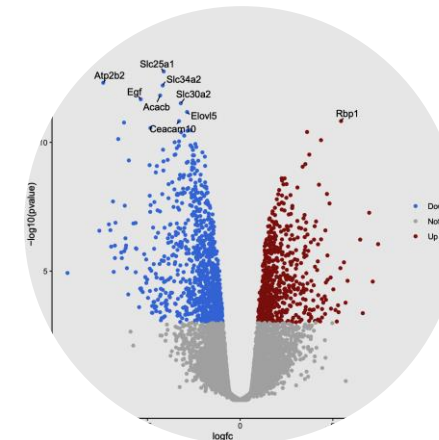
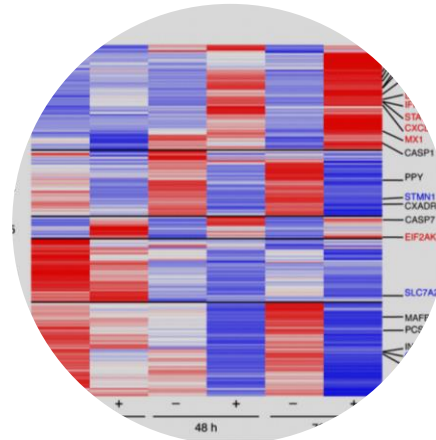
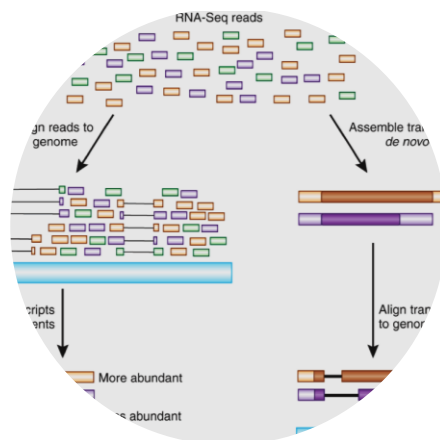


