



RNA-Seq Data Analysis 7-8th June, 2018

Organised and delivered by Bioinformatics Core at WHG:
Helen Lockstone M.Sc.
Ben Wright PhD
Eshita Sharma PhD
Santiago Revale, M.Sc.





What we do when we do RNAseq?



- What it is?
- Scope of RNAseq
- Usual approaches for RNAseq library preparation?
- Considerations for RNAseq experiments
- General methods for RNAseq data analysis.

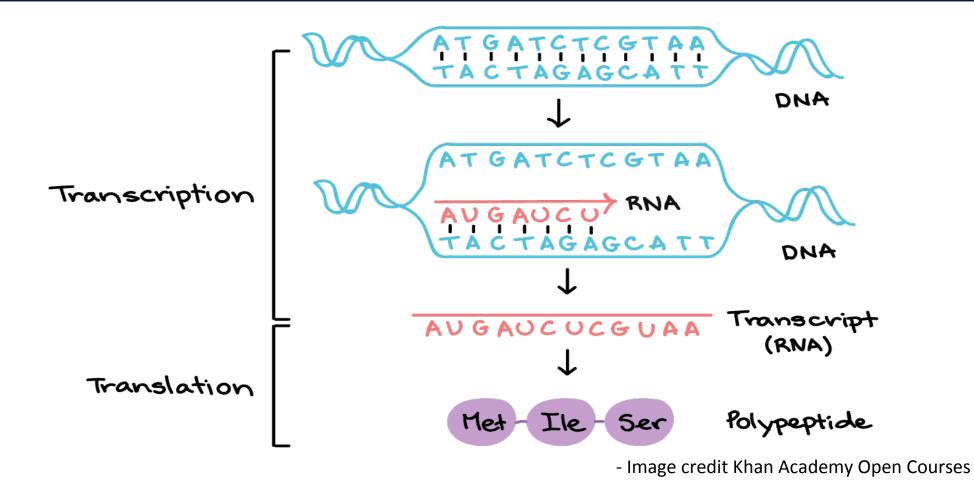
Eshita Sharma, eshita.sharma@well.ox.ac.uk Research Associate in Funtional Genomics, Bioinformatics Core





What it is? RNA - Mid-point of the information cascade





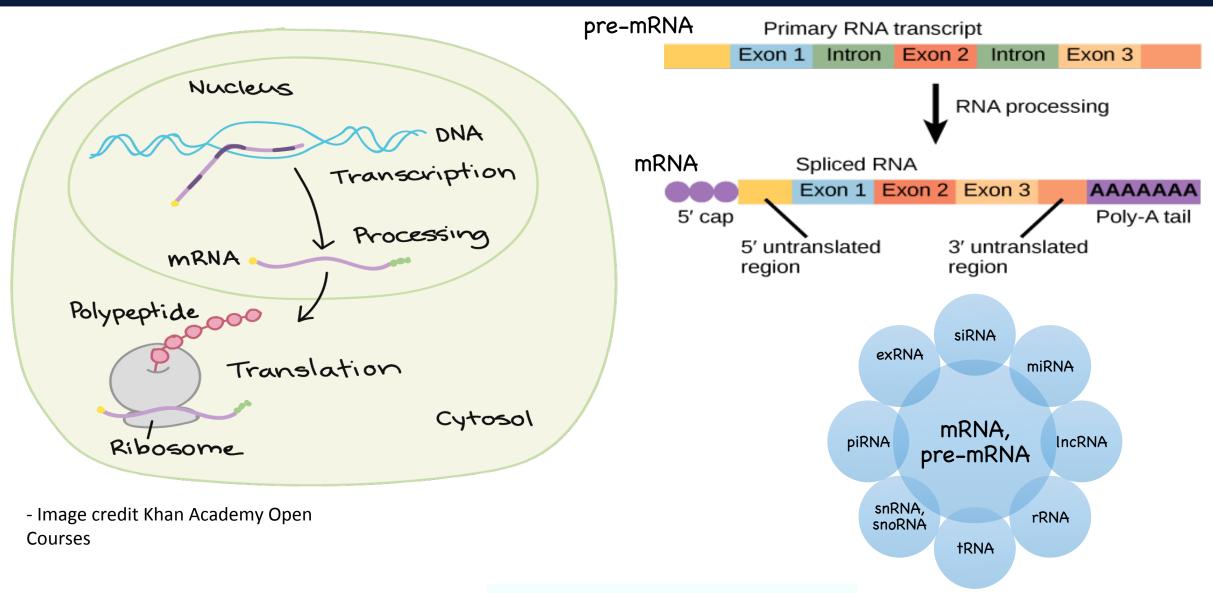
We identify the mRNA molecule and extrapolate the knowledge to say something about the proteins and DNA





The RNA repertoire or Transcriptome sum total of all RNA molecules expressed from the Genome





