

# RNAseq and Annotation

Isheng Jason Tsai

Introduction to NGS Data and Analysis  
Lecture 6 [v2020]



# Lecture outline

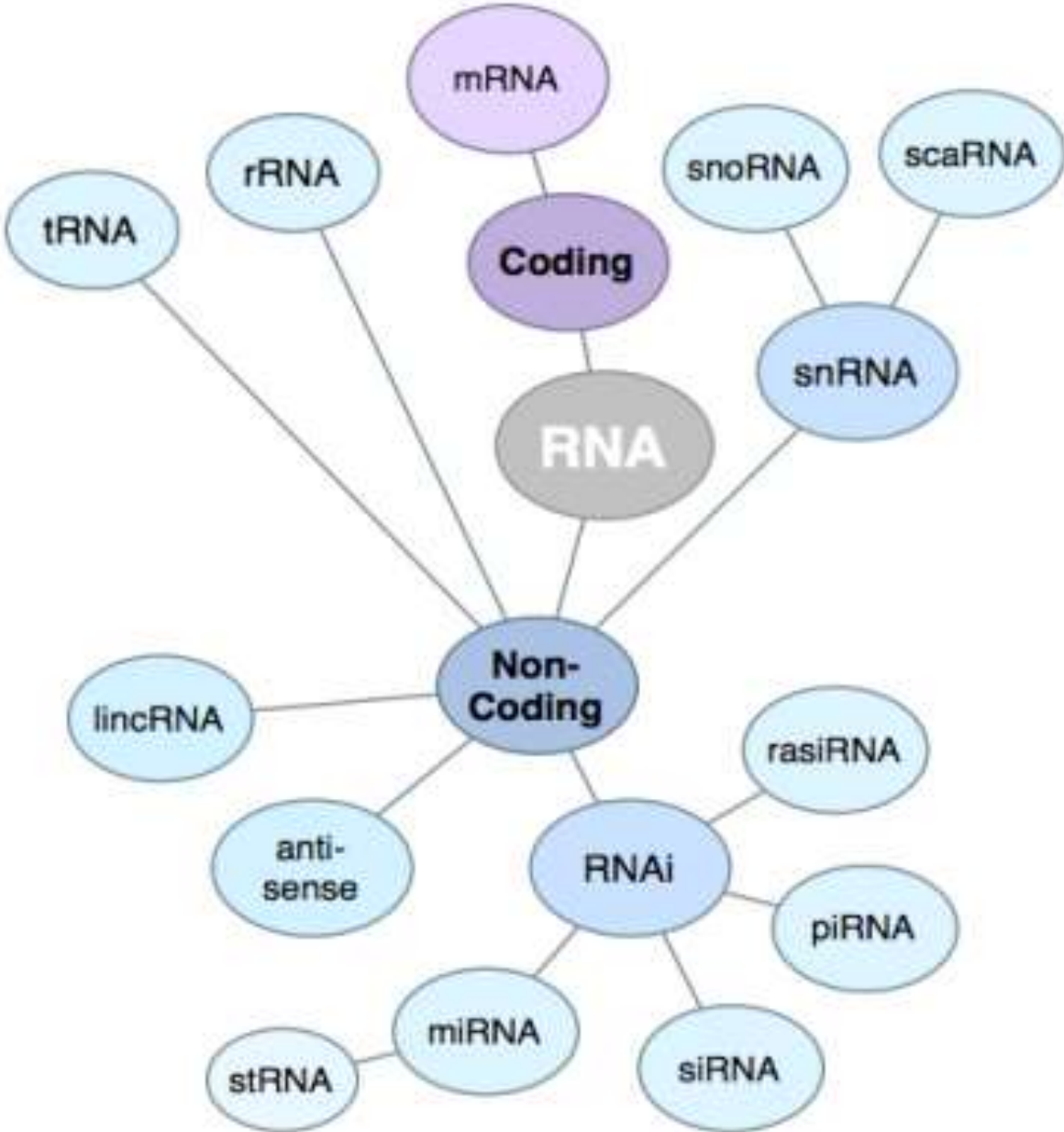
## **What I will cover today**

- mRNAseq (for **bulk** RNAseq)
- mapping
- assembly / reconstruction
- annotation
- experimental design
- differential expression
- single cell genomics

## **What I won't cover due to time constraints but equally important**

- *pseudoalignment* (kallisto, salmon etc)
- noncoding RNAs
- long read mapping

# Types of RNA



Type	Percent of total RNA by mass	Molecules per cell	Average size (kb)	Total weight picograms/cell	Notes	Reference
rRNAs	80 to 90	3–10 × 10 <sup>6</sup> (ribosomes)	6.9	10 to 30		Blobel and Potter (1967), Wolf and Schlessinger (1977), Duncan and Hershey (1983)
tRNA	10 to 15	3–10 × 10 <sup>7</sup>	<0.1	1.5 to 5	About 10 tRNA molecules /ribosome	Waldron and Lacroute (1975)
mRNA	3 to 7	3–10 × 10 <sup>5</sup>	1.7	0.25 to 0.9		Hastie and Bishop (1976), Carter et al. (2005)
hnRNA (pre-mRNA)	0.06 to 0.2	1–10 × 10 <sup>3</sup>	10*	0.004 to 0.03	Estimated at 2–4% of mRNA by weight	Mortazavi et al. (2008), Menet et al. (2012)
Circular RNA	0.002 to 0.03	3–20 × 10 <sup>3</sup>	~0.5	0.0007 to 0.005	Estimated at 0.1–0.2% of mRNA**	Salzman et al. (2012), Guo et al. (2014)
snRNA	0.02 to 0.3	1–5 × 10 <sup>5</sup>	0.1–0.2	0.008 to 0.04		Kiss and Filipowicz (1992), Castle et al. (2010)
snoRNA	0.04 to 0.2	2–3 × 10 <sup>5</sup>	0.2	0.02 to 0.03		Kiss and Filipowicz (1992), Cooper (2000), Castle et al. (2010)
miRNA	0.003 to 0.02	1–3 × 10 <sup>5</sup>	0.02	0.001 to 0.003	About 10 <sup>5</sup> molecules per 10 pg total RNA	Bissels et al. (2009)
7SL	0.01 to 0.2	3–20 × 10 <sup>4</sup>	0.3	0.005 to 0.03	About 1–2 SRP molecules/100 ribosomes	Raue et al. (2007), Castle et al. (2010)
Xist	0.0003 to 0.02	0.1–2 × 10 <sup>3</sup>	2.8	0.0001 to 0.003		Buzin et al. (1994), Castle et al. (2010)
Other lncRNA	0.03 to 0.2	3–50 × 10 <sup>3</sup>	1	0.002 to 0.03	Estimated at 1–4% of mRNA by weight	Mortazavi et al. (2008), Ramsköld et al. (2009), Menet et al. (2012)

\*The size for the average unspliced pre-mRNA is 17 kb; however, most pre-mRNAs are partially spliced at any given time, and the average size of hnRNA is estimated at 10 kb (Salditt-Georgieff et al., 1976).

\*\*Based on the finding that 1–2% of all mRNA species generate circular RNA, which is present at 10% of the level of the parental mRNA.

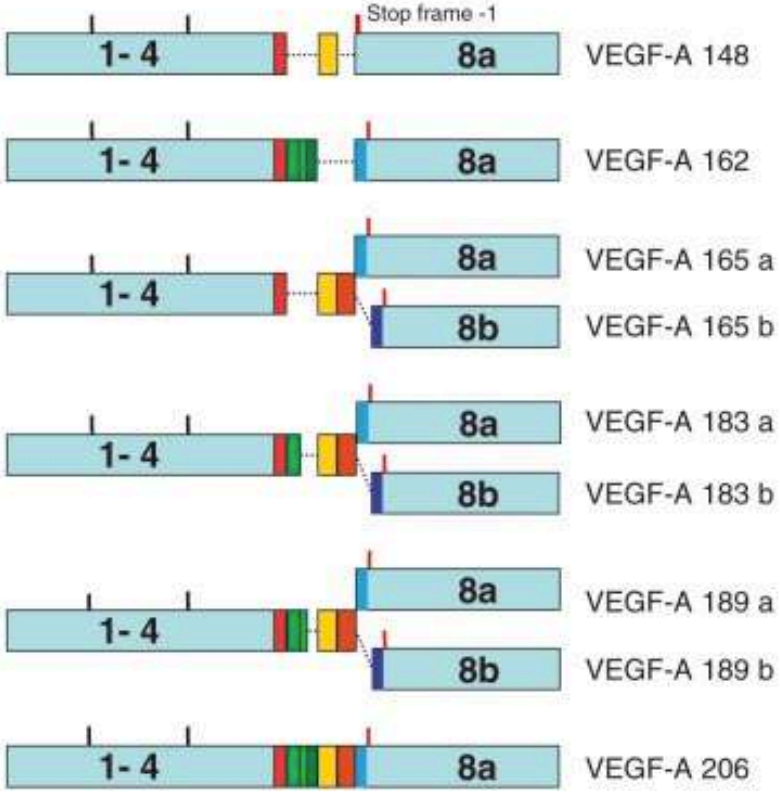
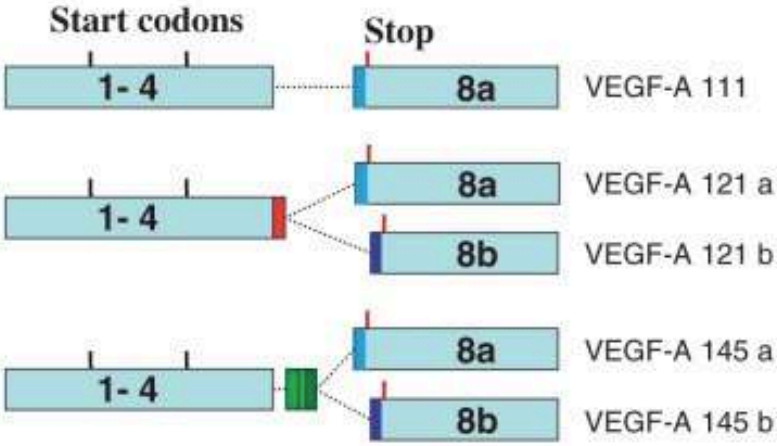
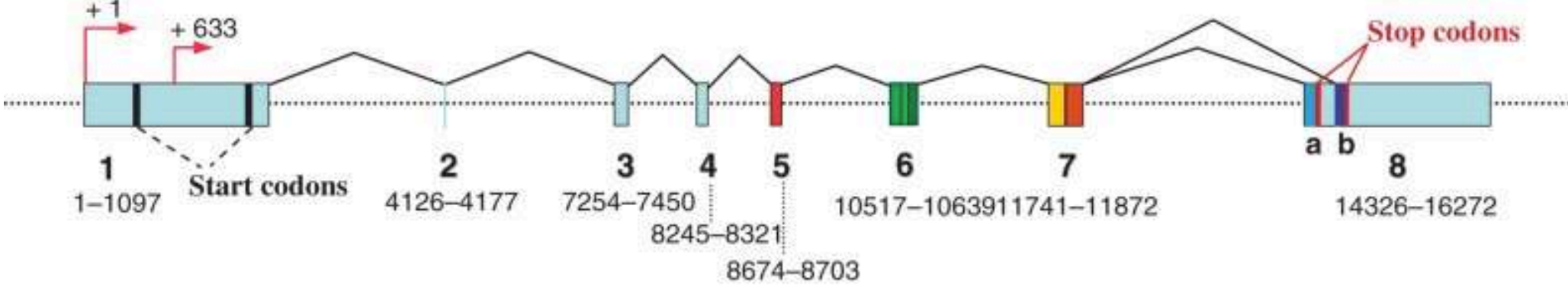


Type	Percent of total RNA by mass	Molecules per cell	Average size (kb)	Total weight picograms/cell	Notes	Reference
rRNAs	80 to 90	3–10 × 10 <sup>6</sup> (ribosomes)	6.9	10 to 30		Blobel and Potter (1967), Wolf and Schlessinger (1977), Duncan and Hershey (1983)
tRNA	10 to 15	3–10 × 10 <sup>7</sup>	<0.1	1.5 to 5	About 10 tRNA molecules /ribosome	Waldron and Lacroute (1975)
mRNA	3 to 7	3–10 × 10 <sup>5</sup>	1.7	0.25 to 0.9		Hastie and Bishop (1976), Carter et al. (2005)
hnRNA (pre-mRNA)	0.06 to 0.2	1–10 × 10 <sup>3</sup>	10*	0.004 to 0.03	Estimated at 2–4% of mRNA by weight	Mortazavi et al. (2008), Menet et al. (2012)
Circular RNA	0.002 to 0.03	3–20 × 10 <sup>3</sup>	~0.5	0.0007 to 0.005	Estimated at 0.1–0.2% of mRNA**	Salzman et al. (2012), Guo et al. (2014)
snRNA	0.02 to 0.3	1–5 × 10 <sup>5</sup>	0.1–0.2	0.008 to 0.04		Kiss and Filipowicz (1992), Castle et al. (2010)
snoRNA	0.04 to 0.2	2–3 × 10 <sup>5</sup>	0.2	0.02 to 0.03		Kiss and Filipowicz (1992), Cooper (2000), Castle et al. (2010)
miRNA	0.003 to 0.02	1–3 × 10 <sup>5</sup>	0.02	0.001 to 0.003	About 10 <sup>5</sup> molecules per 10 pg total RNA	Bissels et al. (2009)
7SL	0.01 to 0.2	3–20 × 10 <sup>4</sup>	0.3	0.005 to 0.03	About 1–2 SRP molecules/100 ribosomes	Raue et al. (2007), Castle et al. (2010)
Xist	0.0003 to 0.02	0.1–2 × 10 <sup>3</sup>	2.8	0.0001 to 0.003		Buzin et al. (1994), Castle et al. (2010)
Other lncRNA	0.03 to 0.2	3–50 × 10 <sup>3</sup>	1	0.002 to 0.03	Estimated at 1–4% of mRNA by weight	Mortazavi et al. (2008), Ramsköld et al. (2009), Menet et al. (2012)

\*The size for the average unspliced pre-mRNA is 17 kb; however, most pre-mRNAs are partially spliced at any given time, and the average size of hnRNA is estimated at 10 kb (Salditt-Georgieff et al., 1976).

\*\*Based on the finding that 1–2% of all mRNA species generate circular RNA, which is present at 10% of the level of the parental mRNA.

# Gene and isoforms



# Typical RNAseq Workflow

1. Experiment / Generate data
2. Map or Assemble
3. Count / Differential expression
4. Analysis

REVIEW

Open Access

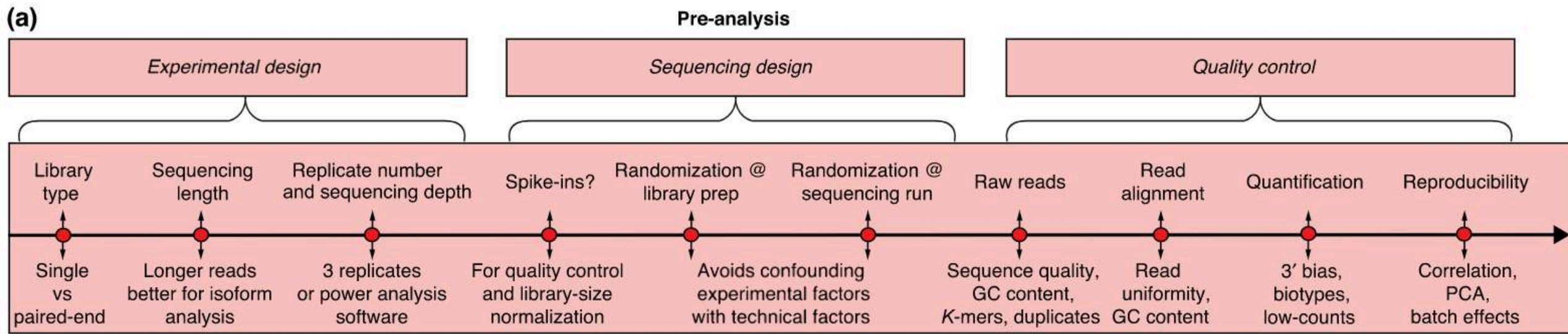
# A survey of best practices for RNA-seq data analysis



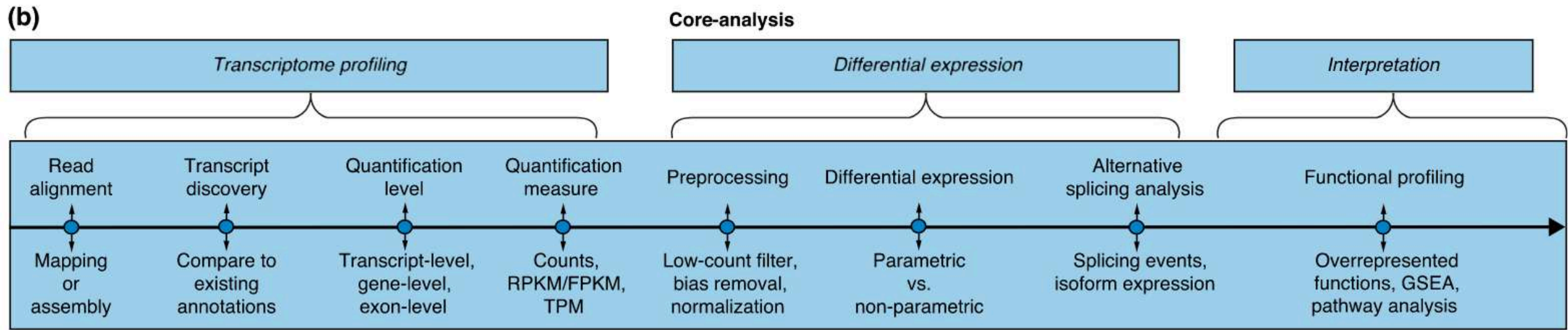
Ana Conesa<sup>1,2\*</sup>, Pedro Madrigal<sup>3,4\*</sup>, Sonia Tarazona<sup>2,5</sup>, David Gomez-Cabrero<sup>6,7,8,9</sup>, Alejandra Cervera<sup>10</sup>, Andrew McPherson<sup>11</sup>, Michał Wojciech Szczęśniak<sup>12</sup>, Daniel J. Gaffney<sup>3</sup>, Laura L. Elo<sup>13</sup>, Xuegong Zhang<sup>14,15</sup> and Ali Mortazavi<sup>16,17\*</sup>



# Pre-analysis



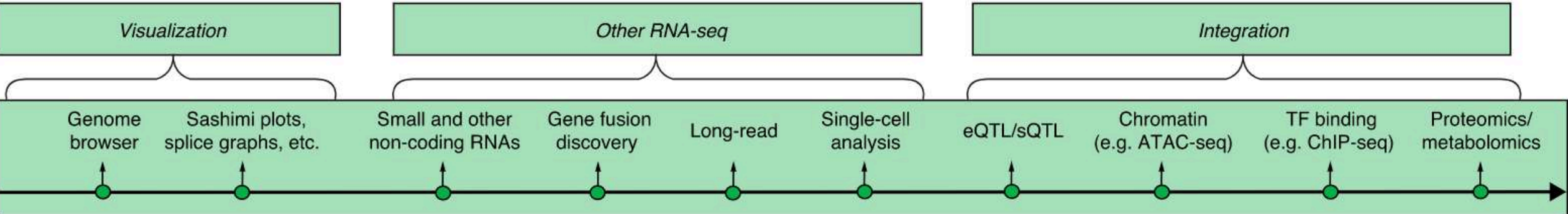
# Core-analysis



# Advanced-analysis (not covered in this lecture but should be mentioned)

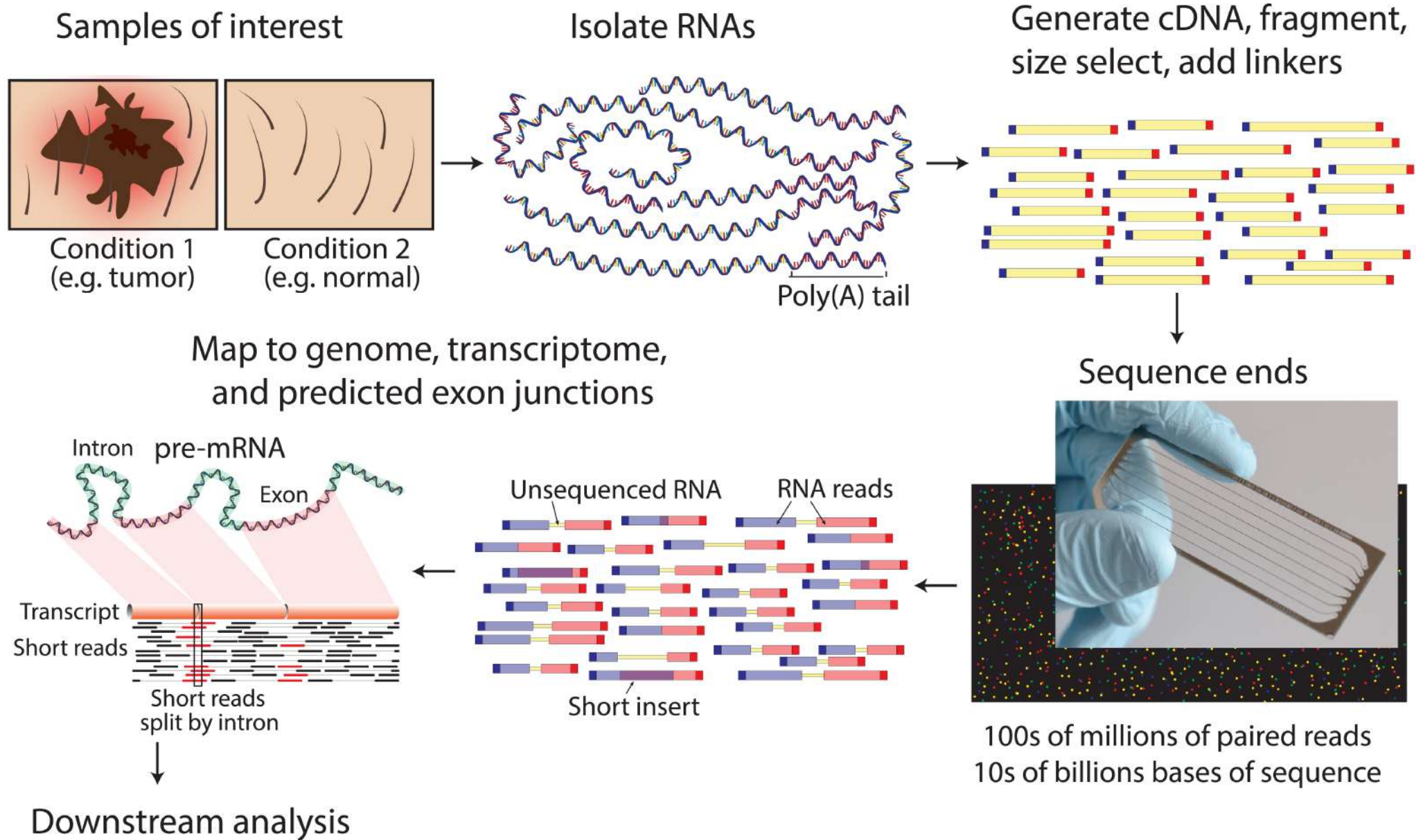
(c)

## Advanced-analysis



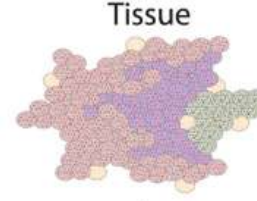
1. Generate data

# RNA-seq data generation



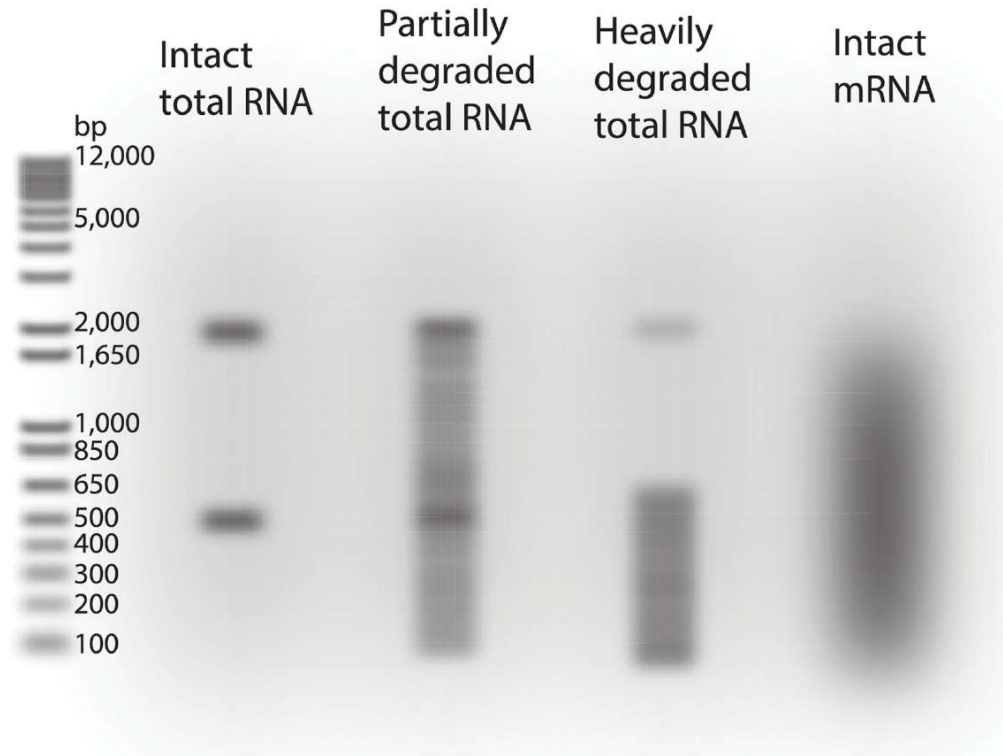


# RNA-seq data generation

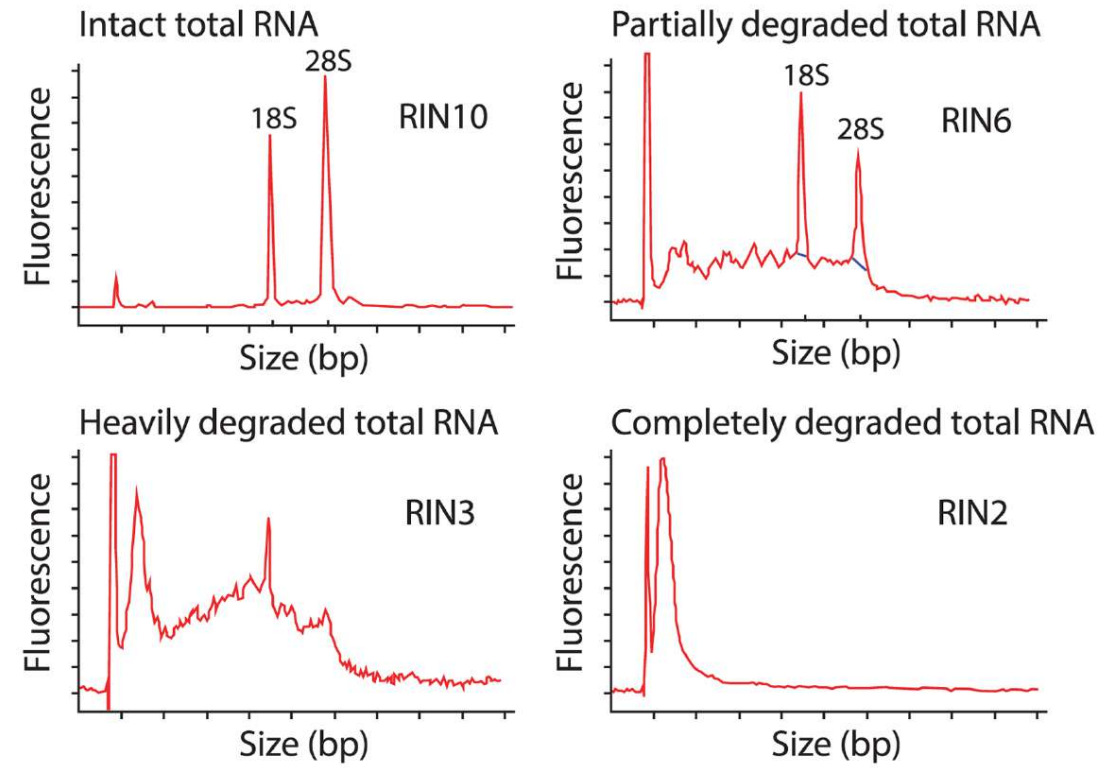


Assess RNA quality ← Isolate total RNA

Gel electrophoresis of RNA

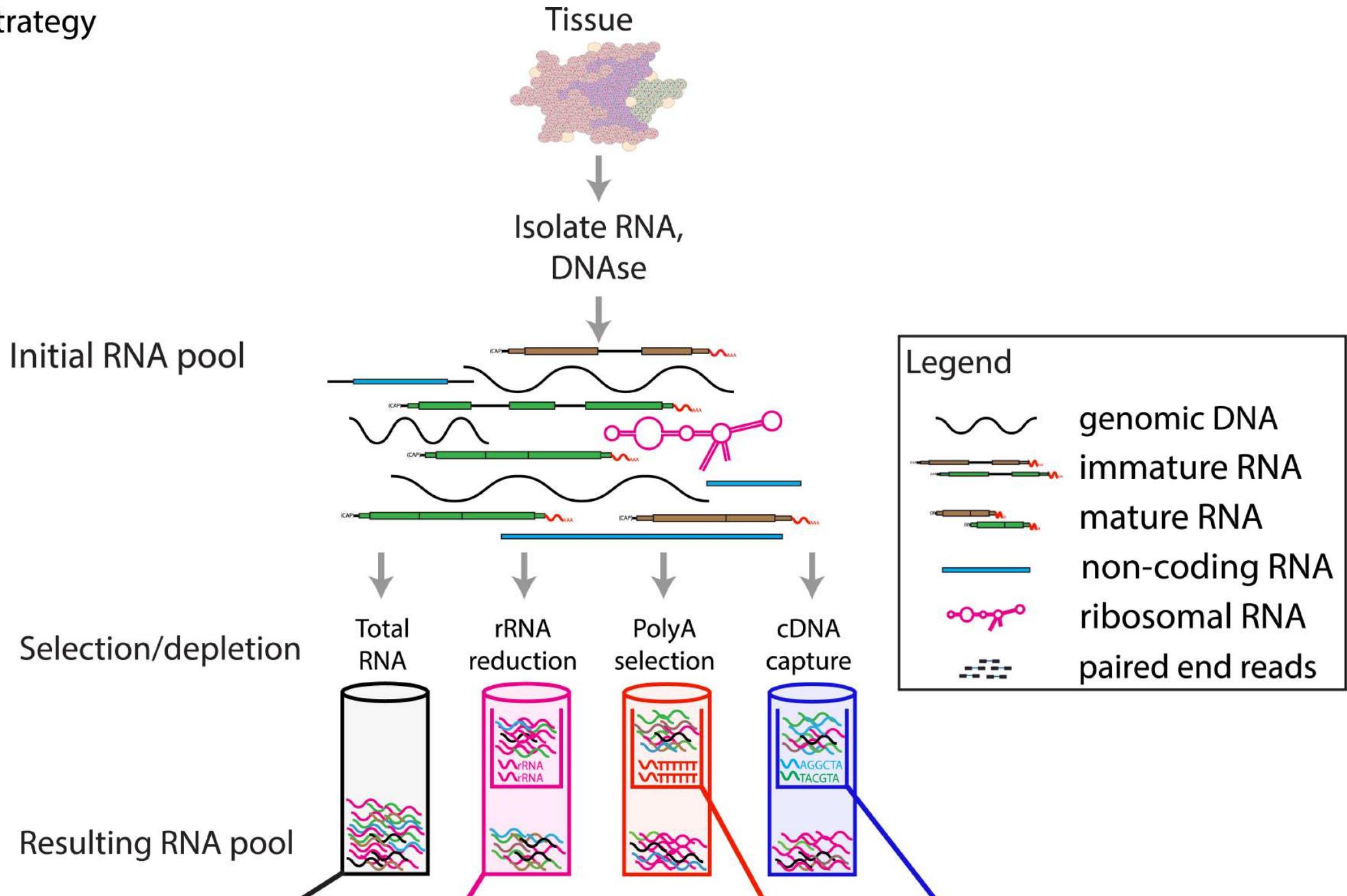


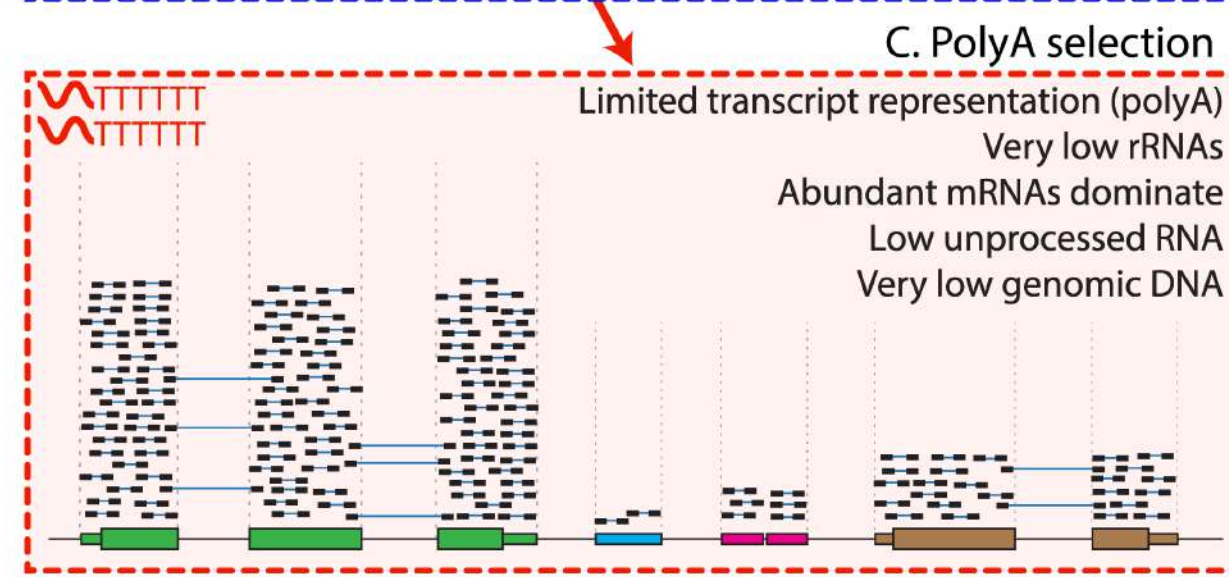
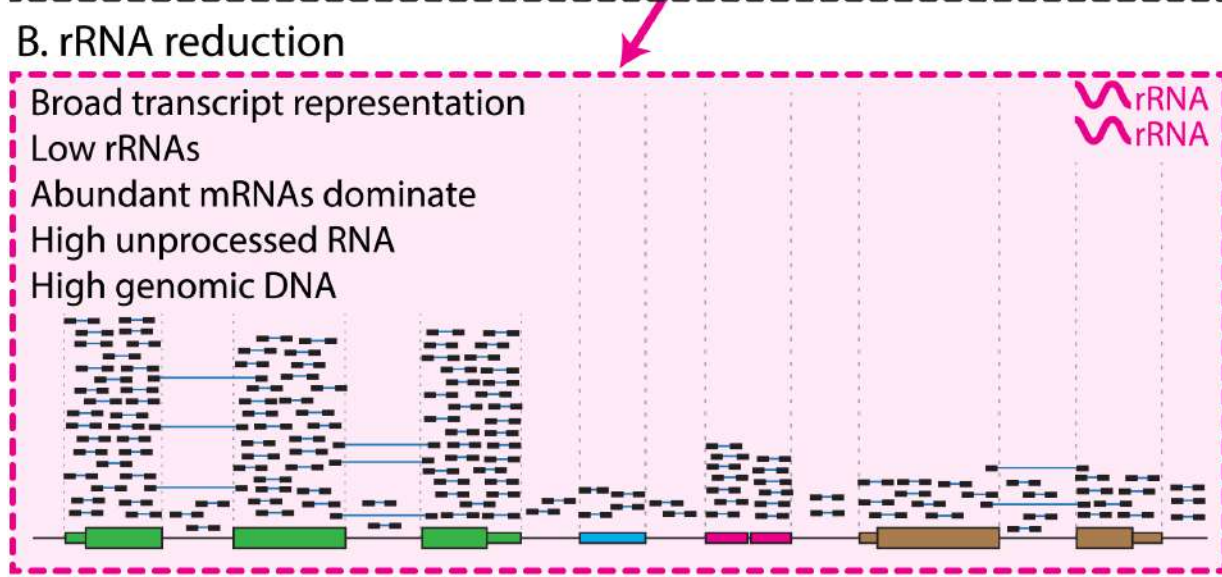
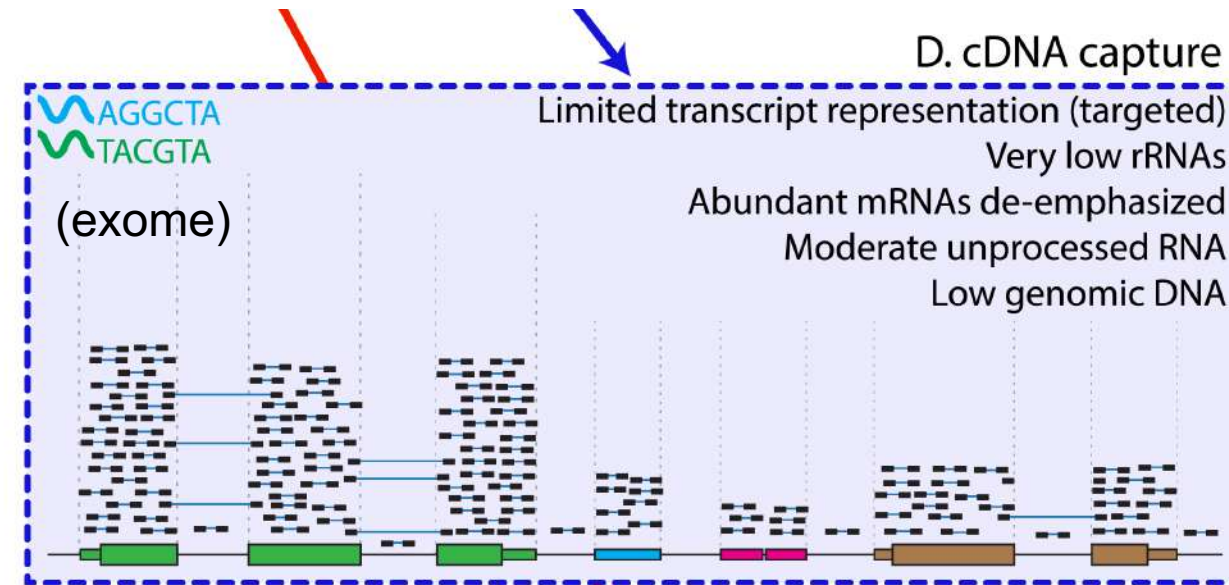
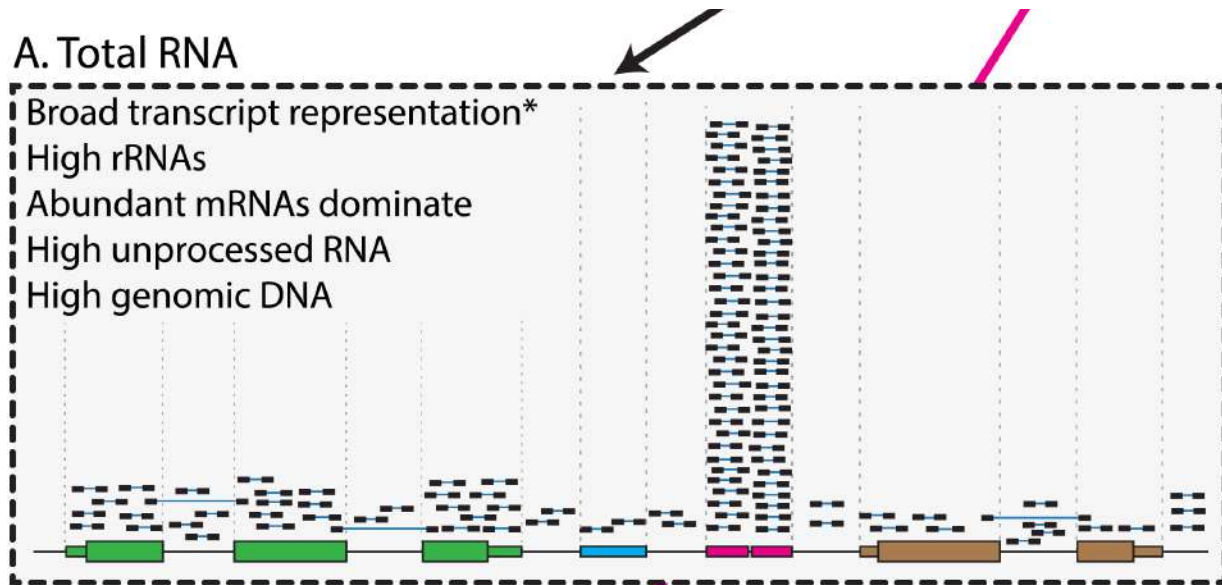
Capillary electrophoresis of total RNA