PMB2025

PATHOGEN MULTIOMICS AND BIOINFORMATICS

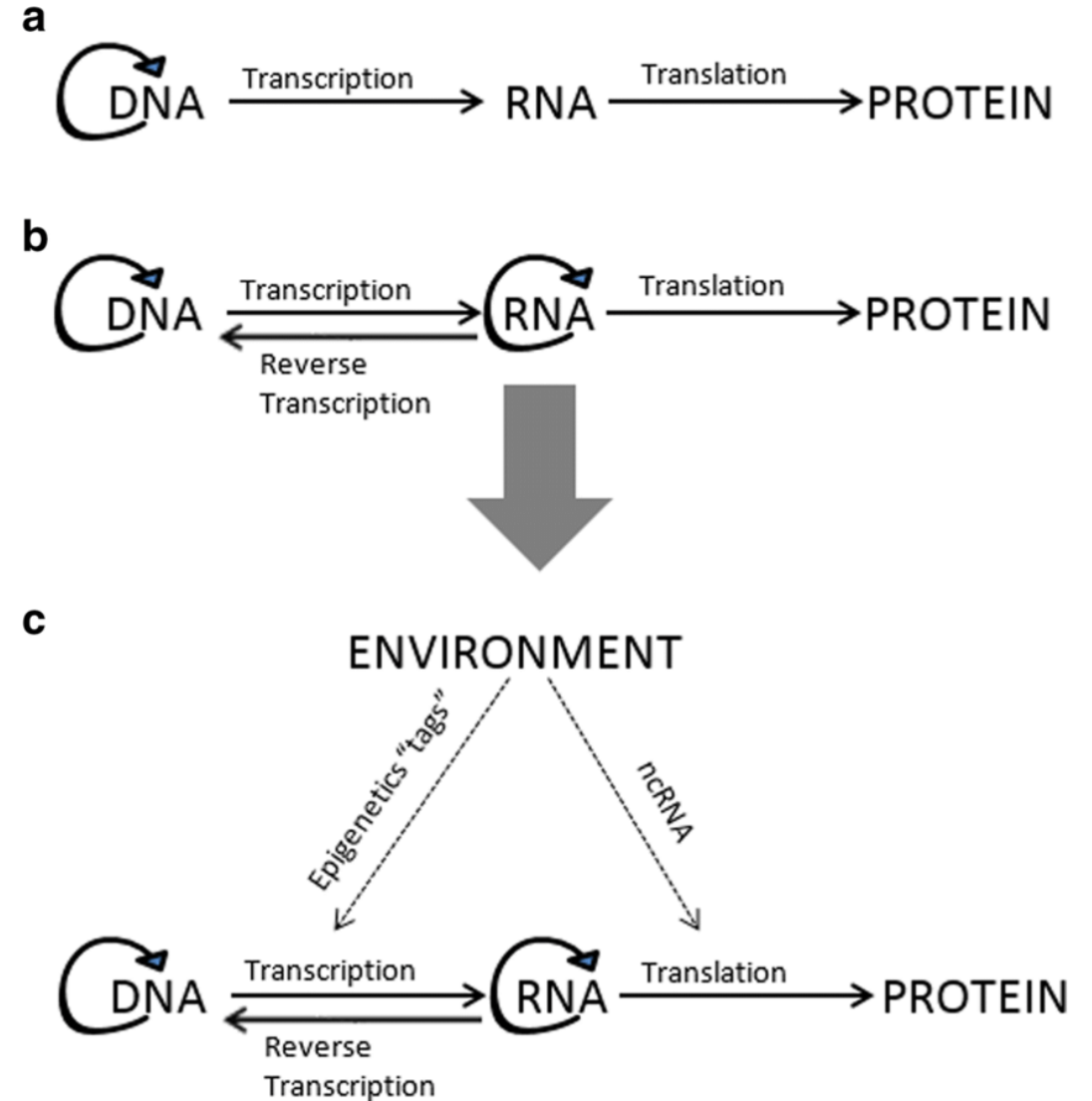
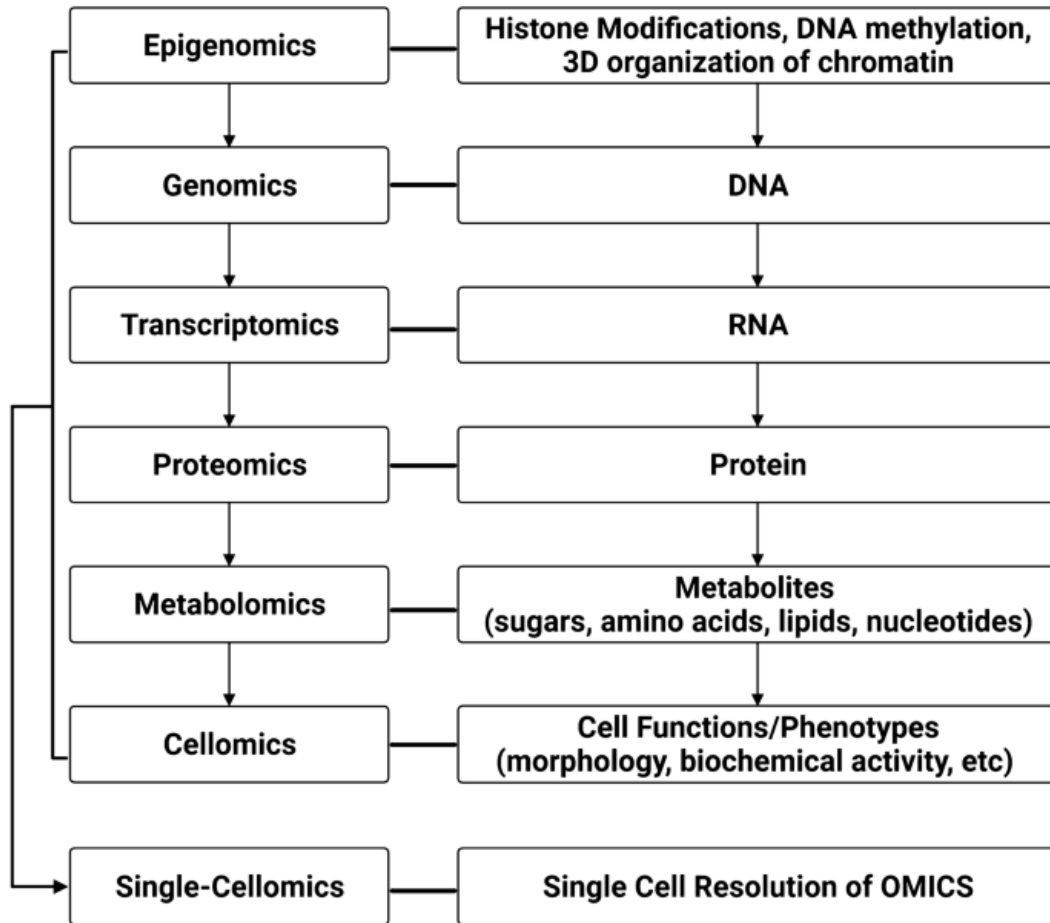
Recife PE 2025

Module 6: RNA-Seq and Transcriptomics



Omics Cascade

The OMICS cascade:



Why using RNA-Seq?

Use of new sequencing technologies to capture and study the transcriptome