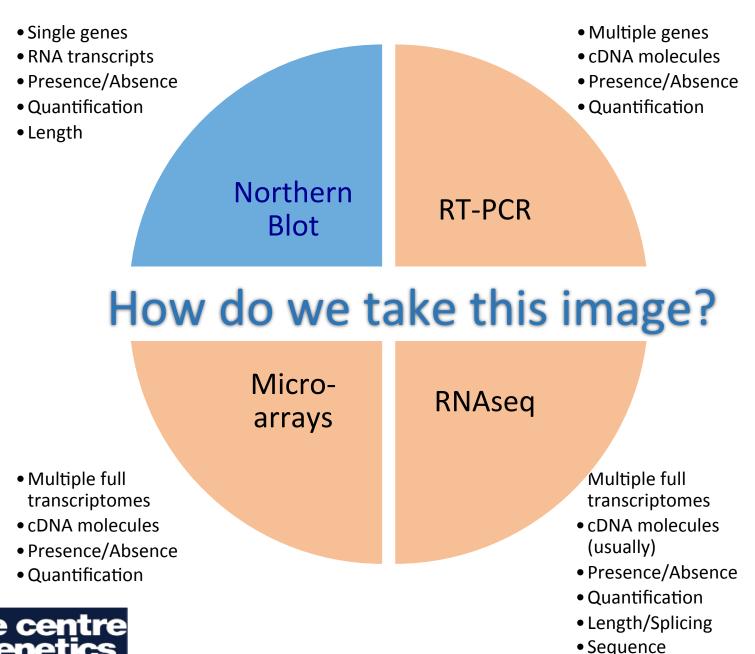


RNA repertoire is dynamic! It varies in time and space.



RNAseq is a method for Transcriptome profiling Image of the transcribed genome at any point of time!









Scope of RNAseq It's always about the goals!



At RNA transcript level, it provides the ability to:

- √ look at alternative gene spliced transcripts,
- ✓ post-transcriptional modifications,
- ✓ gene fusion,
- √ mutations/SNPs,
- ✓ changes in gene expression.

Can look at different populations of RNA to include:

- √ total RNA,
- √ mRNA,
- ✓ small RNA (miRNA, tRNA, ribosomal profiling, etc.)

Can be used to:

- √ determine exon/intron boundaries,
- ✓ verify or amend previously annotated 5' and 3' gene boundaries.





Common main goals



- Catalog all species of transcripts, e.g. messengers, non-coding, small, etc.
- Determine the transcriptional structure of genes, in terms of their starting sites, 5' and 3' ends, splicing patterns and other post-transcriptional modifications.
- Quantify the changes in the expression levels of each transcript during development and/or in different conditions.





Usual approaches for RNAseq library preparation

samples/lane HiSeq4000

lane HiSeq4000

• •	_	1 51	. '	
polyA enriched mRNAseq	Ribodepleted Total RNAseq	SMARTer/Smartseq2 mRNAseq	Small RNAseq <200nt ncRNAseq	3' mRNAseq mRNAseq
.	+	+	.	
Oligo dT ₂₅ selection	Biotinylated ribosomal RNA probes	Oligo dT primer used for Reverse transcription	3' adapter ligation	Oligo dT ₂₅ primed RT (First strand synthesis)
#	+	Ţ.		₽
Fragmentation of mRNA	Bind ribosomal RNA	Template switching by RT	5' adapter ligation	Removal of RNA template
4st I Sand I I S				
1 st strand → 2 nd strand → cDNA synthesis of fragments	extracted with streptavidin beads	PCR pre-amplification of full-length cDNA	1 st strand cDNA synthesis	Random priming and second strand synthesis
Adapter ligation → PCR amplification	Fragmentation → cDNA synthesis → adapters → PCR	Tn5 transposase tagmentation & library prep.	PCR enrichment → Size selection	Bead purification of tagged cDNA library -> PCR
	•	Propri	Facus on 21 25nt miDNA	Ţ.
Mature mRNA	Mature mRNA, nascent RNA,non-coding transcripts	Full-length mature mRNA	Focus on 21-25nt miRNA/ siRNA involved in gene regulation	200-300bp insert libraries of 3' ends of mRNA
•	•	•	Size specificity,	•
Directionality	Works with low-quality RNA, e.g. FFPE samples	Pre-amplification of low- input RNA	Directional, Low-input and Low depth	Lower sequencing depth; poor-quality RNA works
	<u> </u>	Poguiros good guality		
Requires good quality total RNA	Requires high sequencing depth	Requires good quality total RNA; no directionality	Requires good quality total RNA	Mainly useful for expression quantification
	•	•	1 up total DNIA 20 c	0.5 mg 2 mg 40 as mg 1 a
100ng-1ug; 8-10 samples/	100ng to 1ug; 4-6 samples/lane HiSeg4000	10pg to 10ng; 10-30 samples/lane HiSeg4000	1ug total RNA; 20+ samples/lane HiSeq2500;	0.5 ng – 2 ug; 48 samples per lane HiSeq4000;

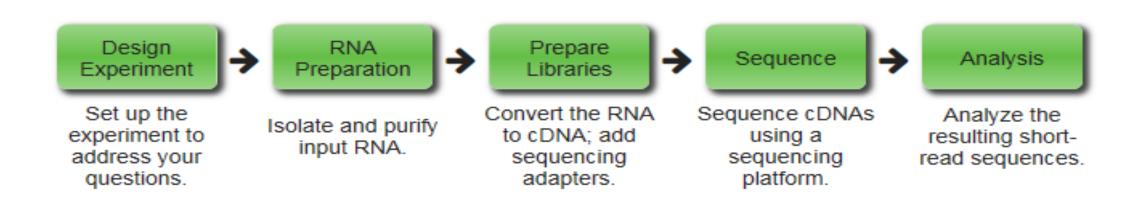
samples/lane HiSeq4000

50bp SE

50-75bp SE

Typical RNAseq experiment









Key considerations for Experimental Design

How is an experiment designed?



