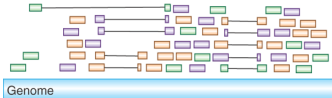
Ita designa

RNA-Seq reads



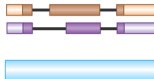
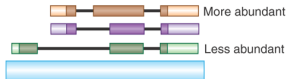
Align reads to genome

Assemble transcripts *de novo*



Assemble transcripts from spliced alignments

Align transcripts to genome



RNA-seq analysis

Gene expression from RNA-seq

routes for analysis

RNA-seq analysis

de-novo assembly

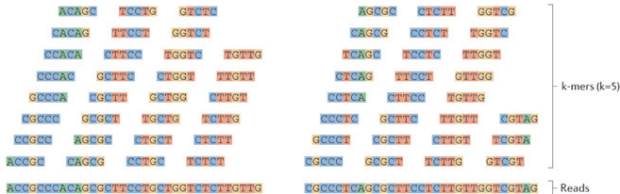
Challenges in relation to genome assembly

RNA-seq analysis

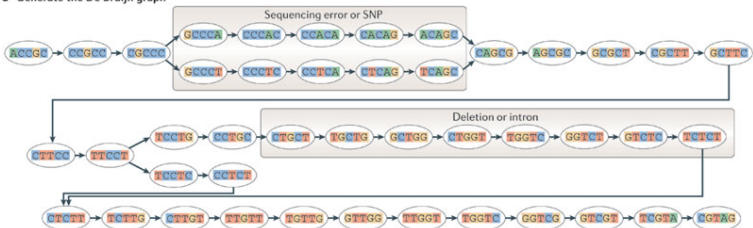
de-novo assembly

programs available

Generate all substrings of length k from the reads



b Generate the De Bruijn graph



c Collapse the De Bruijn graph



d Traverse the graph



e Assembled isoforms



RNA-seq analysis

de-novo assembly

a



Read set

...
G
T
A
G
T
T
T
A

Extend in k -mer space and break ties

>a121:len = 5,845

>a122:len = 2,560

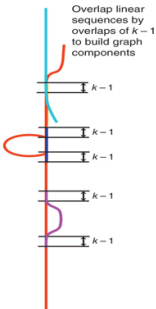
>a123:len = 4,443

>a124:len = 48

>a126:len = 66

Linear sequences

b



c

De Bruijn graph ($k = 5$)

Compacting



Compact graph

Finding paths



Compact graph with reads

Extracting sequences

...CTTCGCAA...TGATCGGAT...
...ATTTCGCAA...TCATCGGAT...

Transcripts

RNA-seq analysis

de-novo assembly

RNA-seq analysis

de-novo assembly

- with ref.

RNA-seq analysis

de-novo assembly

- without ref.