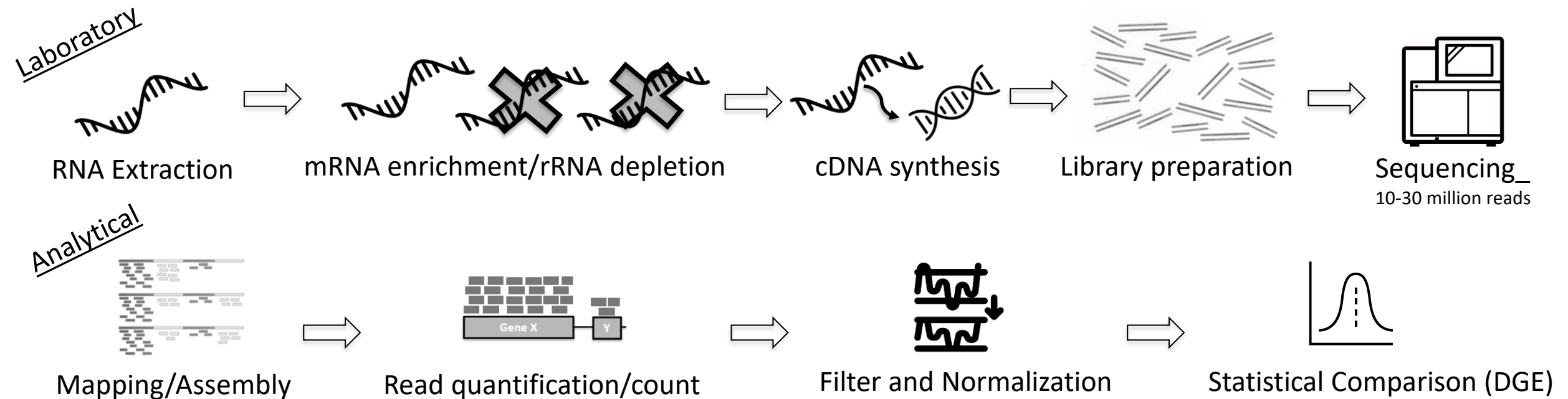
- **Identify novel transcripts**
- **Exon/transcript boundaries**
- **Splice junctions/alternative splicing**
- **Measure transcript abundance**
- **Gene expression differences across multiple samples (i.e. differential expression)**

RNA-Seq: General Workflow

What is RNA-Seq?