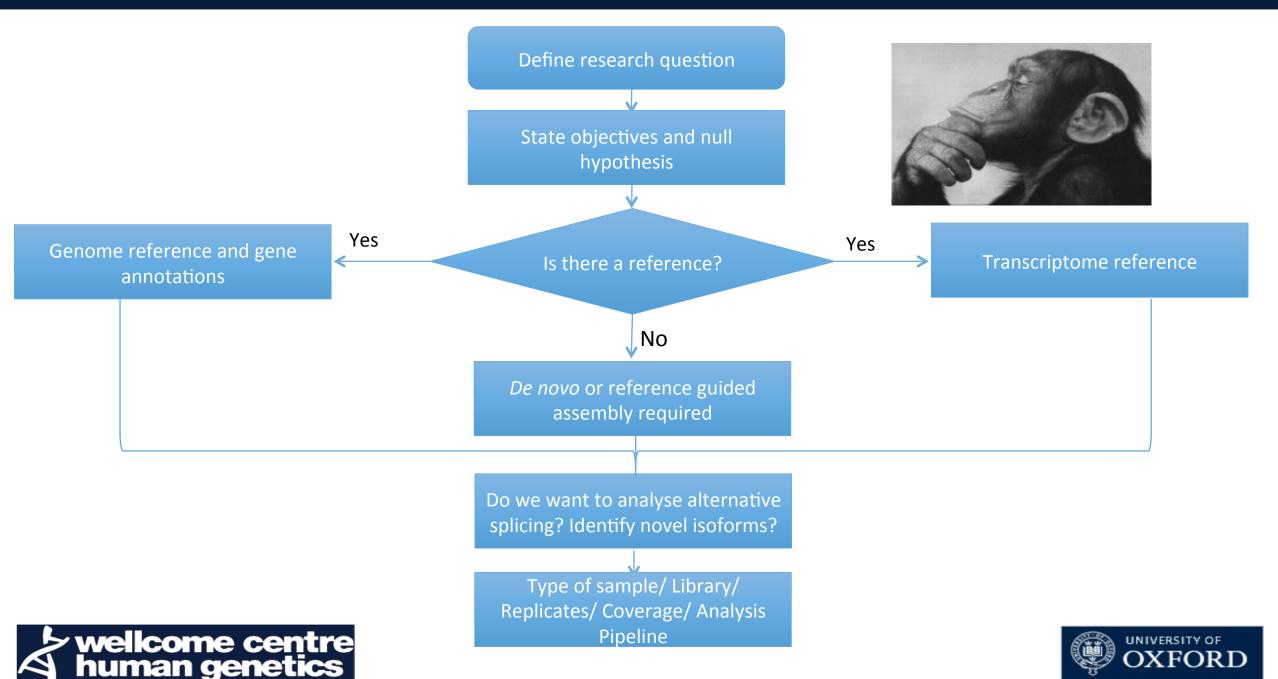
While a good design does not guarantee a successful experiment, a suitably bad design guarantees failure.





Before we begin...





Sample selection



SKIN

0031424 keratinization 2.9 \times 10⁻¹⁴ 0006955 immune response 3.5 \times 10⁻¹³ 0031069 hair follicle morphogenesis 4.1 \times 10⁻⁷

LUNG

0030324 lung development 6.2 \times 10⁻¹⁶ 0006954 inflammatory response 2.1 \times 10⁻¹⁵ 0043330 response to exogenous dsRNA 6.2 \times 10⁻⁶

ADRENAL

0006700 C21-steroid hormone biosynthesis 4.6×10^{-8} 0017157 regulation of exocytosis 4.2×10^{-4} 0006584 catecholamine metabolism 1.4×10^{-3}

KIDNEY

0001822 kidney development 1.4×10^{-6} 0007588 excretion 1.3×10^{-3} 0001736 establishment of planar polarity 2.9×10^{-3}

MUSCLE

0006941 striated muscle contraction 7.7×10^{-11} 0005977 glycogen metabolism 1.8×10^{-9} 0045445 myoblast differentiation 8.0×10^{-7}

TESTIS

0007059 chromosome segregation 9.1×10^{-15} 0007276 gametogenesis 8.1×10^{-4} 0006349 imprinting 1.5×10^{-3}

BRAIN

0007268 synaptic transmission 8.9 \times 10⁻⁴¹ 0016358 dendrite morphogenesis 1.2 \times 10⁻¹⁰ 7.9 \times 10⁻⁶

THYMUS

0019882 antigen presentation 7.1 \times 10⁻²¹ 0045059 positive thymic T cell selection 9.8 \times 10⁻⁸ 0045060 negative thymic T cell selection 2.6 \times 10⁻⁷

HEART

0006099 tricarboxylic acid cycle 2.5×10^{-15} 0045214 sarcomere organization 7.5×10^{-12} 0008016 regulation of heart contraction rate 8.3×10^{-7}

LIVER

0008203 cholesterol metabolism 2.6×10^{-8} 0007596 blood coagulation 2.0×10^{-7} 0000050 urea cycle 5.0×10^{-5}

SPLEEN

0050766 positive regulation of phagocytosis 4.5×10^{-9} 0030183 B cell differentiation 1.5×10^{-7} 0030217 T cell differentiation 2.6×10^{-7}

INTESTINE

0006955 immune response 7.0×10^{-13} 0007586 digestion 9.3×10^{-5} 0050892 intestinal absorption 4.6×10^{-4}