RIN = 28S:18S ratio

# RNA-seq Strategy

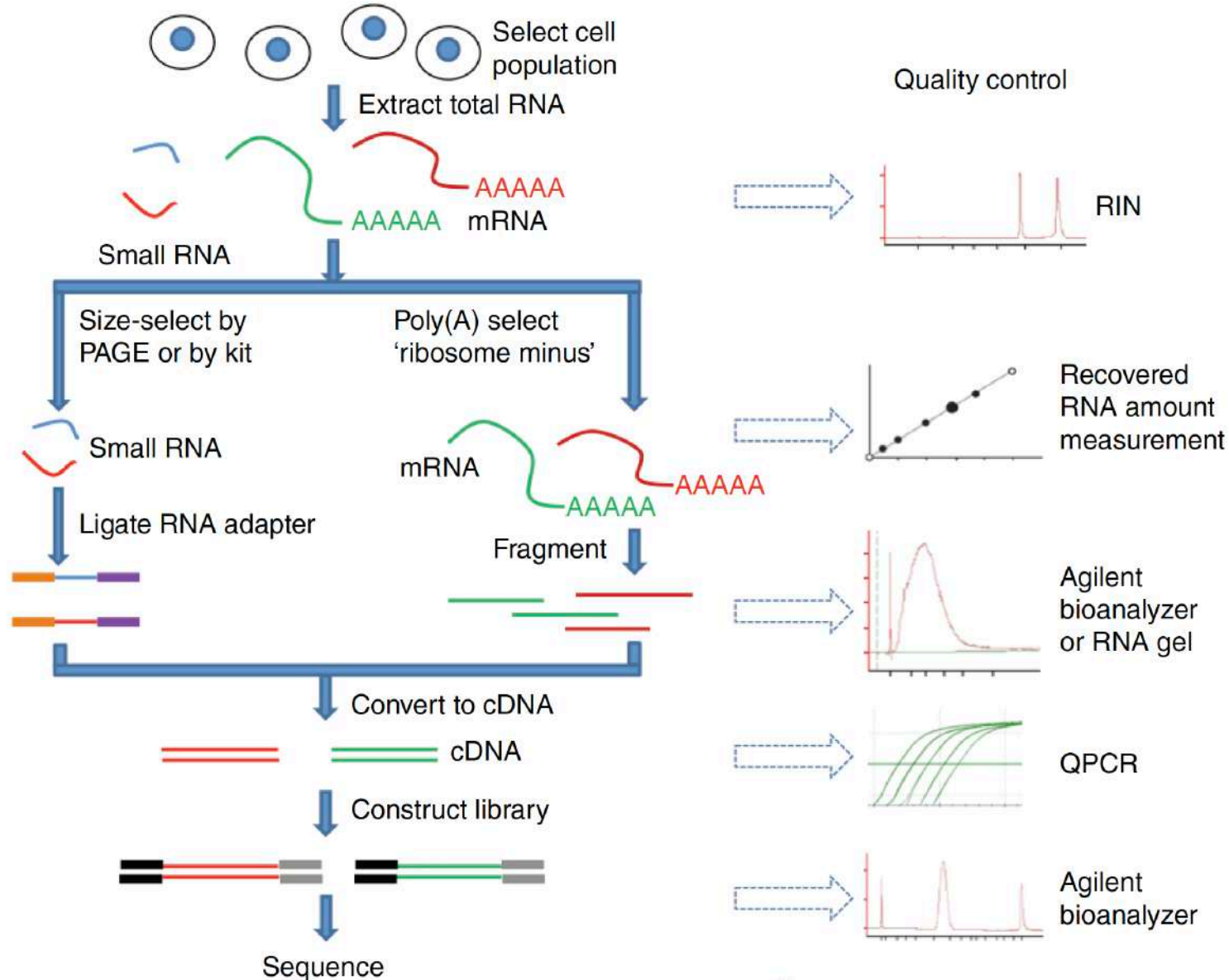




Expected Alignments



# General RNA library preparation workflow



i) RNA extraction and measuring its integrity,

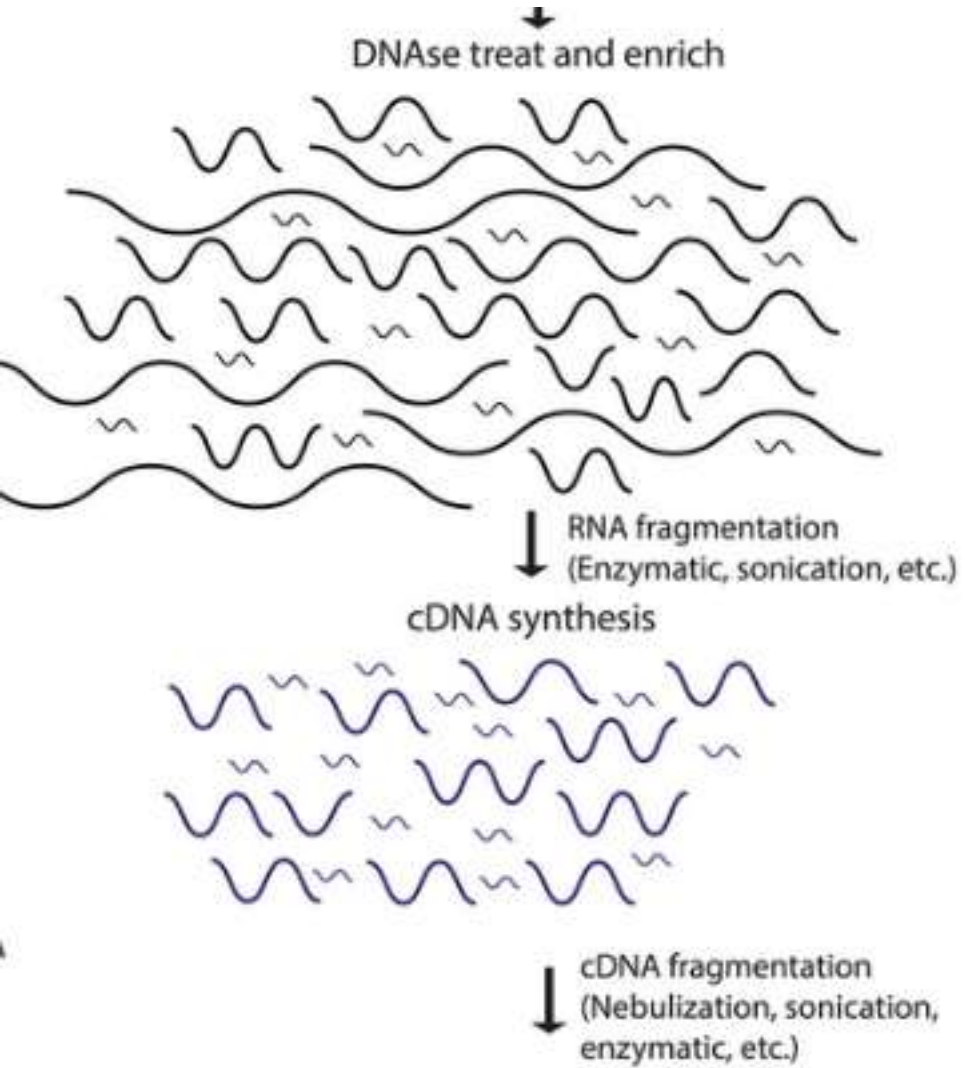
ii) rRNA is depleted (either using poly(A)-selection or rRNA depletion)

iii) the remaining RNA molecules are fragmented, ideally achieving a uniform size distribution.

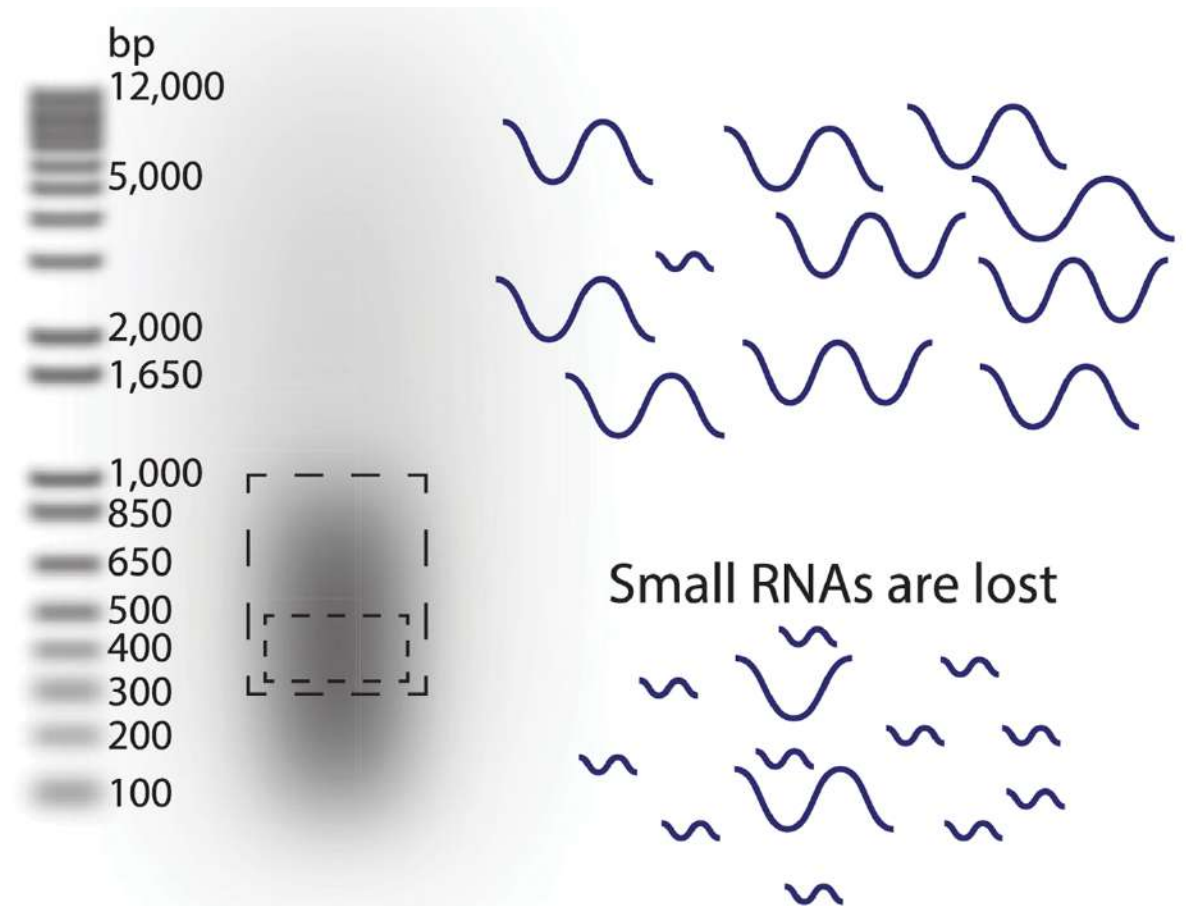
iv) Double-stranded cDNA is synthesized and the adapters for sequencing are added to construct the final library whose fragment size distribution should be unimodal and well-defined.

# RNA-seq data generation ; cDNA

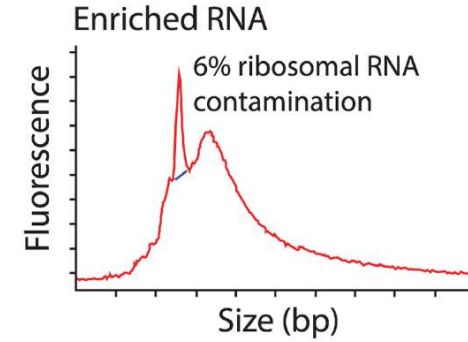
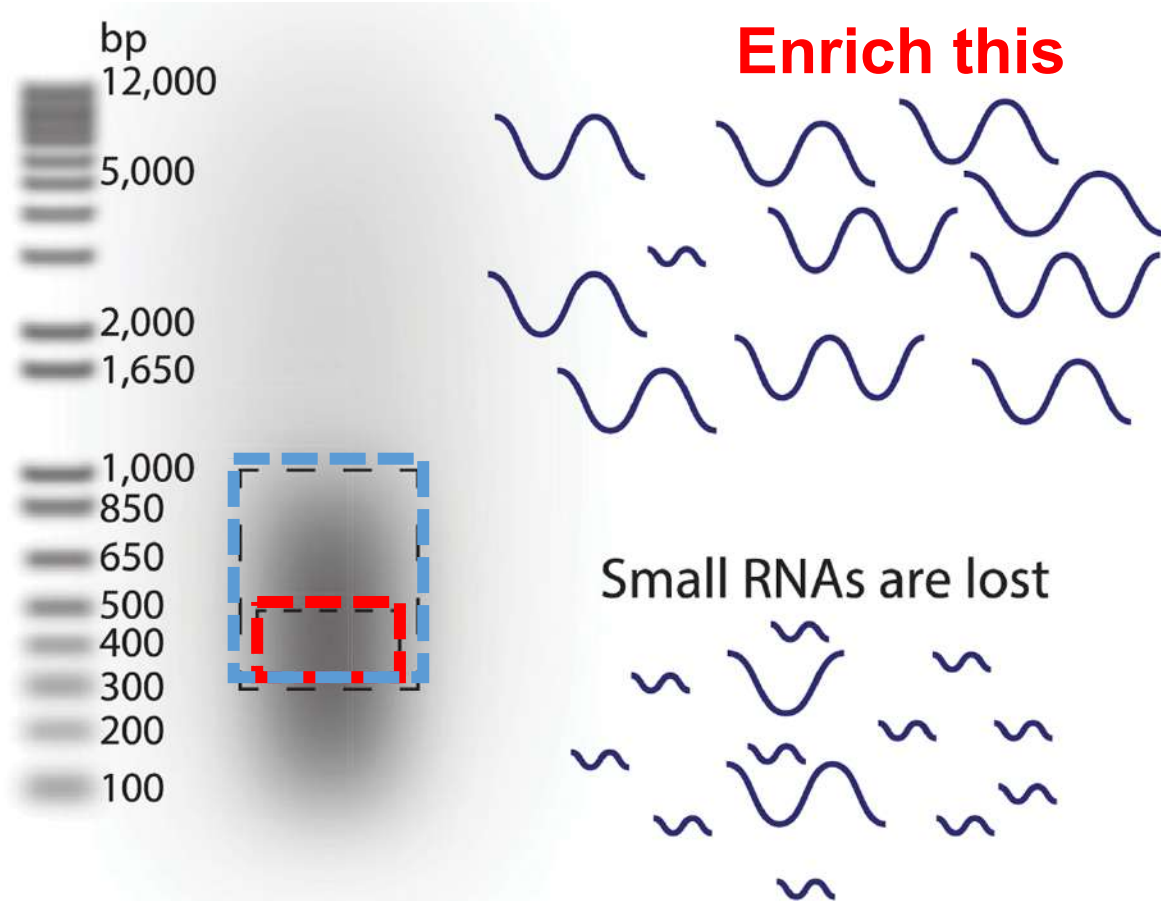
## Total RNA



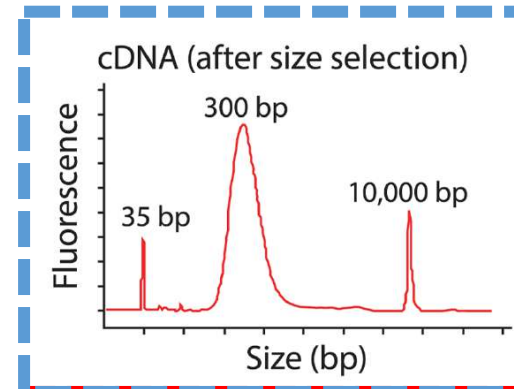
## Size selection



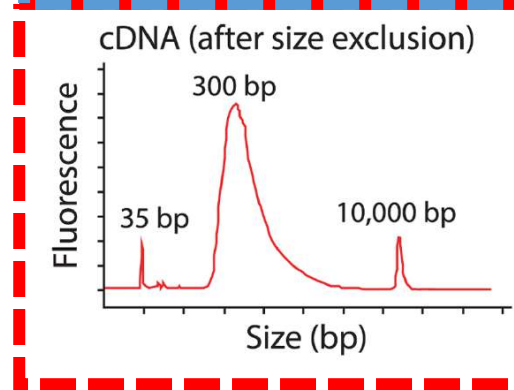
# RNA-seq data generation



.)



Column selection

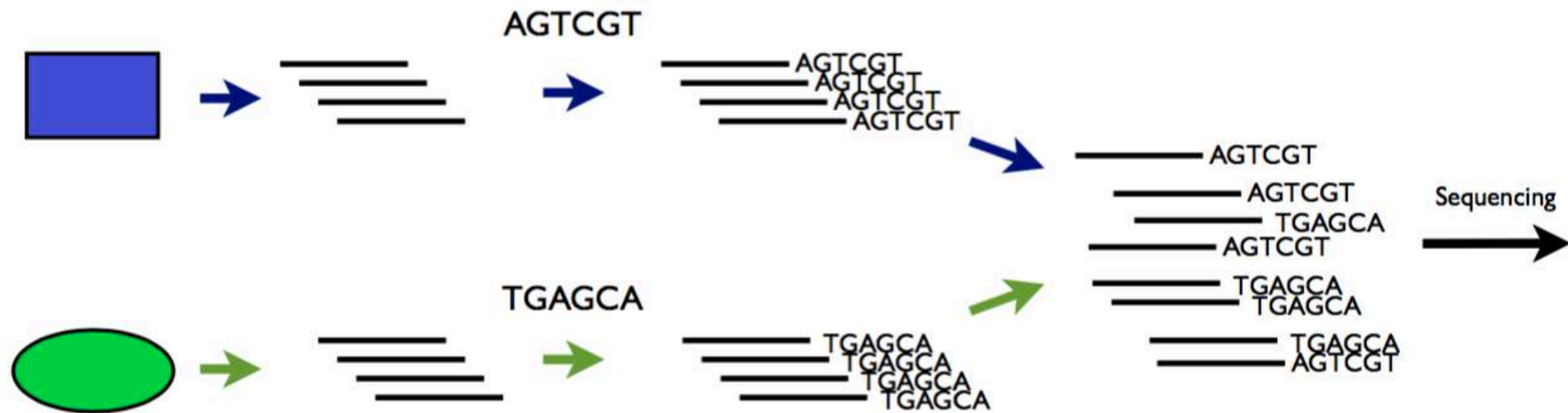


Gel selection

Enriched cDNA -> Added sequencing adaptors -> Sequencing

High coverage of Illumina allows multiplexing:

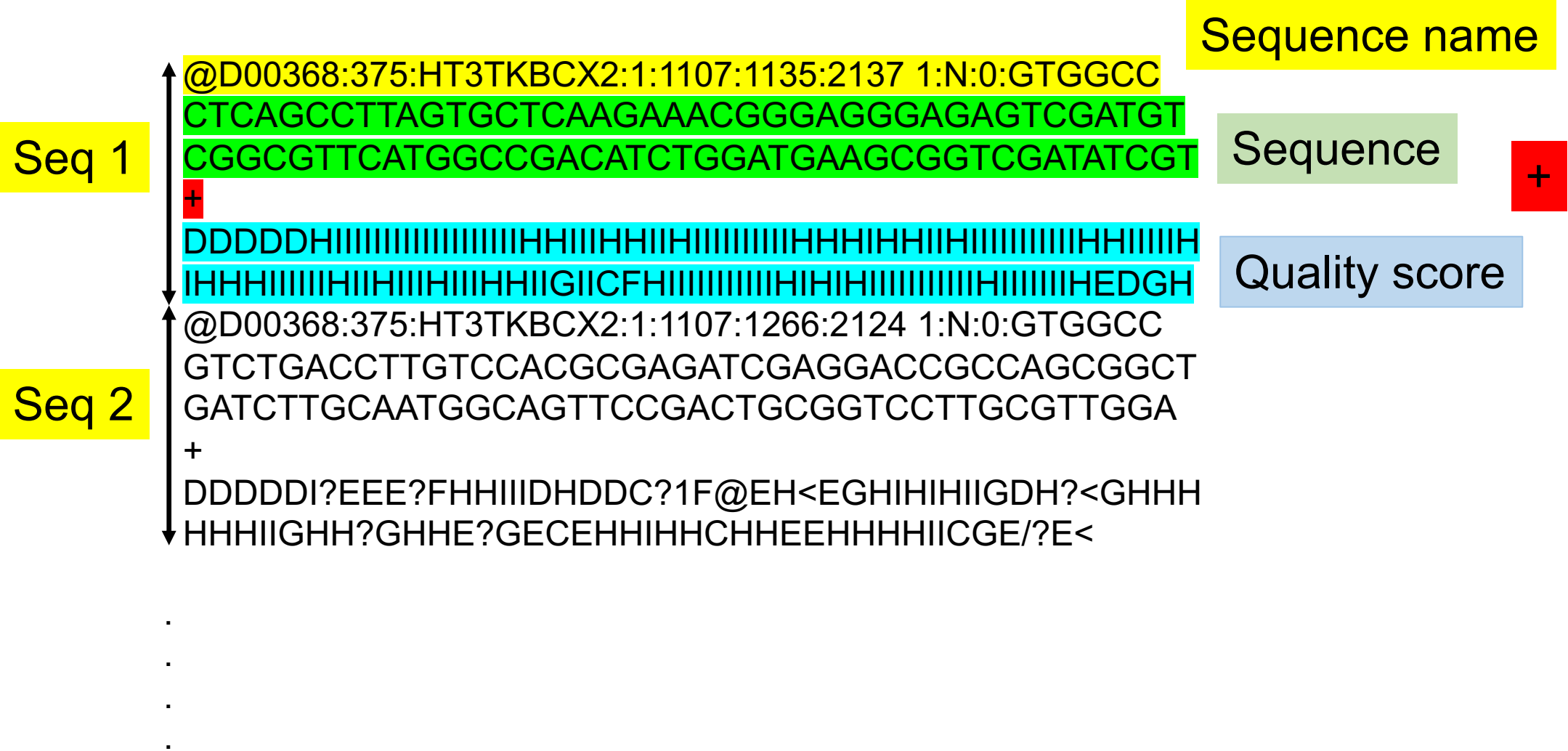
(Use of 4-6 nucleotides to identify different samples in the same run)



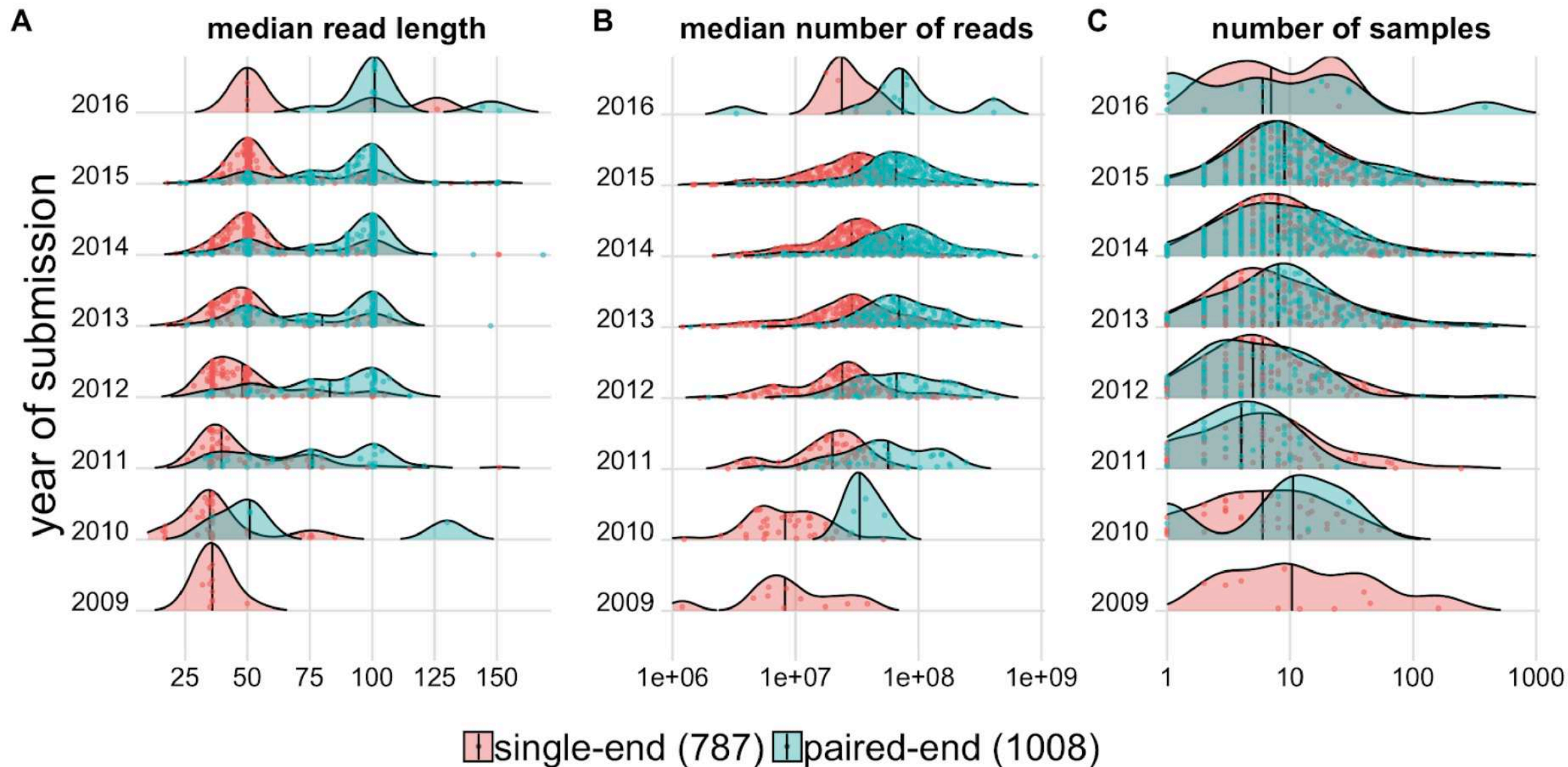
## 1.1 Data type



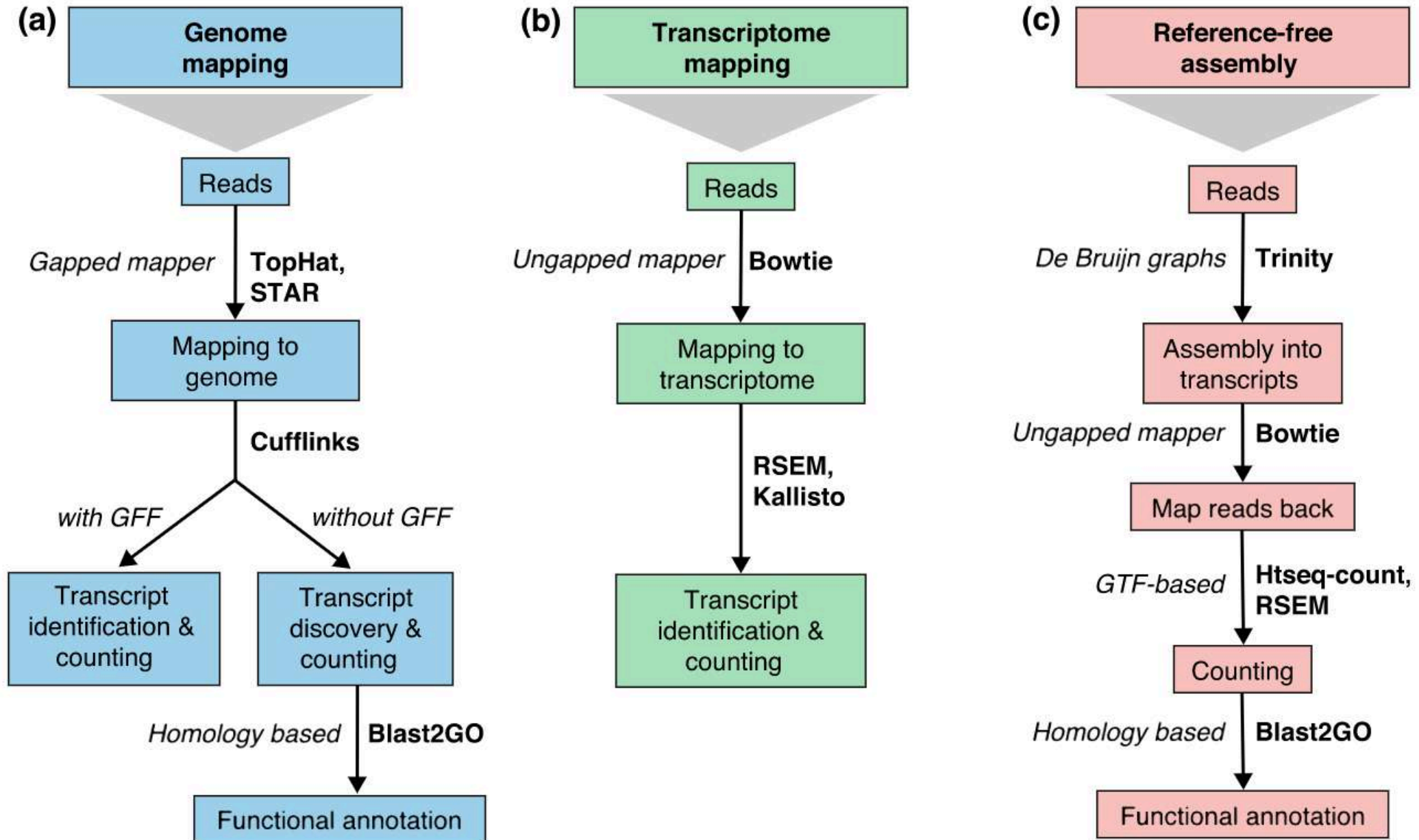
# Fastq file



# Evolution of RNAseq over time (from SRA)



# Read mapping and transcript identification strategies

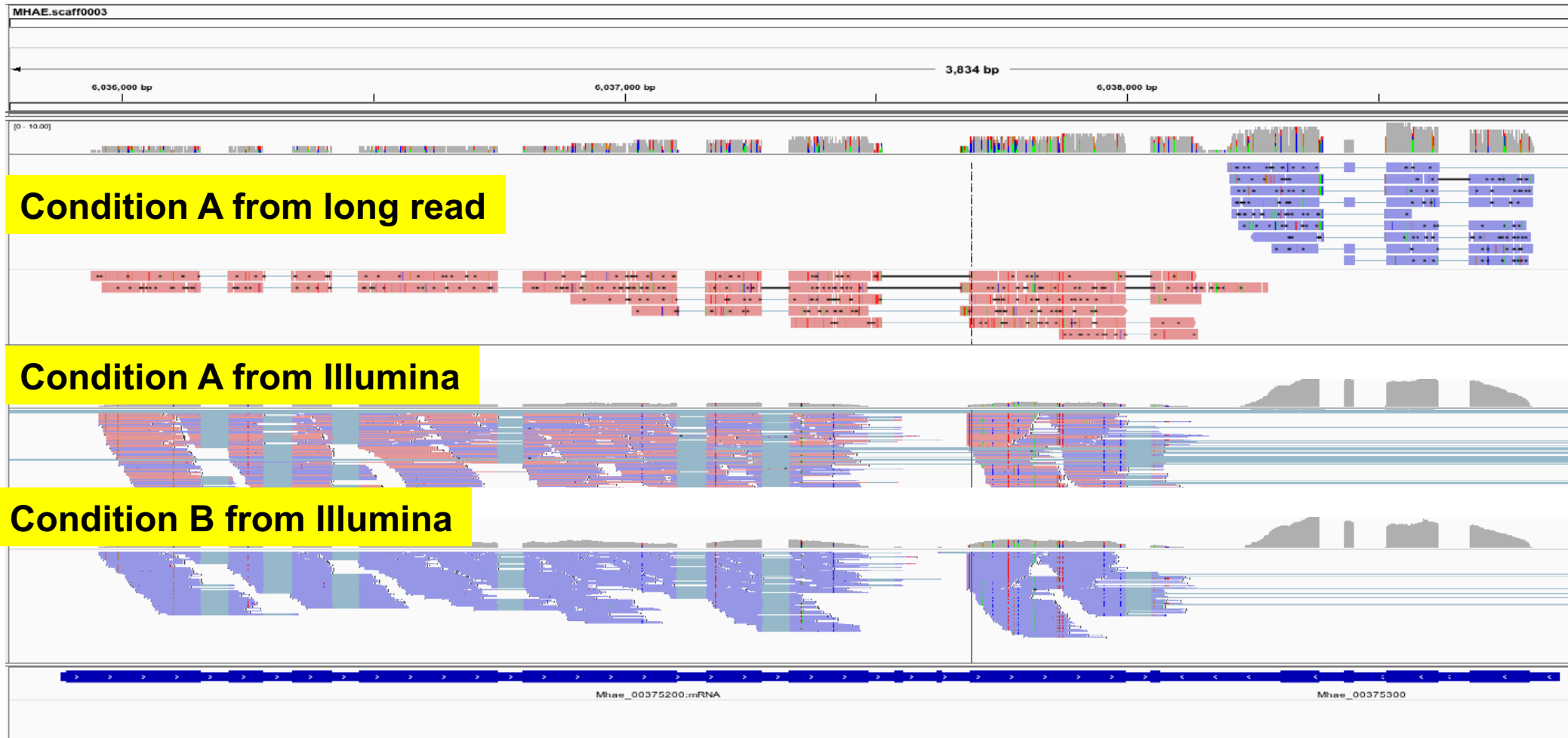


# Read visualisation

Load reference, annotation and bam into a program

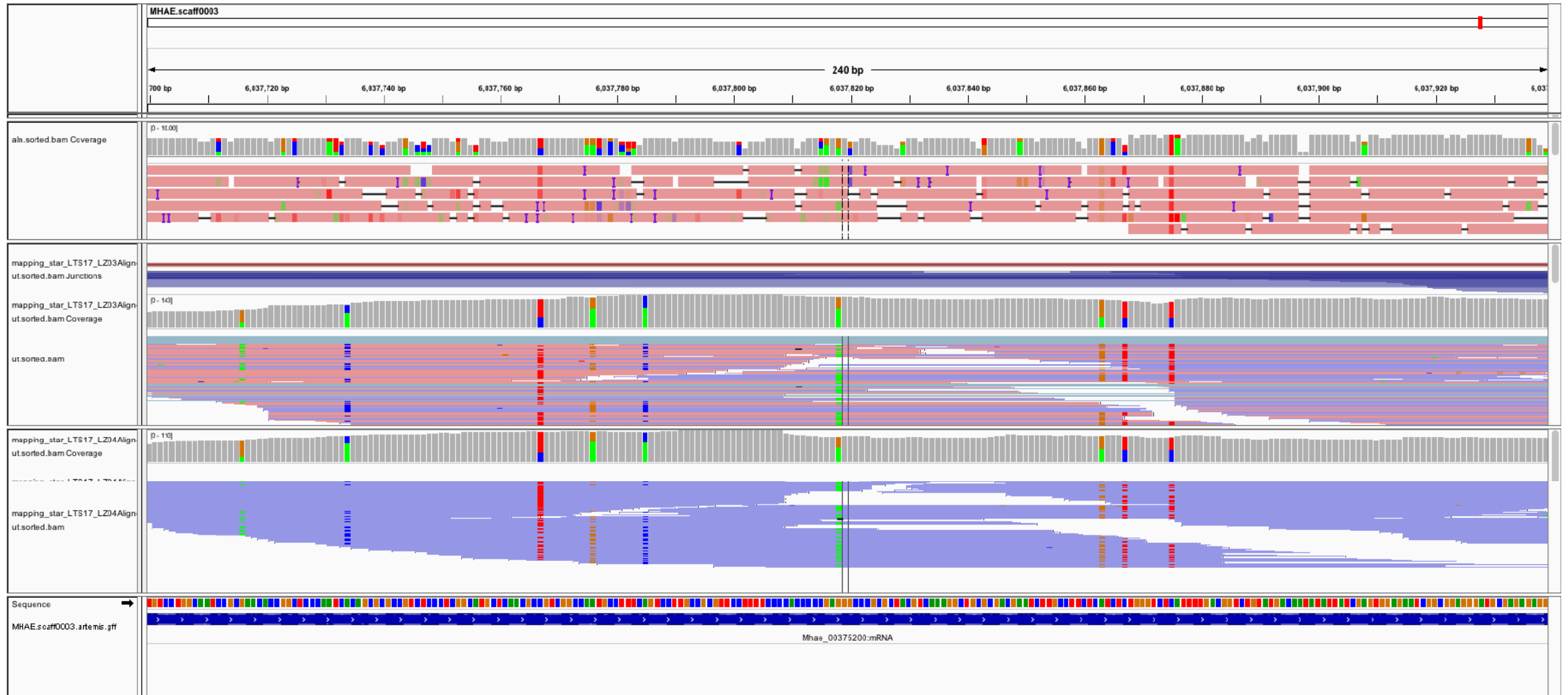
- Artemis <http://www.sanger.ac.uk/science/tools/artemis>
- IGV <http://software.broadinstitute.org/software/igv/>

# Scenario 1



Annotation: two genes of two orientations

# Long reads have more errors



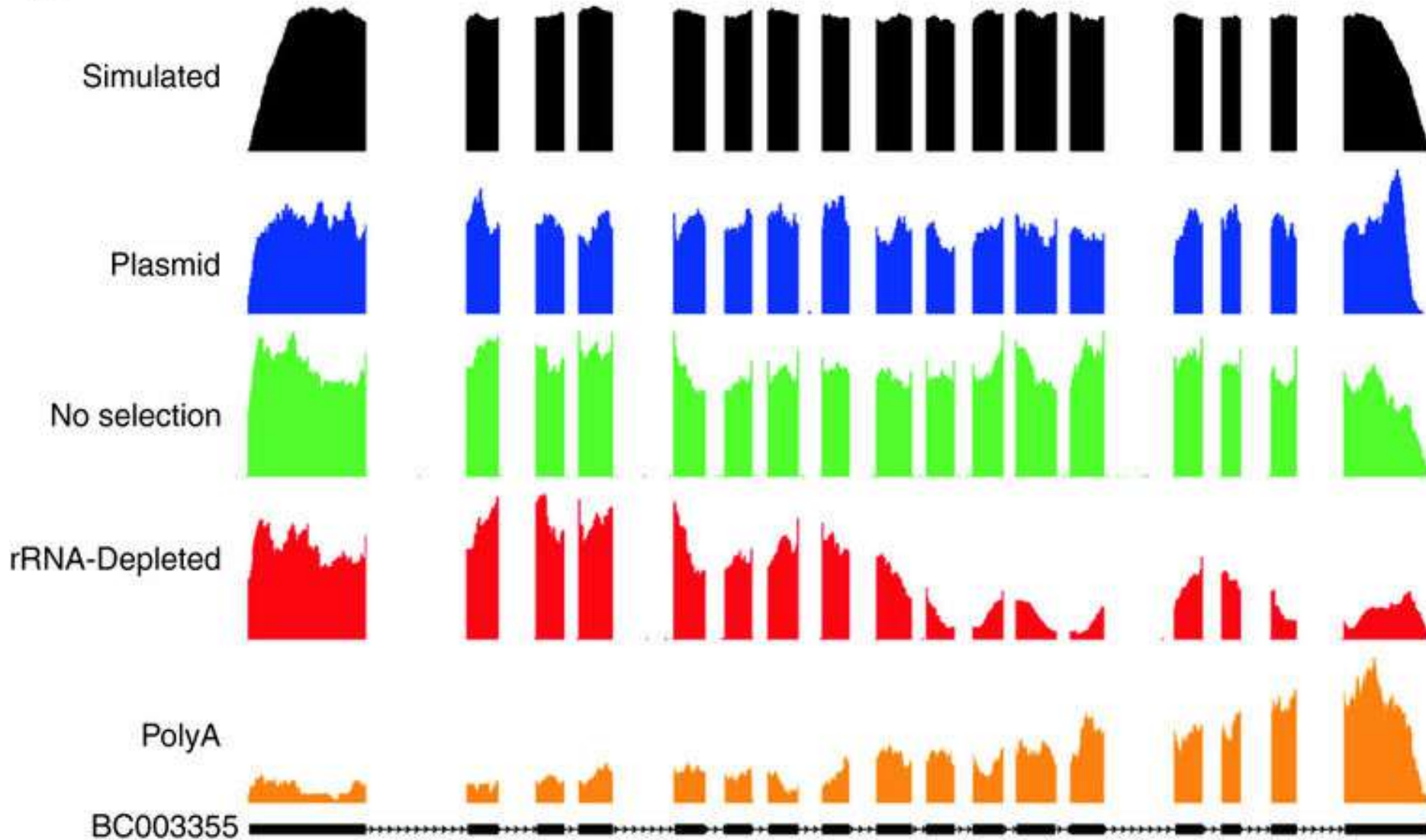
# Wrong annotation (wrong gene fusion / wrong exons)