RNA-Seq consists of a method to analyze the transcriptomics of thousands of features in a single assay and, hence, evaluate and compare gene expression in a genome-wide manner.

Two main stages:



RNA-Seq vs cDNA/EST Seq vs Microarrays

Technology	Tiling microarray	cDNA or EST sequencing	RNA-Seq
<i>Technology specifications</i>			
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
<i>Application</i>			
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
<i>Practical issues</i>			
Required amount of RNA	High	High	Low
Cost for mapping transcriptomes of large genomes	High	High	Relatively low

Wang *et al* 2009

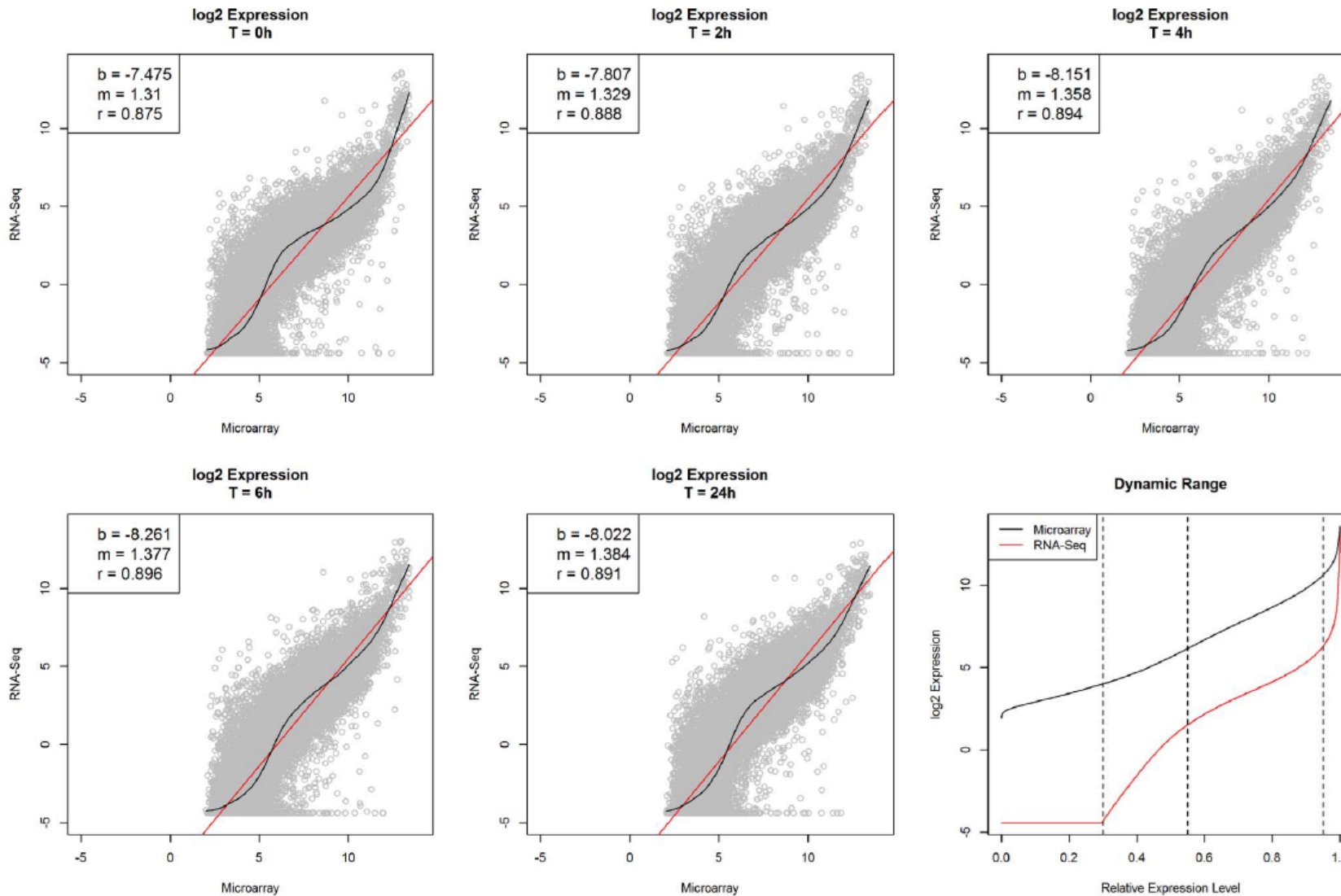
Ability to detect novel transcripts

Wider dynamic range

Higher specificity and sensitivity

Simple detection of rare and low-abundance transcripts

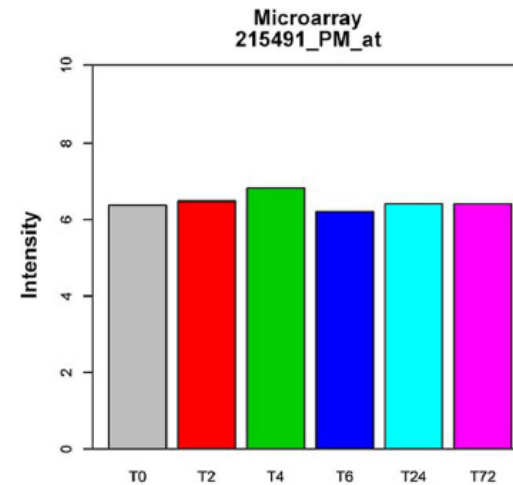
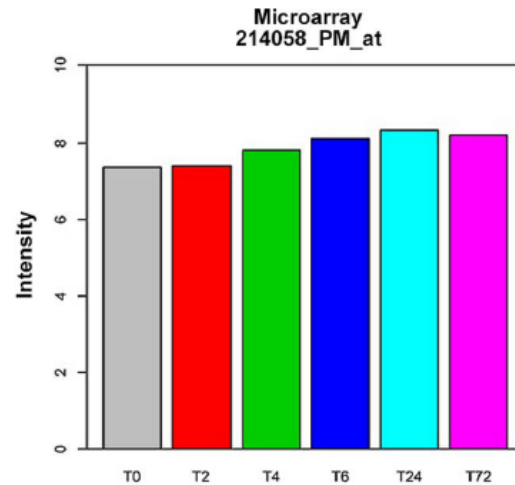
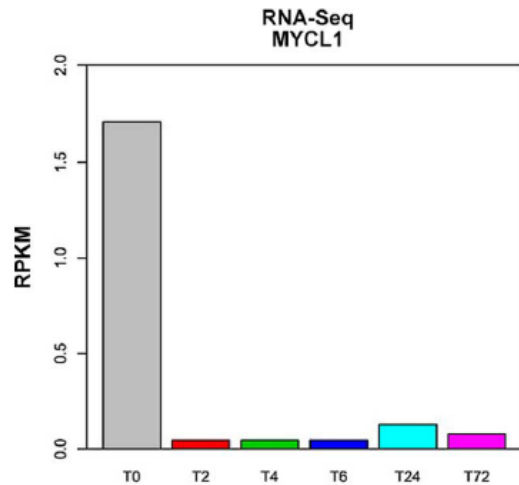
RNA-Seq vs Microarrays



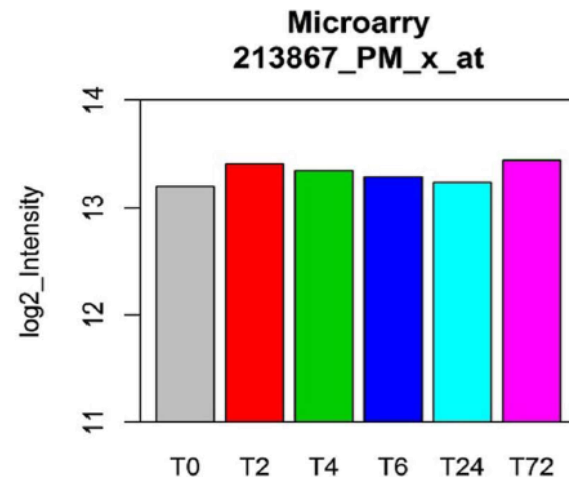
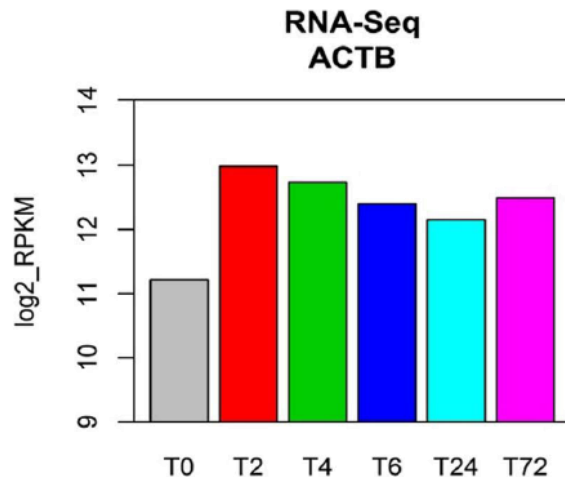
The plots show that the overall dynamic range of the 18,306 common genes generated by the two platforms is much broader in RNA-Seq (2.66105) than in microarray (3.66103).

RNA-Seq vs Microarrays

RNA-Seq is able to detect subtle changes to the level of genes with low expression levels whereas microarrays are not