OVARY

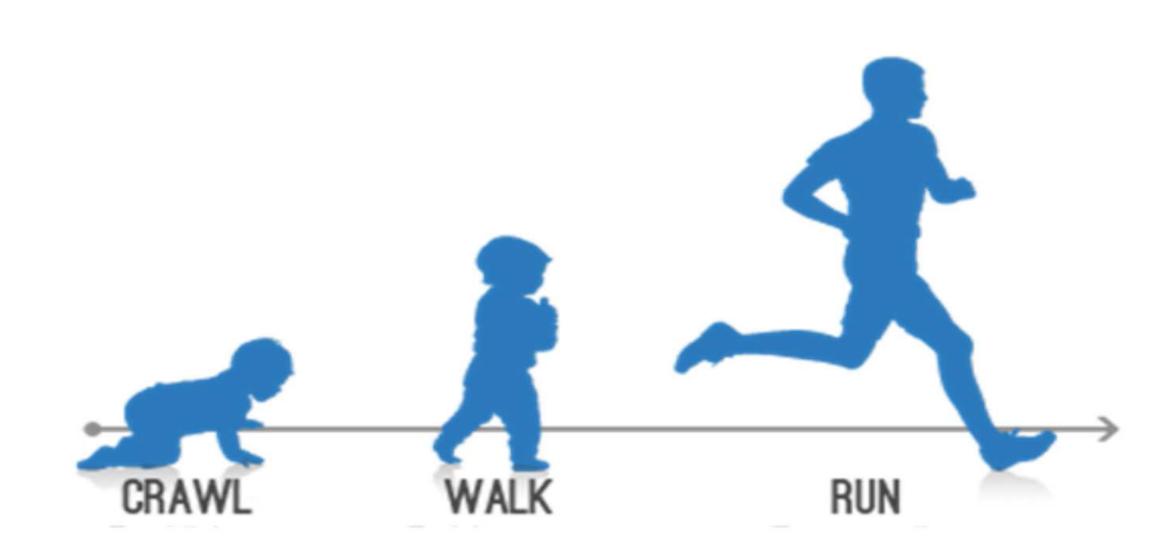
0007059 chromosome segregation 1.0×10^{-12} 0007276 gametogenesis 8.6×10^{-8} 0006349 imprinting 3.5×10^{-5}





Time of sampling



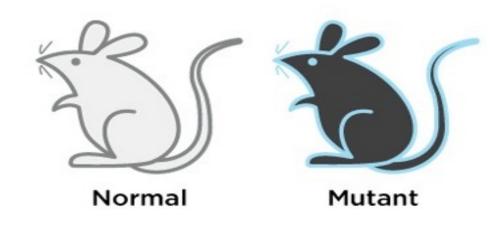


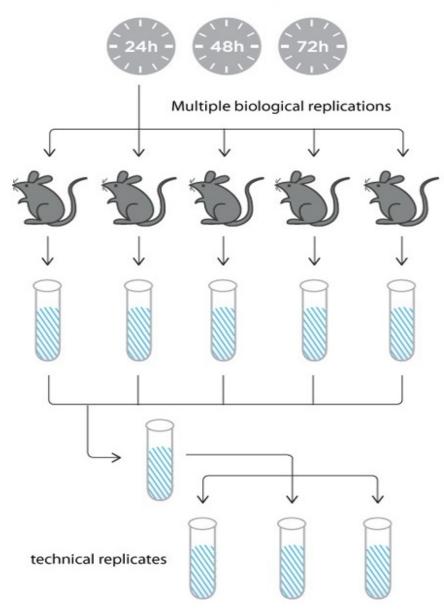




Replicates (technical / biological)









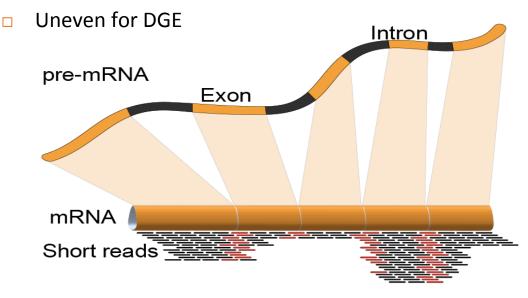


Coverage

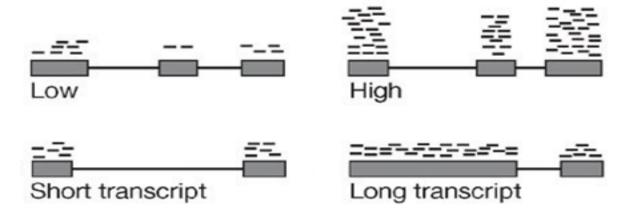
Short reads



Even for annotation
 pre-mRNA
 Exon



Target transcript properties (low abundant vs high abundant transcripts)



