RNA Sequencing

MMG 835, SPRING 2016 Eukaryotic Molecular Genetics

George I. Mias

Department of Biochemistry and Molecular Biology gmias@msu.edu

RNA-Sequencing

- High throughput method
- Count the number of reads (transcripts)
- Can probe entire genome for expression (non-targeted approach)
- Splicing information
- Alternative splicing and novel isoforms
- Possible to phase (algorithmically/new longer read technology)
- Can get variant information
- Now cheaper (almost comparative to microarray, cost mainly library prep)

RNA-Sequencing RNA-Seq Advantages

Technology	Tiling microarray	cDNA or EST sequencing	RNA-Seq
Technology specifications			
Principle	Hybridization	Sanger sequencing	High-throughput sequencing
Resolution	From several to 100 bp	Single base	Single base
Throughput	High	Low	High
Reliance on genomic sequence	Yes	No	In some cases
Background noise	High	Low	Low
Application			
Simultaneously map transcribed regions and gene expression	Yes	Limited for gene expression	Yes
Dynamic range to quantify gene expression level	Up to a few-hundredfold	Not practical	>8,000-fold
Ability to distinguish different isoforms	Limited	Yes	Yes
Ability to distinguish allelic expression	Limited	Yes	Yes
Practical issues			
Required amount of RNA	High	High	Low
Cost for mapping transcriptomes of large genomes	High	High	Relatively low

Wang, Gerstein & Snyder, Nature Review Genetics 10:57-63 (2009)

RNA-Sequencing Keep in mind

- Quantitation Issues
- Biases
- Transcribed region Complexity
- Sample quality (degradation)
- Aligner dependent alignment

Data Considerations

Steps

- Experiment Design
- Library Prep
- Sequencing
- · Data
- Quality Control (QC)
- · Assembly
- Annotation
- Quantification
- Differential Expression
- Biological Significance



illumina.com



Possibilities

mRNA Poly(A)-selection

- hybridization to oligo-dT beads
- N.B. 3' bias
- rRNA depletion
- · small RNA
 - size selection
- Paired end
- Low RNA amounts amplification

Normalization

- Long transcripts
- Coverage differences

RPKM, Reads **p**er **K**ilobase per **M**illion mapped reads

FPKM, Fragments **p**er **K**ilobase per **M**illion mapped reads (paired-end)

Mortazavi A., et al. Nature Methods 5, 621 - 628 (2008)



Mismatches

- biological variation
 - genomic DNA
 - RNA-editing
- errors in library preparation
 - hexamer mispriming
 - PCR errors
- sequencing errors [id of bases]

van Gurp, McIntyre & Verhoeven, PLoS ONE 8(12): e85583 (2013)

 mismatches mainly in the first seven nucleotides of first strand cDNA





Jiang L et al., Genome Res 21: 1543–1551 (2011) Pickrell, Gilad & Pritchard, Science 335: 1302–authorreply1302 (2012).

Strand specific positional error rates in ERCC control RNAseq data

random hexamer mispriming





van Gurp, McIntyre & Verhoeven, PLoS ONE 8(12): e85583 (2013)

Rna-Seq Considerations



Li, Tighe et al., Nature Biotechnology 12(9), p. 915 (2014)

Data Considerations

Mapping Reference

- Genome
- Transcriptome

Alignment

- Short read alignment
- Long read alignment
- Spliced Vs. Non Spliced Alignment
- Mismatches & Duplications
- Non-Uniqueness
- Variant Calls

De Novo Assembly

- Genome based
- non-genome based

.fastq format

Four lines repeating:

- 1. @title and optional description / Sequence Identifier
- 2. Sequence
- 3. + (and *optional* repeat of title line)
- 4. Quality scores corresponding to Sequence (2)

Pearson & Lipman Proc. Natl Acad. Sci. USA, 85, 2444–2448 (1988) Cock et al., Nucleic Acids Research, Vol. 38, No. 6 1767–1771 (2010) <u>http://support.illumina.com/content/dam/illumina-support/help/BaseSpaceHelp_v2/</u> <u>Content/Vault/Informatics/Sequencing_Analysis/BS/swSEQ_mBS_FASTQFiles.htm</u>