# Library enrichment result in sequencing bias

**A**





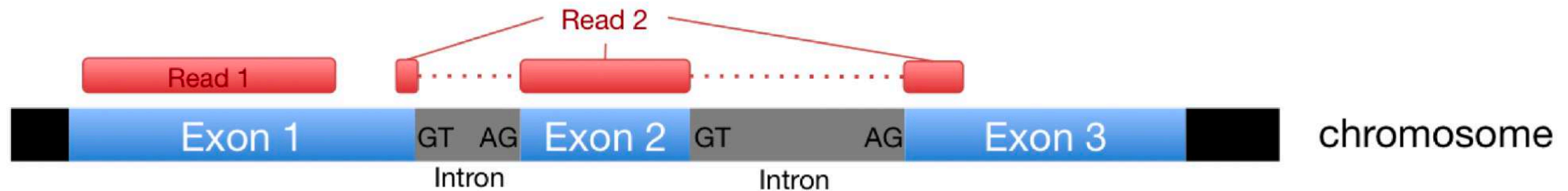
## 2. Mapping

# Read mapping and transcript identification strategies

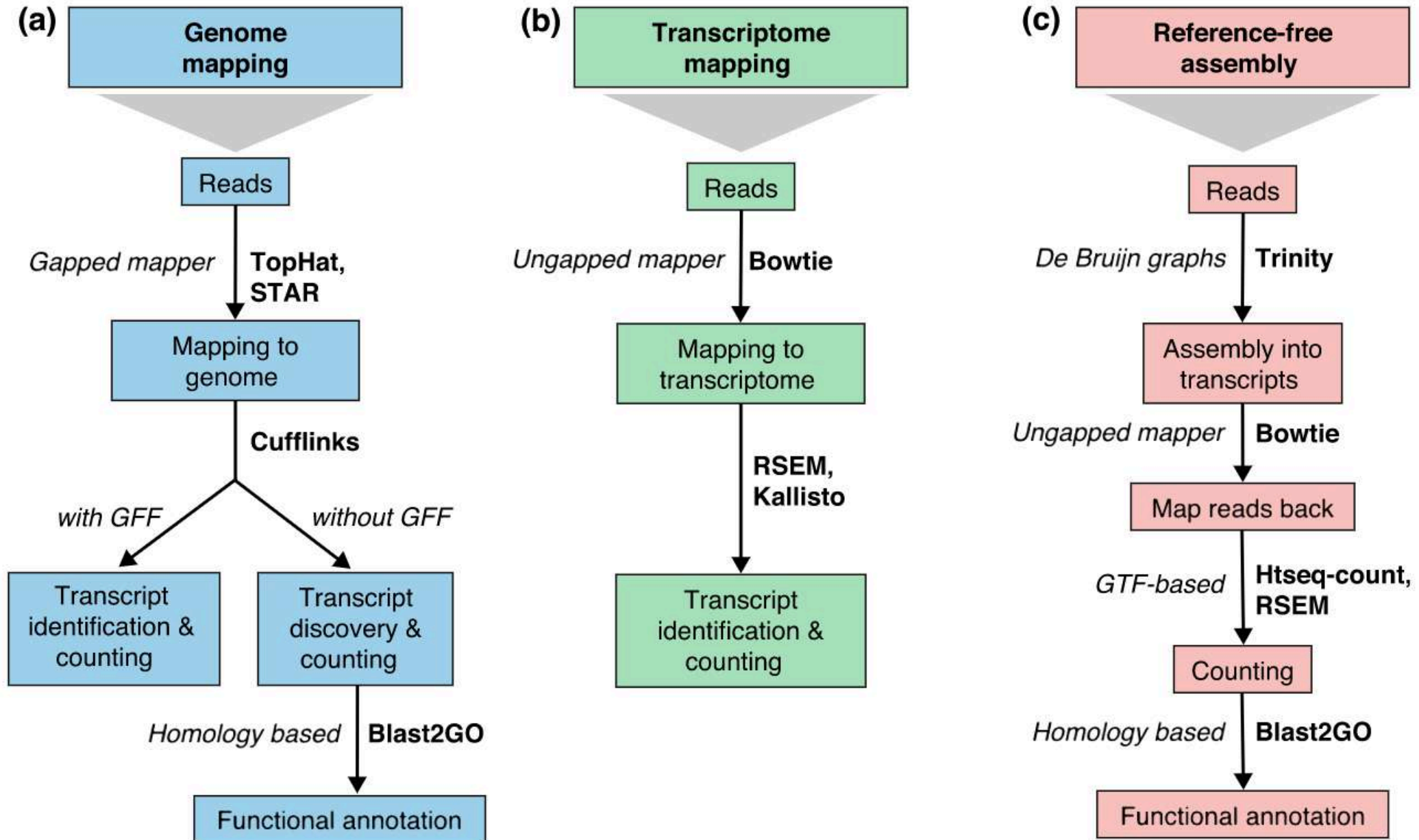
(a) Aligning to the transcriptome



(b) Aligning to the genome



# Read mapping and transcript identification strategies

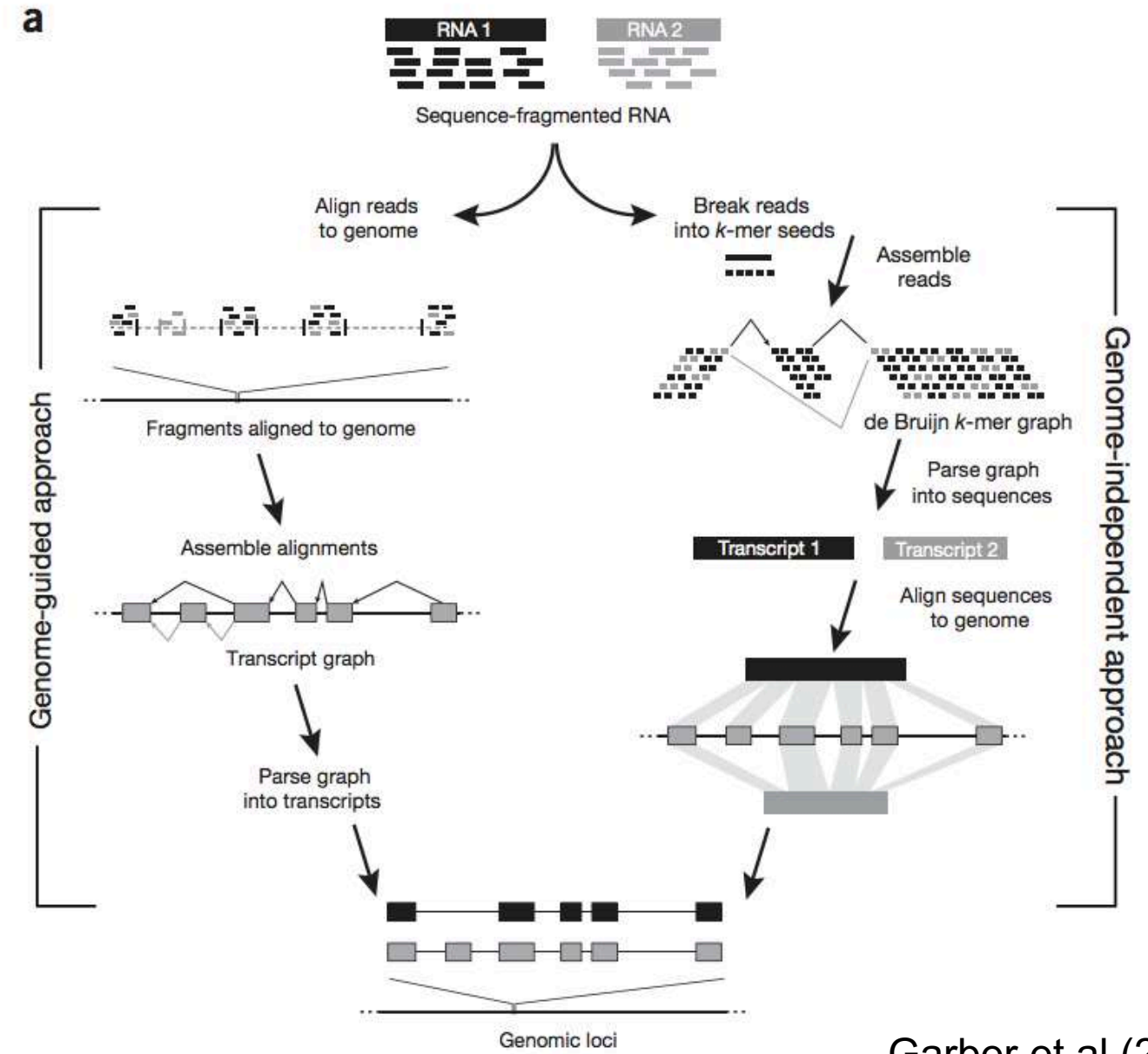


# General workflow for RNAseq to produce annotation

Options:

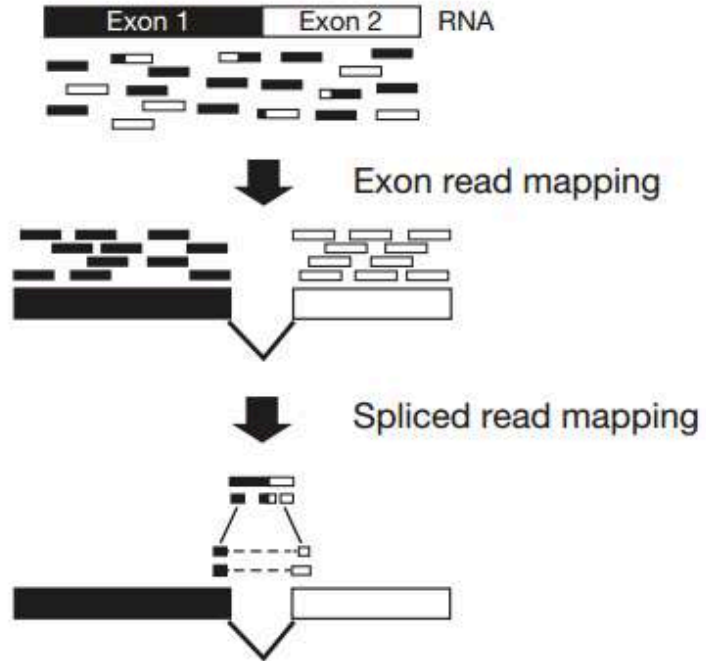
- Align and then assemble
- Assemble and then align

Align to  
Genome  
Transcriptome (if no genome)

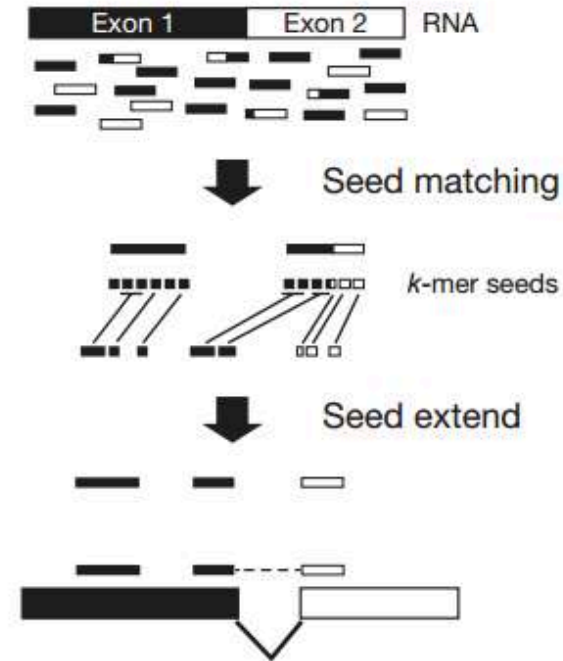


# Strategies for gapped alignments of RNAseq reads

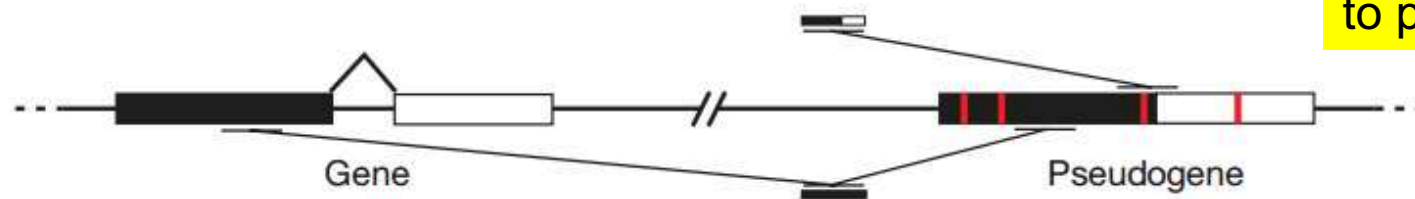
## a Exon-first approach



## b Seed-extend approach



## c Potential limitations of exon-first approaches

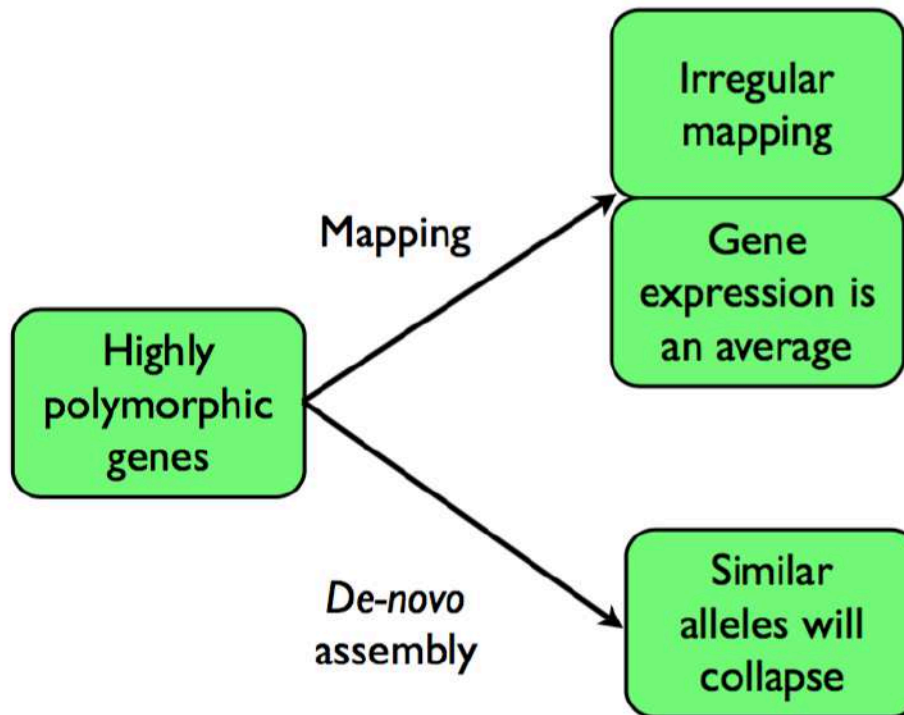


Preferential alignment  
(mismatch rather than split)  
to pseudogene

# A potential mapping problem

## High heterozygosity/Polyploid problem:

mRNA from species with a high heterozygosity or a polyploid genome can produce highly polymorphic reads for the same gene.



**Reference Gene I**

**ATGCGCGCTAGACGACATGACGACA**

**CACTTGACGACATGACG**      **Gene I A**

**CTTGACGACATGACGAC**

**CCCTTGACGACATGACG**      **Gene I B**

**CGCCCTTGACGACATGA**

**Expression Gene I = A + B**

**CACTTGACGACATGACG**      **Gene I A**

**CTTGACGACATGACGAC**

**CCCTTGACGACATGACG**      **Gene I B**

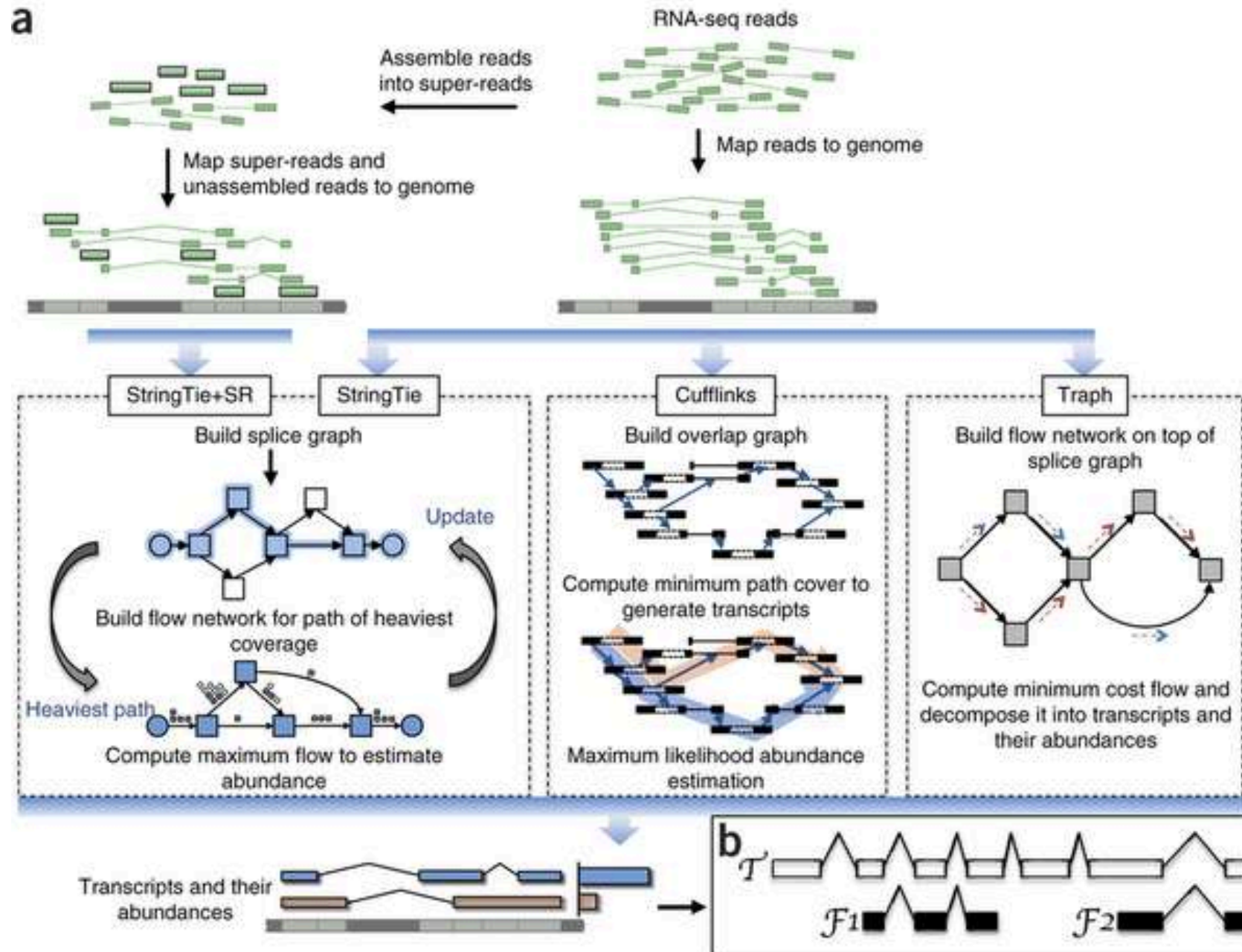
**CGCCCTTGACGACATGA**

**CGCCCTTGACGACATGACGACA**

**Collapsed consensus Gene A + Gene B**



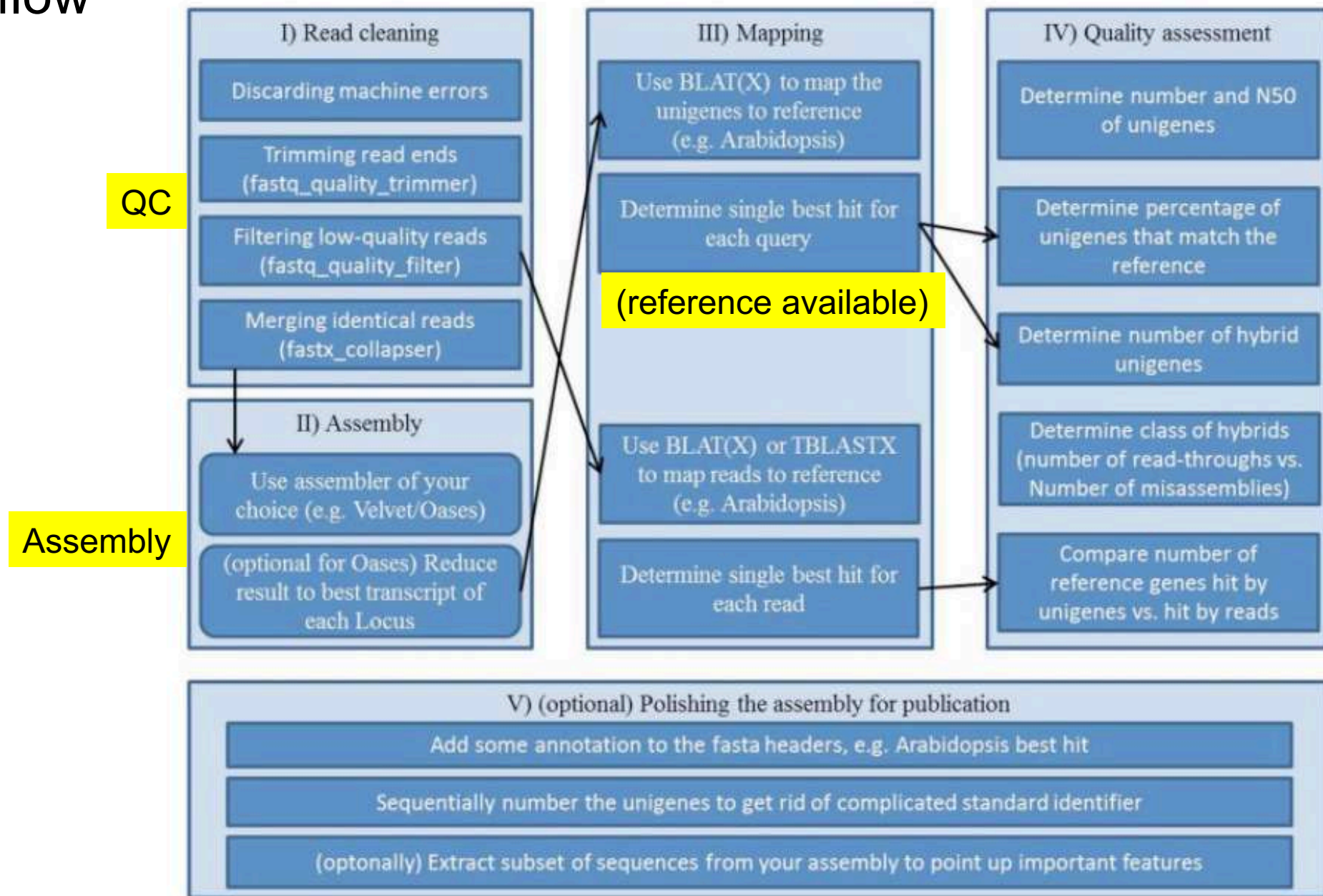
# Transcript reconstruction: Cufflinks and StringTie



*De novo* assembly of transcriptomes



# Workflow



# Transcriptome assembly benchmarks

Table 1: Benchmarking of Trans-**ABySS**, **IDBA**-tran, **SOAP**denovo-Trans, **Trinity**, and **SPAdes** on *M. musculus* RNA-seq dataset (accession number SRX648736, 11 million Illumina 100 bp long paired-end reads). The annotated transcriptome of *M. musculus* consists of 38924 genes and 94545 isoforms. The best values for each metric are highlighted with bold.

Assembler	ABySS	IDBA	SOAP	Trinity	SPAdes
Assembled transcripts	63871	38304	61564	47717	48876
Unaligned transcripts	232	98	273	160	817
Misassemblies	156	272	<b>35</b>	247	456
Database coverage, %	17.7	16.9	17.1	<b>18.4</b>	17.9
Duplication ratio	1.09	<b>1.004</b>	1.013	1.155	1.015
50%-assembled genes	6368	6562	6383	6695	<b>6972</b>
95%-assembled genes	1763	1572	1804	2251	<b>2391</b>
50%-assembled isoforms	6984	6795	6592	<b>7461</b>	7140
95%-assembled isoforms	1815	1572	1818	<b>2388</b>	2391

What is the best protocol for RNAseq analysis?

(Quick answer: no quick answer)

# What analysis combinations should we do?

ARTICLE

DOI: [10.1038/s41467-017-00050-4](https://doi.org/10.1038/s41467-017-00050-4)

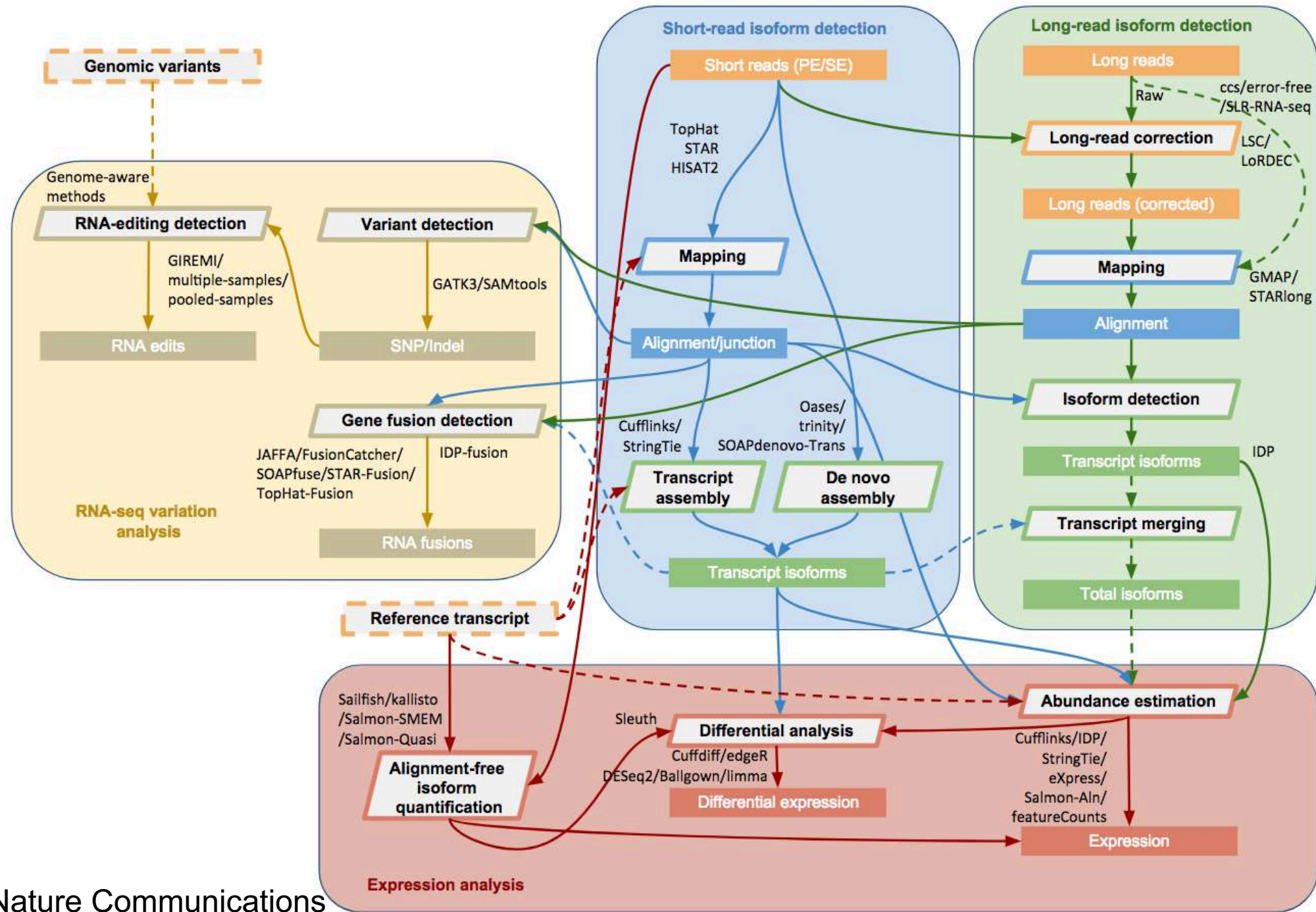
OPEN

## Gaining comprehensive biological insight into the transcriptome by performing a broad-spectrum RNA-seq analysis

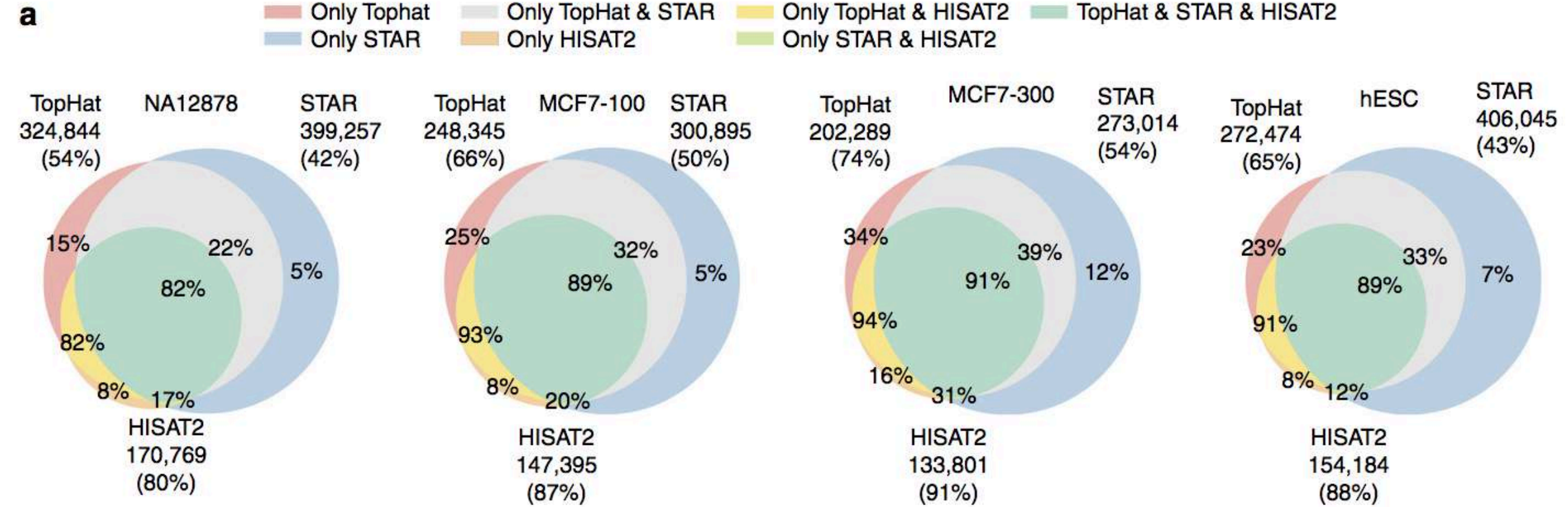
Sayed Mohammad Ebrahim Sahraeian<sup>1</sup>, Marghoob Mohiyuddin<sup>1</sup>, Robert Sebra<sup>2</sup>, Hagen Tilgner<sup>3</sup>, Pegah T. Afshar<sup>4</sup>, Kin Fai Au<sup>5</sup>, Narges Bani Asadi<sup>1</sup>, Mark B. Gerstein<sup>6</sup>, Wing Hung Wong<sup>7</sup>, Michael P. Snyder<sup>3</sup>, Eric Schadt<sup>2</sup> & Hugo Y.K. Lam<sup>1</sup>

... Here we conduct an extensive study analysing a broad spectrum of RNA-seq workflows. Surpassing the expression analysis scope, our work also includes **assessment of RNA variant-calling, RNA editing and RNA fusion detection techniques**. Specifically, **we examine both short- and long-read RNA-seq technologies, 39 analysis tools resulting in ~120 combinations**, and ~490 analyses involving 15 samples with a variety of germline, cancer and stem cell data sets.





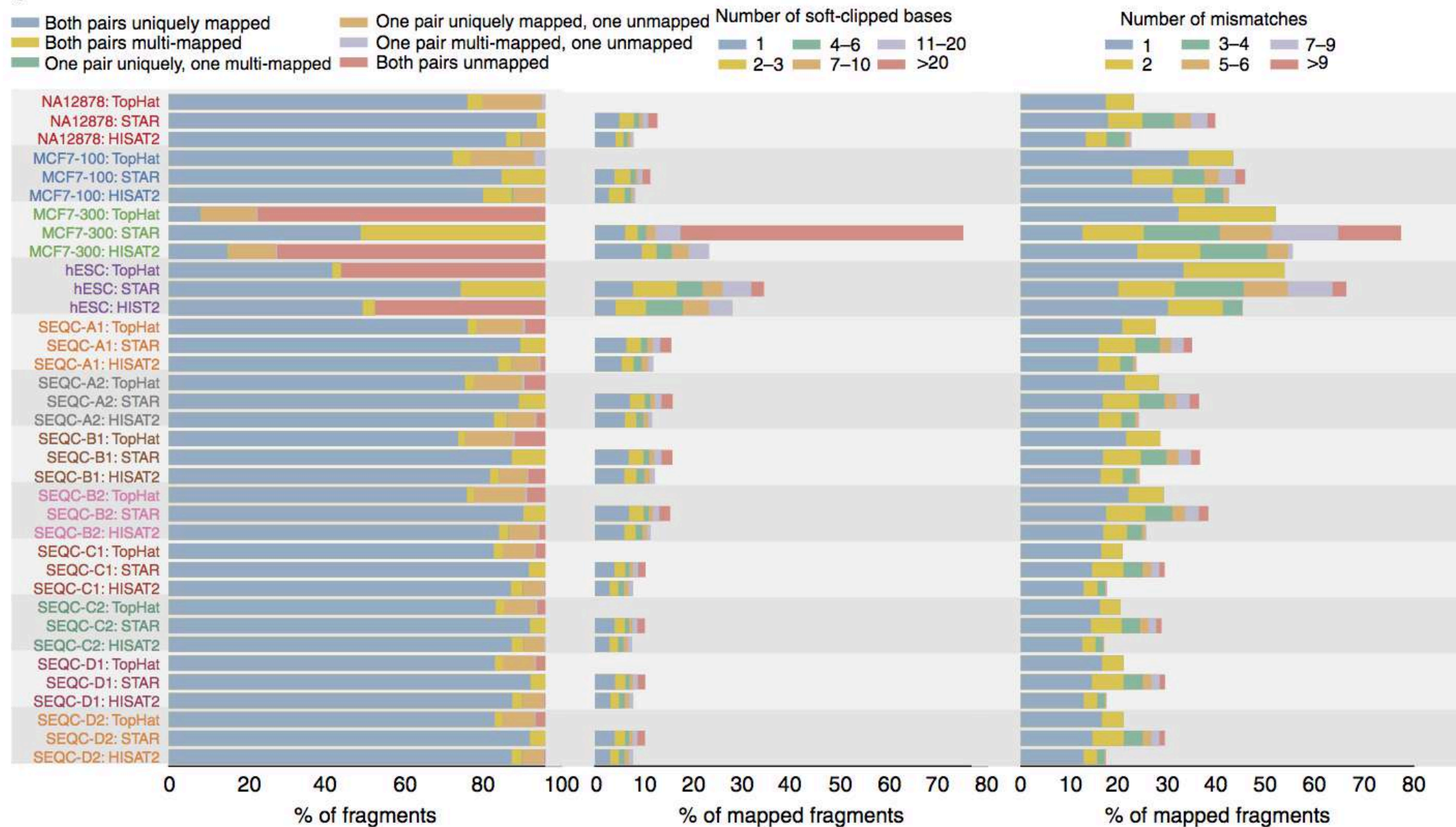
# Can be ~10% difference in mapping





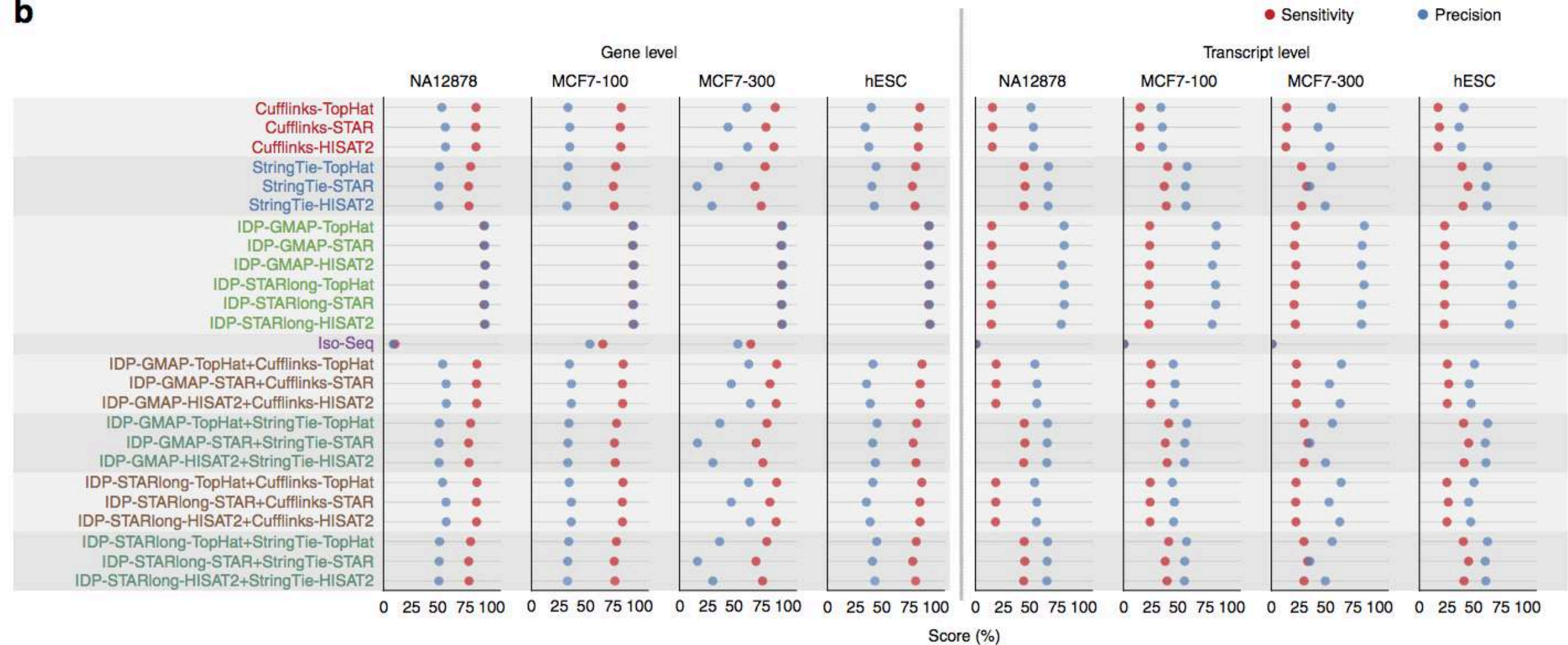
# Read mapping different amongst tools

**b**



# Performance of different transcriptome reconstruction schemes

**b**



# Typical RNAseq Workflow

## 2.3 Annotation

(focus only on gene annotation)

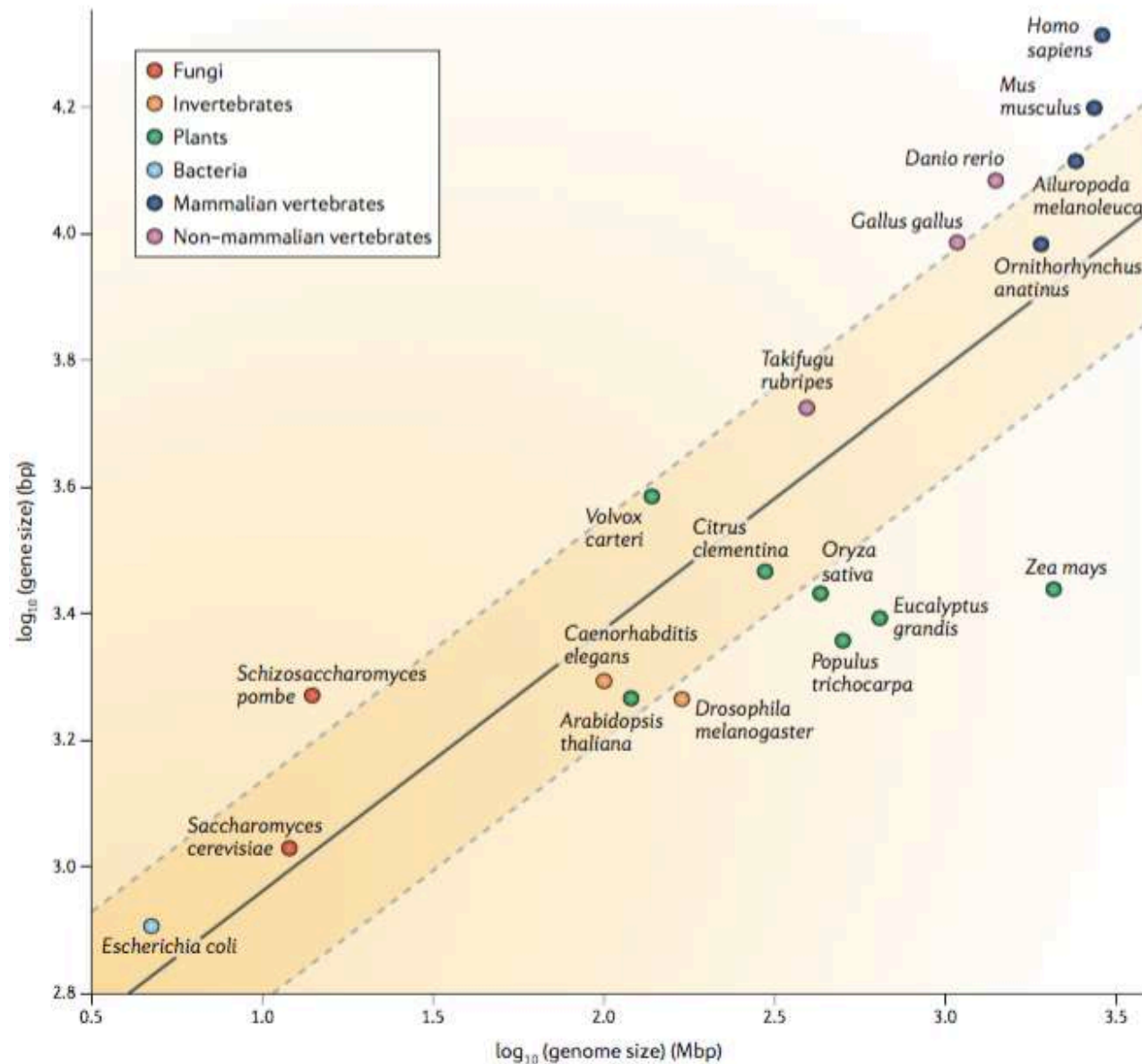


# A beginner's guide to eukaryotic genome annotation

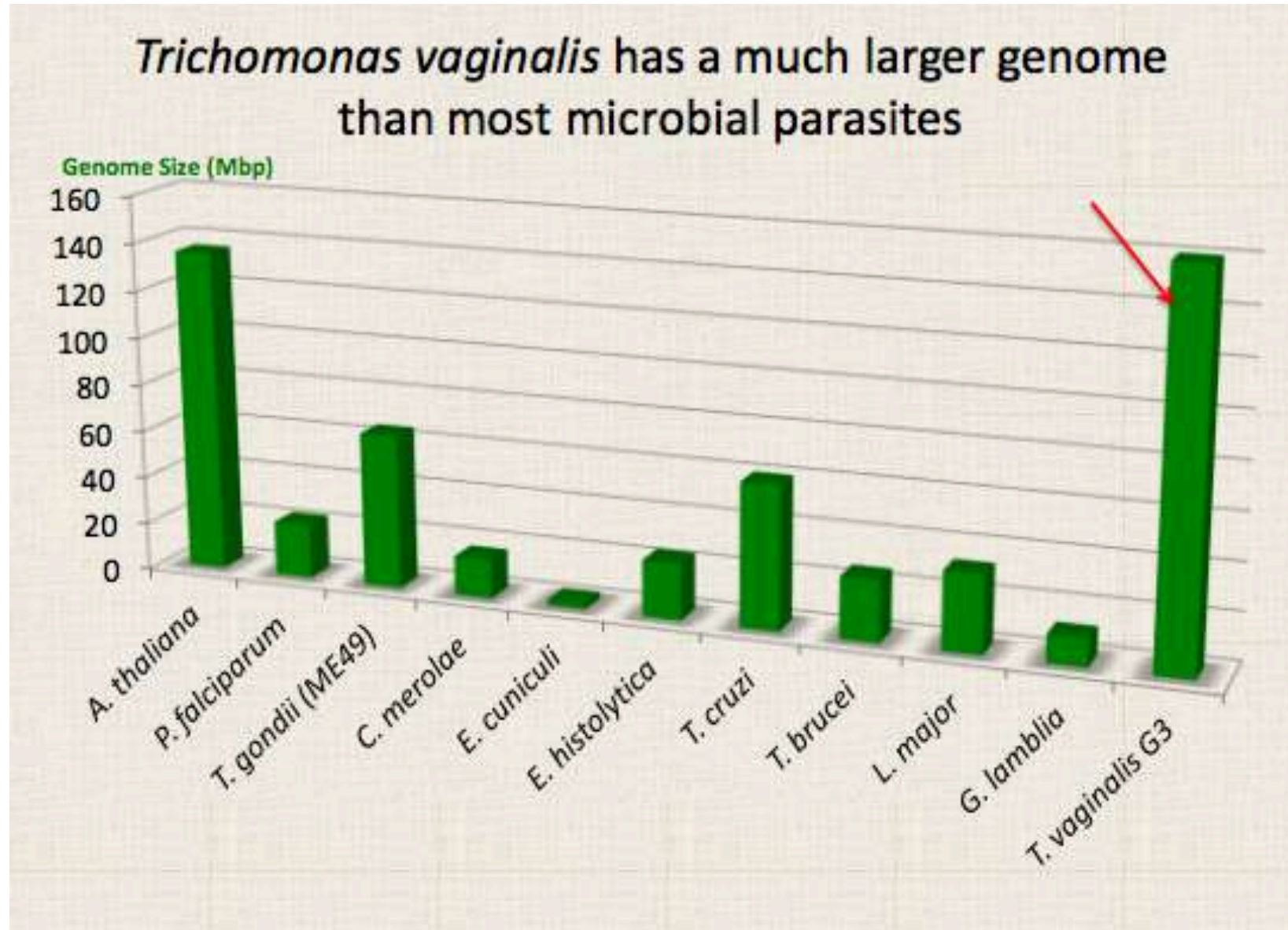
*Mark Yandell and Daniel Ence*

Abstract | The falling cost of genome sequencing is having a marked impact on the research community with respect to which genomes are sequenced and how and where they are annotated. Genome annotation projects have generally become small-scale affairs that are often carried out by an individual laboratory. Although annotating a eukaryotic genome assembly is now within the reach of non-experts, it remains a challenging task. Here we provide an overview of the genome annotation process and the available tools and describe some best-practice approaches.