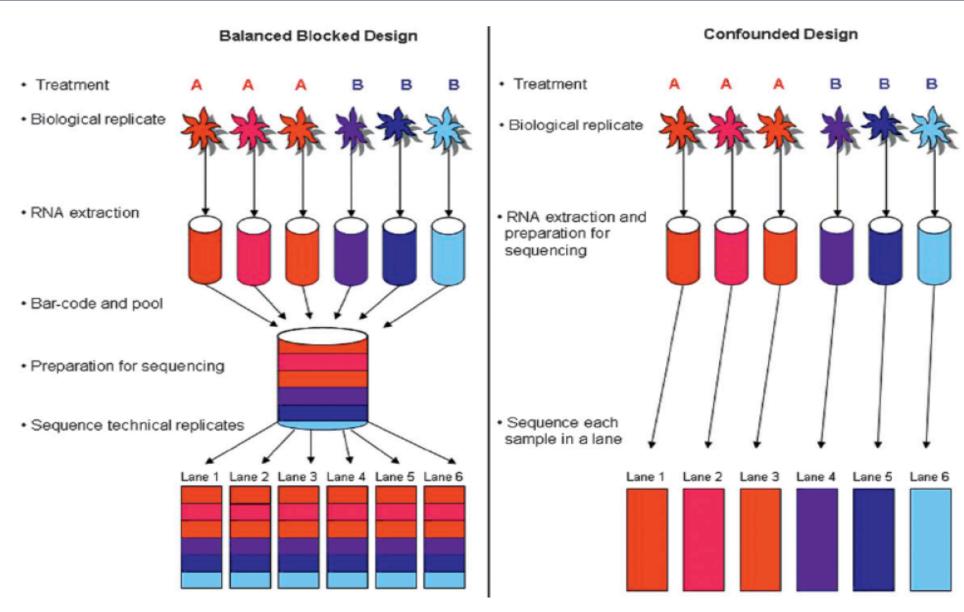
- Allele might not be detected (not in the genome/not being expressed)
- Estimate expression of each allele





Randomization and Blocking









Data management & Downstream analysis and interpretation of the data



- ✓ Several Gigabytes (70-75 Gb Avg)
- ✓ Different layers of interpretations have to be considered (e.g. biological, clinical, regulatory functions, etc.)







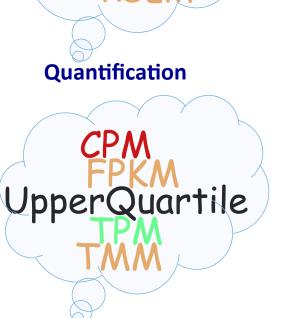
Subjectivity of the analysis

Oxford Genomics CENTRE

- ✓ Multitude of algorithms and pipelines available.
- ✓ Most approaches correct, but have to be tailored to the needs of the investigators in order to better capture the desired effect.







Count normalization

Salmon

featureCounts



Alternative splicing/ Isoform level analysis

voom Limma baySeq DESeq2 NOIseq edgeR

Differential expression analysis





Adapter trimming





Oases

Data Analysis

Typical RNA-seq analysis pipeline



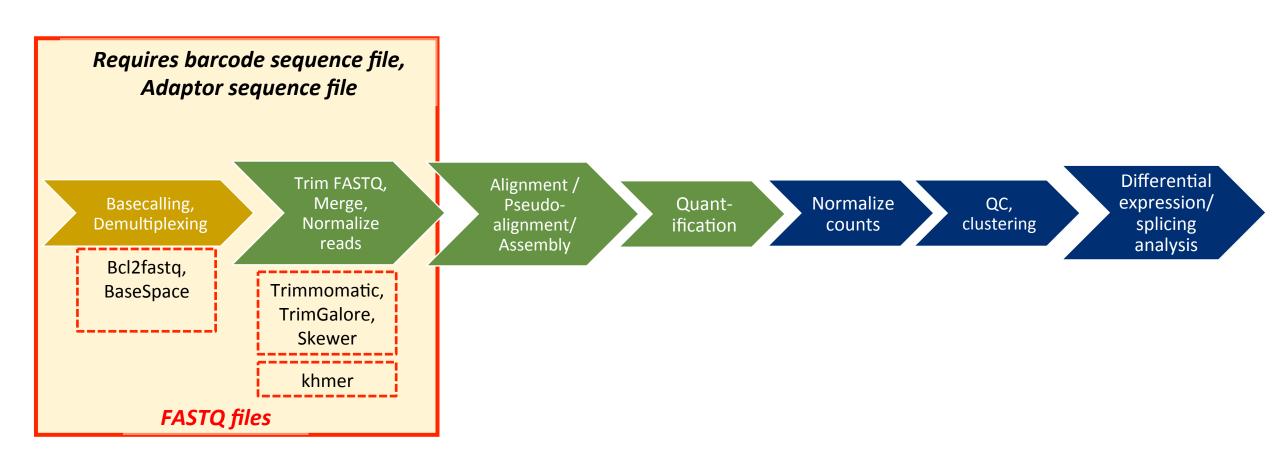
- 1. Demultiplex, filter, and trim sequencing reads.
- 2. Normalize sequencing reads (if performing *de novo* assembly)
- 3. de novo assembly of transcripts (if ref. genome is not available)
- 4. Map sequencing reads to reference genome or transcriptome
- 5. Annotate transcripts assembled or to which reads have been mapped
- 6. "Count" mapped reads to estimate transcript abundance
- 7. Perform statistical analysis to identify differential expression (or differential splicing) among samples or treatments
- 8. Perform multivariate statistical analysis/visualization to assess transcriptome-wide differences among samples





Read processing









Read processing



For de novo assembly

Software	De- mulitplexing	Adaptor Trimming	Quality Filtering/ Trimming	K-mer Filtering	K-mer Normalization	
ASTX-Toolkit	~	~	~			
Goby	~	~				
chmer				~	~	
NGS_backbone		~	~			
Stacks	~	~	✓	~	~	
rimmomatic		~	✓			
piopieces	~	_	~			





Read processing, alignment/assembly+alignment, quantification



Genome; type of analysis – Assembly or just Alignment

Basecalling, Demultiplexing Trim FASTQ, Merge, Normalize reads Alignment /
Pseudoalignment/
Assembly

Quantification Normalize counts

QC, clustering Differential expression/ splicing analysis

HISAT2, STAR, Subread

Salmon, Kallisto

Trinity, StringTie, Cufflinks2, STAR

Alignment bam files or direct quantification in alignment free approach

FASTQ file if de novo assembly, GTF file if reference-guided assembly

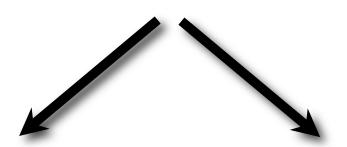




de novo assembly or reference mapping?