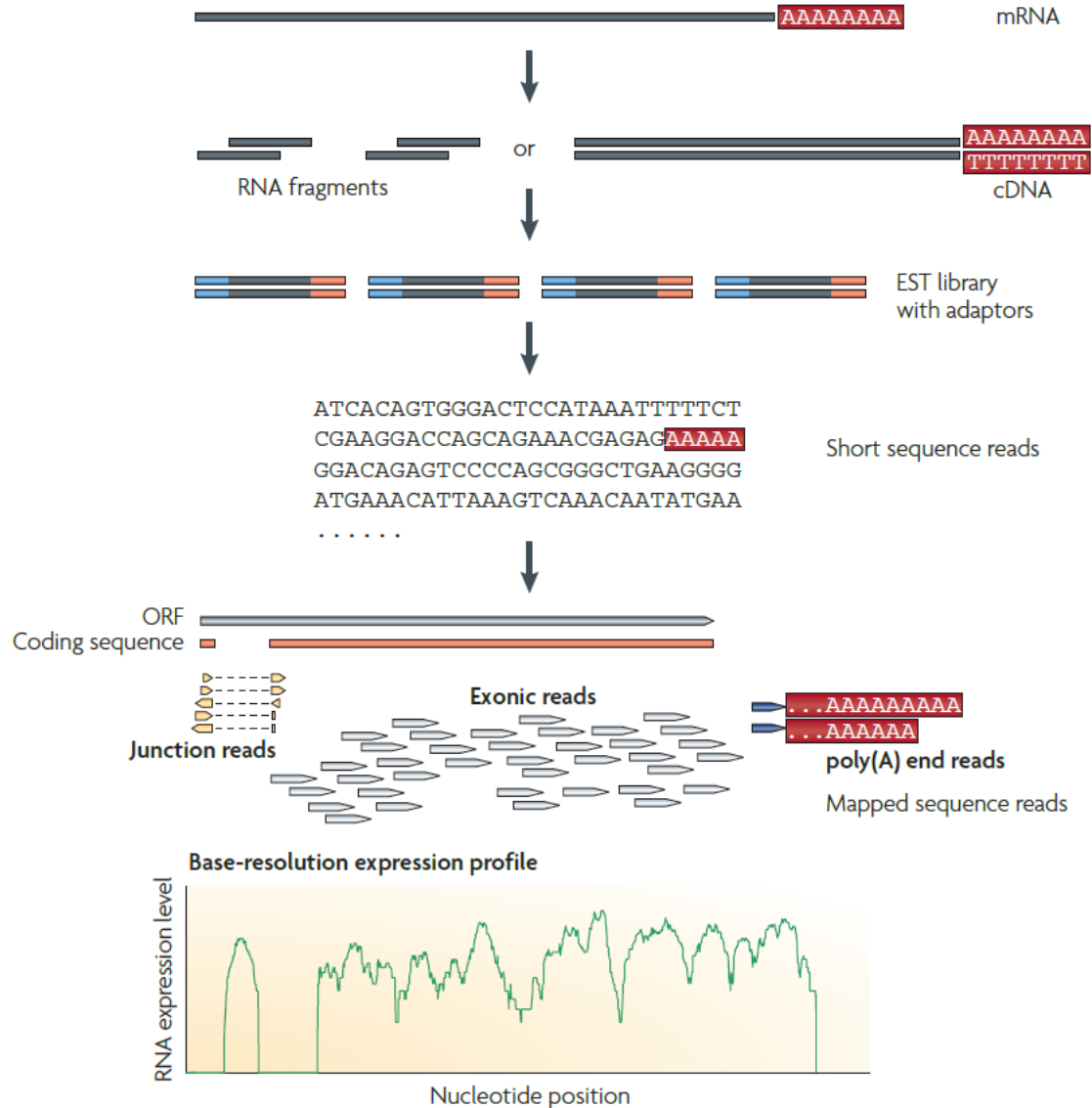
Zhao et al 2014



Zhao et al 2014

... Similarly RNA-Seq is able to detect expression level changes to highly expressed genes and microarrays are not (saturation).

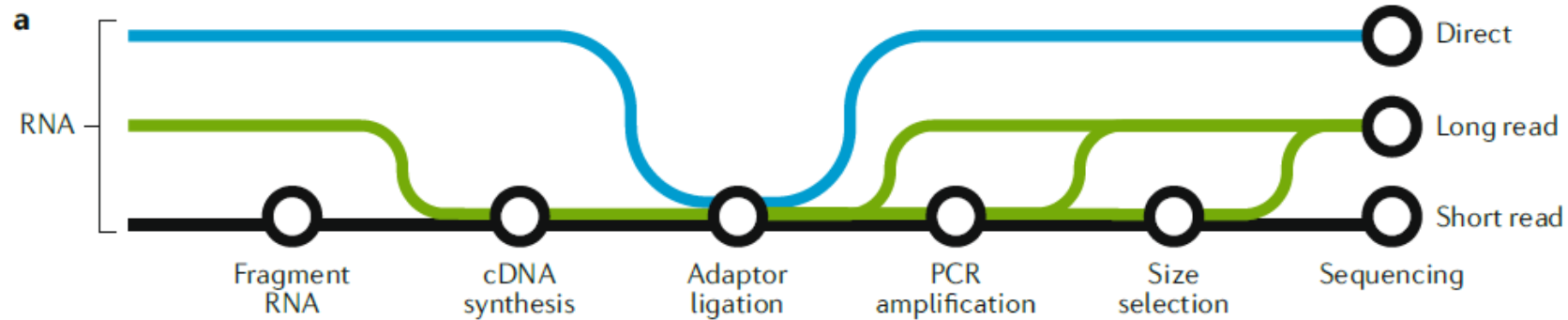
Library Preparation



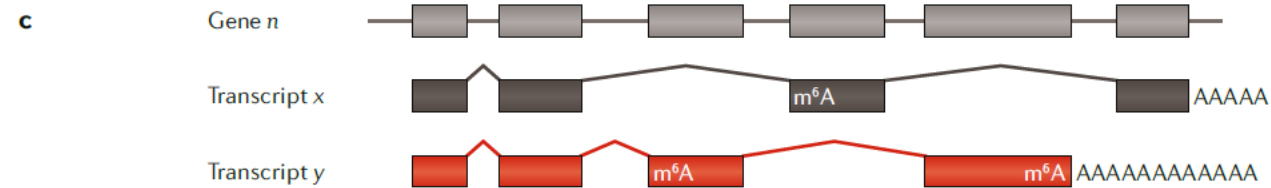
Aspects and factors to consider:

- RNA Source: Total RNA, mRNA, depletion of rRNA?
- Strand specific?
- Replicates?
 - Technical (multiple libraries from the same sample)
 - Biological (multiple samples from the same condition)
- Which platform?
- Multiple samples/multiplexing

Library Preparation: comparison between technologies and limitations



All long-read and short-read approaches require adapter ligation



Short-read cDNA	Ambiguous to exon	
	Unambiguous to exon	
	Ambiguous to isoform	
	Unambiguous to isoform	

Short-reads can be ambiguously mapped to different isoforms

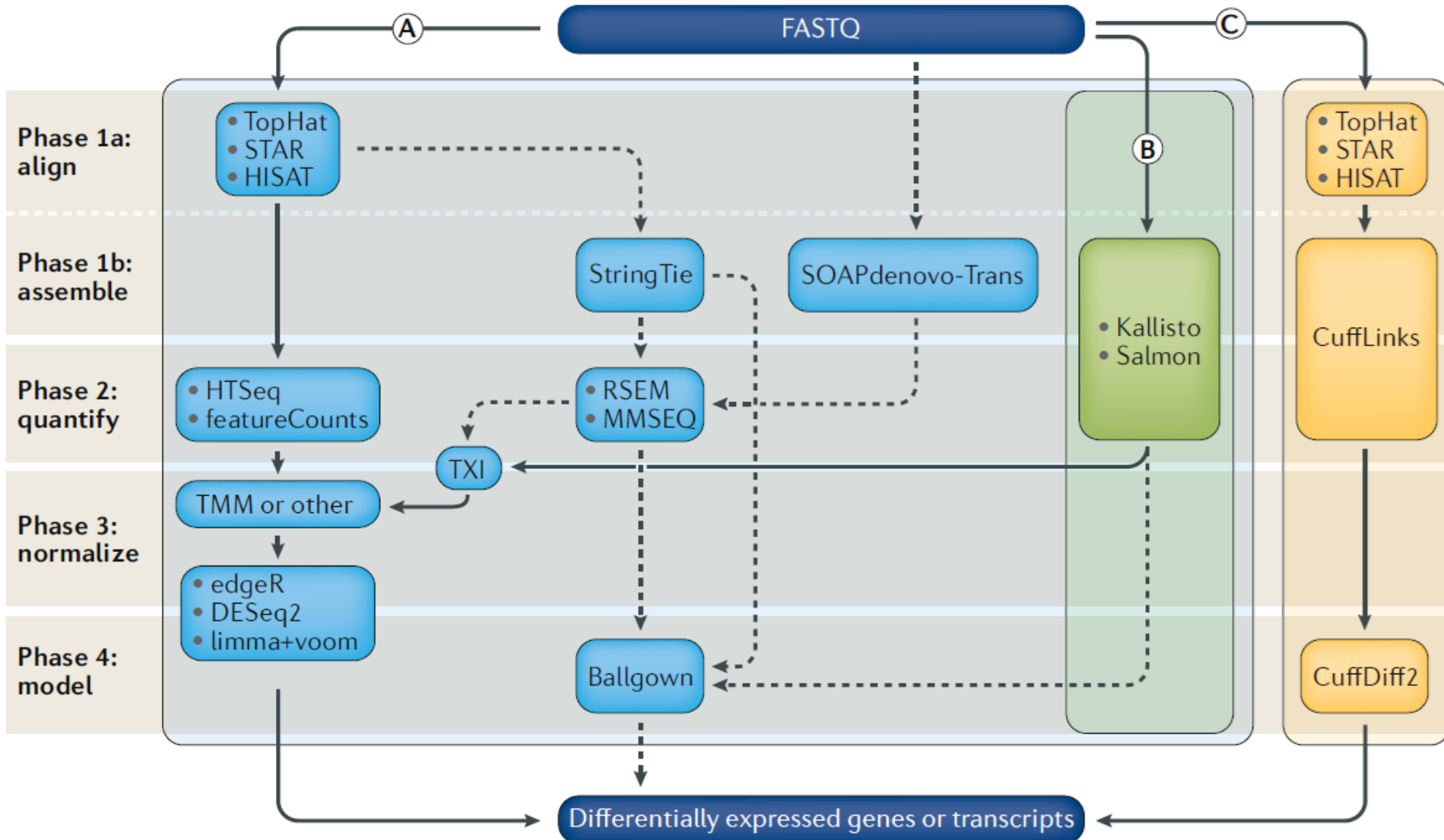
Long-read cDNA	Unambiguous to isoform	
----------------	------------------------	--

Direct RNA-Seq	Unambiguous to isoform	
----------------	------------------------	--