.fastq format

Four lines repeating:

- 1. @title and optional description / Sequence Identifier
- 2. Sequence
- 3. + (and *optional* repeat of title line)
- 4. Quality scores corresponding to Sequence (2)

Example:

@SIM:1:FCX:1:15:6329:1045 1:N:0:2 TCGCACTCAACGCCCTGCATATGACAAGACAGAATC

╋

Pearson & Lipman Proc. Natl Acad. Sci. USA, 85, 2444–2448 (1988) Cock et al., Nucleic Acids Research, Vol. 38, No. 6 1767–1771 (2010) <u>http://support.illumina.com/content/dam/illumina-support/help/BaseSpaceHelp_v2/</u> <u>Content/Vault/Informatics/Sequencing_Analysis/BS/swSEQ_mBS_FASTQFiles.htm</u>

.fastq format

@<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<x-pos> <read>:<is filtered>:<control number>:<sample number>

Element	Requirements	Description
@	@	Each sequence identifier line starts with @
<instrument></instrument>	Characters allowed: a–z, A–Z, 0–9 and underscore	Instrument ID
<run number=""></run>	Numerical	Run number on instrument
<flowcell id=""></flowcell>	Characters allowed: a–z, A–Z, 0–9	
<lane></lane>	Numerical	Lane number
<tile></tile>	Numerical	Tile number
<x_pos></x_pos>	Numerical	Run number on instrument
<y_pos></y_pos>	Numerical	X coordinate of cluster
<read></read>	Numerical	Read number. 1 can be single read or Read 2 of paired-end
<is filtered=""></is>	Y or N	Y if the read is filtered (did not pass), N otherwise
<control number=""></control>	Numerical	0 when none of the control bits are on, otherwise it is an even number. On HiSeq X systems, control specification is not performed and this number is always 0.
<sample number=""></sample>	Numerical	Sample number from sample sheet
http://support.ill	umina.com/conte	nt/dam/illumina-support/help/BaseSpaceHelp_v2/
Content/Vault/In	formatics/Sequer	ncing Analysis/BS/swSEQ mBS FASTQFiles htm

Phred Scores

- Phred Algorithm
- Used with Sanger Data.
- Algorithm assigns probability of error in calling a base

```
Quality score (Q-Score)
Q= -10 x log(error probability)
P=10^{-Q/10}
```

Ewing et al., Genome Research 8: 175-185 (1998). Ewing & Green, Genome Research 8: 186-194 (1998).

Q-Scores

_	Quality Score	Error Probability				
	Q40	0.0001 (1 in 10,000)				
	Q30	0.001 (1 in 1,000)				
	Q20	0.01 (1 in 100)				
	Q10	0.1 (1 in 10)				

Quality score (Q-Score) Q= -10 x log(error probability) P= $10^{-Q/10}$

http://www.illumina.com/content/dam/illumina-marketing/documents/products/ technotes/technote_understanding_quality_scores.pdf

Q-Scores Ascii table

Encode Q-Score to characters Using ASCII

Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	@	96	60	`
1	1	[START OF HEADING]	33	21	1	65	41	Α	97	61	а
2	2	[START OF TEXT]	34	22		66	42	В	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	С	99	63	С
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	е
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	1.00	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	н	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49	1.00	105	69	i.
10	Α	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	В	[VERTICAL TAB]	43	2B	+	75	4B	Κ	107	6B	k
12	С	[FORM FEED]	44	2C	,	76	4C	L	108	6C	1
13	D	[CARRIAGE RETURN]	45	2D	÷	77	4D	Μ	109	6D	m
14	E	[SHIFT OUT]	46	2E		78	4E	Ν	110	6E	n
15	F	[SHIFT IN]	47	2F	1	79	4F	0	111	6F	0
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	Ρ	112	70	р
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	S
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	т	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	V
23	17	[ENG OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Υ	121	79	У
26	1A	[SUBSTITUTE]	58	3A	1.00	90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	{
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	١	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D	1	125	7D	}
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	~
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	[DEL]