When to use each?



de novo

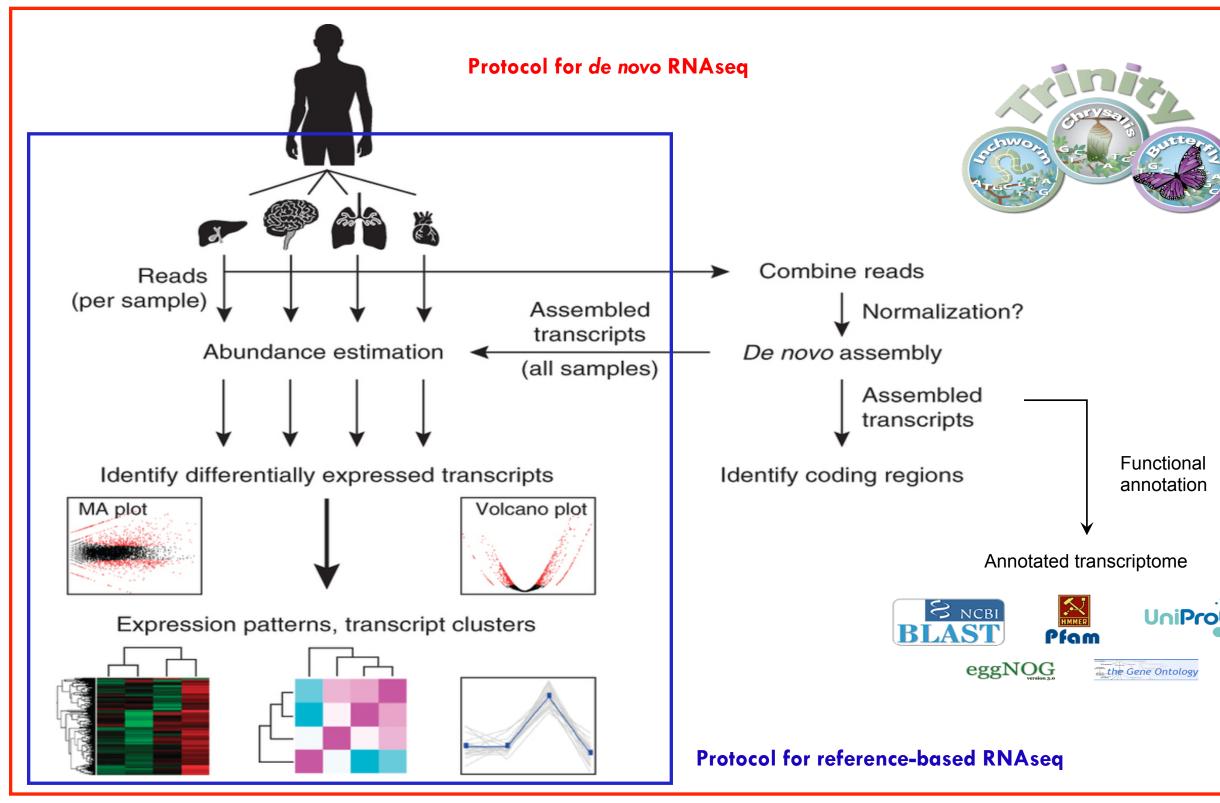
(do not know the transcriptome) (main goal is to discover NOT to quantify)

reference

(do know the transcriptome)
(main goal is to quantify NOT to discover)