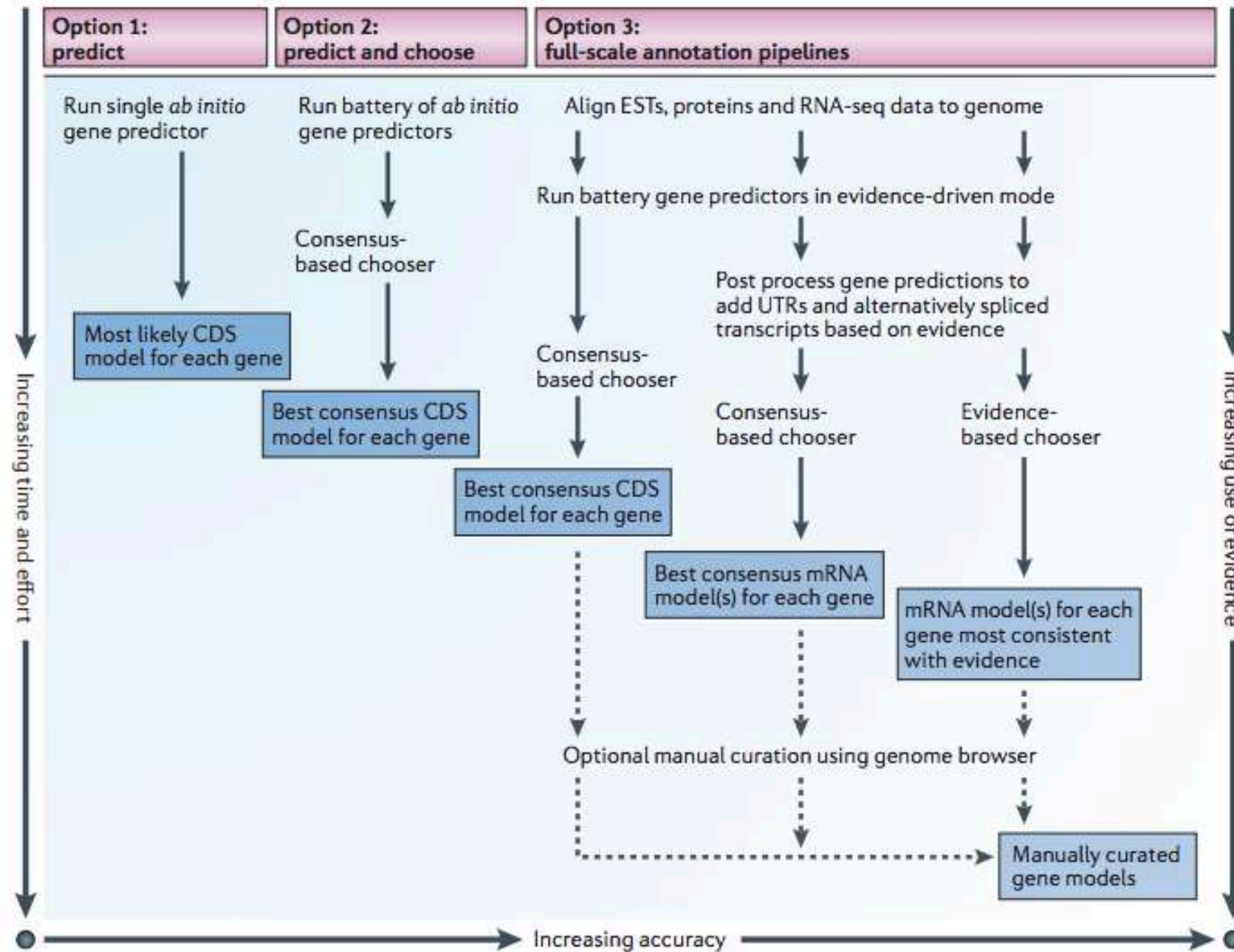
# Know your genome size (and gene numbers)



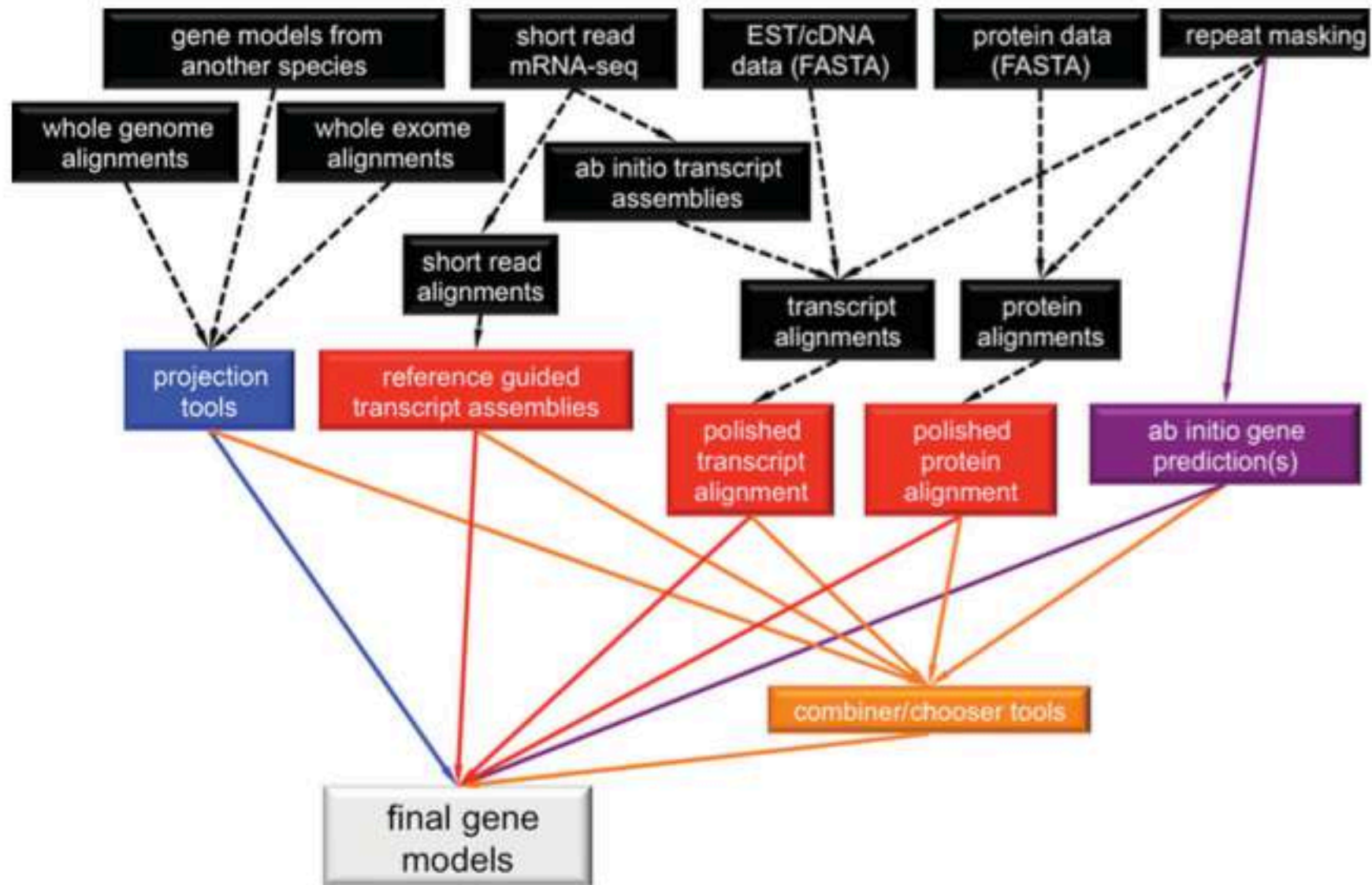
There's always exceptions (due to TE *Maverick* expansion)







# Multiple evidences; Update



# Basic rule of thumb

Just genome with no closely related species

Different *de novo* predictors, and combine them with combiners

Genome + closely related species + RNAseq

*de novo* predictors + evidence + combiners

**Genome + closely related species available + RNAseq**

*de novo* predictors + evidence + RNAseq evidence + combiners

**Genome + closely related species available + RNAseq + manual efforts**

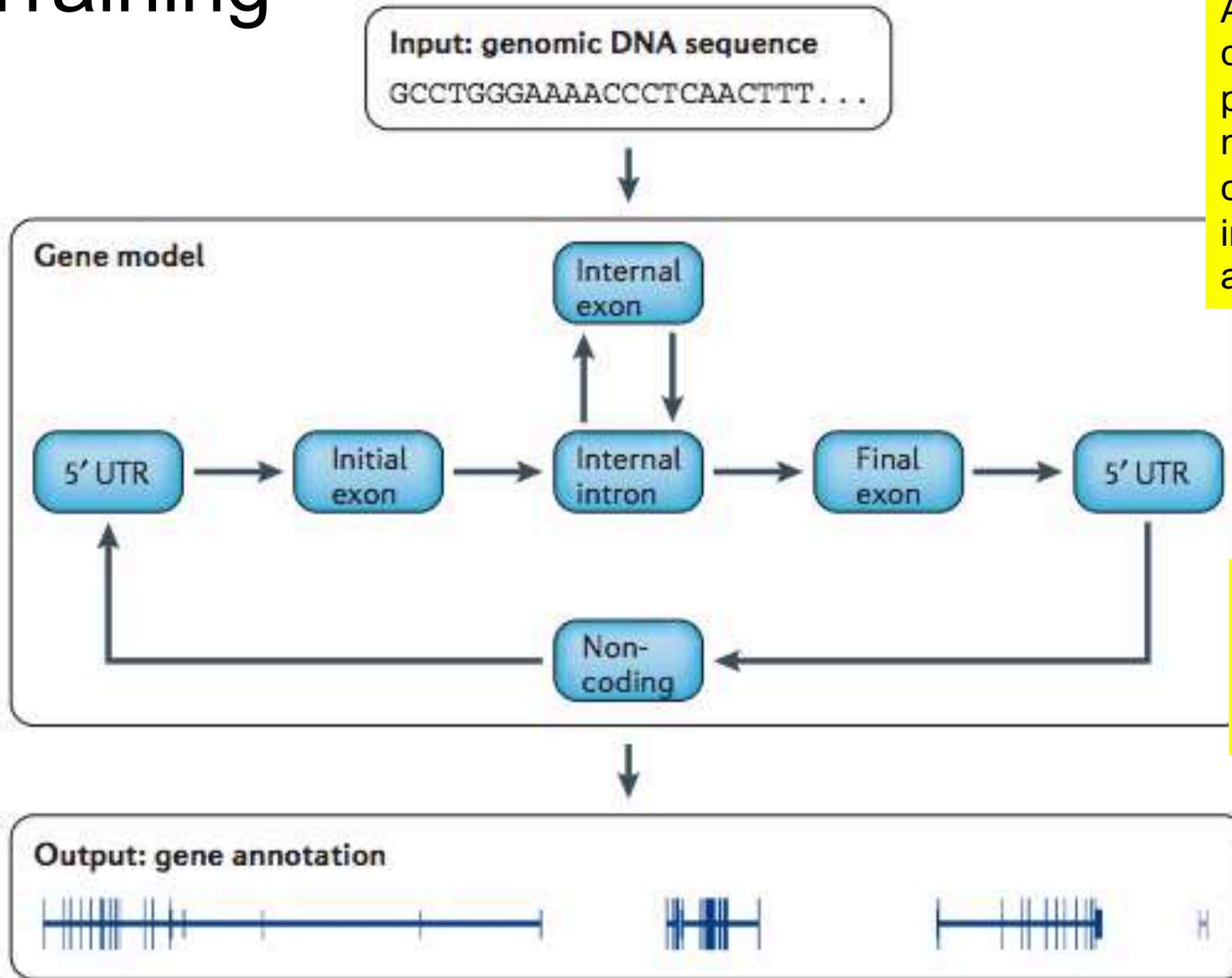
manual curation to train *de novo* predictors

Trained predictors + protein evidence + RNAseq evidence + combiners

Genome + initial annotations + RNAseq

protein evidence -> Trying to improve existing annotations

# Training



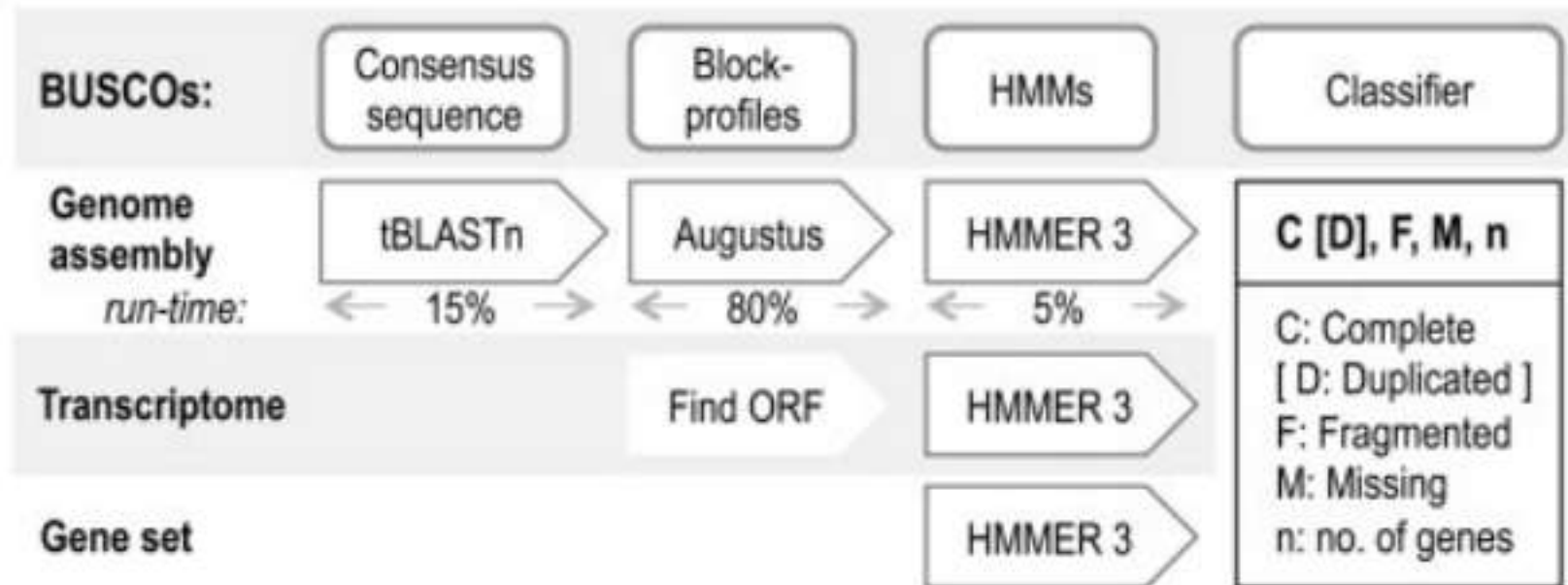
A simplified gene-finding model that captures the basic properties of a protein-coding gene is shown. The model takes the DNA sequence of a chromosome, or a portion thereof, as input and produces detailed gene annotations as output.

Note that this simplified model is incapable of identifying overlapping genes or multiple isoforms of the same gene. UTR, untranslated region.

# Where to find initial “correct” genes

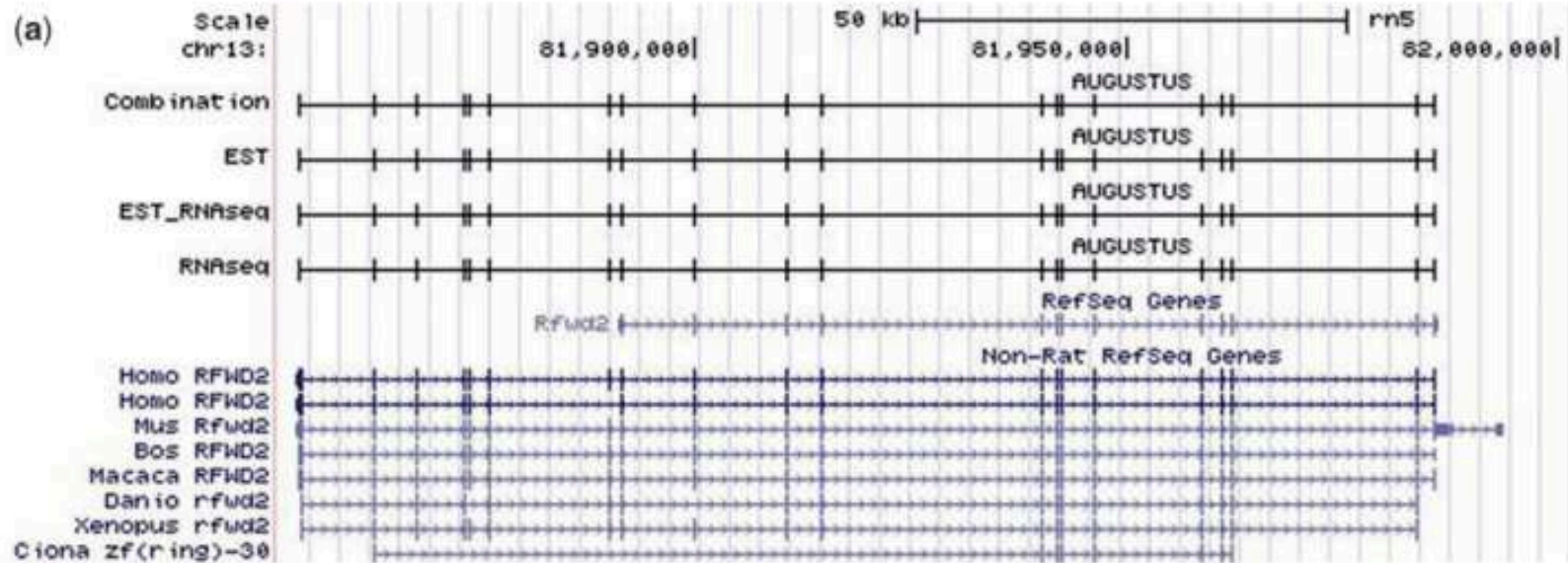
## **BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs**

Felipe A. Simão<sup>†</sup>, Robert M. Waterhouse<sup>†</sup>, Panagiotis Ioannidis, Evgenia V. Kriventseva and Evgeny M. Zdobnov\*



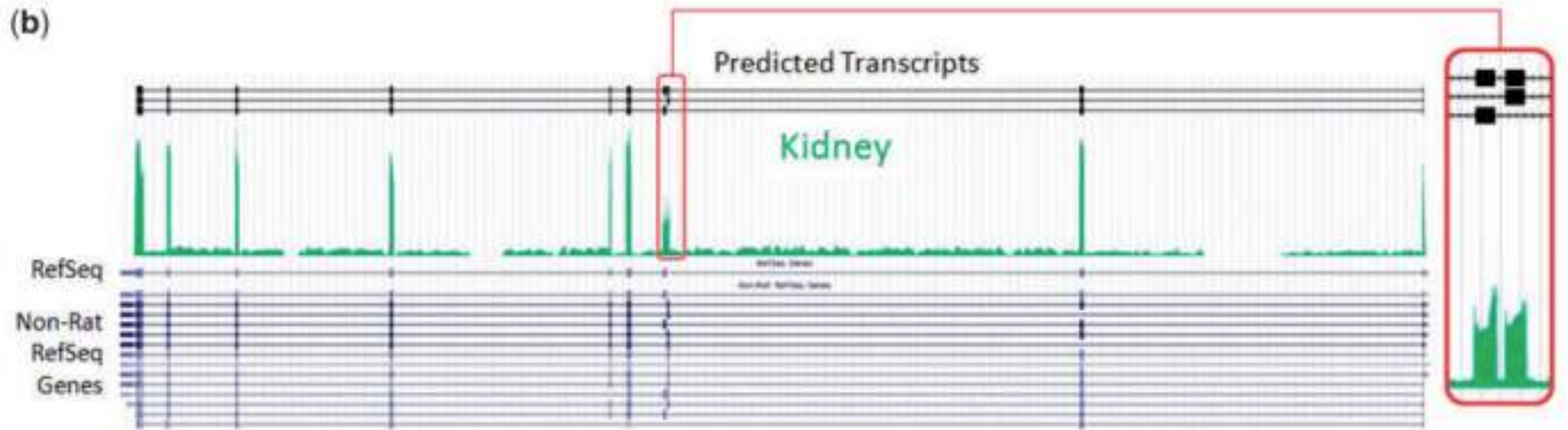


# Combine multiple evidence will improve annotation

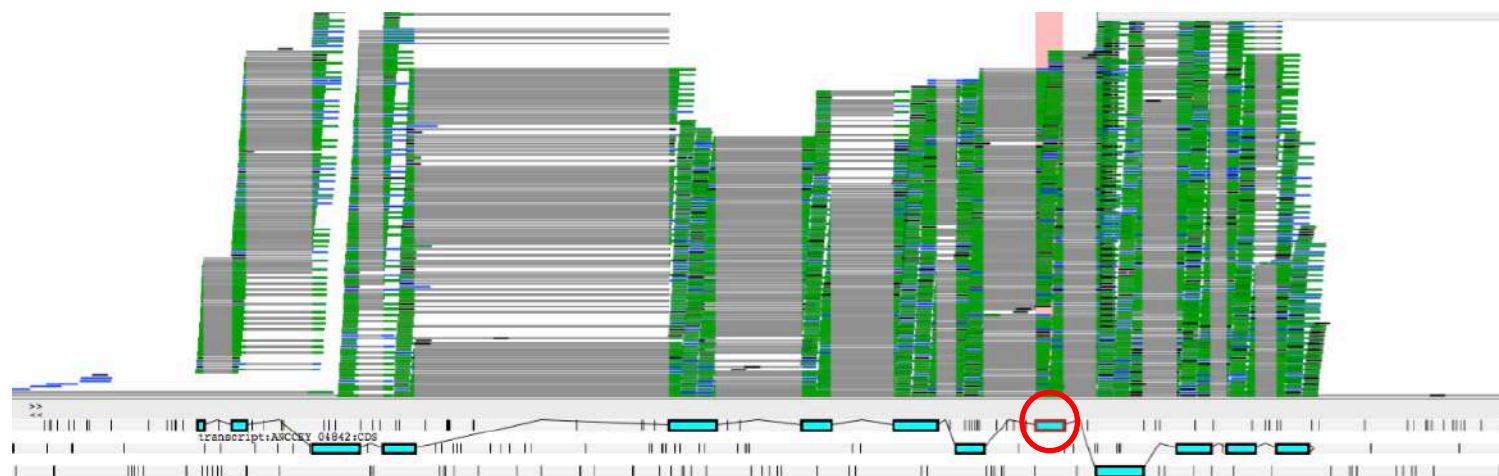
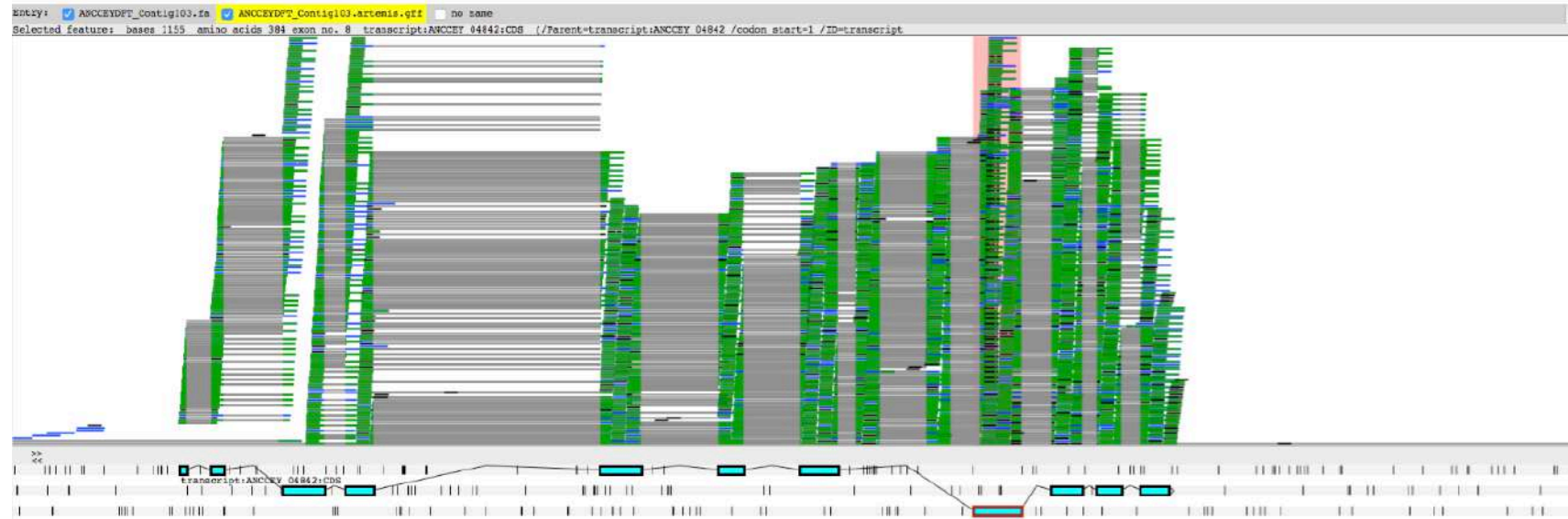




# Novel isoforms

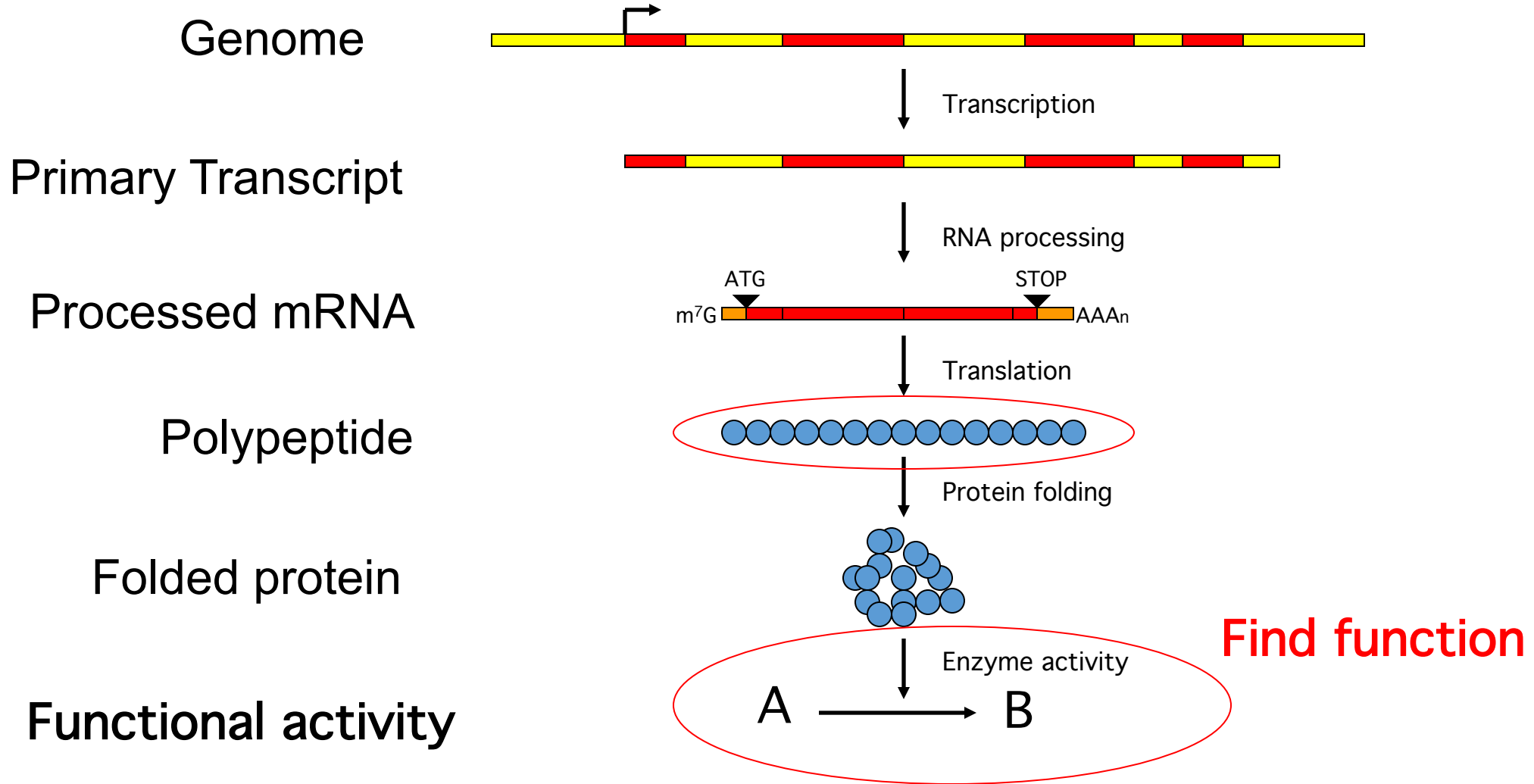


# Manual curation using artemis



Functional annotation

# Functional annotation



# Functional annotation

**Name** the protein correctly

Attaching biological information to genomic elements

- Biochemical function
  - Biological function
  - Involved regulation and interactions
  - Expression
- Utilize known **structural annotation to predicted protein sequence**

# Functional annotation – Homology Based

Most common way

Predicted Exons/CDS/ORF are searched against the non-redundant protein database (NCBI, SwissProt) to search for similarities

Visually assess the **top 5-10 hits** to identify whether these have been assigned a function

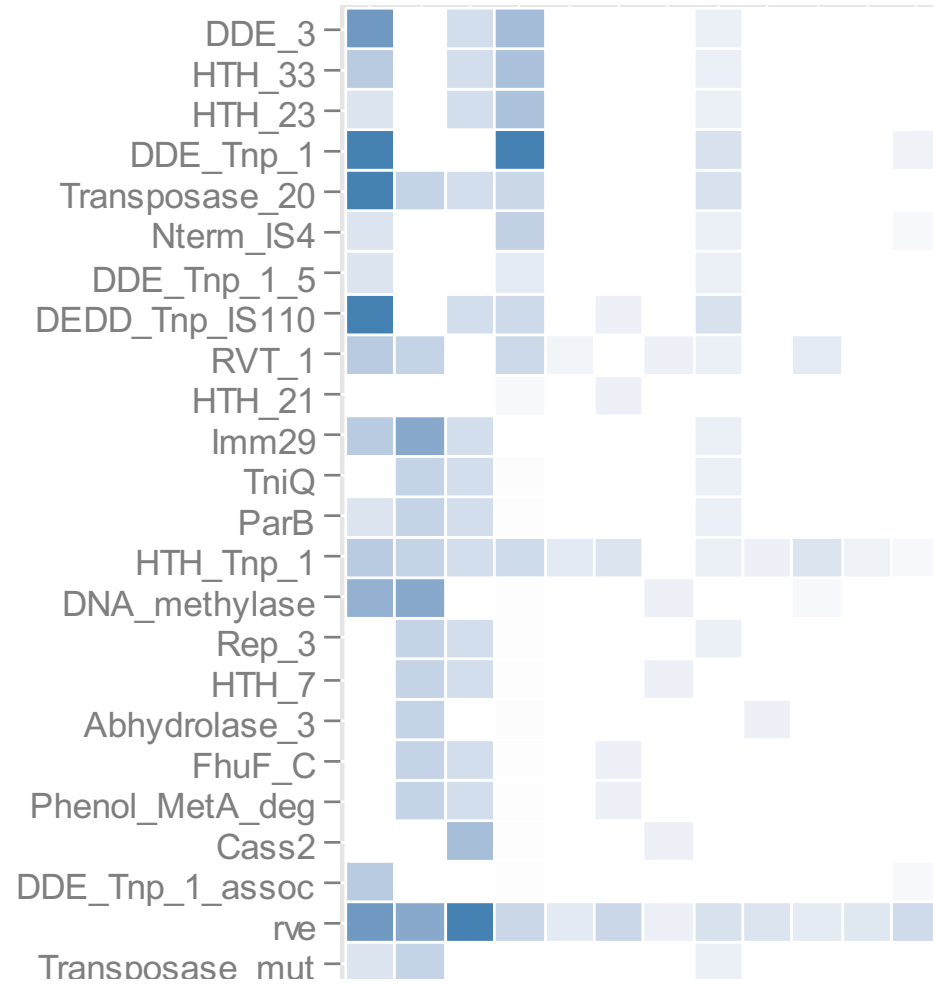
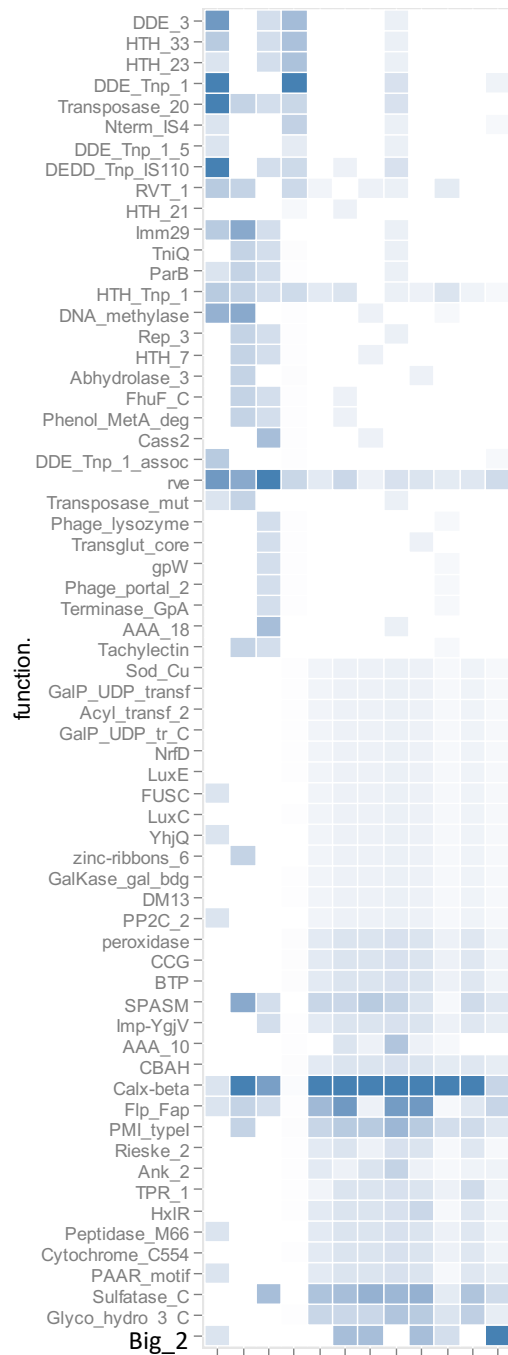
Functions (**and names**) are assigned



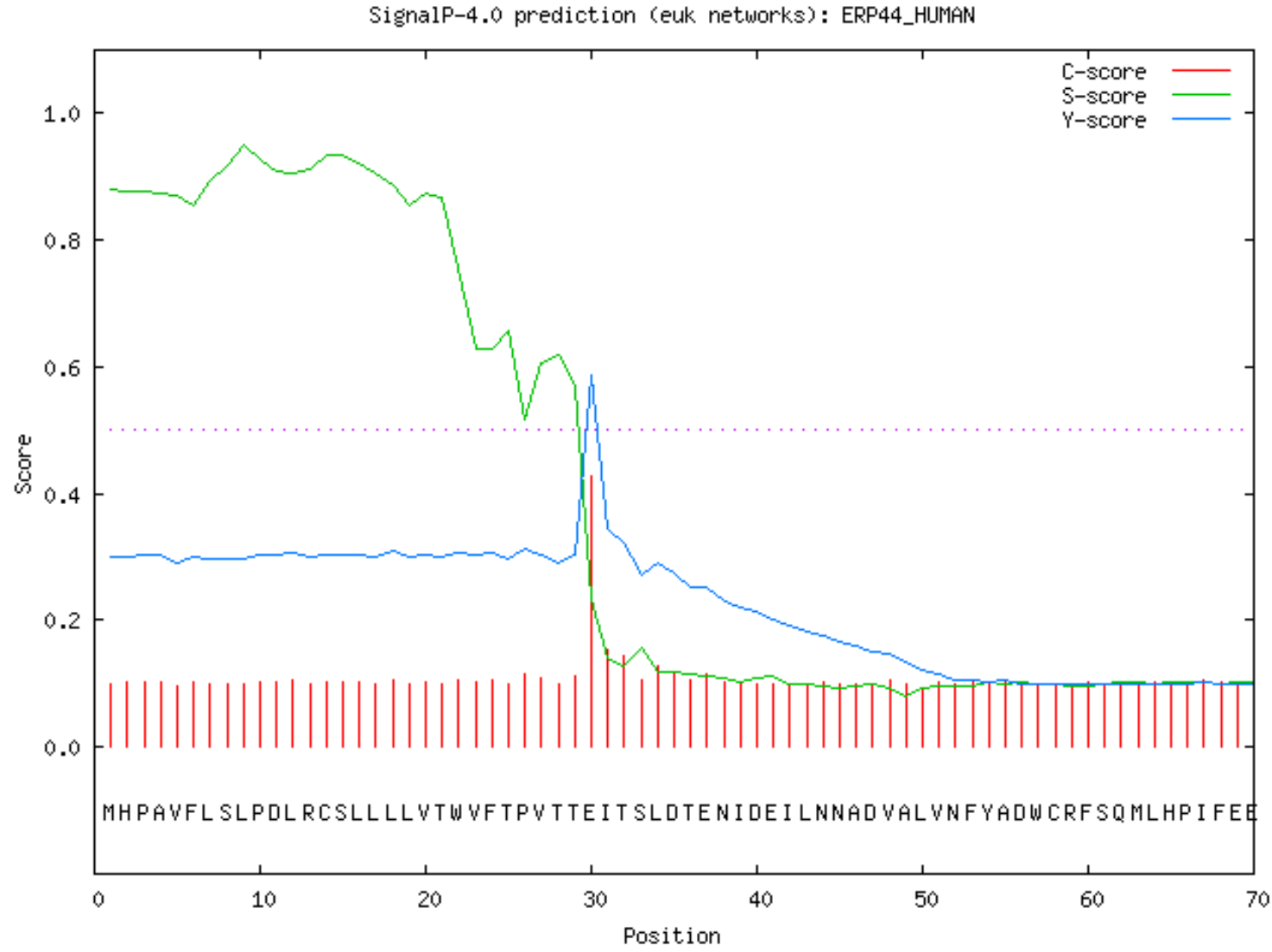
# Other features which can be determined