Library Preparation: comparison between technologies and limitations

Sequencing technology	Platform	Advantages	Disadvantages	Key applications
Short-read cDNA	Illumina, Ion Torrent	<ul style="list-style-type: none"> • Technology features very high throughput: currently 100–1,000 times more reads per run than long-read platforms • Biases and error profiles are well understood (homopolymers are still an issue for Ion Torrent) • A huge catalogue of compatible methods and computational workflows are available • Analysis works with degraded RNA 	<ul style="list-style-type: none"> • Sample preparation includes reverse transcription, PCR and size selection adding biases to all methods • Isoform detection and quantitation can be limited • Transcript discovery methods require a de novo transcriptome alignment and/or assembly step 	Nearly all RNA-seq methods have been developed for short-read cDNA sequencing: DGE, WTA, small RNA, single-cell, spatialomics, nascent RNA, translatoe, structural and RNA–protein interaction analysis, and more are all possible
Long-read cDNA	PacBio, ONT	<ul style="list-style-type: none"> • Long reads of 1–50 kb capture many full-length transcripts • Computational methods for de novo transcriptome analysis are simplified 	<ul style="list-style-type: none"> • Technology features low-to-medium throughput: currently only 500,000 to 10 million reads per run • Sample preparation includes reverse transcription, PCR and size selection (for some protocols), adding biases to many methods • Degraded RNA analysis is not recommended 	Sequencing is particularly suited to isoform discovery, de novo transcriptome analysis, fusion transcript discovery, and MHC, HLA or other complex transcript analysis
Long-read RNA	ONT	<ul style="list-style-type: none"> • Long reads of 1–50 kb capture many full-length transcripts • Computational methods for de novo transcriptome analysis are simplified • Sample preparation does not require reverse transcription or PCR-reducing biases • RNA base modifications can be detected • Poly(A) tail lengths can be directly estimated from single-molecule sequencing 	<ul style="list-style-type: none"> • Technology features low throughput: currently only 500,000 to 1 million reads per run • Sample preparation and sequencing biases are not well understood • Degraded RNA analysis is not recommended 	<ul style="list-style-type: none"> • Sequencing is particularly suited to isoform discovery, de novo transcriptome analysis, fusion transcript discovery, and MHC, HLA or other complex transcript analysis • Ribonucleotide modifications can be detected

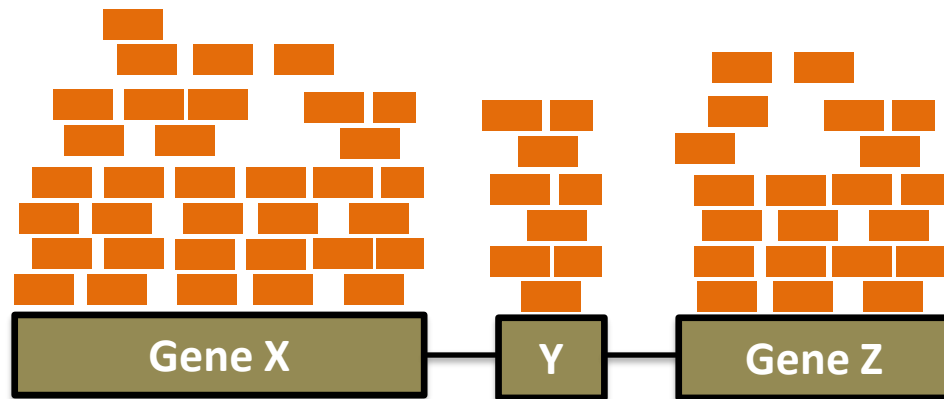
Analytical Pipeline Overview



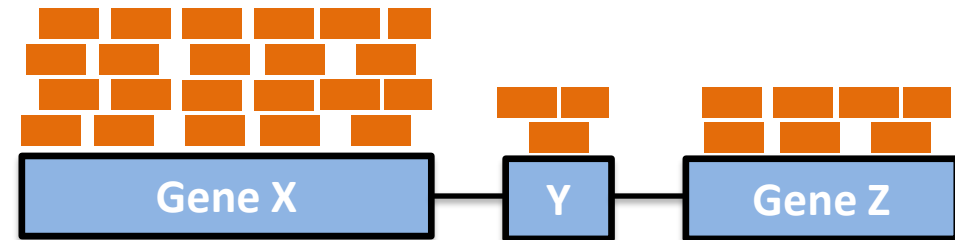
Searching for Differentially Expressed Genes

For differential expression the number of reads mapping to each gene (read count) is used to evaluate expression levels...

In which sample is Gene X overexpressed?



Sample A



Sample B

Raw counts cannot be used to evaluate or compare expression between samples! And within a sample?

Factors to consider:

- *Sequencing Depth*
- *Gene Length*
- *RNA Composition*

Normalizing Raw Read Counts: CPM/RPM, RPKM/FKPM and TPM

The answer: Normalization - this is required for differential expression analysis vizualization, etc.

Some Normalization Methods:

CPM/RPM – Counts/Reads per million

$$CPM = \frac{\text{No. reads mapped to gene} \times 10^6}{\text{Total number of mapped reads}}$$

RPKM/FPKM – Reads/fragments per kilobase million

$$\text{Scaling Factor (SF)} = \frac{\text{Total number of mapped reads}}{10^6}$$

$$RPM = \frac{\text{No. reads mapped to gene}}{SF}$$

$$RPMK = \frac{RPM}{\text{gene length(Kbp)}}$$

$$RPMK = \frac{\text{No. reads mapped to gene} \times 10^3 \times 10^6}{\text{Total number of mapped reads} \times \text{gene length(bp)}}$$

TPM – Transcripts per kilobase million

$$RPK = \frac{\text{No. reads mapped to gene}}{\text{gene length(Kbp)}}$$

$$\text{Scaling Factor (SF)} = \frac{\sum RPK}{10^6}$$

$$TPM = \frac{RPK}{SF}$$

Normalizing Raw Read Counts: Median of Ratios (DESeq2)

DESeq2 – Median of Ratios Method Normalization

Accounts for sequencing depth and RNA composition... but not gene length

1. Starting on raw counts, calculate the geometric mean for each gene across all sample – pseudo-reference;



2. Calculate the ratio of each sample to the pseudo-reference;



3. Calculate the normalization factor for each sample (**size factor**) by taking the median of all ratios;



4. Normalized counts are obtained by dividing the raw count of each gene by the normalization factor;

Anders and Huber Genome Biology 2010, 11:R106
<http://genomebiology.com/2010/11/10/R106>