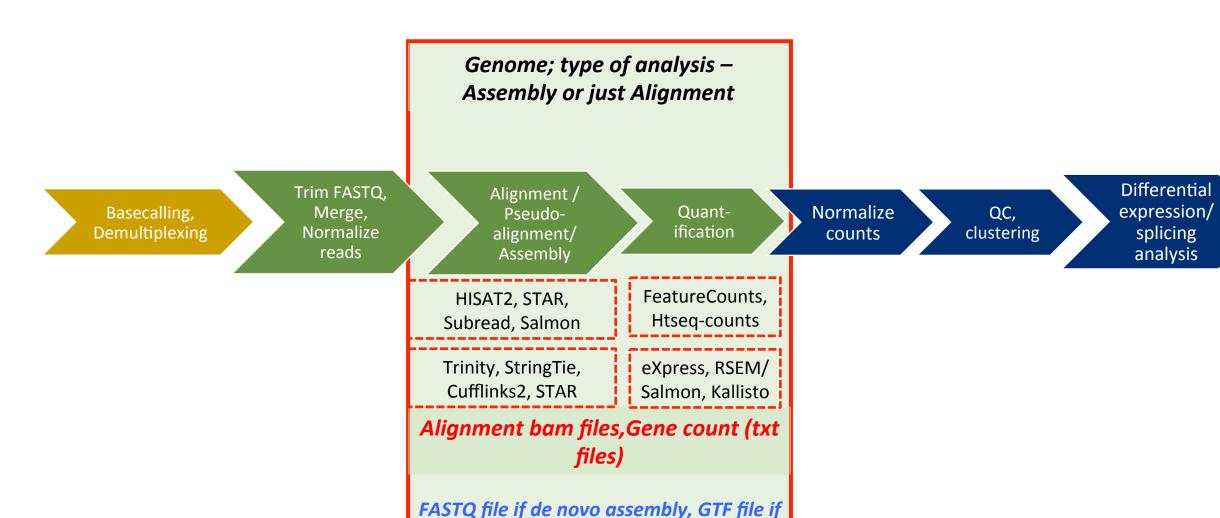






Read processing, alignment/assembly+alignment, quantification





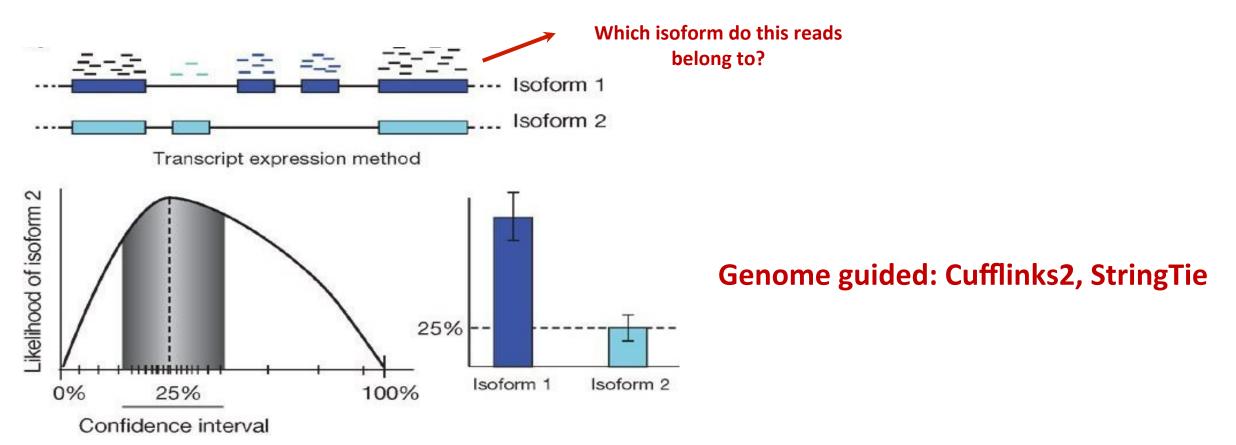
reference-guided assembly



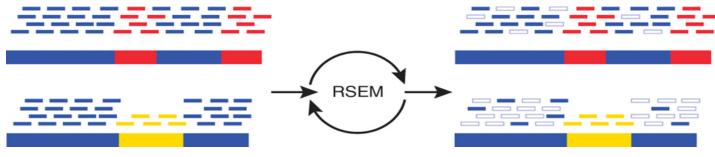


Expression quantification





Trancriptome guided: RSEM, eXpress, Salmon, Kallisto



Normalization, Expression quantification



Basecalling, Demultiplexing Trim FASTQ, Merge, Normalize sequencing data

Alignment /
Pseudoalignment/
Assembly

Quantification Samplegroups; Pairing; Candidate genes; Other technical variables

Normalize counts

QC, clustering

Differential expression/ splicing analysis

CPM, RPKM, FPKM,
TPM, Upper
Quartile, TMM,
SizeFactor

Empirical analysis, Hierarchical clustering, PCA, MDS

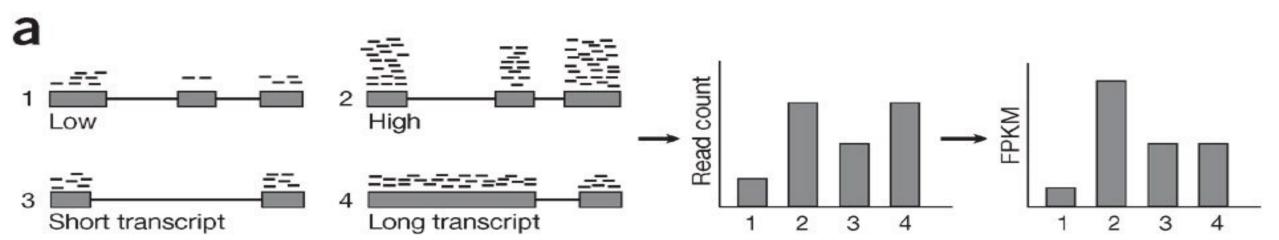
normalized counts (txt file)
PCA/tSNE plot, heatmap,
QC metrics





Count normalization





Influence of length: Counts are proportional to the transcript length times the mRNA expression level.

Influence of sequencing depth: The higher sequencing depth, the higher counts.

"Gene counts" should be corrected in order to minimize these biases: normalization.

Statistical model should take into account "length" and "sequencing depth".