- Signal peptides
- Transmembrane domains
- Low complexity regions
- Various binding sites, glycosylation sites etc.
- Protein Domain
- Secretome

# PFAM

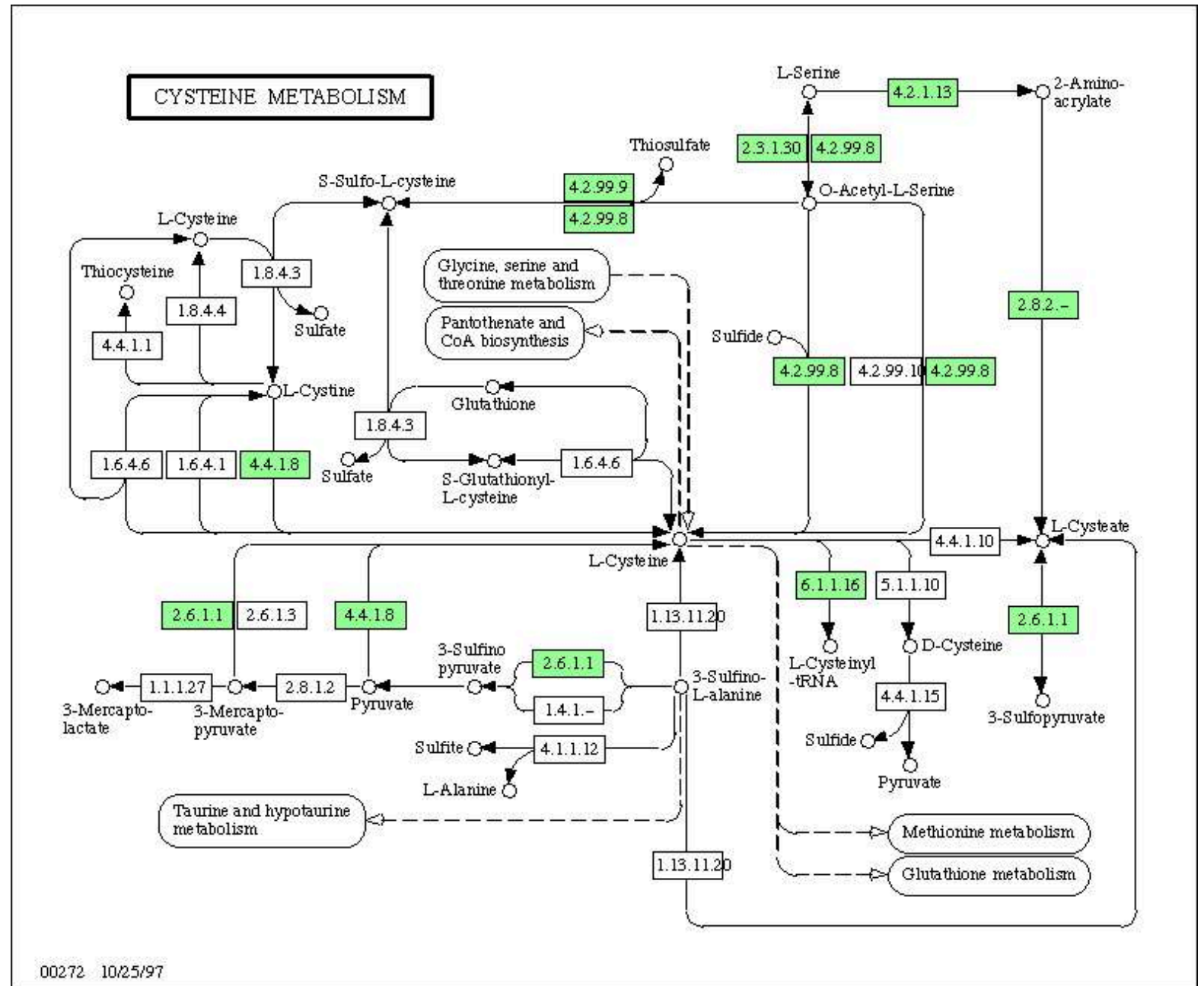


# SignalP: predicts the presence and location of signal peptide



# KEGG

Help improve annotation by showing missing genes in essential pathways



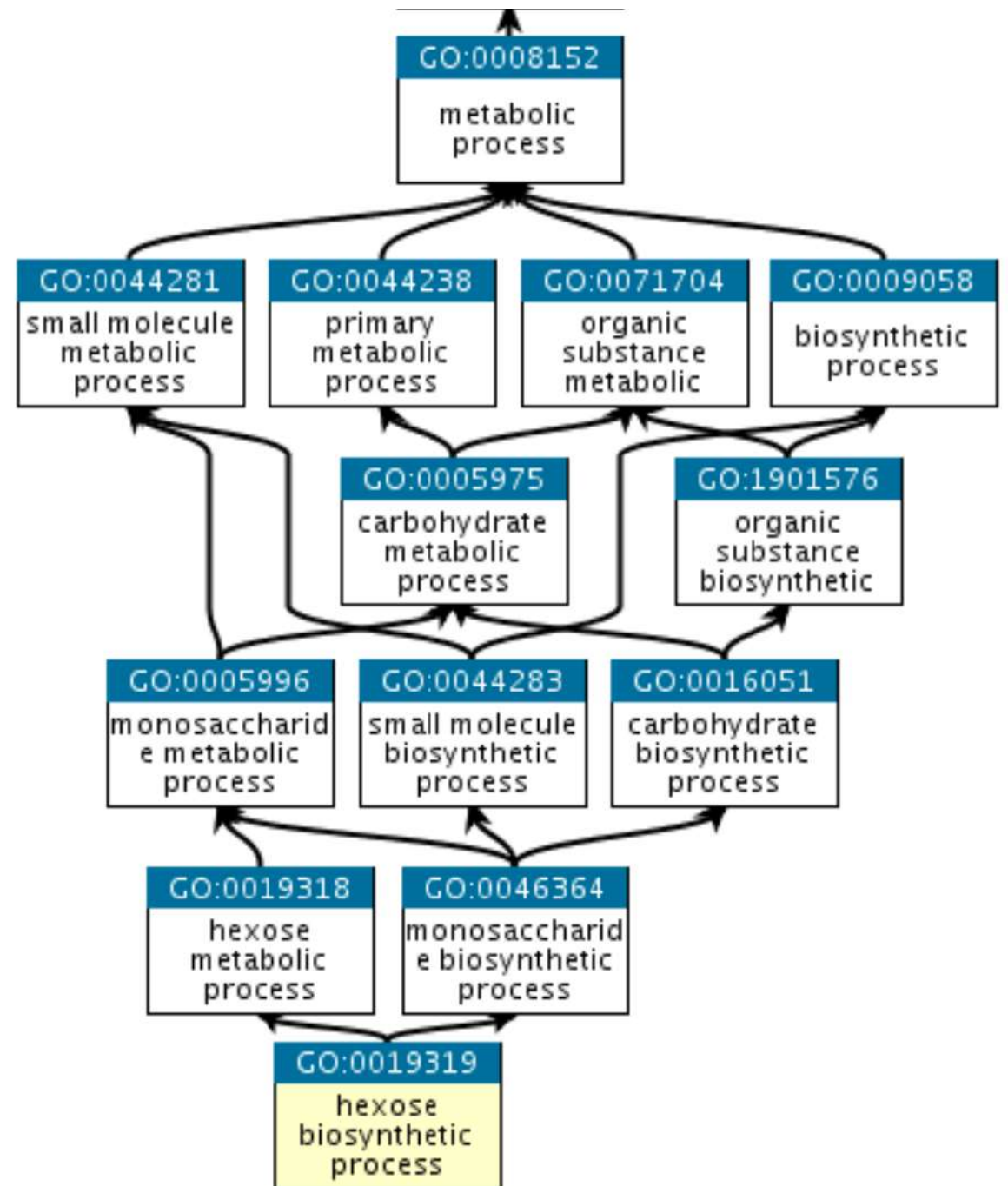
# Gene Ontology

- A controlled vocabulary for annotating three aspects of a gene product's biology:
- **Biological Process** (BP) – the molecular, cellular, and organismal level processes in which a gene product is involved
- **Molecular Function** (MF) – the molecular activity of a gene product
- **Cellular Component** (CC) – the subcellular localization of a gene product

# Gene Ontology

“An ontology is a formal representation of a body of knowledge within a given domain. Ontologies usually consist of a set of classes (or terms or concepts) with relations that operate between them.”

“GO is loosely hierarchical, with ‘child’ terms being more specialized than their ‘parent’ terms, but unlike a strict hierarchy, a term may have more than one parent term.”





# BLAST2GO

The screenshot displays the Blast2GO PRO application window. At the top, there is a navigation bar with icons for 'start', 'blast', 'interpro', 'mapping', 'annot', 'charts', 'graphs', and 'select'. Below this is a table listing sequences with columns for 'nr', 'SeqName', 'Description', 'Length', '#Hits', 'e-Value', 'sim mean', '#GO', 'GO list', 'Enzyme list', and 'InterPro Scan'. A 'Run Blast' dialog box is open in the foreground, titled 'Run Blast' and 'Blast Options'. It asks the user to choose one option: 'CloudBlast', 'NCBI Blast', 'AWS Blast', or 'Local Blast'. Each option has a brief description and a logo. The 'NCBI Blast' option is currently selected. At the bottom of the dialog are buttons for 'Default', '< Back', 'Next >', 'Cancel', and 'Run'. Below the dialog, a 'Welcome Message' tab is active, showing 'Blast Result of C04018A12'. It displays 'Query Name (Length): C04018A12 (715)', 'E-Value Cutoff: 0.001', and 'Annotation: -'. Below this is a table titled 'Sequences Producing Significant Alignments' with columns for 'Gene Name' and 'Sequence'. The table lists several sequences with their accession numbers and descriptions, including 'chitinase CHI1 [Citrus sinensis]', 'class IV chitinase [Galega orientalis]', and several 'PREDICTED: hypothetical protein [Vitis vinifera]' entries.

nr	SeqName	Description	Length	#Hits	e-Value	sim mean	#GO	GO list	Enzyme list	InterPro Scan
1	C0401...	mpk3_arath ame: full-mitogen-activated prote...	717	20	5.3E-144	87.3%	0	-	-	-
2	C0401...	protein	708							
3	C0401...	protein	620							
4	C0401...	class iv chitinase	715							
5	C0401...	cytl_vigun ame: full-cysteine proteinase inhibi...	863							
6	C0401...	protein phosphatase 2c	863							
7	C0401...	protein	578							
8	C0401...	lgul_orysj ame: full-lactoylglutathione lyase a...	800							
9	C0401...	mt2_actde ame: full-metallothionein-like prote...	625							
10	C0401...	protein	612							
11	C0401...	protein phosphatase	645							

Gene Name	Sequence
	gi 3608477 gb AAC35981.1 chitinase CHI1 [Citrus sinensis]
	gi 33414046 gb AAP03085.1 class IV chitinase [Galega orientalis]
	gi 225434068 ref XP_002275122.1 PREDICTED: hypothetical protein [Vitis vinifera]
	gi 157353727 emb CAO46259.1 unnamed protein product [Vitis vinifera]
LOC100250948	gi 225434052 ref XP_002274620.1 PREDICTED: hypothetical protein [Vitis vinifera]
	gi 157353719 emb CAO46251.1 unnamed protein product [Vitis vinifera]

Europe, Germany: DE2 Version: b2g\_sep14 /Users/sgoetz/b2gWorkspace/blast2go\_project.20141104\_2248.dat

# Case study: eukaryote annotation (2018)

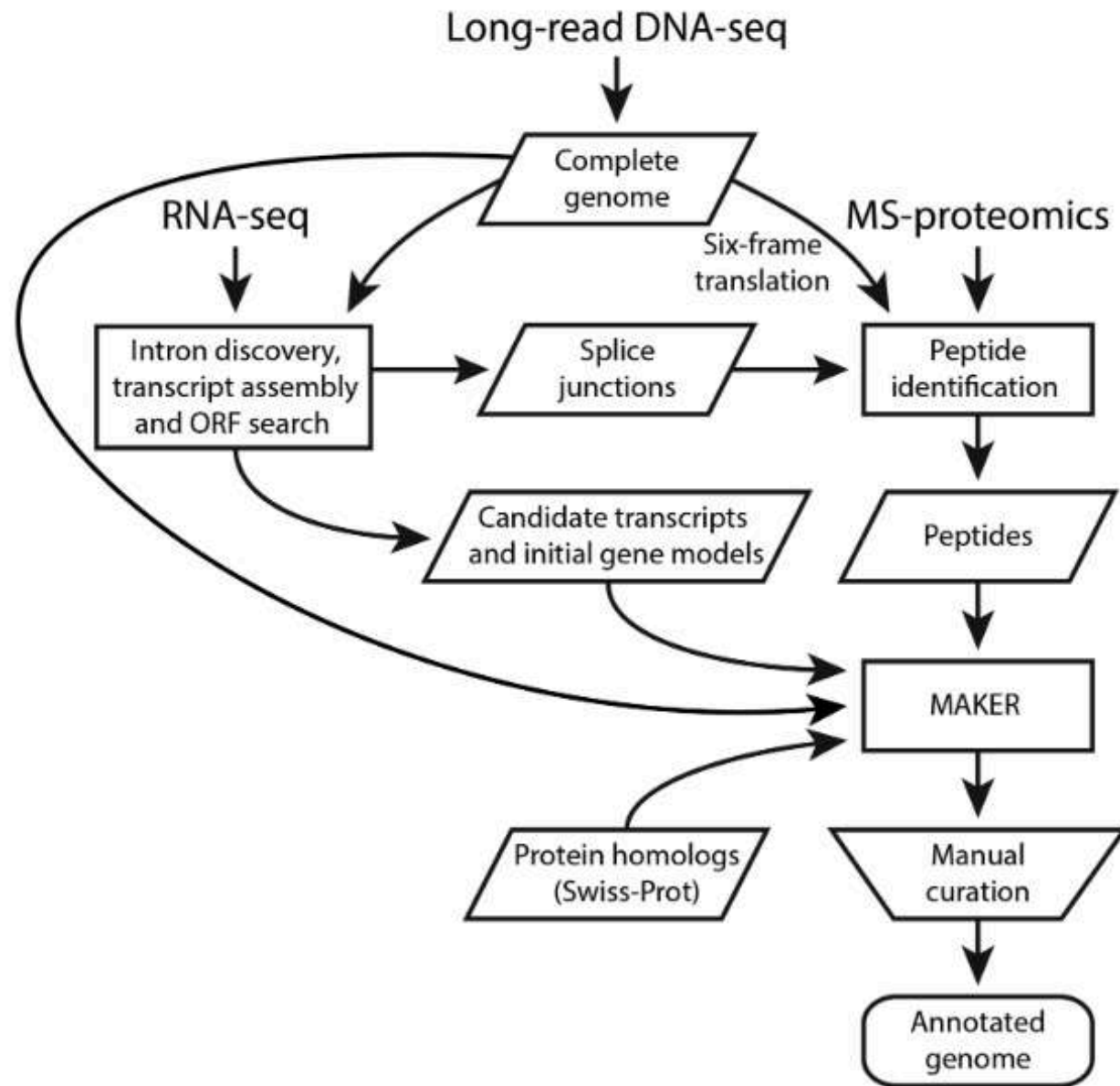
Published online 18 January 2017

*Nucleic Acids Research*, 2017, Vol. 45, No. 5 2629–2643

doi: 10.1093/nar/gkx006

## **Proteogenomics produces comprehensive and highly accurate protein-coding gene annotation in a complete genome assembly of *Malassezia sympodialis***

**Yafeng Zhu<sup>1,†</sup>, Pär G. Engström<sup>2,†</sup>, Christian Tellgren-Roth<sup>3</sup>, Charles D. Baudo<sup>4</sup>, John C. Kennell<sup>4</sup>, Sheng Sun<sup>5</sup>, R. Blake Billmyre<sup>5</sup>, Markus S. Schröder<sup>6</sup>, Anna Andersson<sup>7</sup>, Tina Holm<sup>7</sup>, Benjamin Sigurgeirsson<sup>8</sup>, Guangxi Wu<sup>9</sup>, Sundar Ram Sankaranarayanan<sup>10</sup>, Rahul Siddharthan<sup>11</sup>, Kaustuv Sanyal<sup>10</sup>, Joakim Lundeberg<sup>8</sup>, Björn Nystedt<sup>12</sup>, Teun Boekhout<sup>13</sup>, Thomas L. Dawson, Jr.<sup>14</sup>, Joseph Heitman<sup>5</sup>, Annika Scheynius<sup>15,\*</sup>,<sup>‡</sup> and Janne Lehtiö<sup>1,\*</sup>,<sup>‡</sup>**

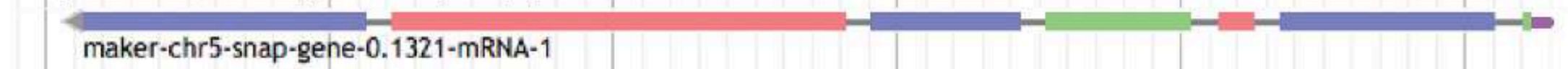




Current annotation



MAKER prediction (homology, RNA-seq and peptides)



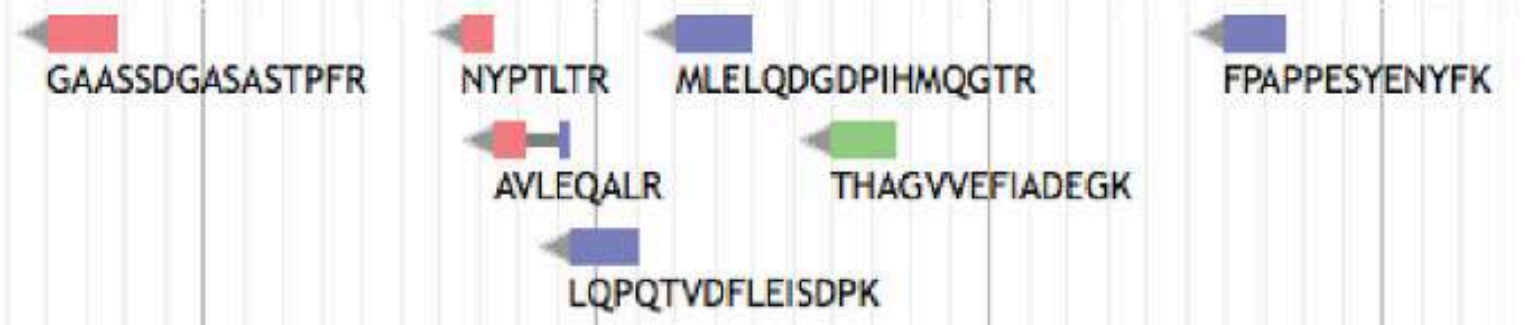
Gioti *et al* annotation



RNA-seq coverage



Peptide evidence



**Table 2.** Characteristics of *M. sympodialis* gene sets

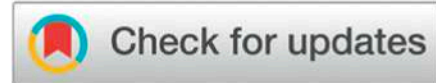
	Published (MAKER with homology evidence) (5)	MAKER with homology and RNA-seq evidence	MAKER with homology, RNA-seq and peptide evidence	Manually curated annotation
Protein-coding genes	3536	3612	4113	4493
Gene density (genes/kb) <sup>1</sup>	0.46	0.46	0.53	0.58
Coding sequence (Mb)	5.40	5.35	6.14	6.72
Coding exons	6995	8453	9212	9793
Introns	3462	5030	5267	5350
Mean exon size (bp) <sup>2</sup>	772	635	669	687
Mean intron size (bp)	65	52	50	30
Genes supported by peptides	3176	3176	3674	3891
Introns supported by RNA-seq	1661 (48%)	4246 (84%)	4275 (81%)	5271 (99%)
Out-frame peptides	4658 (13%)	5453 (15%)	1796 (5%)	338 (1%)

<sup>1</sup>Gene density was computed relative to the size of the corresponding genome assembly (7.71 Mb for the draft assembly of Gioti *et al.* (5) and 7.79 Mb for the current assembly).

<sup>2</sup>Excluding untranslated regions.



# Case study: Annotation using long reads



**Breakthrough Technologies**

---

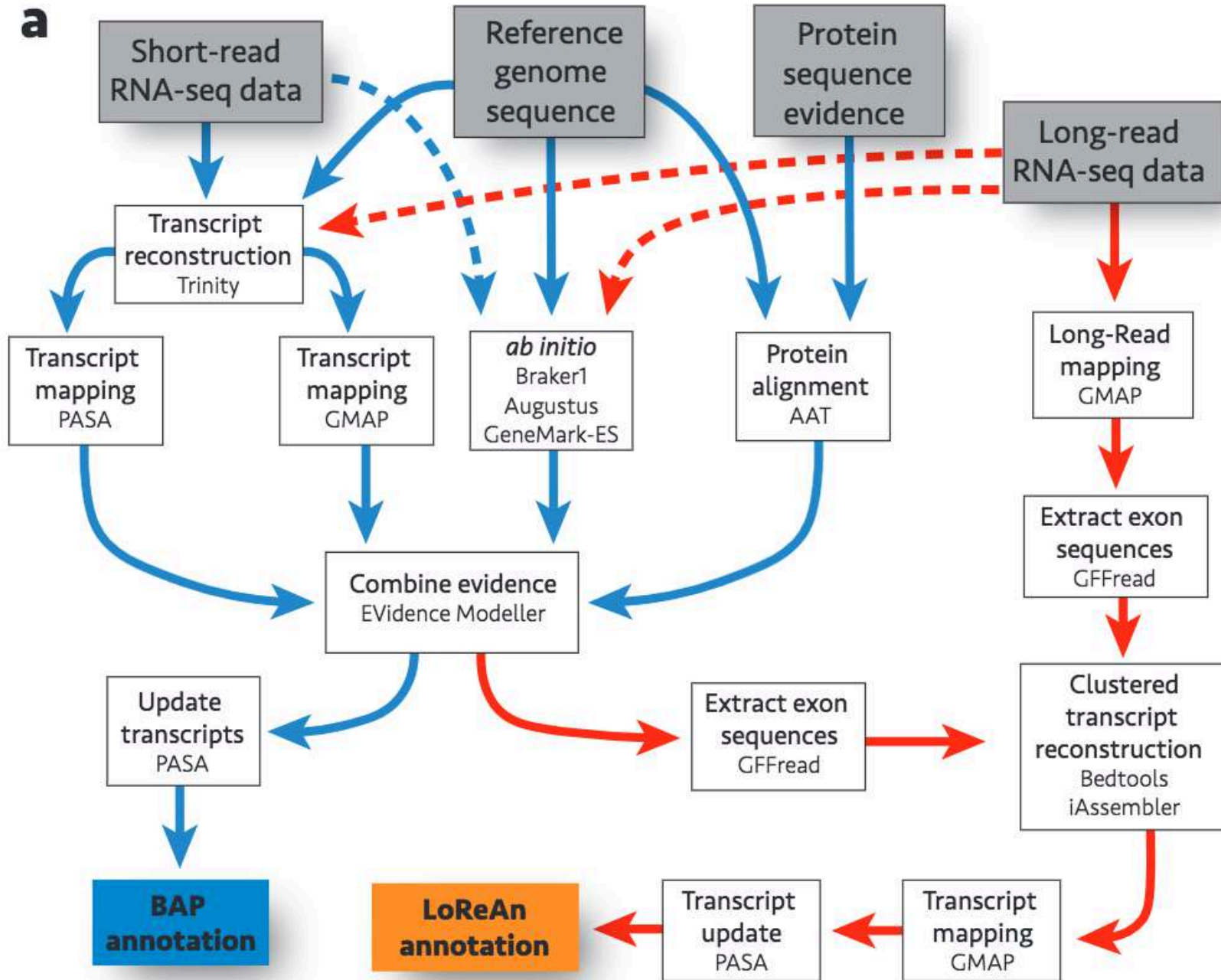
## **Long-Read Annotation: Automated Eukaryotic Genome Annotation Based on Long-Read cDNA Sequencing<sup>1</sup>[OPEN]**

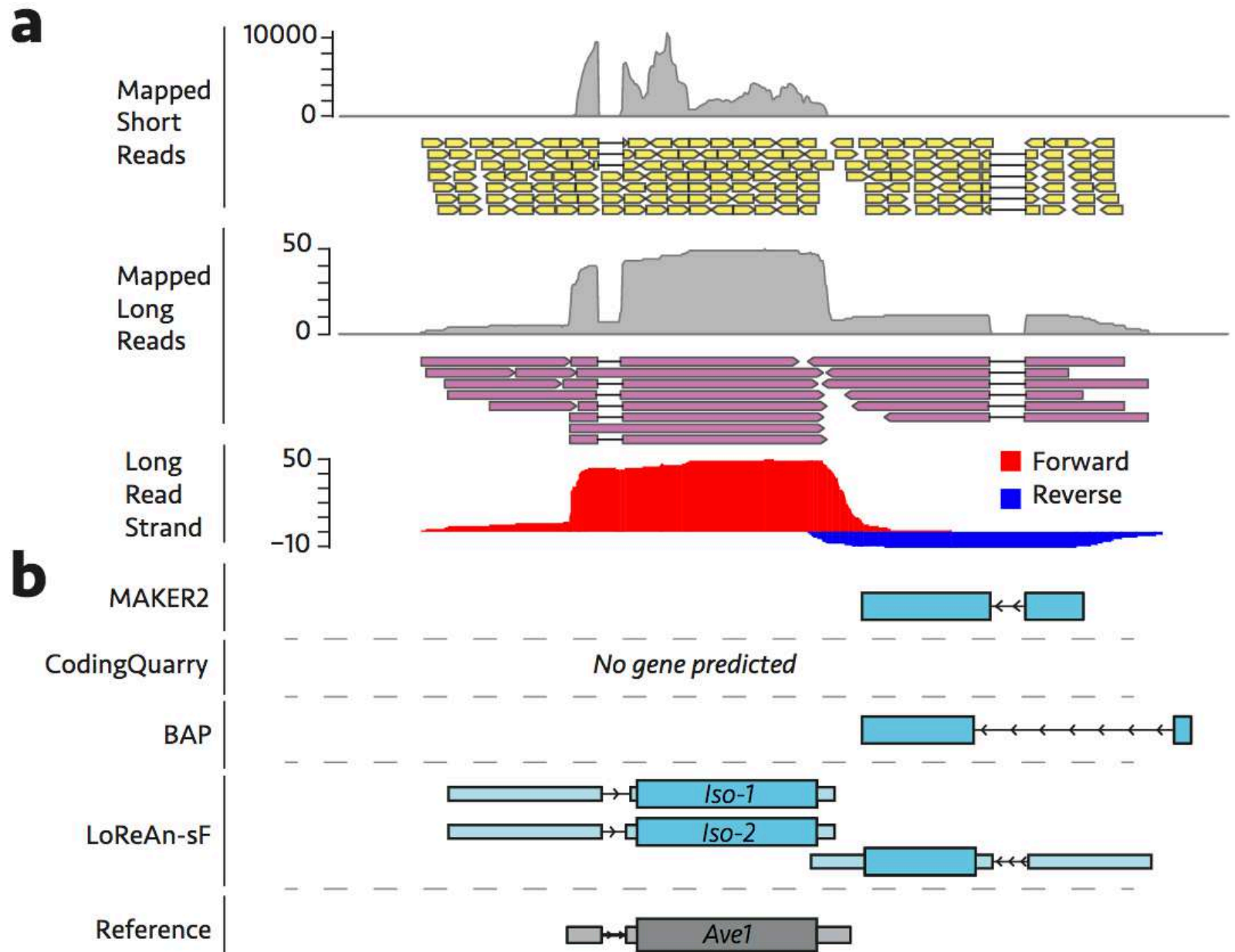
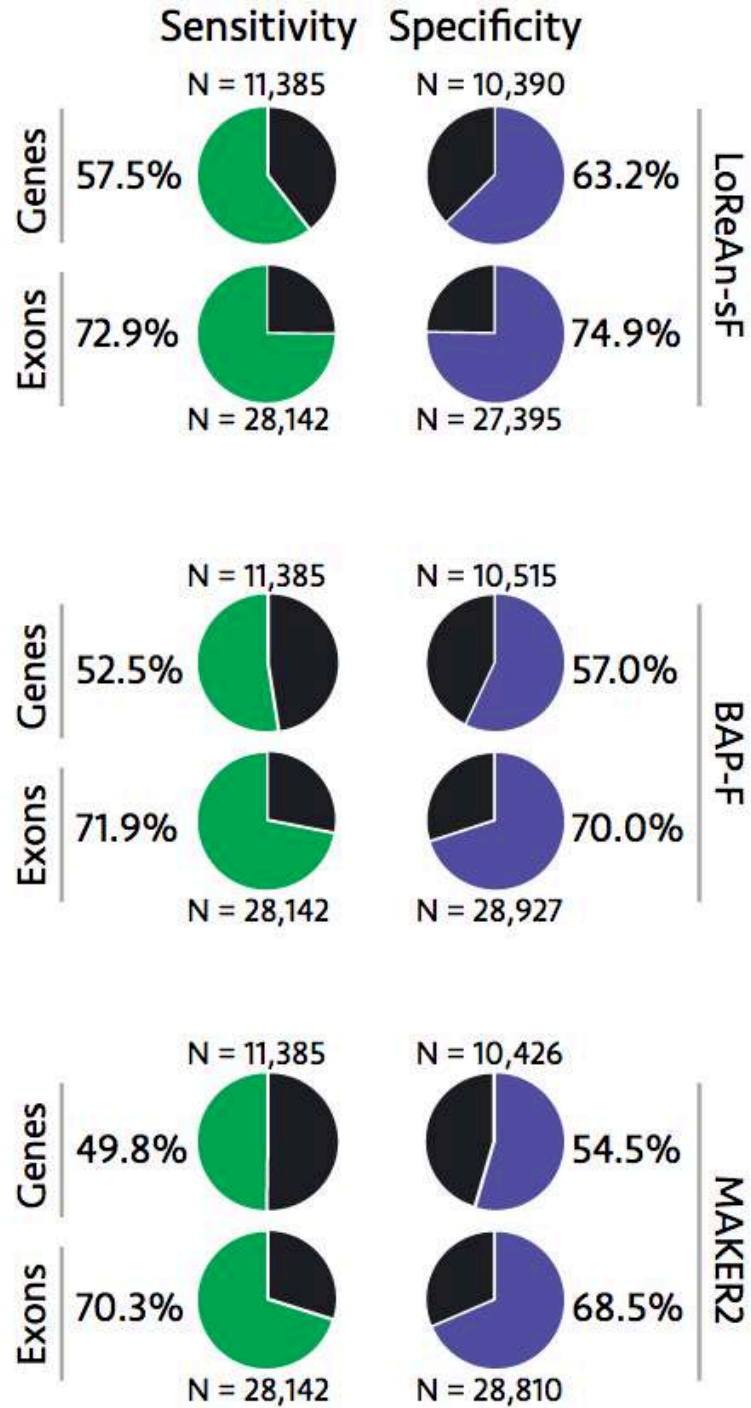
**David E. Cook,<sup>a,2,3</sup> Jose Espejo Valle-Inclan,<sup>a,2,4</sup> Alice Pajoro,<sup>b,5</sup> Hanna Rovenich,<sup>a,6</sup>  
Bart P. H. J. Thomma,<sup>a,7,8,9</sup> and Luigi Faino<sup>a,10,7</sup>**

<sup>a</sup>Laboratory of Phytopathology, Wageningen University and Research, Droevendaalsesteeg 1, 6708 PB Wageningen, the Netherlands

<sup>b</sup>Laboratory of Molecular Biology, Wageningen University and Research, Droevendaalsesteeg 1, 6708 PB Wageningen, the Netherlands

ORCID IDs: 0000-0002-2719-4701 (D.E.C.); 0000-0002-4857-5984 (J.E.V.); 0000-0003-4125-4181 (B.P.H.J.T.); 0000-0002-6807-4191 (L.F.).





Break

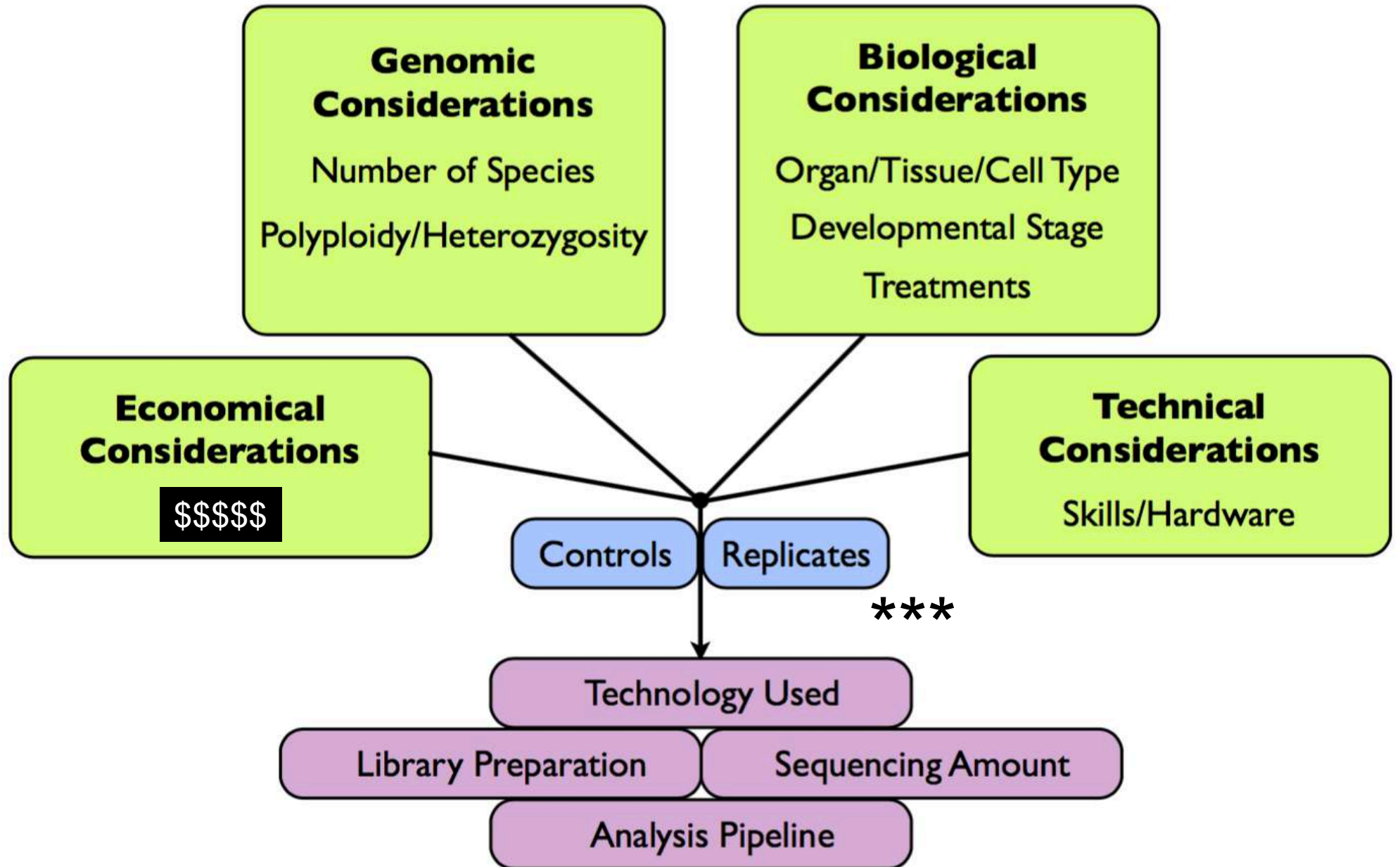
### 3. Differential expression



# Types of experiments

- One genome or multiple genomes (Host / pathogens)
- Multiple alleles
  - High heterozygosity
  - Polyploidy
  - Gene families
- Isoforms?
- Organ / Tissue / Cell type specific
  - Laser Capture Microdissection
  - single cell transcriptomics [not discussed]
- Time points
  - Development
  - Response to treatment (before, during, after)

# Experimental design



# How many reads are enough?

## Genohub

	Small (bacteria)	Intermediate (fruit fly, worm)	Large (mouse, human)
No. of reads for DGE ( $\times 10^6$ )	5 SR	10 SR	20–50 SR
No. of reads for <i>de novo</i> transcriptome assembly ( $\times 10^6$ )	30–65 PE	70–130 PE	100–200 PE
Read length (bp)	50	50–100	>100

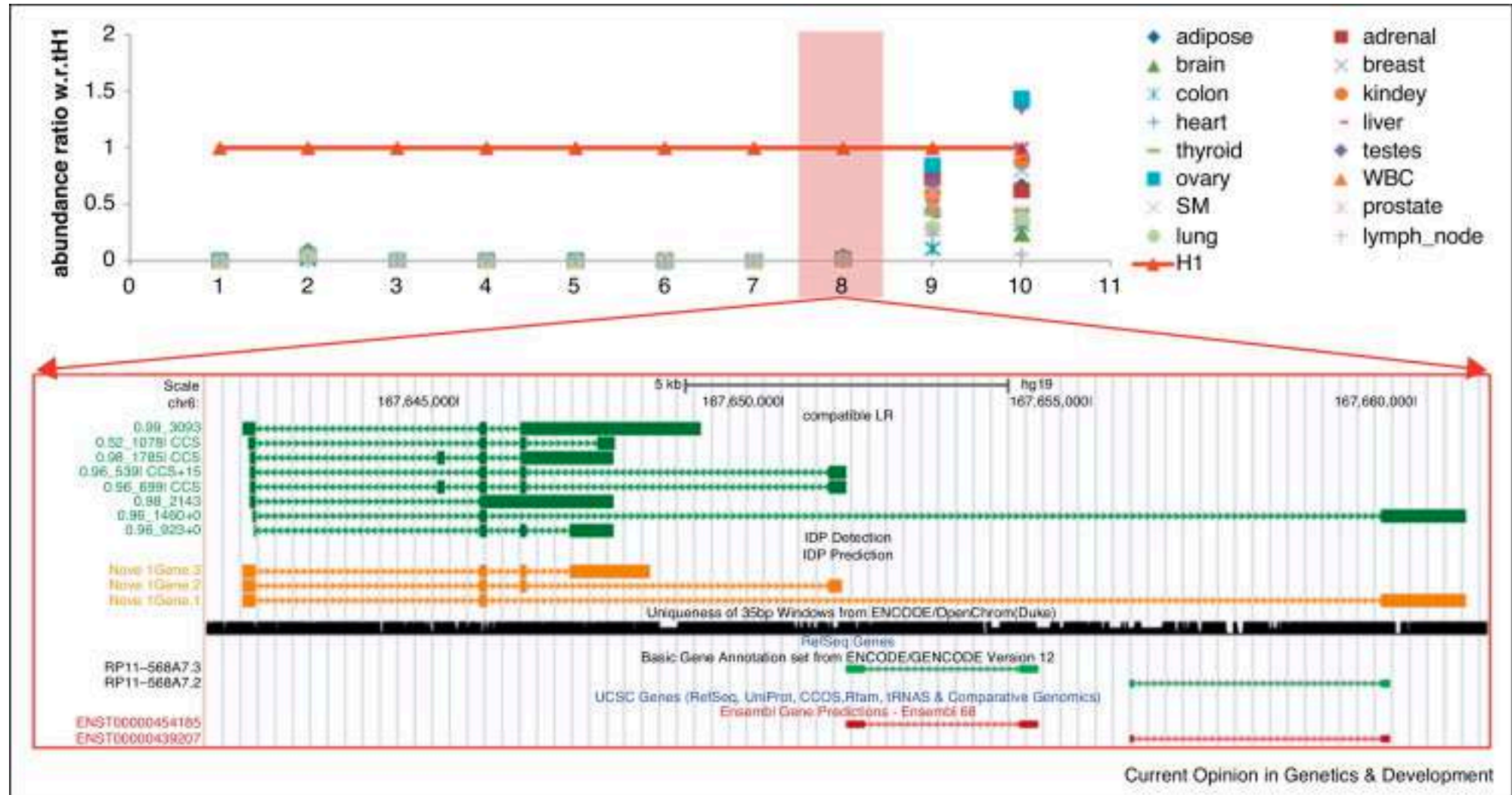
## Illumina

- Gene expression profiling experiments that are looking for a quick snapshot of highly expressed genes may only need 5–25 million reads per sample. In these cases, consider pooling multiple RNA-Seq samples into one lane of a sequencing run. This allows for high multiplexing of samples.
- Experiments looking for a more global view of gene expression, and some information on alternative splicing, typically require 30–60 million reads per sample. This range encompasses most published RNA-Seq experiments for mRNA/whole transcriptome sequencing.