METHOD

Open Access

Differential expression analysis for sequence count data

Simon Anders*, Wolfgang Huber

Gene	Sample A	Sample B	Pseudo-reference
<i>rpoB</i>	1100	750	$\sqrt[2]{1100 \times 750} = 908,30$
<i>eis</i>	15	10	$\sqrt[2]{15 \times 10} = 12,25$

Gene	Ratio Sample A	Ratio Sample B
<i>rpoB</i>	$1100/908,30 = 1,21$	$750/908,30 = 0,83$
<i>eis</i>	$15/12,25 = 1,22$	$10/12,25 = 0,82$

Normalization Factors:

Sample A – Median(1,21; 1,22)= 1,215

Sample B - Median(0,83; 0,82)= 0,825

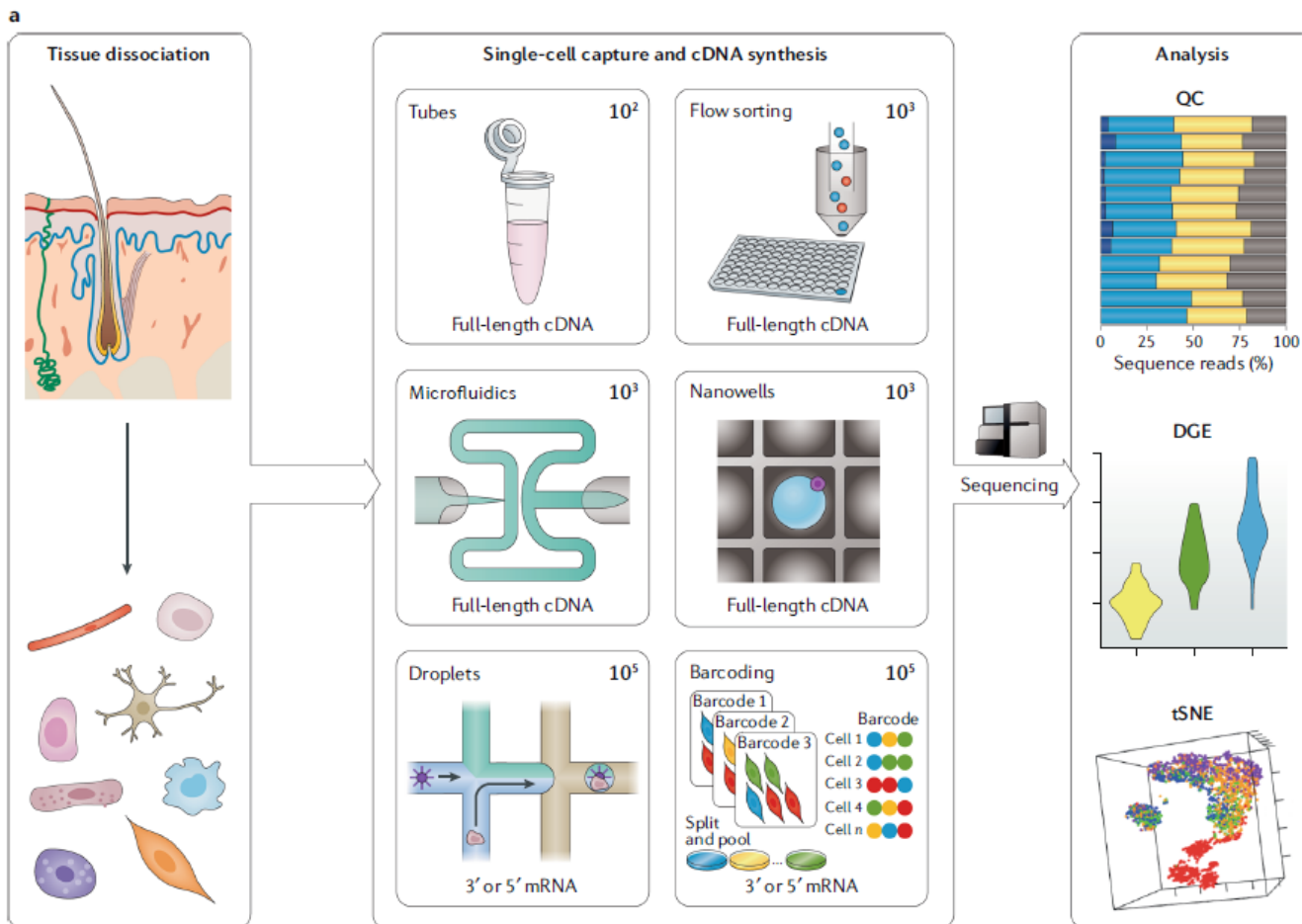
Gene	Normalized Count A	Normalized Count B
<i>rpoB</i>	$1100/1,215 = 905,35$	$750/0,825 = 909,09$
<i>eis</i>	$15/1,215 = 12,35$	$10/0,825 = 12,12$

Comparison of Normalization Methods

Method	Factors Accounted			Applications		
	Sequencing Depth	Gene Length	RNA Composition	Within sample	Comparisons between samples	DE Analysis
CPM	✓	✗	✗	✗	✓	✗
RPKM/FPKM	✓	✓	✗	✓	✗	✗
TPM	✓	✓	✗	✓	✓	✗
Median of Ratios (<i>DESeq2</i>)	✓	✗	✓	✗	✓	✓
Trimmed Mean of M Values (<i>EdgeR</i>)	✓	✓	✓	✓	✓	✓

Other than bulk RNA-Seq ...

Single-cell RNA-Seq



Spatialomics