(ASCII, American Standard Code for Information Interchange)

https://simple.wikipedia.org/wiki/ASCII

Q-Scores

SS	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS	SSSSSSSSSSSSS	SSSSSSSSSS.				
• •		XXXXX	xxxxxxxxxxx	XXXXXXXXXXXX	XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX	xxxxxx	
			IIIIIIIIIIII	IIIIIIIIIIIII	IIIIIIIIIIIII	IIIIII. .	
т.т.1							
- " -	#\$%\$\$'()*+ /0123	156790 • • < ->2			wvv7[\]^ `ab	adofahijklm	nongratuwww.
• 7 	πφοα () ,/0125			KLIMOPQK510V	WXIZ[\] _ aD		
33		59 6	4 73			104	126
0.			1				
		-5	09.	•••••	•••••	40	
			09.			40	
			39.			40	
0.2	2		1				
s –	Sanger	Phred+33,	raw reads	typically	(0, 40)		
х –	Solexa	Solexa+64,	raw reads	typically	(-5, 40)		
I -	Illumina 1.3+	Phred+64,	raw reads	typically	(0, 40)		
л –	Tllumina 1 5+	Phred+64	raw reads	typically	(3 40)		
	with O-unused	1	2-Road Co	erprearry		Indicator	(hold)
	with U=unused	, 1=unusea,	z-Read Se	gment Quall	ty control	Indicator	(DOTA)
L -	Illumina 1.8+	Phred+33,	raw reads	typically	(0, 41)		

ASCII Conversion for different schemes

https://en.wikipedia.org/wiki/FASTQ_format#Encoding

Aligners

- Unspliced
 - Fast
 - Known exons and splice junctions
 - Reference required
 - c.f. DNA alignment
 - Borrows Wheeler Transform based (BWT)
 - BWA
 - Bowtie
 - SHRiMP
 - Stampy

Aligners

- Align to whole genome, including intron-spanning reads that allow large gaps
- Exon first
 - TopHat, MapSplice
 - Unspliced alignment first
 - Unmapped reads after
- Seed-extend
 - ► GSNAP, QPALMA
 - Short seeds mapped first
- K-mer mappers

Exon-first approach

Example 2 Exons



Exon-first approach

RNA

Example 2 Exons Exon 2 Exon 1 Exon read mapping

map full, unspliced reads (exonic reads)

Exon-first approach



Seed-extend approach



store a map of all small words (k-mers) of similar size in the genome in an efficient lookup data structure

Seed-extend approach



- each read is divided into k-mers
- k-mers mapped to the genome via the lookup structure.

Seed-extend approach



- Mapped k-mers extended into larger alignments
- May include gaps flanked by splice sites.

Potential limitations of exon-first approaches



Example : gene & associated retrotransposed pseudogene.

Pseudogenes:

defunct genomic loci with sequence similarity to functional genes

but *lacking coding potential* due to the presence of disruptive mutations such as frame shifts and premature stop codons.

Pei et al., Genome Biology 13:R51 (2012). Garber et al., Nature Methods 8, 469–477 (2011).

Potential limitations of exon-first approaches



Example : gene & associated retrotransposed pseudogene.

Categories of Pseudogenes

- 1. Processed pseudogenes Created by retrotransposition of mRNA from functional protein-coding loci back into the genome
- **2. Duplicated (also referred to as unprocessed) pseudogenes** Derived from duplication of functional genes
- **3. Unitary pseudogenes** Arise through in situ mutations in previously functional protein-coding genes

Pei et al., Genome Biology 13:R51 (2012). Garber et al., Nature Methods 8, 469–477 (2011).



Example : gene & associated retrotransposed pseudogene.

- Exonic reads will map to both the gene and its pseudogene
- Gene placement preferred owing to lack of mutations
- Spliced read could be incorrectly assigned to the pseudogene as it appears to be exonic, preventing higher-scoring spliced alignments from being pursued.



Sequence-fragmented RNA

Reads originating from two different isoforms of the same genes colored **black** and **gray**.

Genome Guided Approach



Genome Guided Approach



Genome Guided Approach



Genome Guided Approach



Genome-Independent Approach Break reads into k-mer seeds



K-mer arranged into a de Bruijn graph structure.