[1] <https://genohub.com/next-generation-sequencing-guide/#depth2>

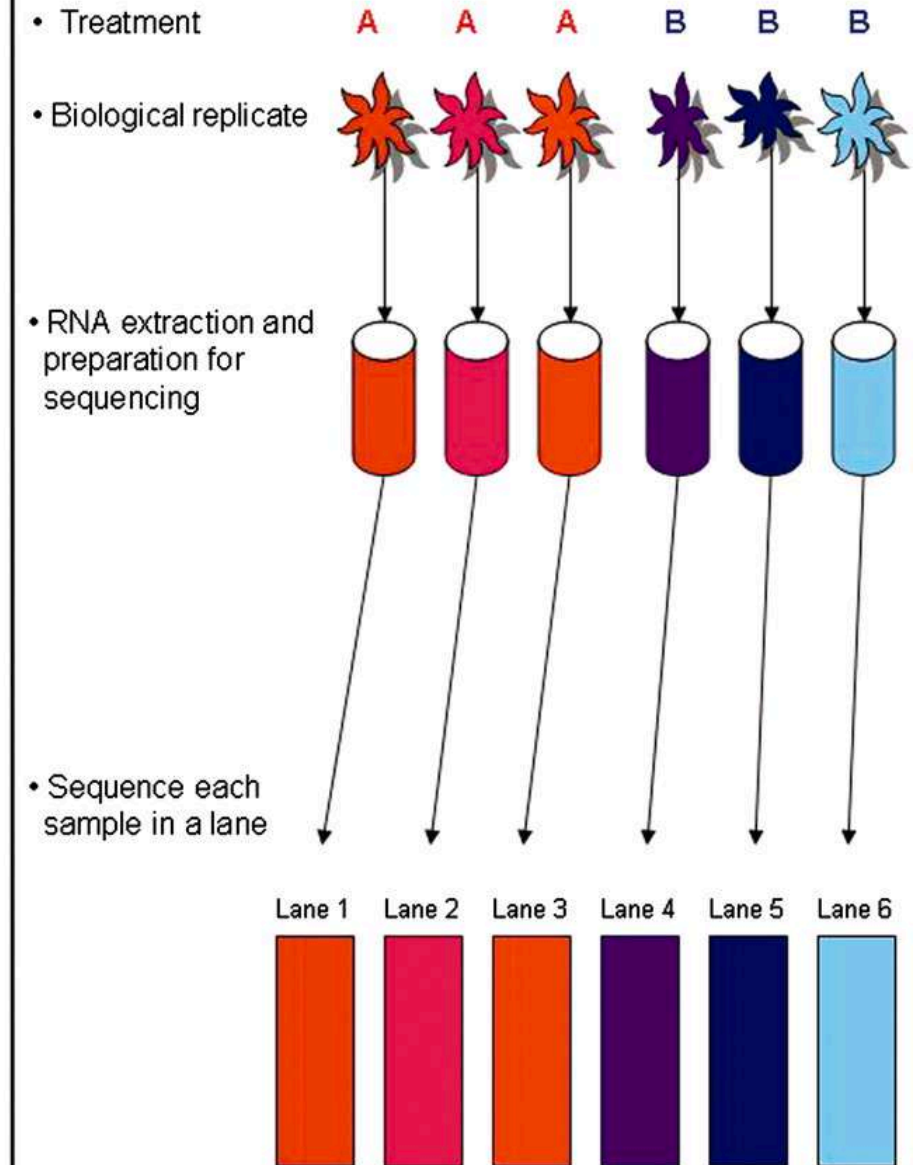
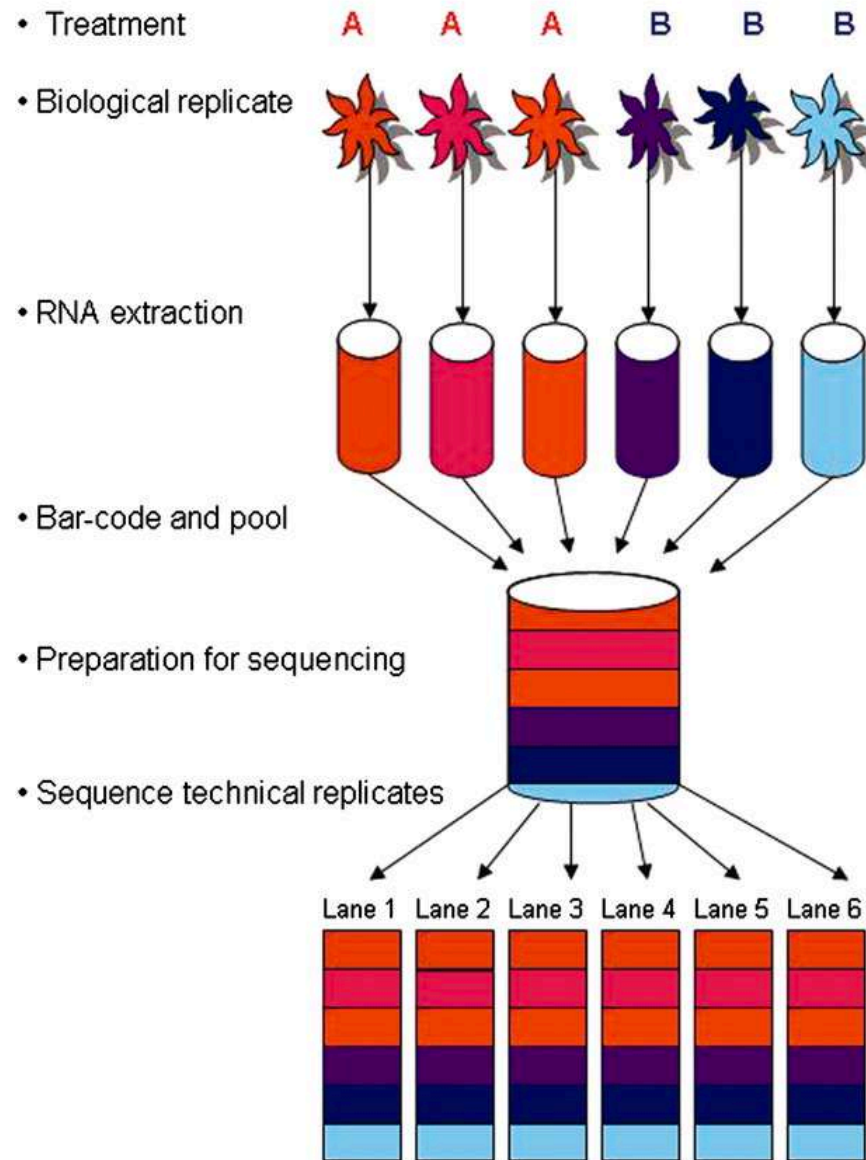
[2] <https://support.illumina.com/bulletins/2017/04/considerations-for-rna-seq-read-length-and-coverage-.html>

For isoform discovery, longer sequences are better



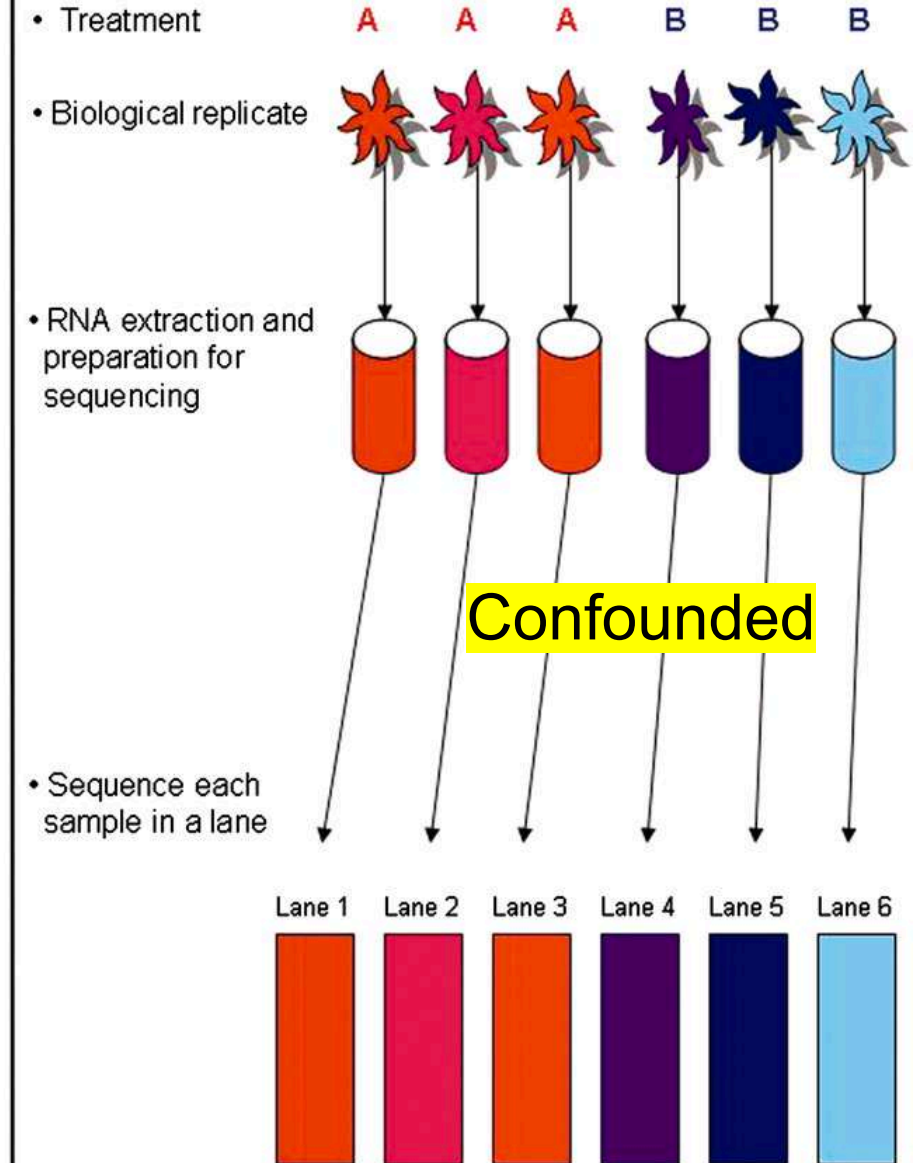
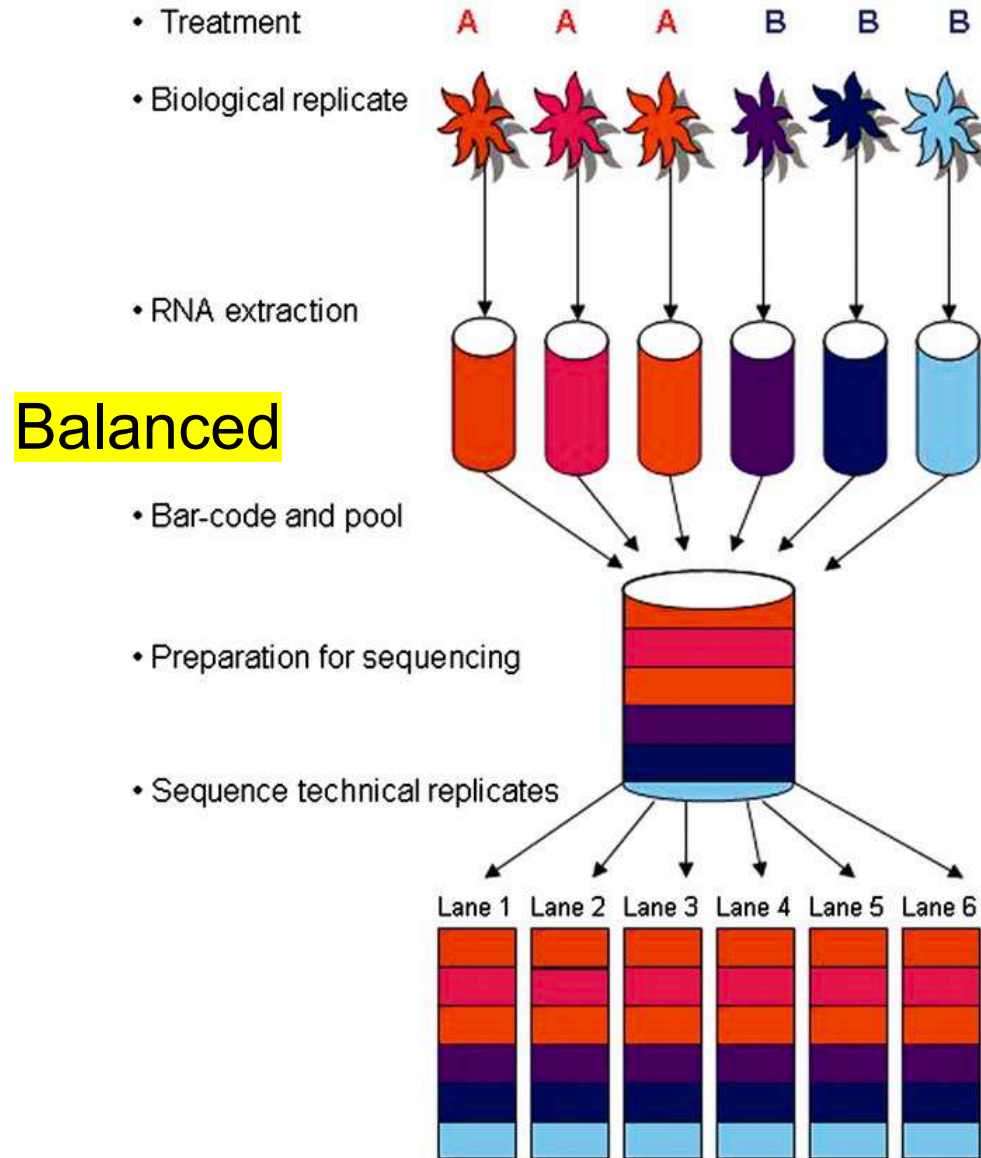


# Which of the following designs is correct?

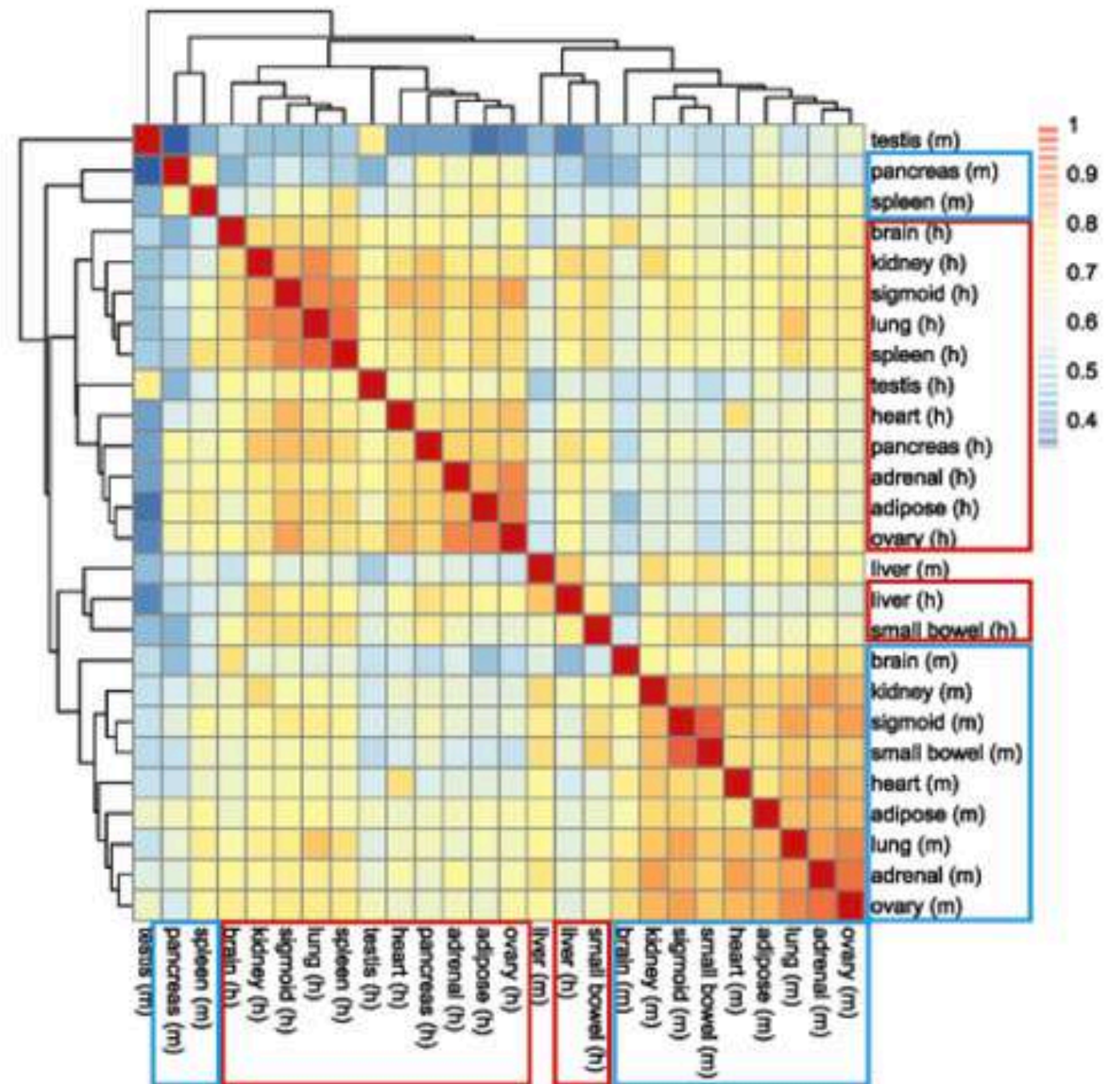




# Which of the following designs is correct?



# Example of batch effect:

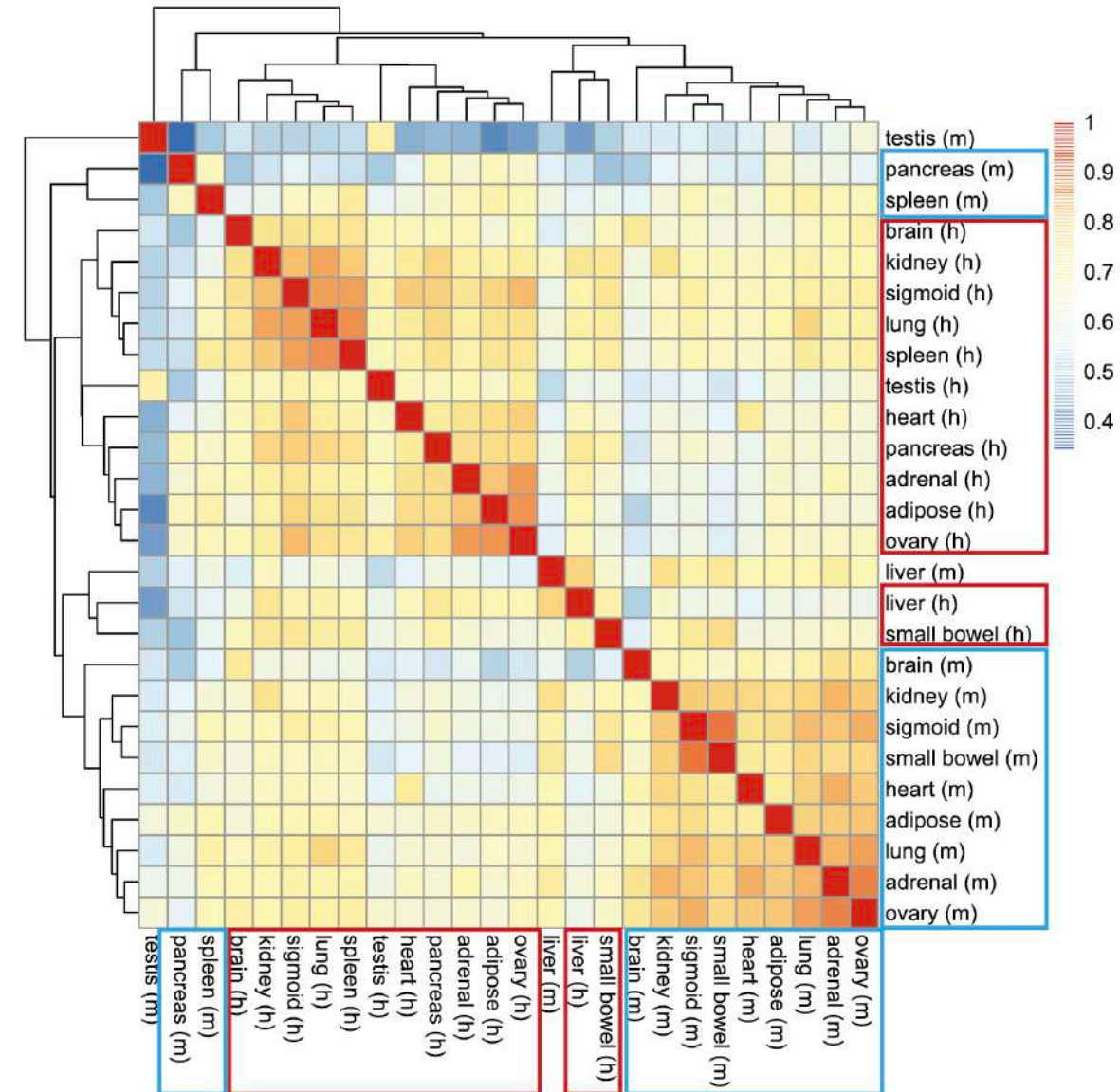
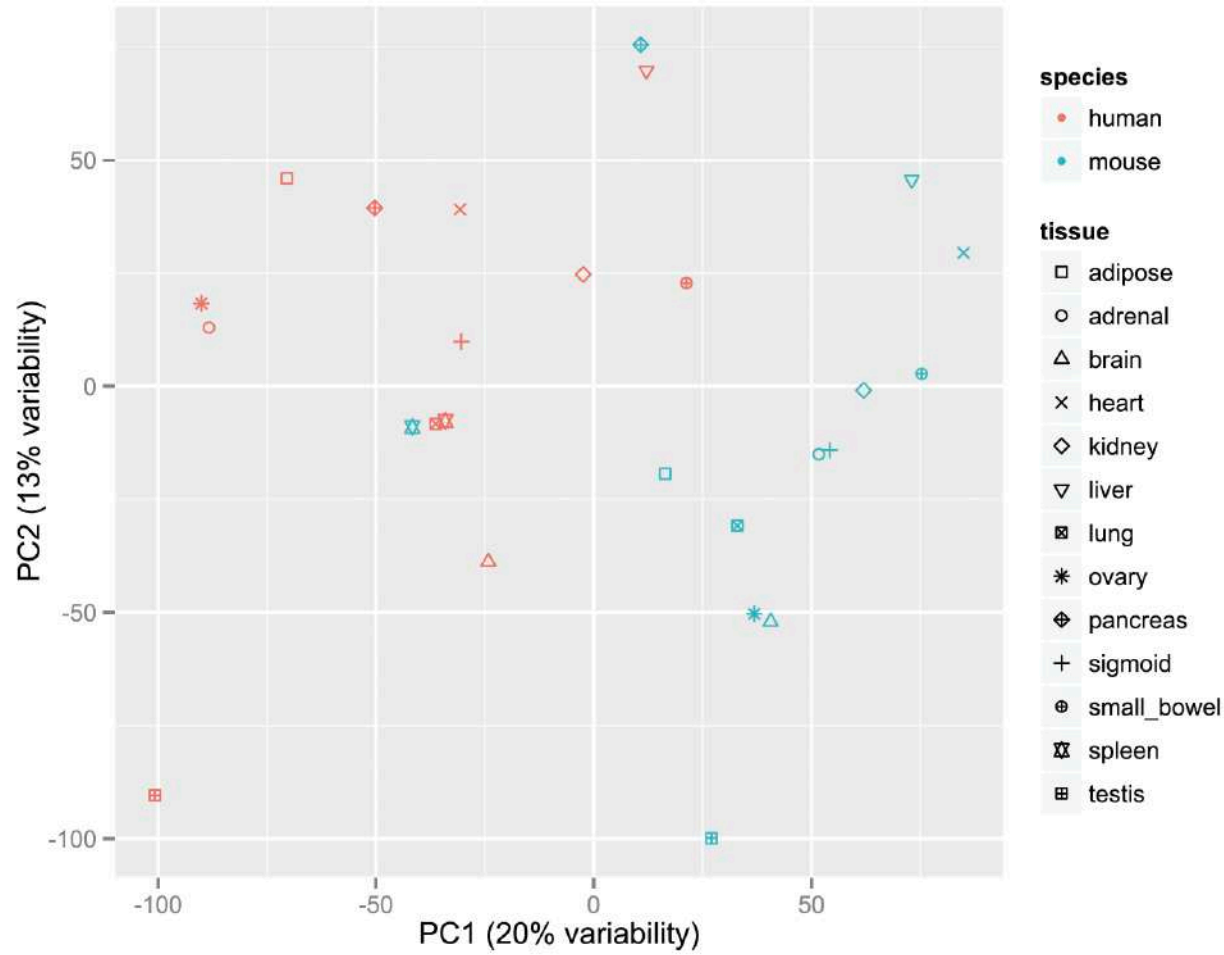


# Example of batch effect:

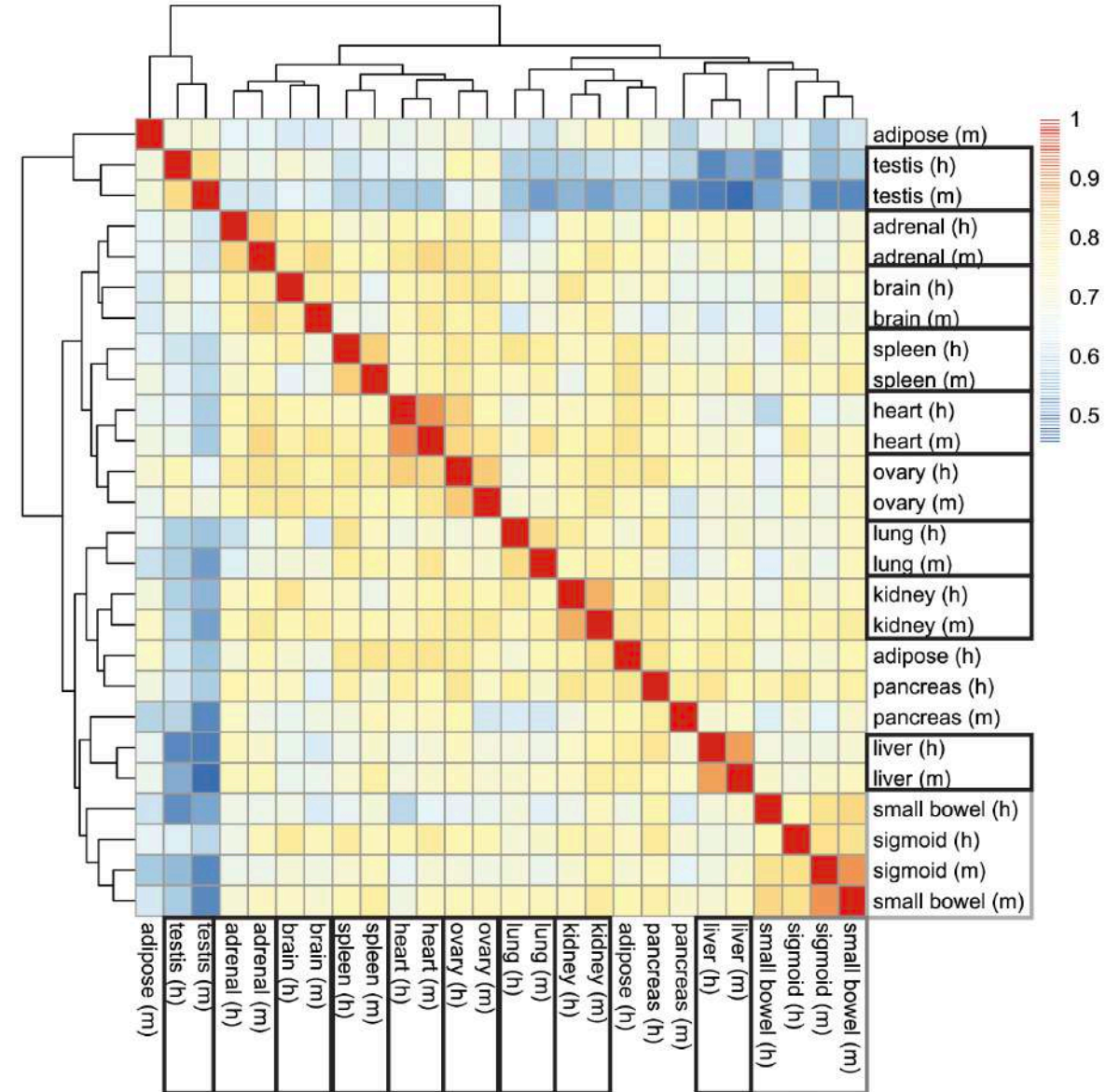
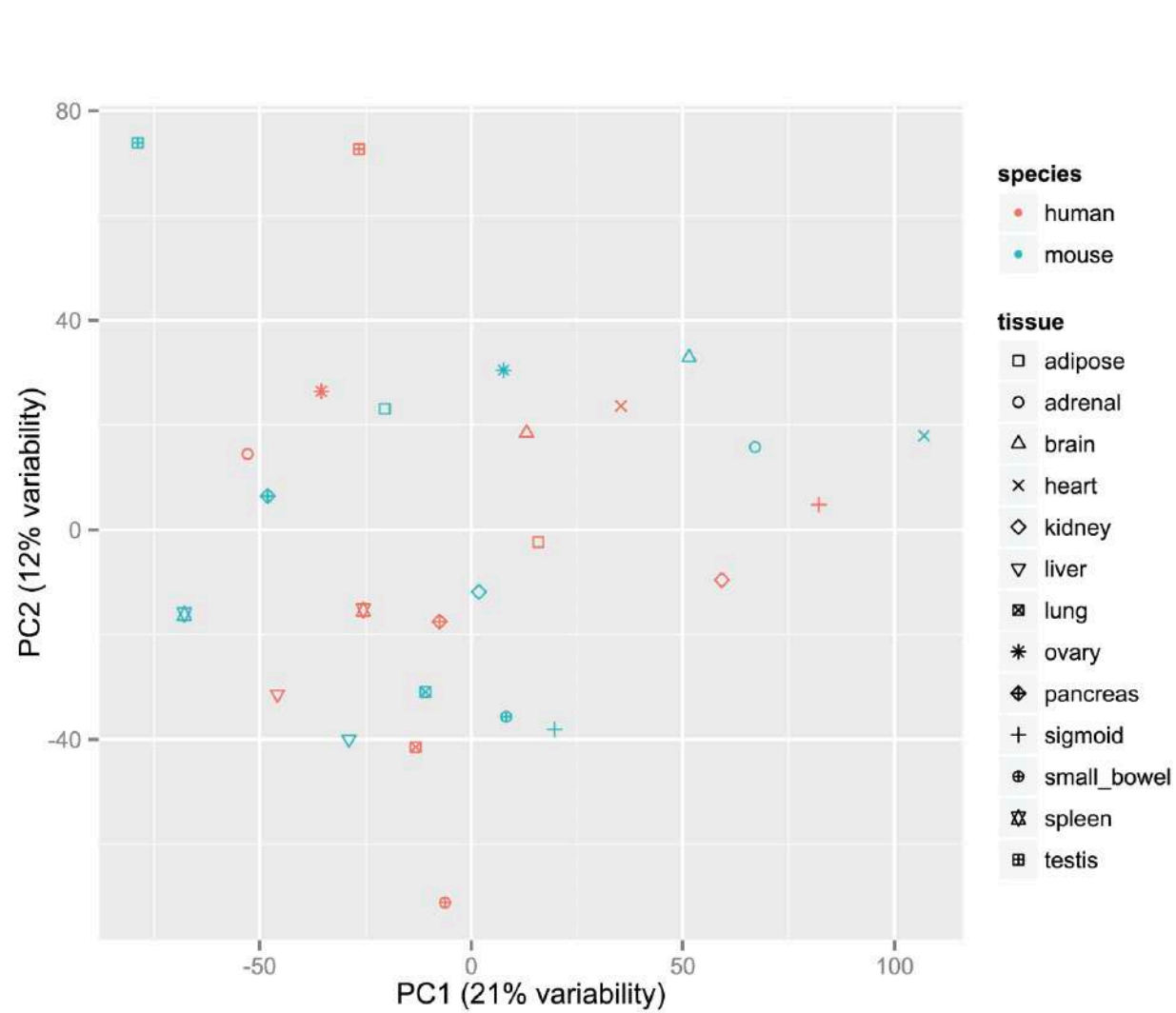
D87PMJN1 (run 253, flow cell D2GUAACXX, lane 7)	D87PMJN1 (run 253, flow cell D2GUAACXX , lane 8)	D4LHBFN1 (run 276, flow cell C2HKJACXX , lane 4)	MONK (run 312, flow cell C2GR3ACXX , lane 6)	HWI-ST373 (run 375, flow cell C3172ACXX , lane 7)
heart	adipose	adipose	heart	brain
kidney	adrenal	adrenal	kidney	pancreas
liver	sigmoid colon	sigmoid colon	liver	brain
small bowel	lung	lung	small bowel	spleen
spleen	ovary	ovary	testis	● Human
testis		pancreas		● Mouse



# Recapitulating the patterns reported by the mouse ENCODE papers



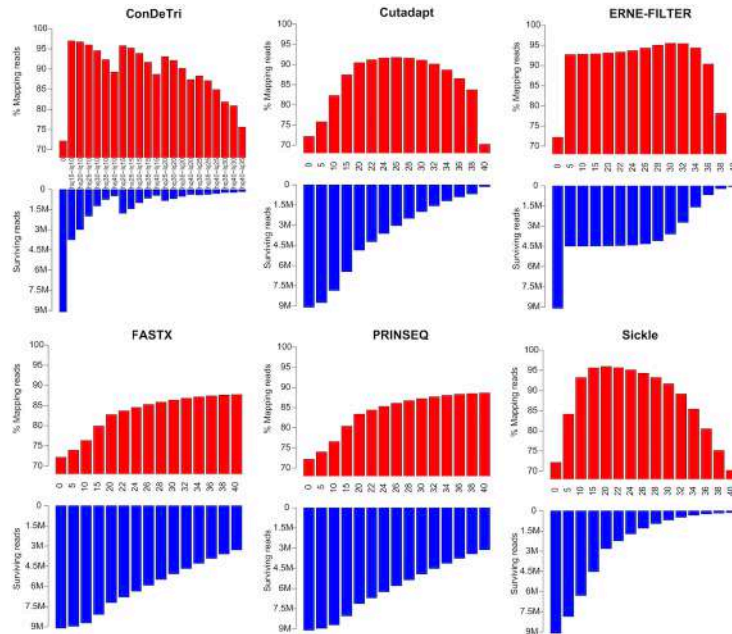
# Clustering of data once batch effects are accounted for





# Is trimming beneficial?

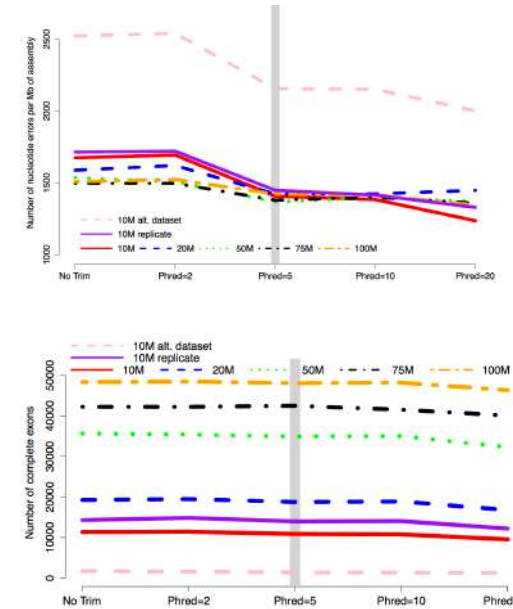
## Software comparison, RNA/DNA-Seq



*“trimming is beneficial in RNA-Seq, SNP identification and genome assembly procedures, with the best effects evident for intermediate quality thresholds (Q between 20 and 30)”*

Del Fabbro C et al (2013) **An Extensive Evaluation of Read Trimming Effects on Illumina NGS Data Analysis**. PLoS ONE 8(12): e85024. doi:10.1371/journal.pone.0085024

## Assembly-oriented, RNA-seq only



Erroneous bases in assembly

# complete exons

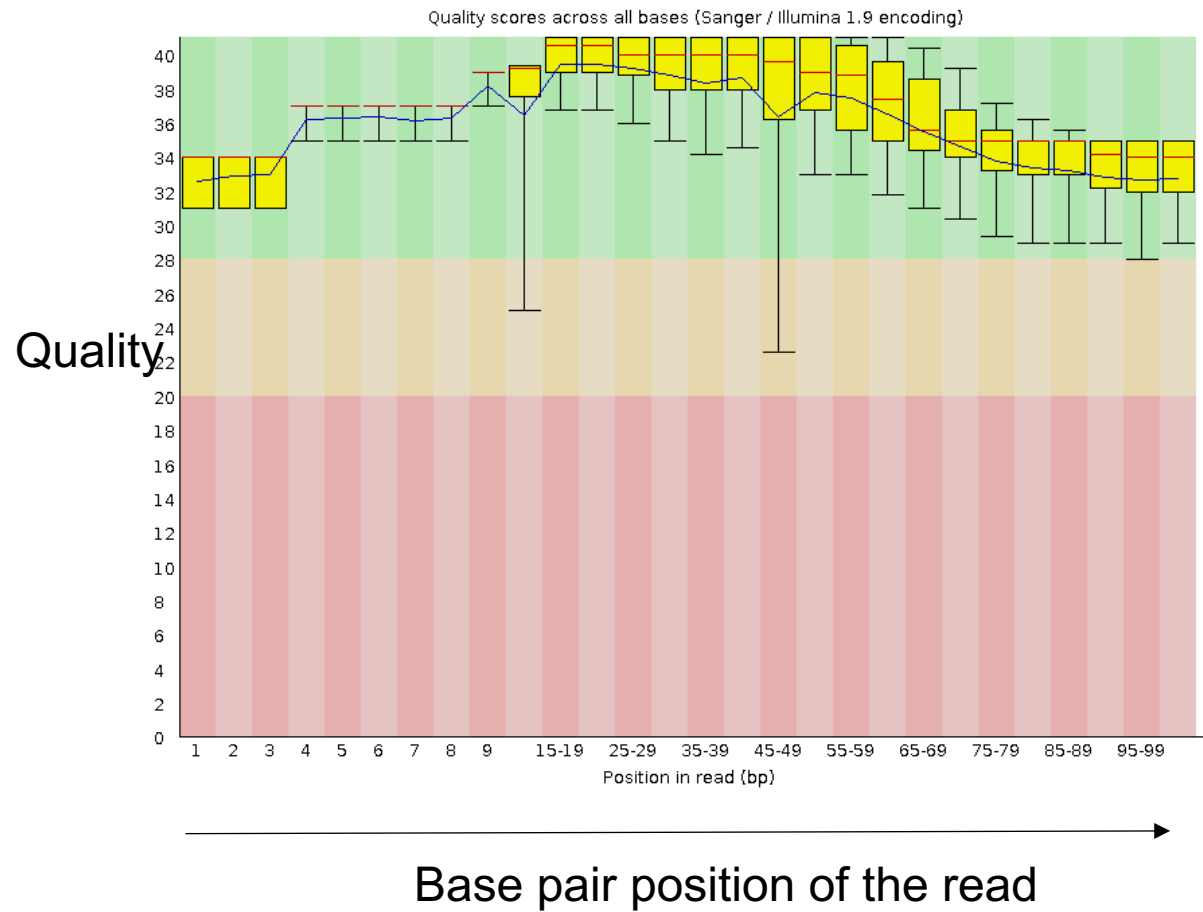
*“Although very aggressive quality trimming is common, this study suggests that a more gentle trimming, specifically of those nucleotides whose Phred score < 2 or < 5, is optimal for most studies across a wide variety of metrics.”*

MacManes MD (2013)