Count normalization



	Counts	СРМ	RPKM/FPKM	TPM
Value	Integer	Fraction	Fraction	Fraction
Depth-bias	X	✓	✓	✓
Length-bias	X	X	✓	✓
Compare same genes across samples	X		(but may have bias)	
Compare different genes in sample	X	X	✓	
Compare different genes across samples and across experiments	X	X	X	
Can be used for barplots/ boxplots of single genes	X		(but may have bias)	
Can be used for heatmaps with multiple genes (log transformed)	X	(as long as we don't compare the colour of different genes)	(but may have bias)	

Differential expression analysis



Basecalling, Demultiplexing Trim FASTQ, Merge, Normalize sequencing data

Alignment /
Pseudoalignment/
Assembly

Quantification Normalize counts

QC, clustering

Types of comparisons

Differential expression/ splicing analysis

edgeR, DESeq2, CuffDiff2, Limma/voom

rMATS, DEXSeq, ballgown

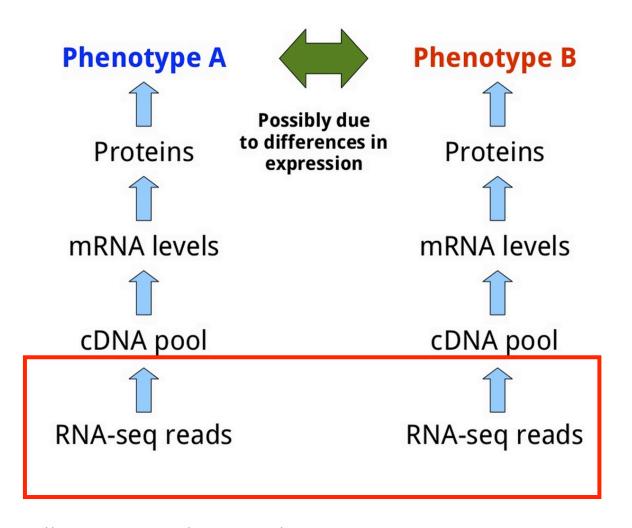
Gene-lists,
Differential
splicing; Pathway
analysis (txt files)

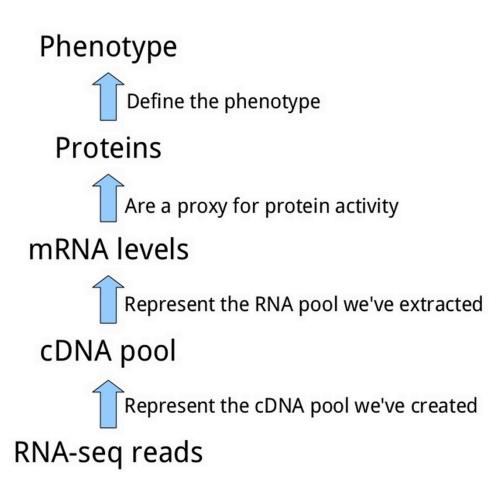




Our assumptions and comparison





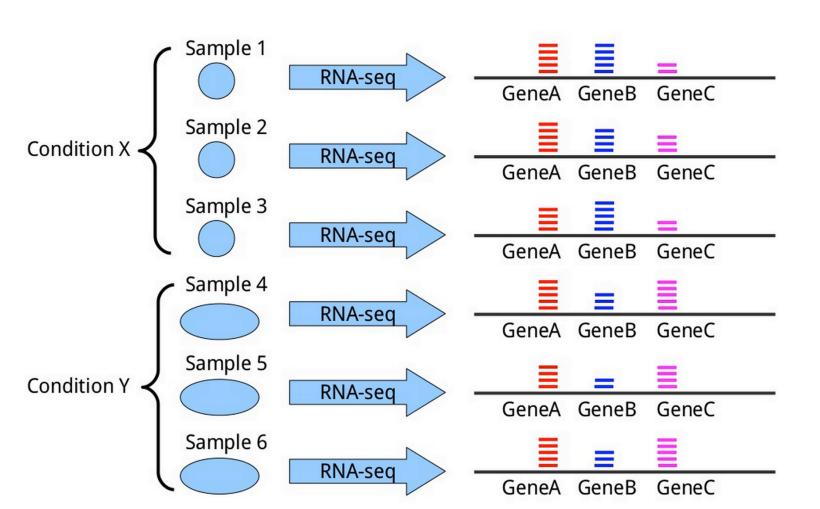


http://www.slideshare.net/joachimjacob/1rna-seqpart1working-tothegoal?related=2





Statistical testing for Differential expression



http://www.slideshare.net/joachimjacob/1rna-seqpart1working-tothegoal?related=2

Read in raw count data

Remove genes with lowexpression (<10 reads per sample across group)

Normalize subset count tables using size factors (e.g. TMM normalization in edgeR)

Unsupervised clustering to identify technical effects and biological effects

Create design matrix with comparison of interest and technical/biological variability

Fit normalized expression matrix to linear models to identify coefficients for each gene

Identify DE genes with predefined statistical criteria (~FDR < 0.05)

What we do when we do RNAseq?



- What it is?
- Scope of RNAseq
- Usual approaches for RNAseq library preparation?
- Considerations for RNAseq experiments
- General methods for RNAseq data analysis.





References



- Cresko Lab, University of Oregon. RNA-seqlopedia. http://rnaseq.uoregon.edu/
- Garber M, Grabherr MG, Guttman M, Trapnell C. Computational methods for transcriptome annotation and quantification using RNA-seq. Nat Met. 2011; 8: 469– 477.
- Grabherr MG, Haas BJ, Yassour M, Levin JZ, et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nature Biotechnology. 2011; 29: 644–652.
- Haas BJ, Papanicolaou A, Yassour M, et al. De novo transcript sequence reconstruction from RNA-seq using the Trinity platform for reference generation and analysis. Nature Protocols. 2013; 8: 1494–1512.
- Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. Nat Rev Genet. 2009 Jan; 10(1): 57–63.
- List of RNAseq bioinformatic tools.
 http://en.wikipedia.org/wiki/List_of_RNA-Seq_bioinformatics_tools
- https://f1000research.com/articles/5-1408/





Thank You!





Eshita.sharma@well.ox.ac.uk