• Examples: Trinity, TransAbyss, Velvet



Branch points and transcript possibilities



Branch points and transcript possibilities



Branch points and transcript possibilities



Garber et al., Nature Methods 8, 469–477 (2011). Select path based on coverage



• Transcripts of different lengths & different read coverage levels



- Transcripts of different lengths & different read coverage levels
- Different total read counts observed for each transcript





- Transcripts of different lengths & different read coverage levels
- Different total read
 counts observed for each transcript
- FPKMnormalized read counts



Reads from alternatively spliced genes may be attributable to a single isoform or more than one isoform.

- Reads are color-coded when their isoform of origin is clear.
- Black reads indicate reads with uncertain origin.









- E.g. gene with 2 expressed isoforms
- Exons colored according to the isoform of origin to all gene isoforms.



- E.g. gene with 2 expressed isoforms
- Exons colored according to the isoform of origin to all gene isoforms.

Gene models used for quantification



- E.g. gene with 2 expressed isoforms
- Exons colored according to the isoform of origin to all gene isoforms.

Exon intersection method

Uses only exons common to all gene isoforms

Gene models used for quantification

Comparison of true versus estimated FPKM values in *simulated* RNA-seq data.



True FPKM

Differential expression analysis



Expression Microarrays Fluorescence based comparison of intensity

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

Differential expression analysis



Differential expression analysis



- Example Gene with 2 isoforms with isoform switch in 2 conditions.
- Similar no. of reads
- Different distribution across isoforms





http://www.labome.com/method/RNA-seq-Using-Next-Generation-Sequencing.html Corney (2013) Mater Methods 3:203

Bowtie Extremely fast, general purpose short read aligner

Trapnell et al., Nature Biotechnology 28(5) p511(2010) Trapnell et al., Nature Protocols 7, 562–578 (2012)



TopHat

fragment sequences to genome

Cufflinks package

Cufflinks

Assembles transcripts

Cuffcompare Compares transcript assemblies to annotation

Cuffmerge Merges two or more transcript assemblies

Cuffdiff

Finds differentially expressed genes and transcripts Detects differential splicing and promoter use

Trapnell et al., Nature Biotechnology 28(5) p511(2010) Trapnell et al., Nature Protocols 7, 562–578 (2012)



to genome TopHat Spliced fragment alignments Cufflinks Abundance estimation Assembly d Mutually incompatible fragments Fragment Transcript coverage length and compatibility distribution Overlap graph С e Maximum likelihood abundances Log-likelihood Minimum path cover Transcript Transcripts and their abundances

fragment sequences

CummeRbund Plots abundance and differential expression results from Cuffdiff

Trapnell et al., Nature Biotechnology 28(5) p511(2010) Trapnell et al., Nature Protocols 7, 562–578 (2012)

ChIP-seq

a DNA-binding protein ChIP-seq

b Histone modification ChIP-seq



Sample fragmentation

Endonuclease (ChIP–exo)

Immunoprecipitate and

then purify DNA

pT7-AA..AA ______

LinDA

Sonication

Exonuclease

AA..AA-pT7



Crosslink proteins and DNA



Sample fragmentation

MNase digestion



Immunoprecipitate and then purify DNA



 Chromatin ImmunoPrecipitation followed by sequencing (ChIP–seq)

- Detect protein–DNA binding
- Detect chemical modifications of histone proteins

Furey, Nature Reviews Genetics 13, 840-852 (2012)



Amplify, if few cells

HI-C & DNA folding



- Cells cross-linked with formaldehyde; covalent links between spatially adjacent chromatin segments
- Chromatin digested with a restriction enzyme (here, HindIII)
- Resulting stickyends are filled in with nucleotides one biotinylated (purple dot)
- DNA is purified and sheared.
- Biotinylated junctions are isolated with streptavidin beads and identified by paired-end sequencing.
- Erez Lieberman-Aiden et al. Science 326, 289 (2009)

HI-C & DNA folding



- Hi-C produces genome-wide contact matrix
- Each pixel represents all interactions between 1-Mb loci
- Intensity as total number of reads

Erez Lieberman-Aiden et al. Science 326, 289 (2009)