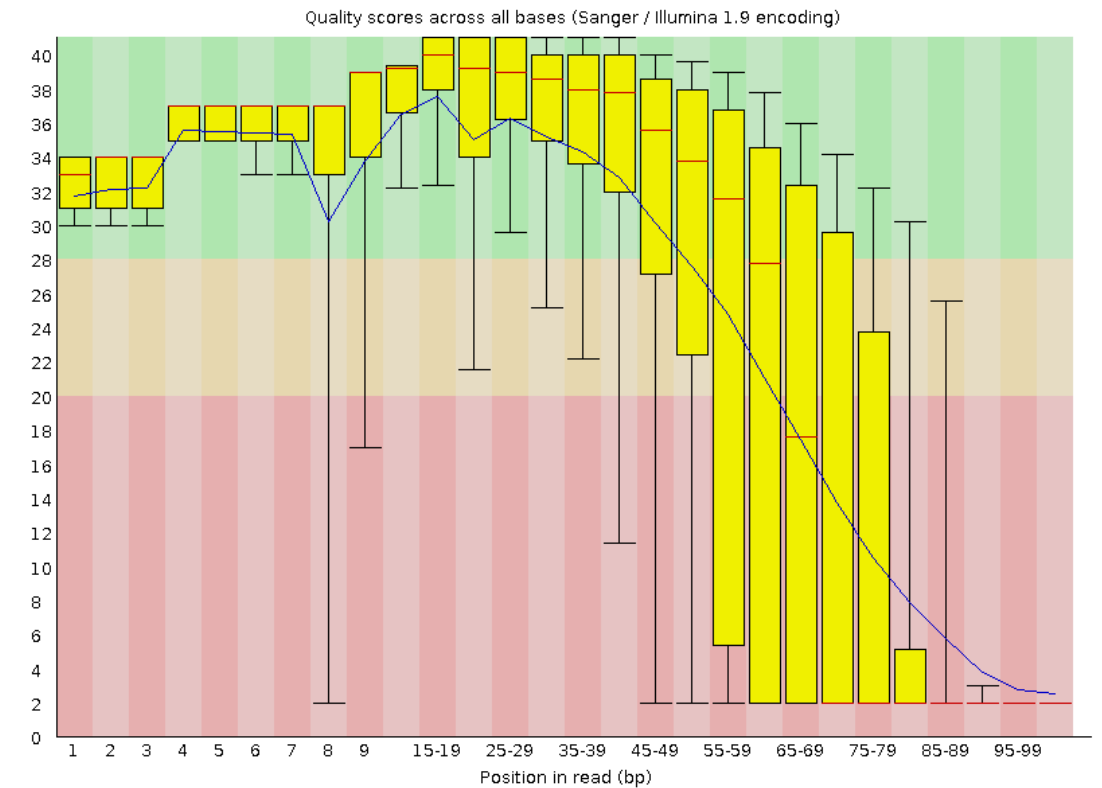
**On the optimal trimming of high-throughput mRNAseq data** doi: 10.1101/000422

# My take: only trim data when you have to

## Good/Trimmed data

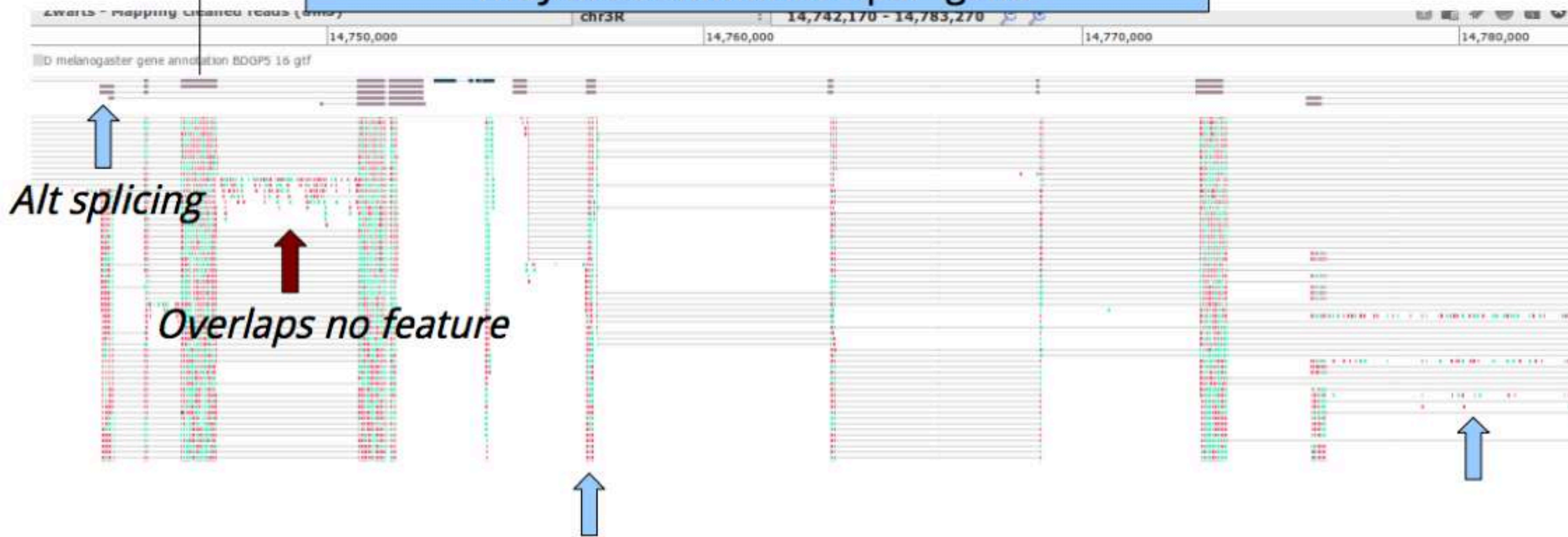


## Poor/Raw data



3.1 Once you have mappings, you can start counting

'Exons' are the type of *features* used here.  
They are summarized per 'gene'

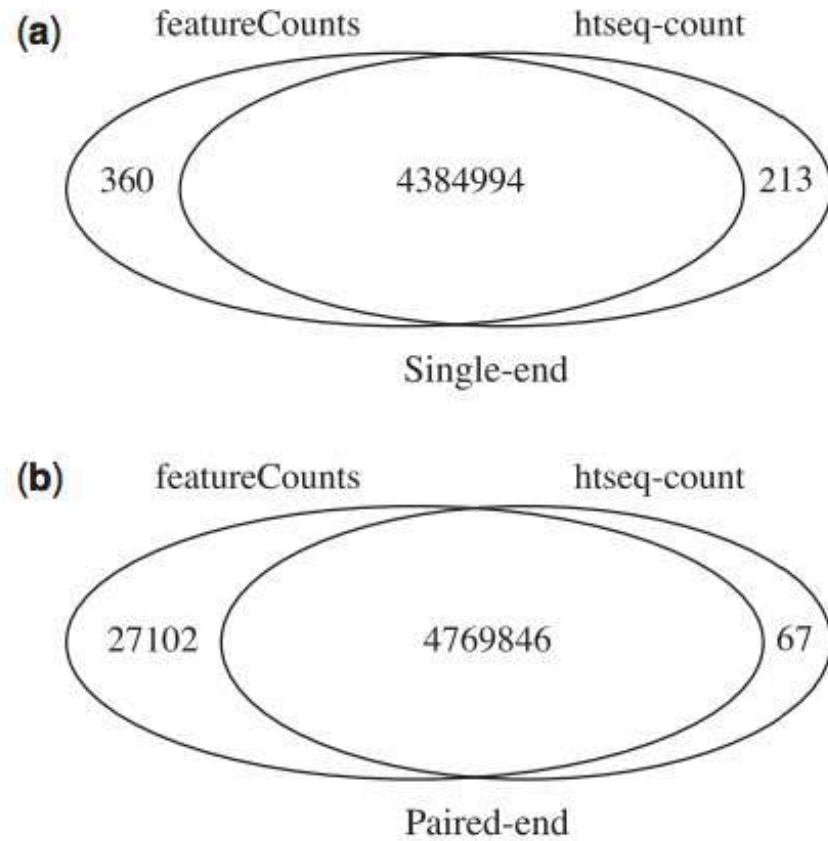


### Concept:

GeneA = exon 1 + exon 2 + exon 3 + exon 4 = 215 reads

GeneB = exon 1 + exon 2 + exon 3 = 180 reads

# Featurecount (much faster!)



**Table 3.** Performance with RNA-seq reads simulated from an annotated assembly of the Budgerigar genome

Methods	Number of reads	Time (mins)	Memory (MB)
<i>featureCounts</i>	7 924 065	0.6	15
<i>summarizeOverlaps</i> (whole genome at once)	7 924 065	12.6	2400
<i>summarizeOverlaps</i> (by scaffold)	7 924 065	53.3	262
<i>htseq-count</i>	7 912 439	12.1	78

*Note:* The annotation includes 16 204 genes located on 2850 scaffolds. *featureCounts* is fastest and uses least memory. Table gives the total number of reads counted, time taken and peak memory used. *htseq-count* was run in 'union' mode.



## Some QC is needed

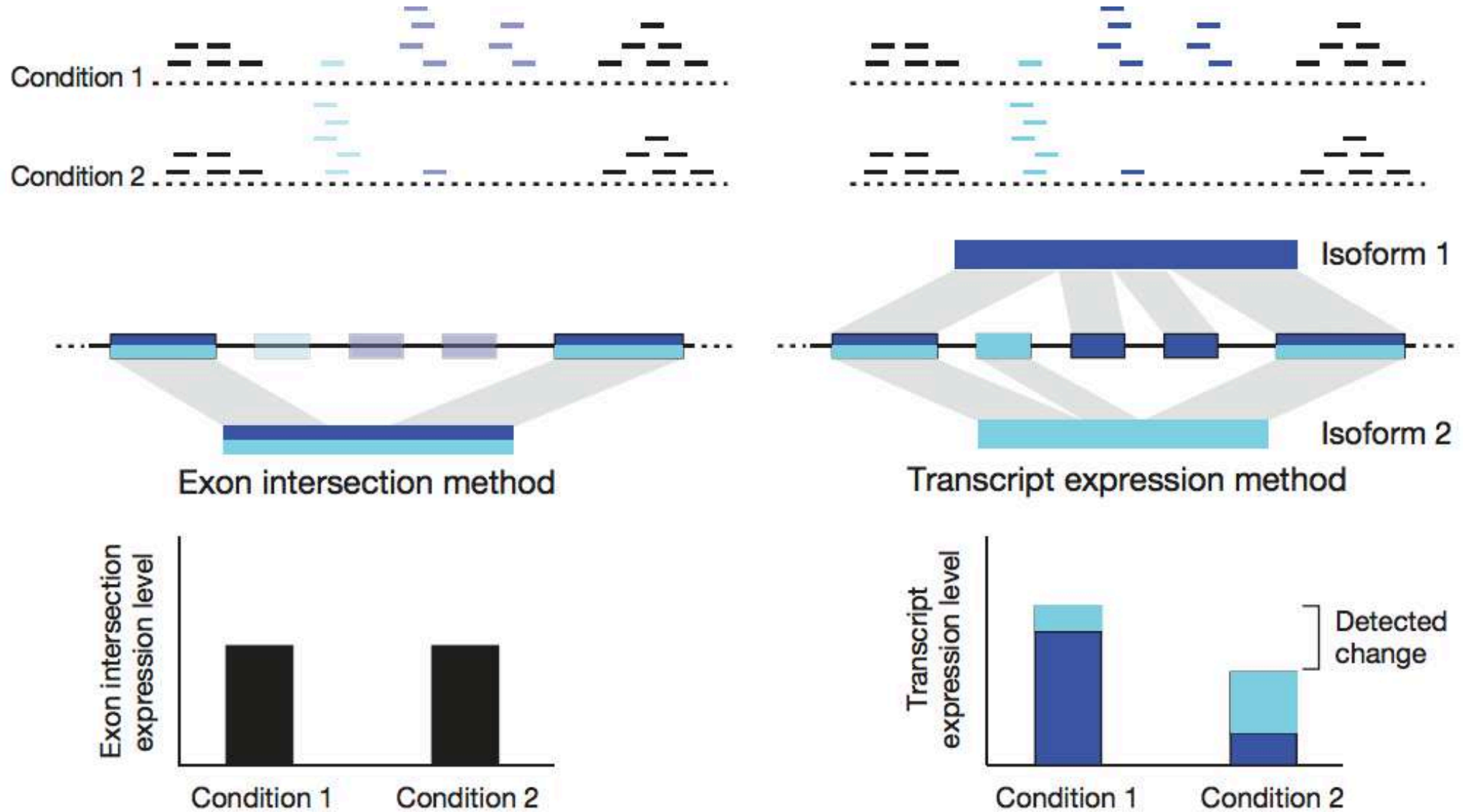
- Which genes are highly counted?
- Any samples with a lot of missing/no count genes?
- Is there any biases on sequencing?
- Anything that may affect sample counts (like batch effect?)

# Ambiguity in counting

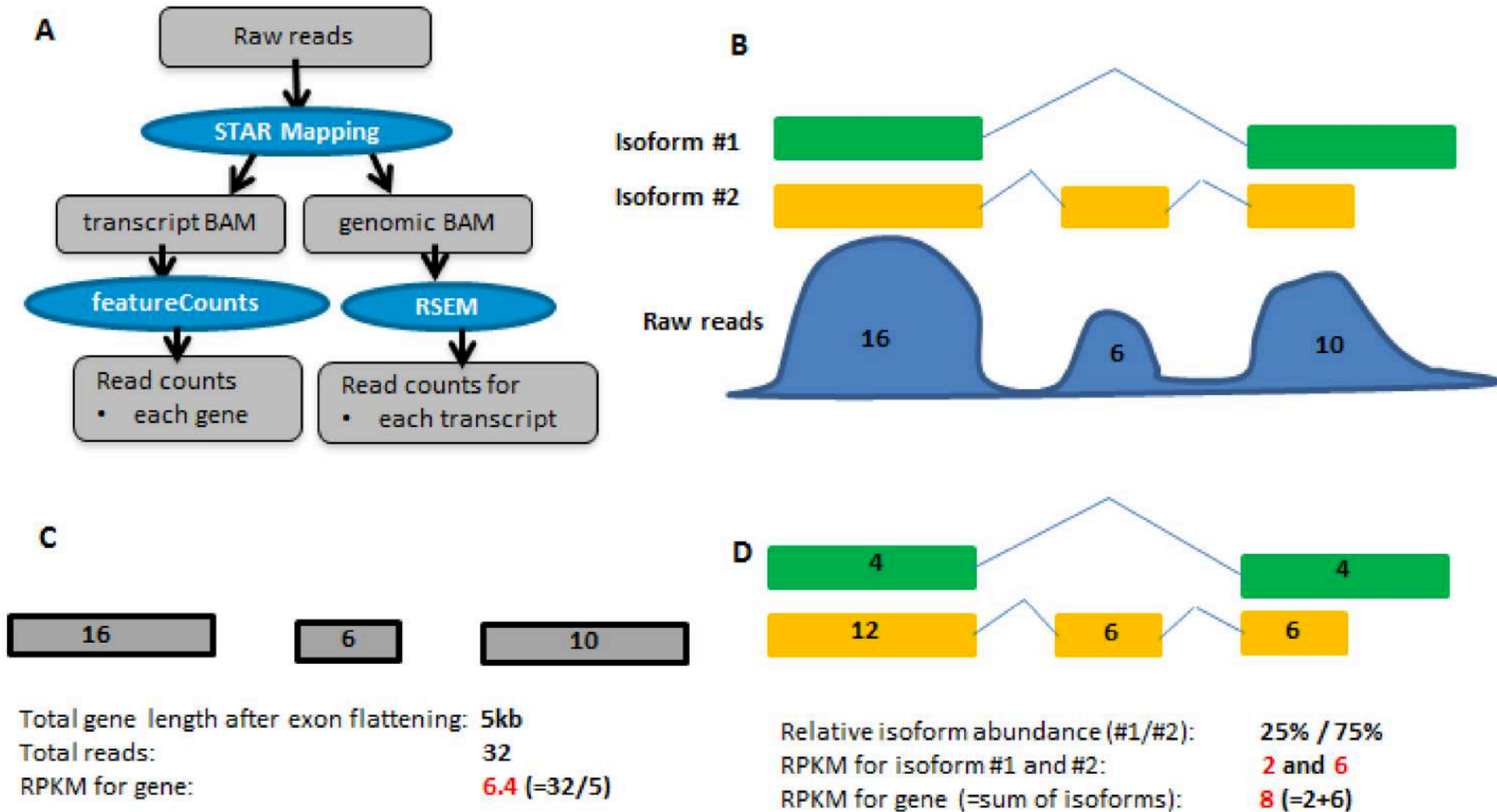
We focus on the **gene level**: merge all counts over different **isoforms** into one, taking into account:

- Reads that do **not overlap** a feature, but appear in introns. Take into account?
- Reads that align to **more than one feature** (exon or transcript). Transcripts can be overlapping - perhaps on different strands. (PE, and strandedness can resolve this partially).
- Reads that **partially** overlap a feature, not following known annotations.

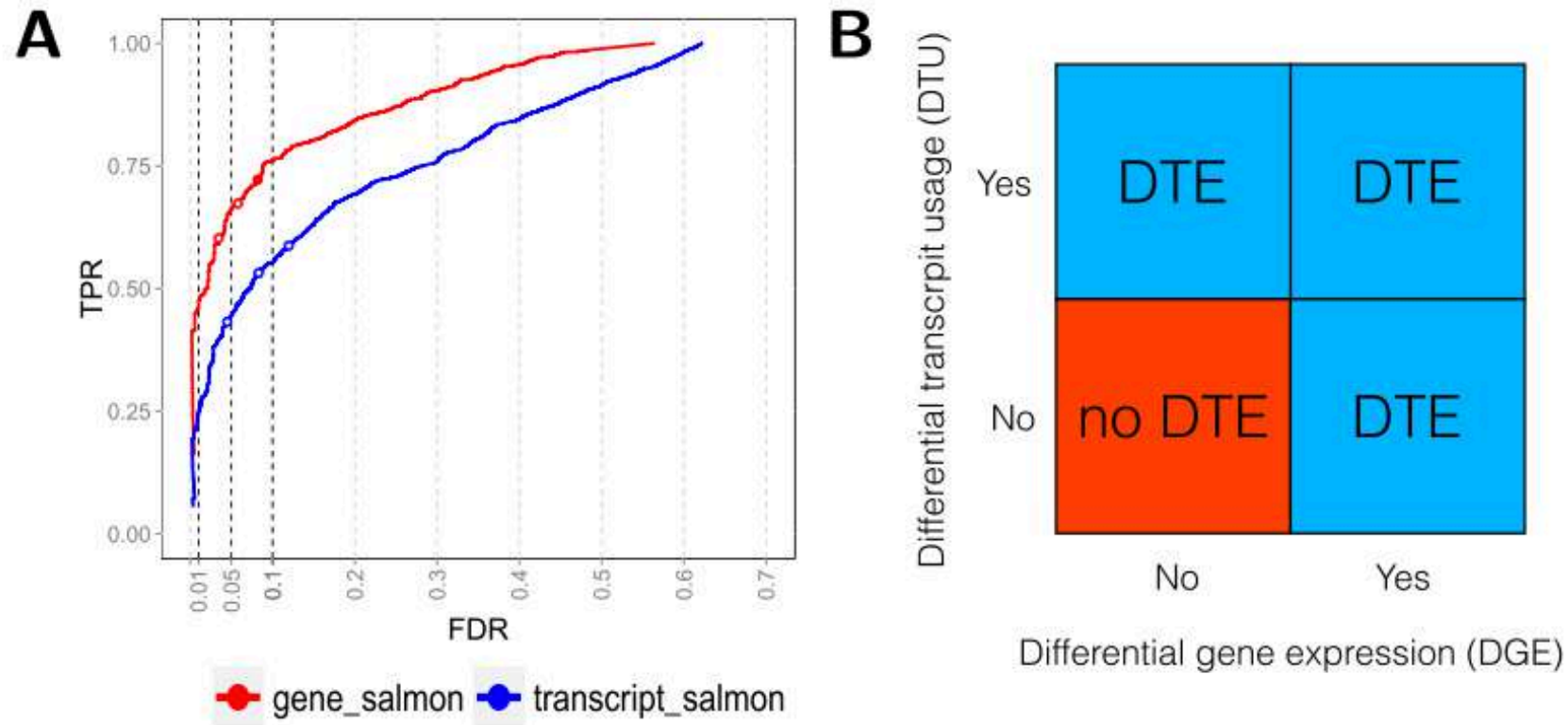
# Transcript counting could be more robust in detecting changes



# Outstanding problems in counting with exon merging model



But Differential transcript expression can lead to inflated false positive rate (and more difficult to interpret biologically)



**Figure 2 (sim2).** **A:** DTE detection performance on transcript- and gene-level, using *edgeR* applied to transcript-level estimated counts from *Salmon*. The statistical analysis was performed on transcript level and aggregated for each gene using the *perGeneQValue* function from the *DEXSeq* R package; aggregated results show higher detection power. The curves trace out the observed FDR and TPR for each significance cutoff value. The three circles mark the performance at adjusted p-value cutoffs of 0.01, 0.05 and 0.1. **B:** Schematic illustration of different ways in which differential transcript expression (DTE) can arise, in terms of absence or presence of differential gene expression (DGE) and differential transcript usage (DTU).

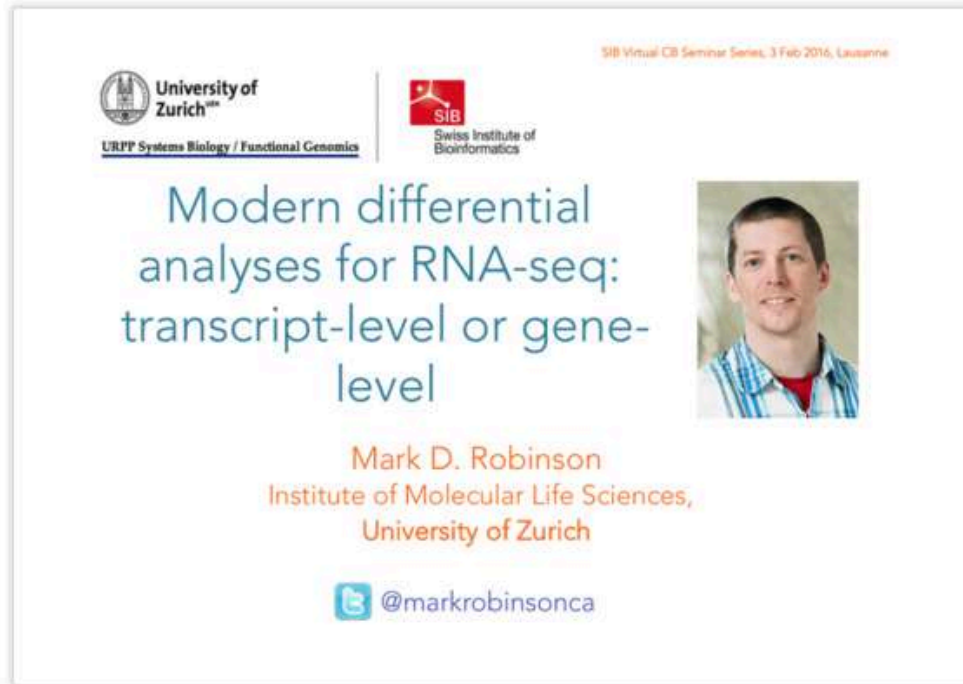


# So use isoform or not?

## Modern RNA-seq differential expression analyses: transcript-level or gene-level

Posted by: RNA-Seq Blog in Presentations February 11, 2016 1,733 Views

Modern RNA-seq differential expression analyses: transcript-level or gene-level



The slide features logos for the University of Zurich (URPP Systems Biology / Functional Genomics) and the Swiss Institute of Bioinformatics (SIB). It includes the text 'SIB Virtual CB Seminar Series, 3 Feb 2016, Lausanne'. The main title is 'Modern differential analyses for RNA-seq: transcript-level or gene-level'. A portrait of Mark D. Robinson is shown. His affiliation is 'Mark D. Robinson, Institute of Molecular Life Sciences, University of Zurich'. A Twitter handle '@markrobinsonca' is also present.

“There is no crisis; the impact of union vs. transcript counting in many datasets is rather small”

“Unless the need dictates, answer the easier questions”

<http://www.rna-seqblog.com/modern-rna-seq-differential-expression-analyses-transcript-level-or-gene-level/>



figshare



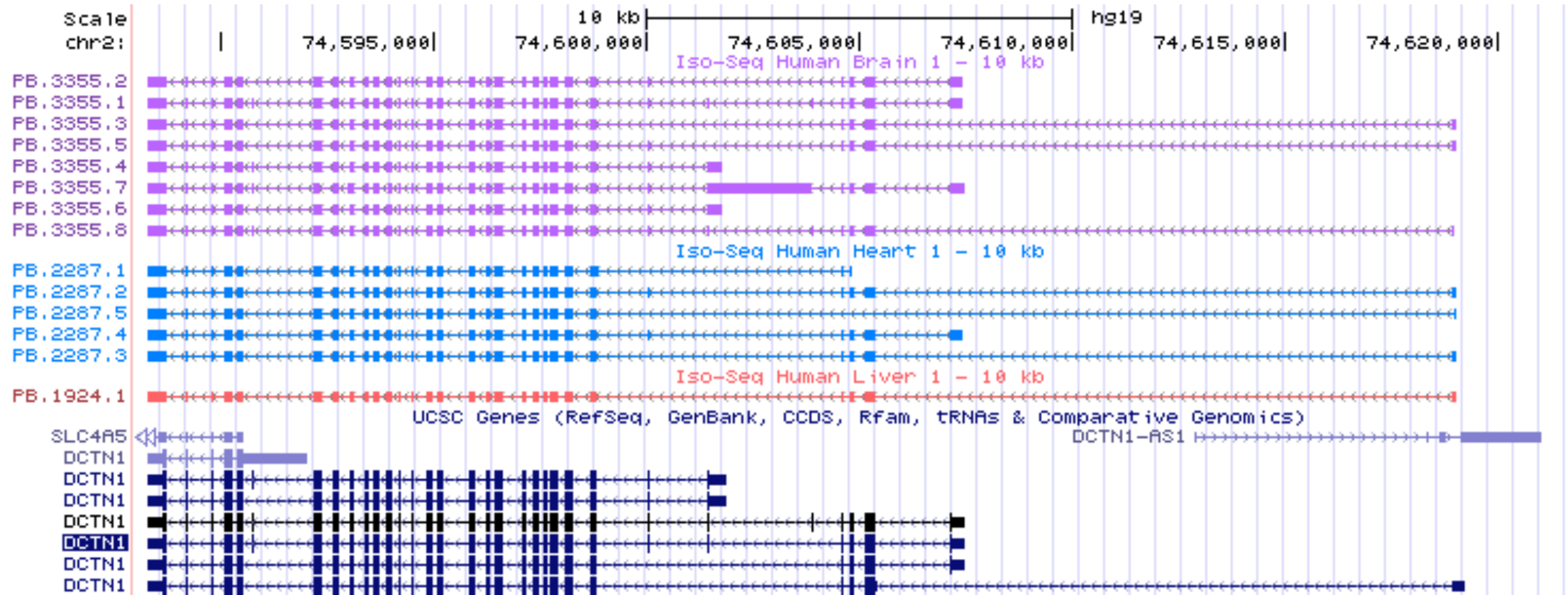
Share



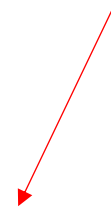
Download (6.53 MB)

# We may end up counting full-length transcripts anyway

## Pacbio IsoSeq

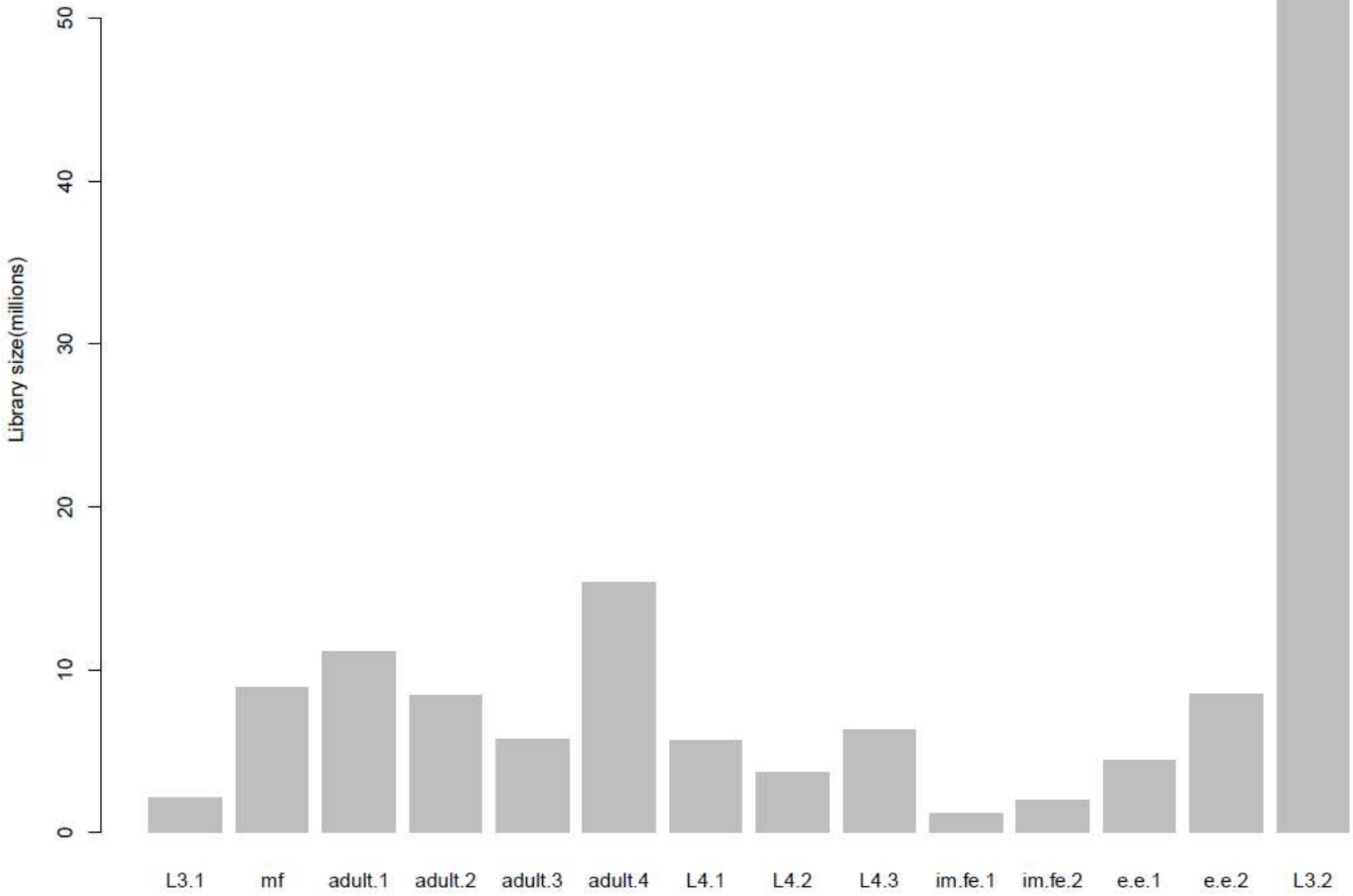


This is the bit we care about!

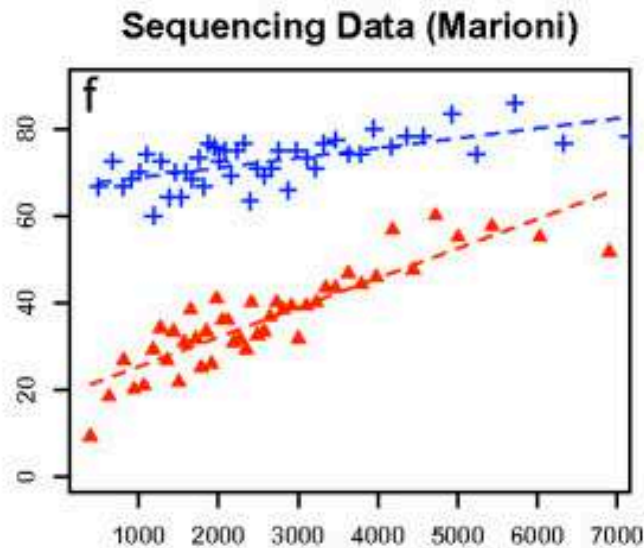
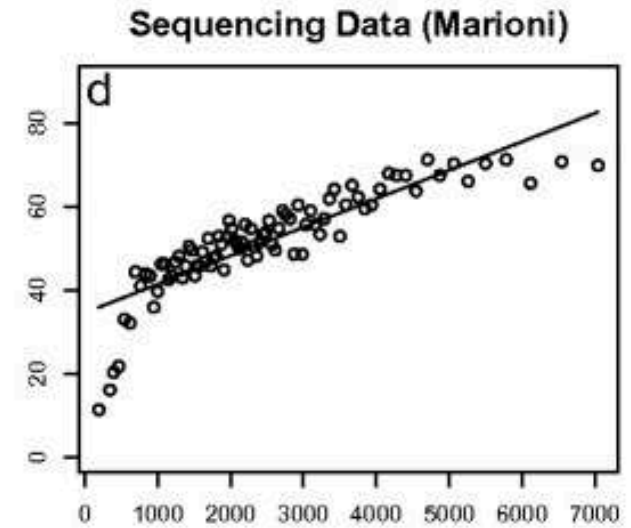
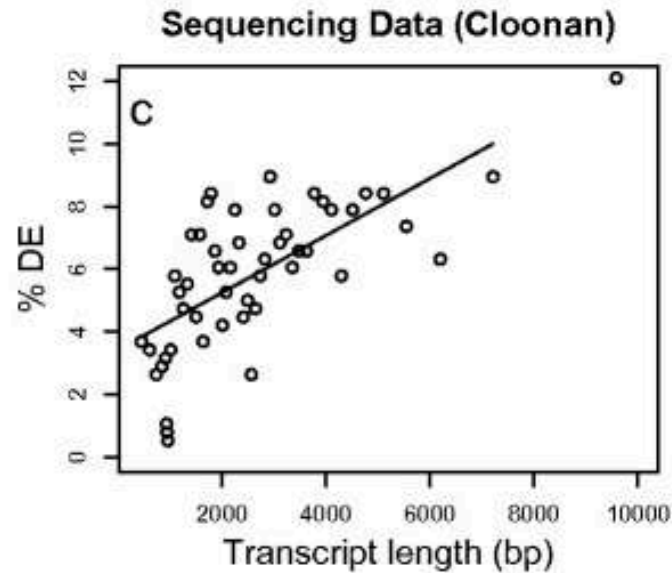
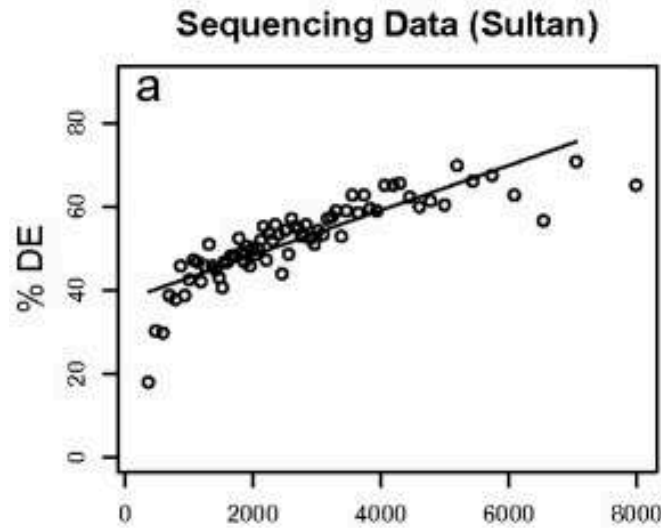


Counts of the gene depends on **expression** ,transcript length  
,sequencing depth and simply chance

# Higher sequencing depth equals more counts



# Counts are proportional to the transcript length x mRNA expression level



33% of highest expressed genes  
33% of lowest expressed genes



# Normalization: different goals

- **Counts per million (CPM)**
- **R/FPKM:** (Mortazavi et al. 2008)
  - **Correct for:** differences in sequencing depth and transcript length
  - **Aiming to:** compare a gene across samples and diff genes within sample
- **TMM:** (Robinson and Oshlack 2010)
  - **Correct for:** differences in transcript pool composition; extreme outliers
  - **Aiming to:** provide better across-sample comparability
- **TPM:** (Li et al 2010, **Wagner** et al 2012)
  - **Correct for:** transcript length distribution in RNA pool
  - **Aiming to:** provide better across-sample comparability
- **Limma voom (logCPM):** (Lawet al 2013)
  - **Aiming to:** stabilize variance; remove dependence of variance on the mean

<https://haroldpimentel.wordpress.com/2014/05/08/what-the-fpkm-a-review-rna-seq-expression-units/>

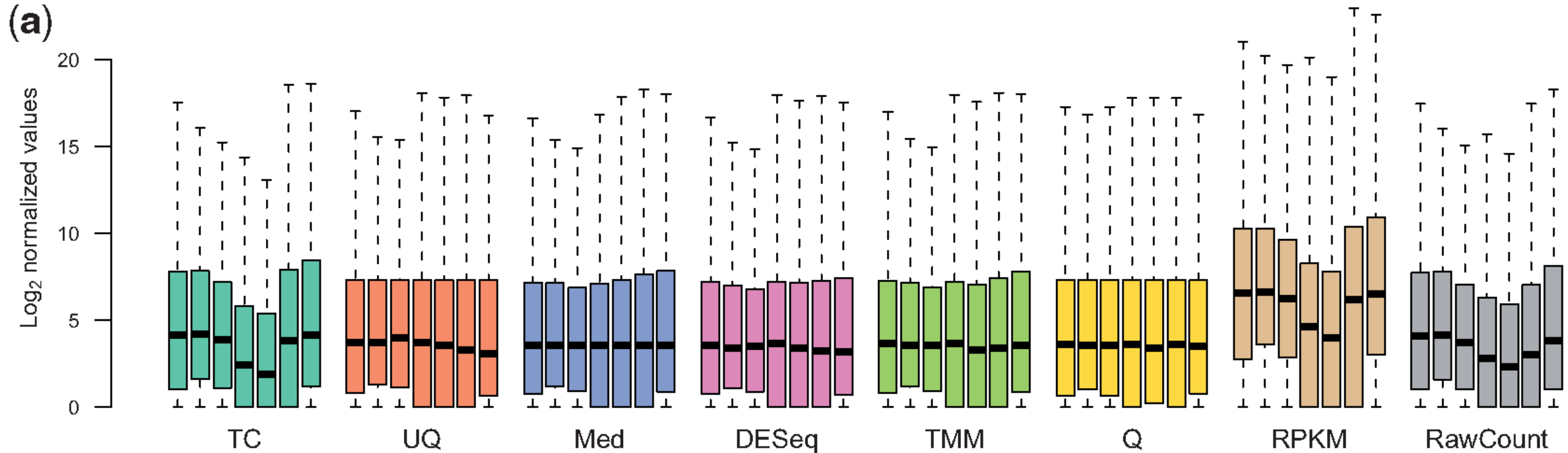
## Optimal Scaling of Digital Transcriptomes

Gustavo Glusman , Juan Caballero, Max Robinson, Burak Kutlu, Leroy Hood

Published: Nov 06, 2013 • DOI: 10.1371/journal.pone.0077885

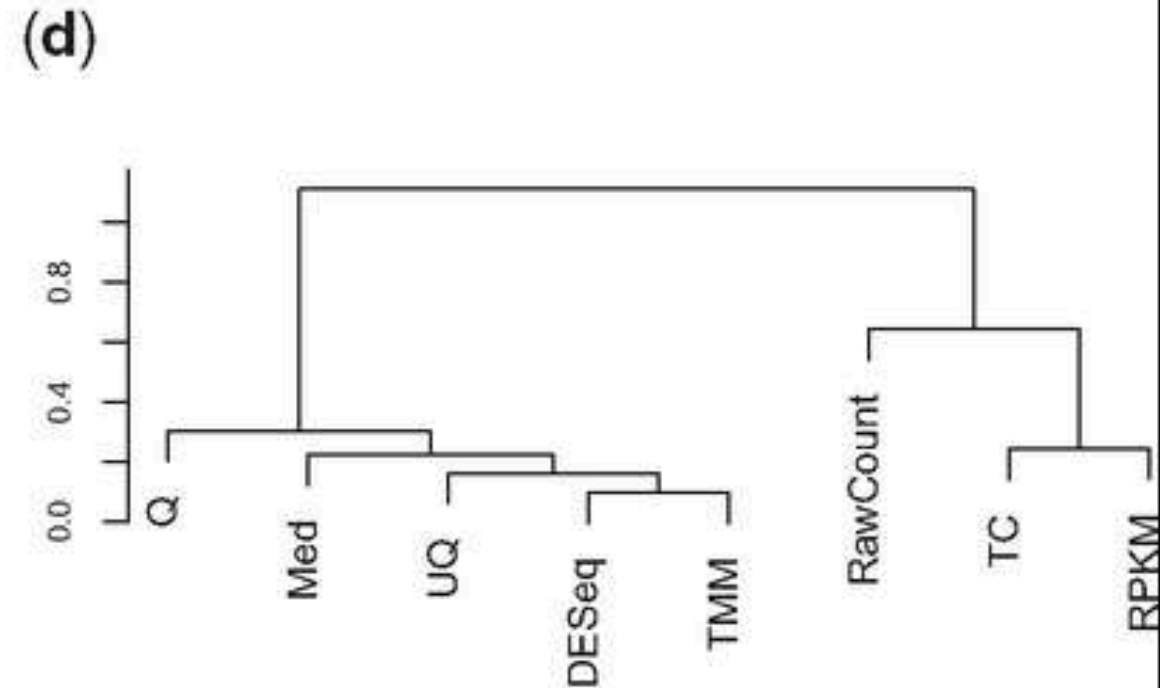
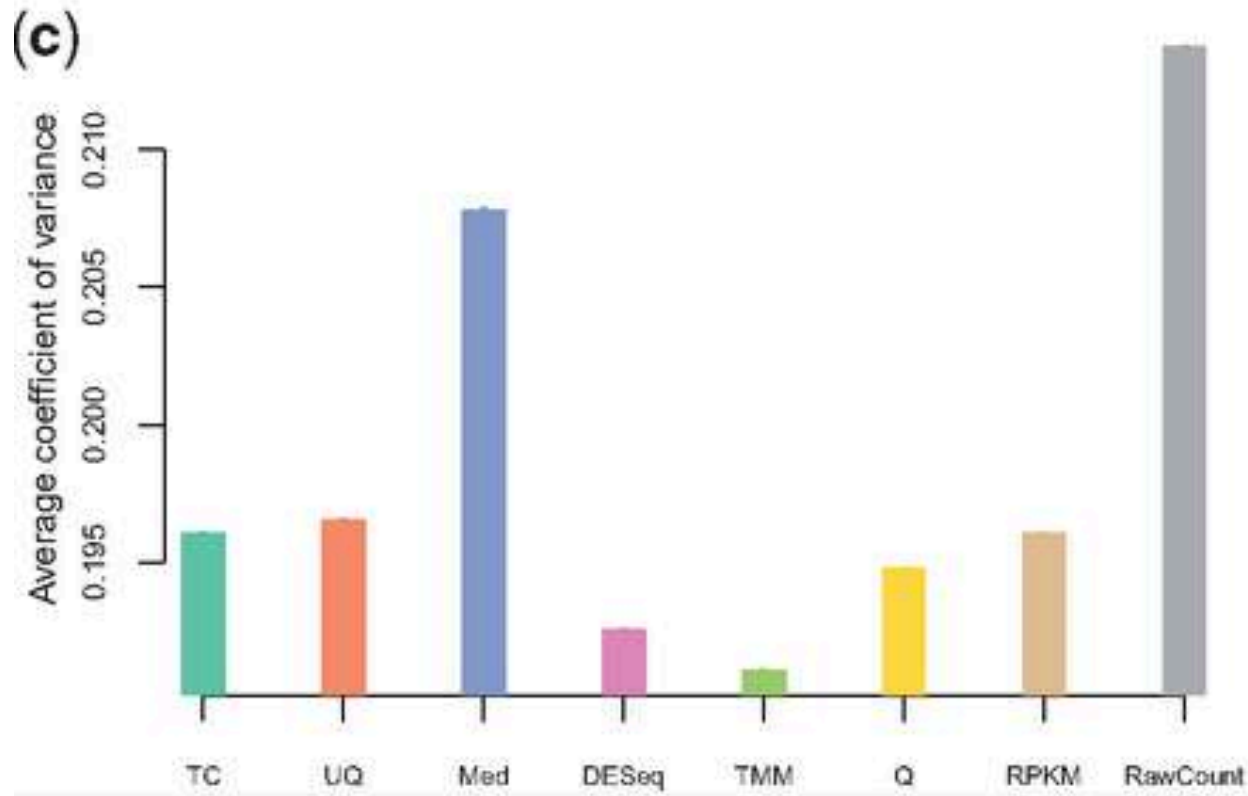
# RPKM shouldn't be used for between sample comparisons

Boxplots of  $\log_2(\text{counts} + 1)$  for **seven** replicates in the *M. musculus* data, by normalization method.



# RPKM shouldn't be used for between sample comparisons

**C)** Analysis of housekeeping genes for the *H. sapiens* data. **(D)** Consensus dendrogram of differential analysis results



# Summary of comparison results for the seven normalization methods under consideration

Method	Distribution	Intra-Variance	Housekeeping	Clustering	False-positive rate
TC	–	+	+	–	–
UQ	++	++	+	++	–
Med	++	++	–	++	–
<b>DESeq</b>	++	++	++	++	++
<b>TMM</b>	++	++	++	++	++
Q	++	–	+	++	–
RPKM	–	+	+	–	–

A '–' indicates that the method provided unsatisfactory results for the given criterion, while a '+' and '++' indicate satisfactory and very satisfactory results for the given criterion.

## **Between sample comparisons**

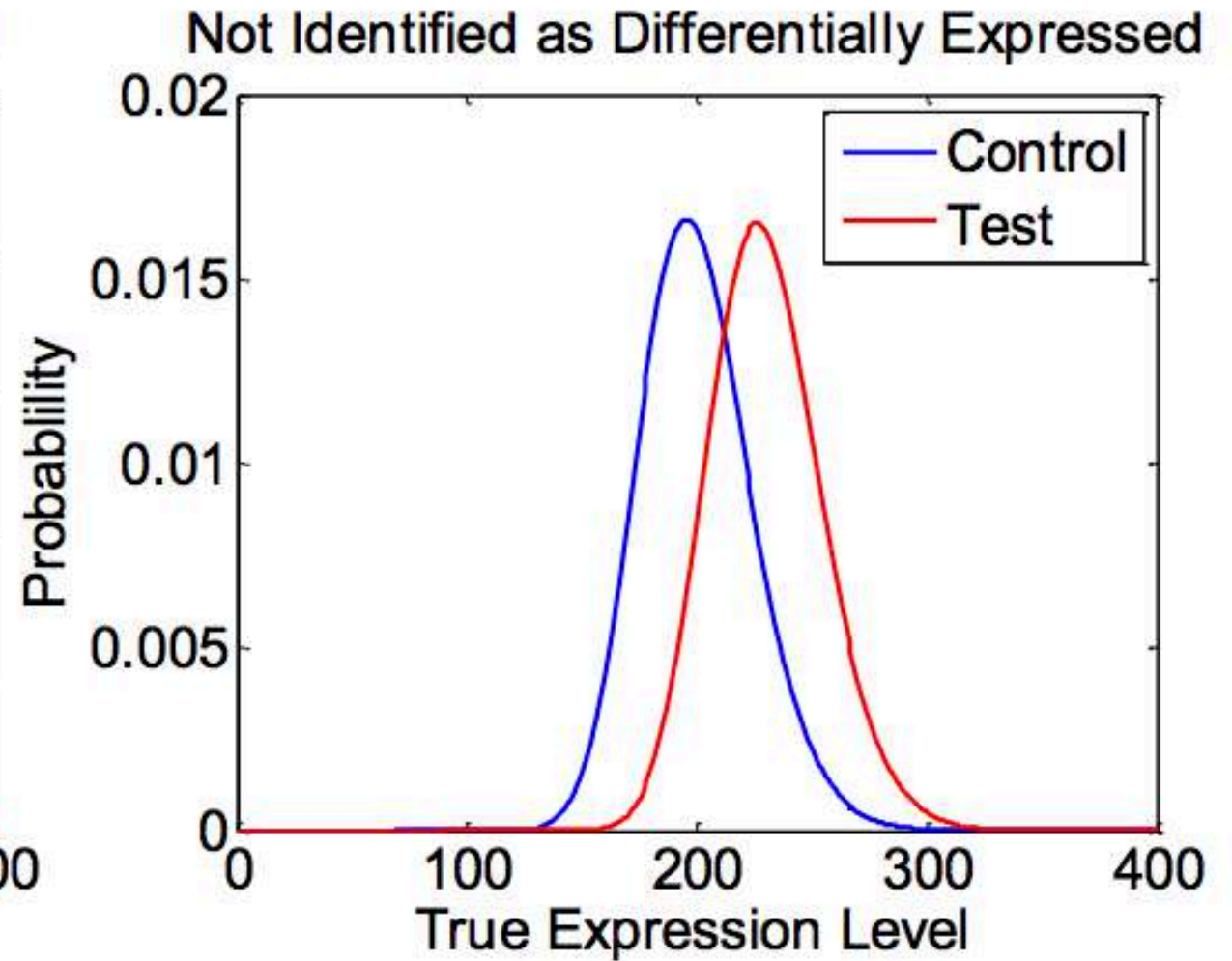
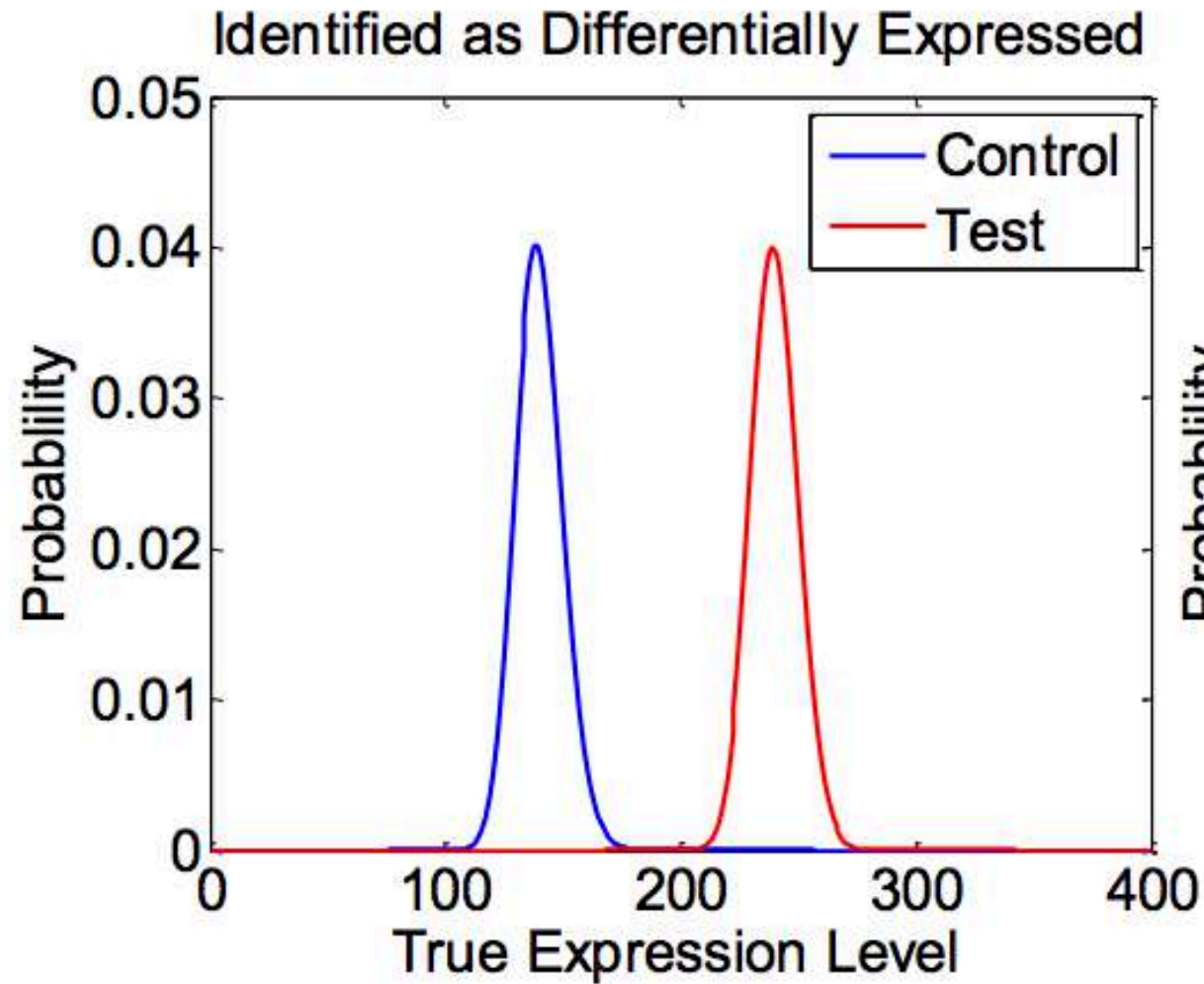
- Differentially expressed genes = counts of genes change between conditions **more systematically** than expected by chance
- Need **biological and technical replicates** to detect differential expression



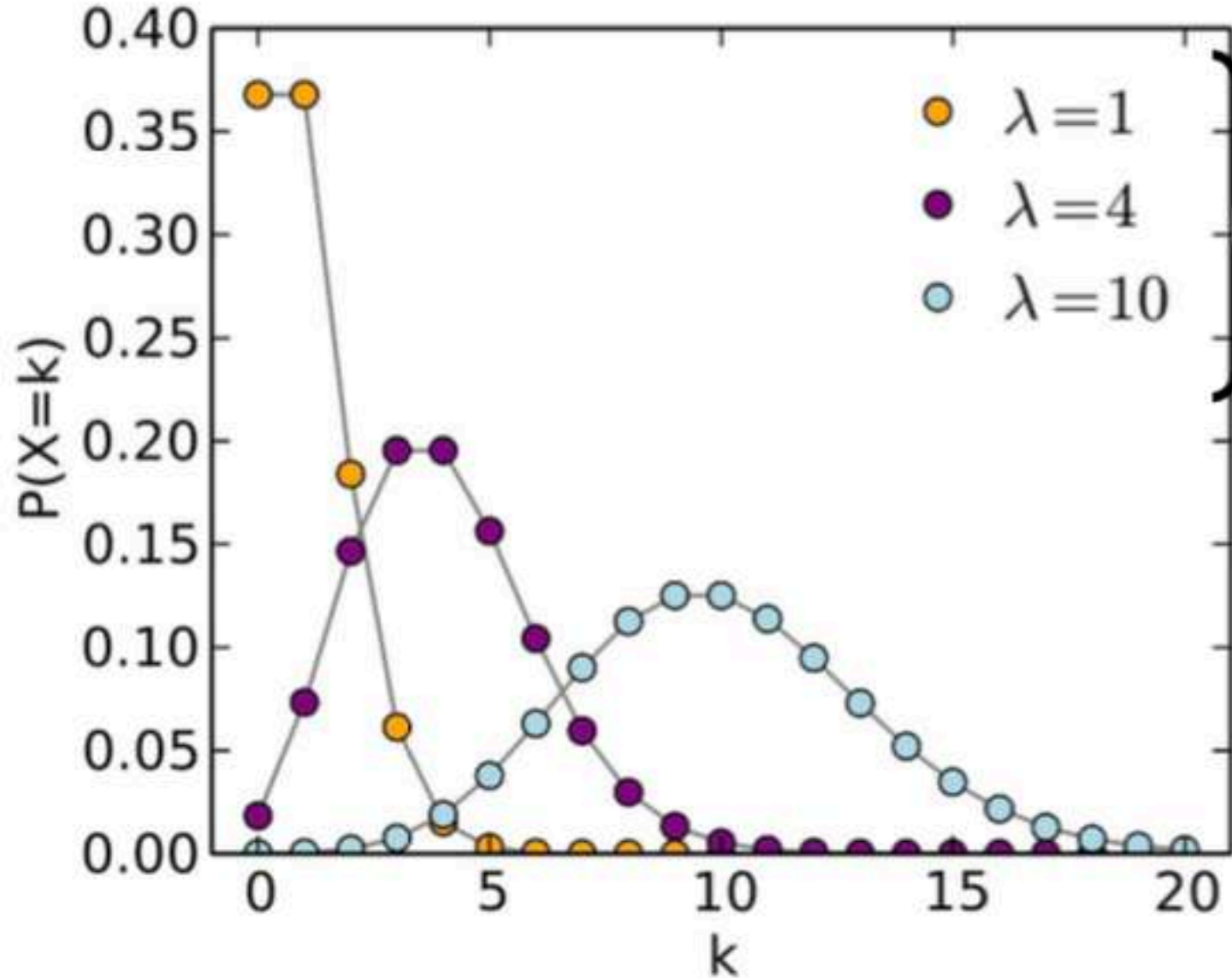
# Replicate categories

	<b>Replicate type</b>	<b>Category</b>
Subjects	Colonies	Biological
	Strains	Biological
	Cohoused groups	Biological
	Gender	Biological
	Individuals	Biological
Sample preparation	Organs from sacrificed animals	Biological
	Methods for dissociating cells from tissue	Technical
	Dissociation runs from given tissue sample	Technical
	Individual cells	Biological
	RNA-seq library construction	Technical
Sequencing	Runs from the library of a given cell	Technical
	Reads from different transcript molecules	Variable
	Reads with unique molecular identifier from a given transcript molecule	Technical

# Fitting a distribution for every gene for **DE**



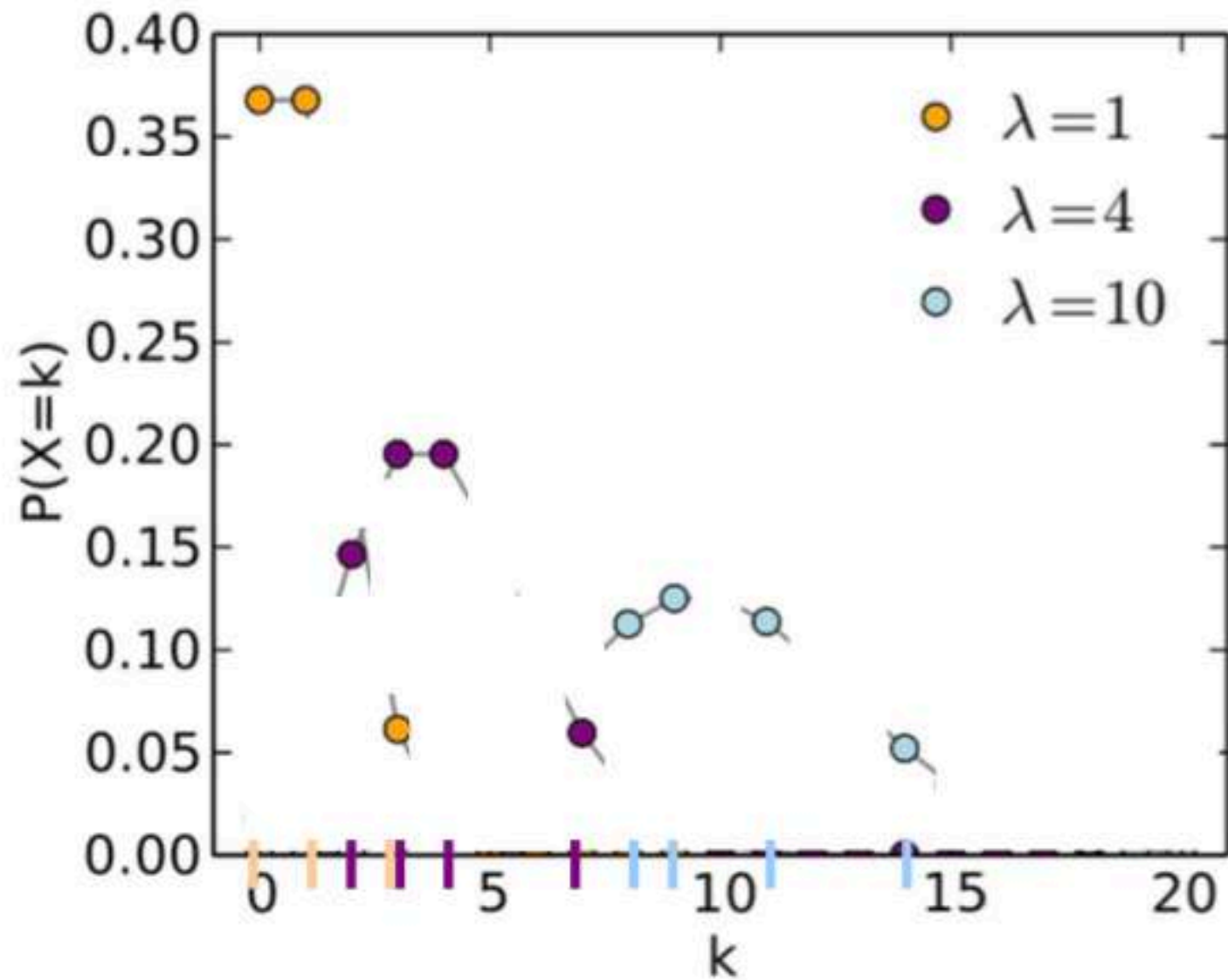
The counts of technical replicates follow a **poisson** distribution (Marioni *et al.*, 2008). So mean = count, variance = count



From Wikipedia. Can be 3 different genes, each with their own poisson distribution. Lambda is the mean of the gene's distribution, with a certain number of reads.

Y-axis: chance to pick that number of reads.

# Four technical replicates



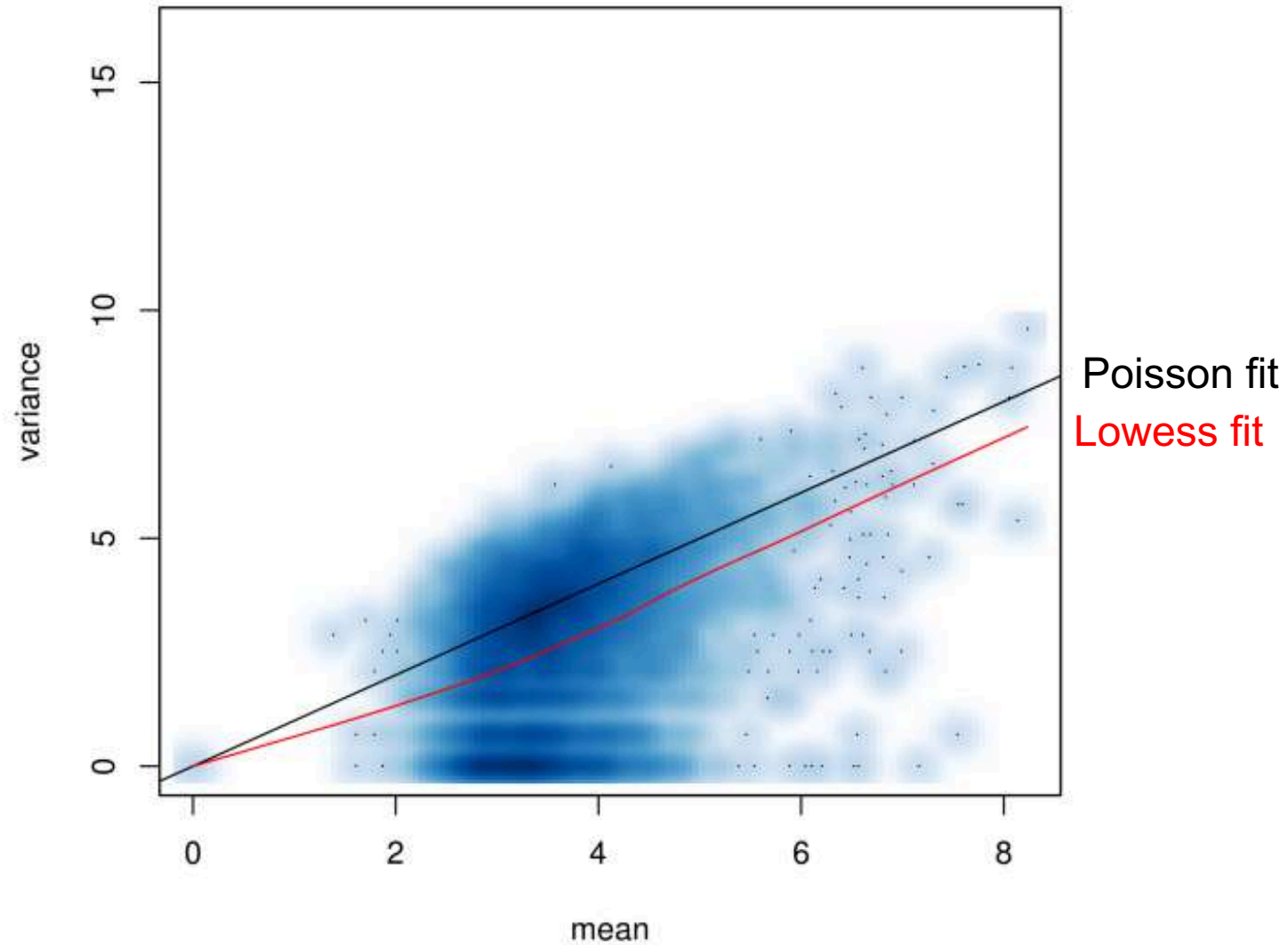
**GeneA** 0, 0, 1, 3

**GeneB** 2, 3, 4, 7

**GeneC** 8, 9, 11, 14



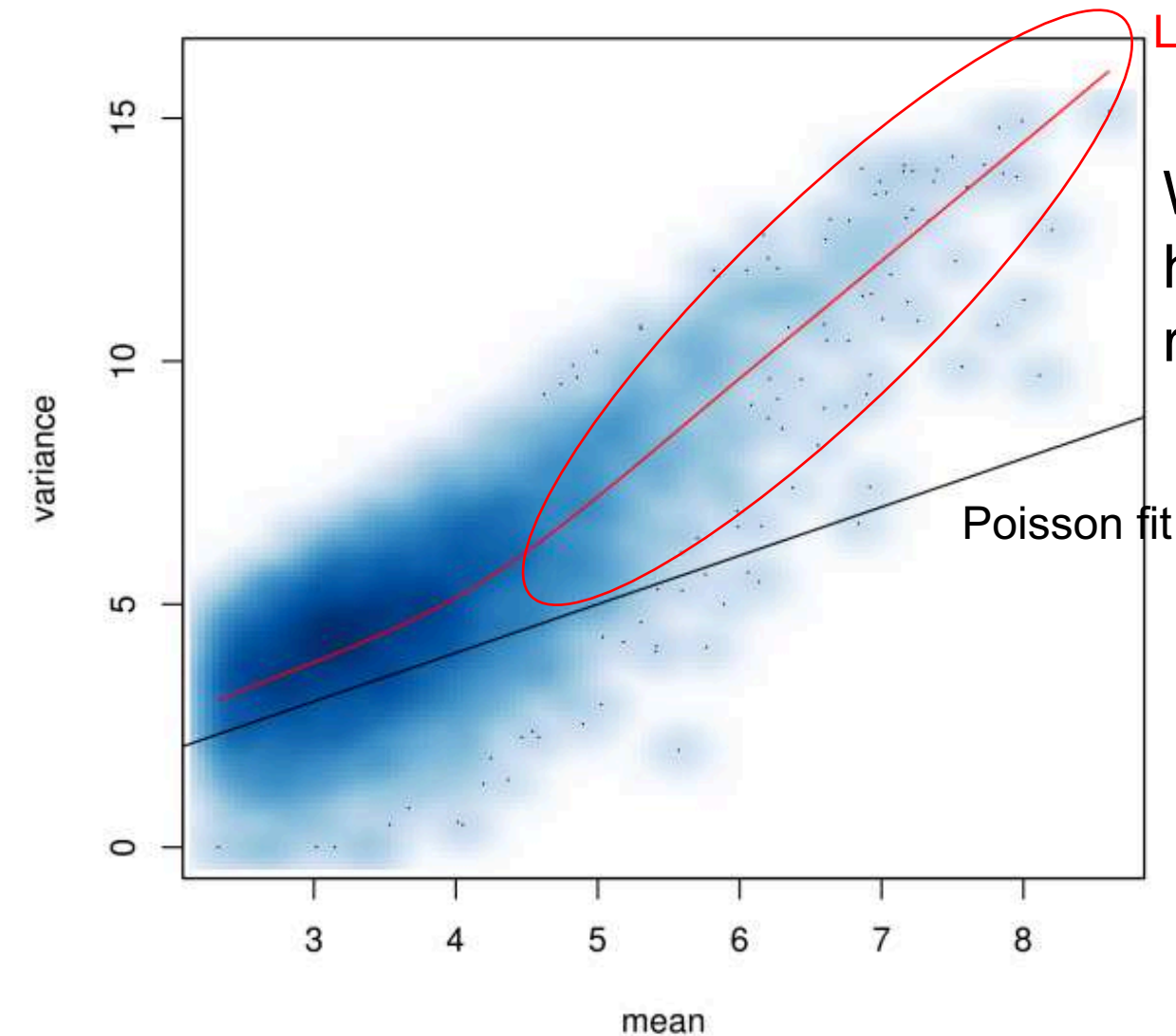
# Poisson model seems good fit in technical replicates



<http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.367.1606&rep=rep1&type=pdf>



# Poisson model seems good fit in technical replicates



Lowess fit

We call this **overdispersion**: the variance is higher for higher counts between biological replicates

Poisson fit

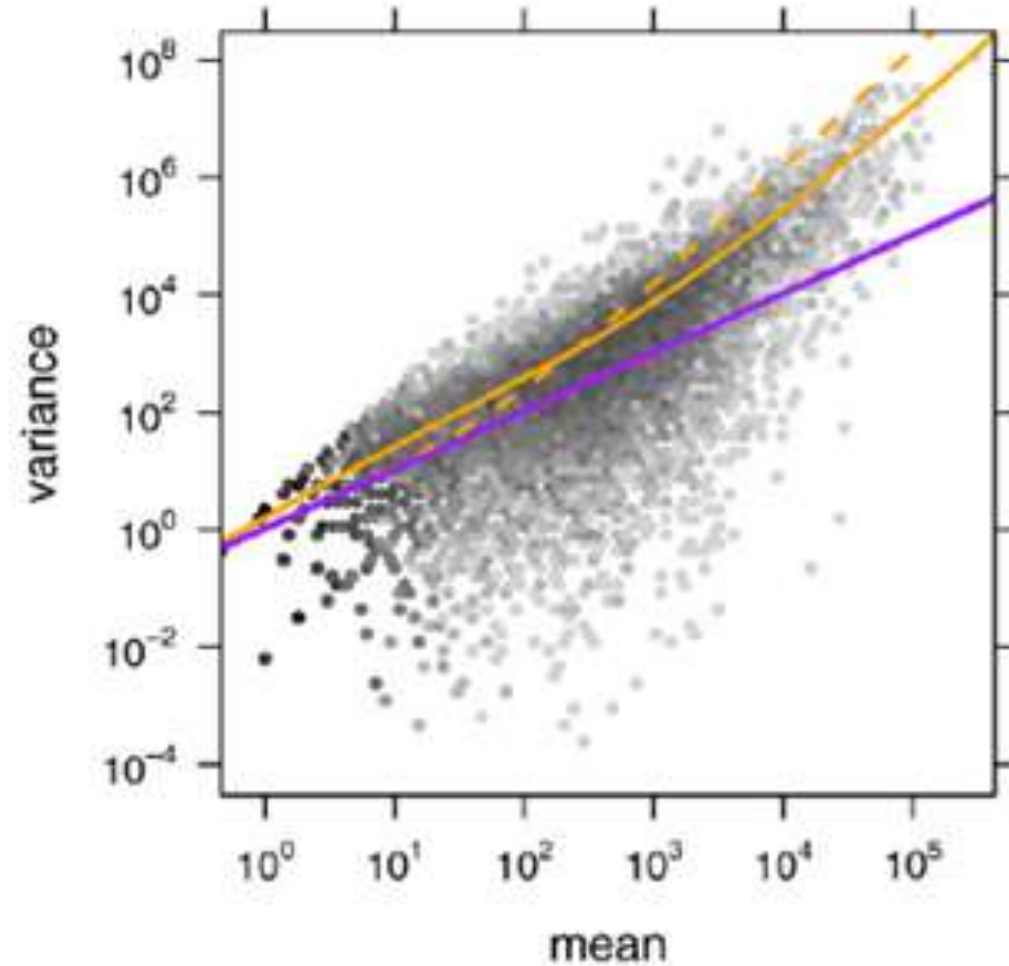
# Variance depends strongly on the mean

Technical replicate: Poisson

Biological replicate: **Negative binomial**

For **low counts**, the Poisson (technical) variation or the measurement error is dominant.

For **higher counts**, the Poisson variation gets smaller, and another source of variation becomes dominant, the **dispersion** or the **biological variation**. Biological variation does not get smaller with higher counts.



- Poisson  $v = \mu$  Poisson distribution
- - - Poisson + constant CV  $v = \mu + \alpha \mu^2$  (edgeR)
- Poisson + local regression  $v = \mu + f(\mu^2)$  (DESeq)

} Negative binomial distribution

# Lots of Differential Gene Expression methods

**Table 1** Methods for calling differentially expressed genes in RNA-seq data analysis. Total citations were based on Google Scholar search result as of 22 September 2015, and normalized by number of years since formal publication. The methods were ranked according to their citations per year.

Method	Total citations	Citations per year	Reference
DESeq <sup>*</sup>	2,987	597	<i>Anders &amp; Huber (2010)</i>
edgeR <sup>*</sup>	2,260	452	<i>Robinson, McCarthy &amp; Smyth (2010)</i>
Cuffdiff2	517	258	<i>Trapnell et al. (2013)</i>
DESeq2 <sup>*</sup>	209	209	<i>Love, Huber &amp; Anders (2014)</i>
voom <sup>*</sup>	143	143	<i>Law et al. (2014)</i>
DEGseq	592	118	<i>Wang et al. (2010)</i>
NOISeq <sup>*,a,b</sup>	324	81	<i>Tarazona et al. (2011)</i>
baySeq	310	62	<i>Hardcastle &amp; Kelly (2010)</i>
SAMSeq <sup>b</sup>	114	57	<i>Li &amp; Tibshirani (2013)</i>
EBSeq	107	53	<i>Leng et al. (2013)</i>
PoissonSeq	99	33	<i>Li et al. (2012)</i>
BitSeq	70	23	<i>Glaus, Honkela &amp; Rattray (2012)</i>
DSS	46	23	<i>Wu, Wang &amp; Wu (2013)</i>
TSPM	70	17	<i>Auer &amp; Doerge (2011)</i>
GPseq	86	17	<i>Srivastava &amp; Chen (2010)</i>
NBPSeq	65	16	<i>Di et al. (2011)</i>
QuasiSeq	47	16	<i>Lund et al. (2012)</i>
GFOLD <sup>*,a</sup>	44	15	<i>Feng et al. (2012)</i>
ShrinkSeq	30	15	<i>Van De Wiel et al. (2013)</i>
NPEBseq <sup>b</sup>	14	7	<i>Bi &amp; Davuluri (2013)</i>
ASC <sup>*,a</sup>	32	6	<i>Wu et al. (2010)</i>
BADGE	2	1	<i>Gu et al. (2014)</i>

**Table 7:** Comparison of programs for differential gene expression identification (Rapaport et al., 2013; Seyednasrollah et al., 2015; Schurch et al., 2015).

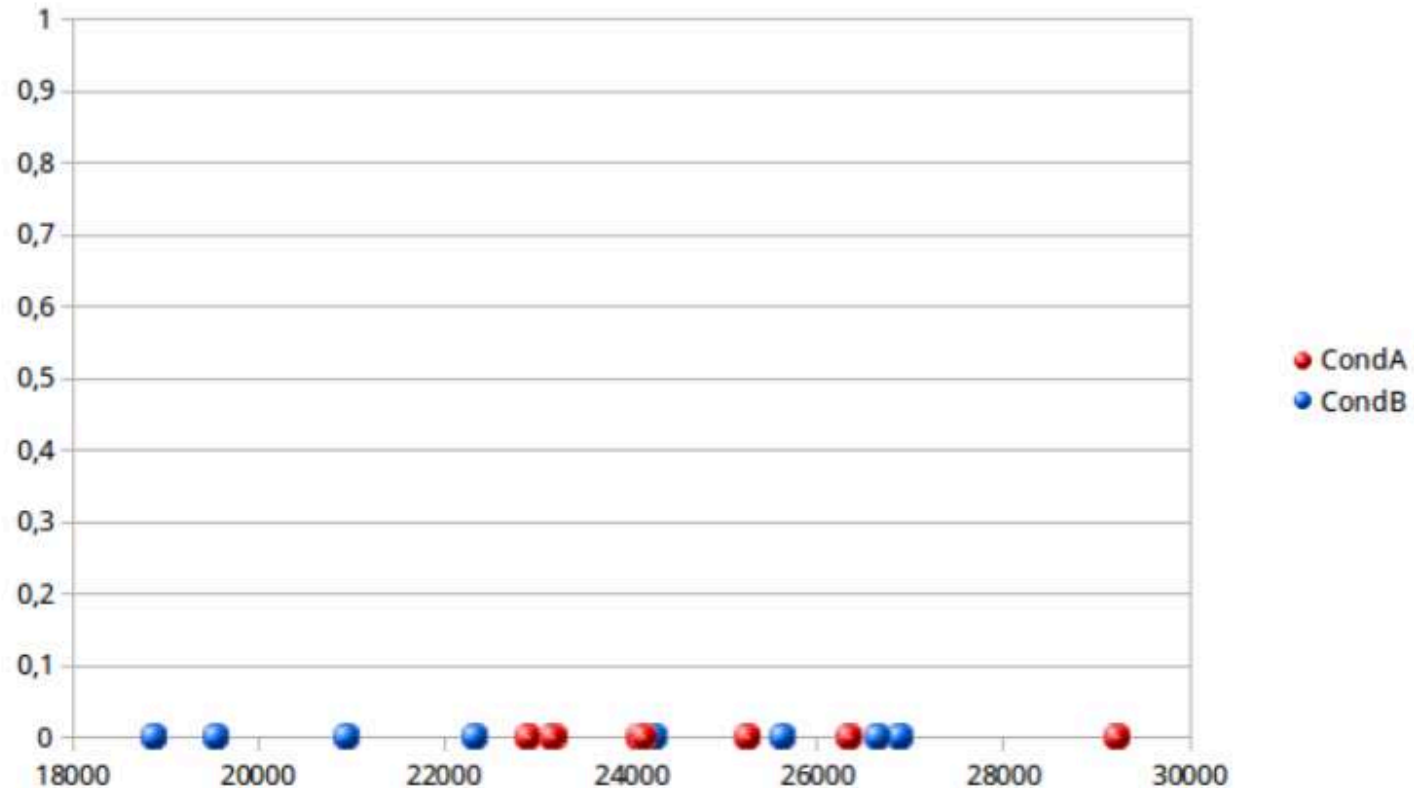
Feature	DESeq2	edgeR	limmaVoom	Cuffdiff
<b>Seq. depth normalization</b>	Sample-wise size factor	Gene-wise trimmed median of means (TMM)	Gene-wise trimmed median of means (TMM)	FPKM-like or DESeq-like
<b>Assumed distribution</b>	Neg. binomial	Neg. binomial	<i>log</i> -normal	Neg. binomial
<b>Test for DE</b>	Exact test (Wald)	Exact test for over-dispersed data	Generalized linear model	<i>t</i> -test
<b>False positives</b>	Low	Low	Low	High
<b>Detection of differential isoforms</b>	No	No	No	Yes
<b>Support for multi-factored experiments</b>	Yes	Yes	Yes	No
<b>Runtime (3-5 replicates)</b>	Seconds to minutes	Seconds to minutes	Seconds to minutes	Hours



# Scenario

gene\_id CAF0006876

<b>Condition A</b>	sample1 23171	sample2 22903	sample3 29227	sample4 24072	sample5 23151	sample6 26336	sample7 25252	sample8 24122
<b>Condition B</b>	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629

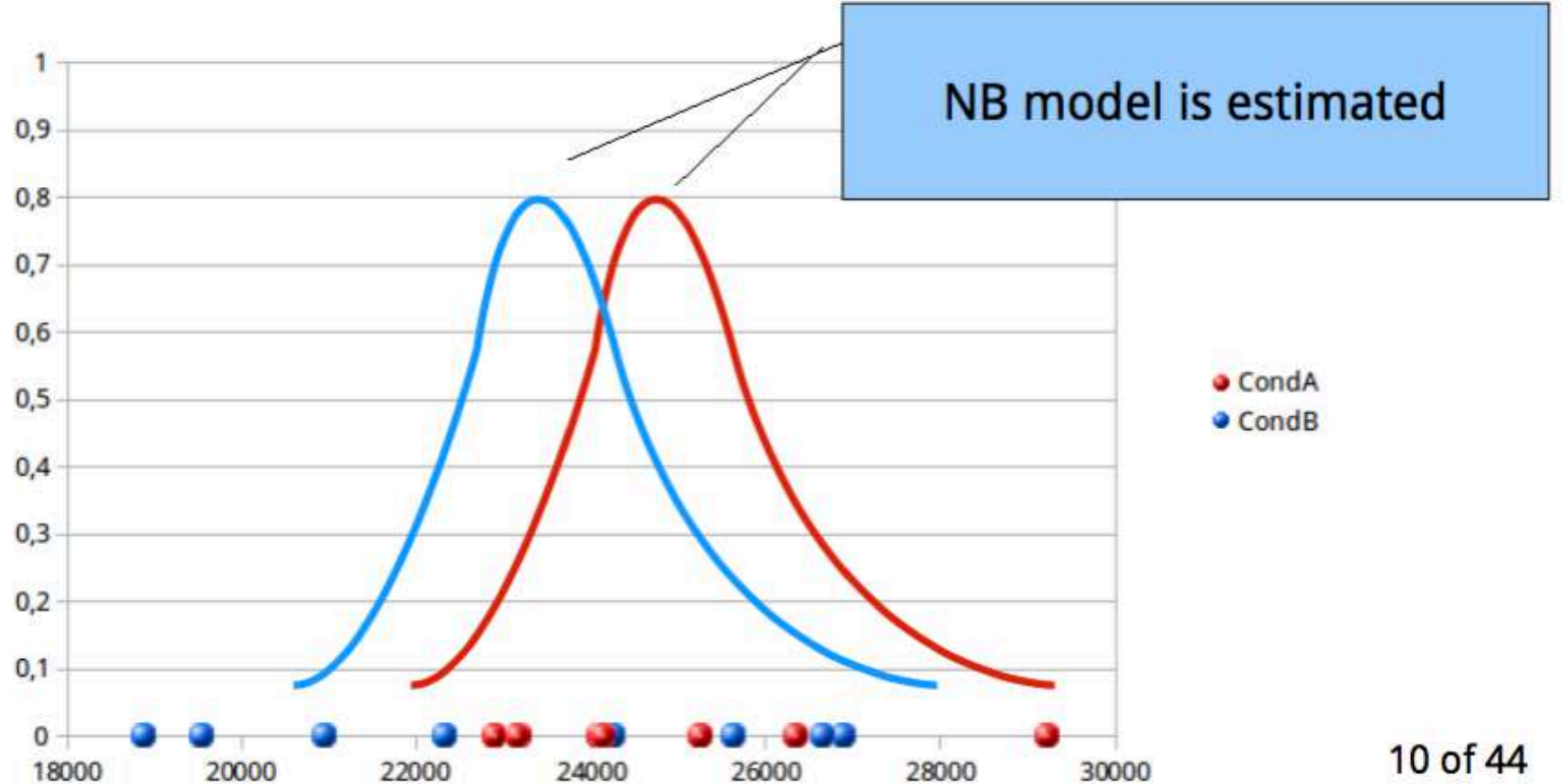




# Scenario

gene\_id CAF0006876

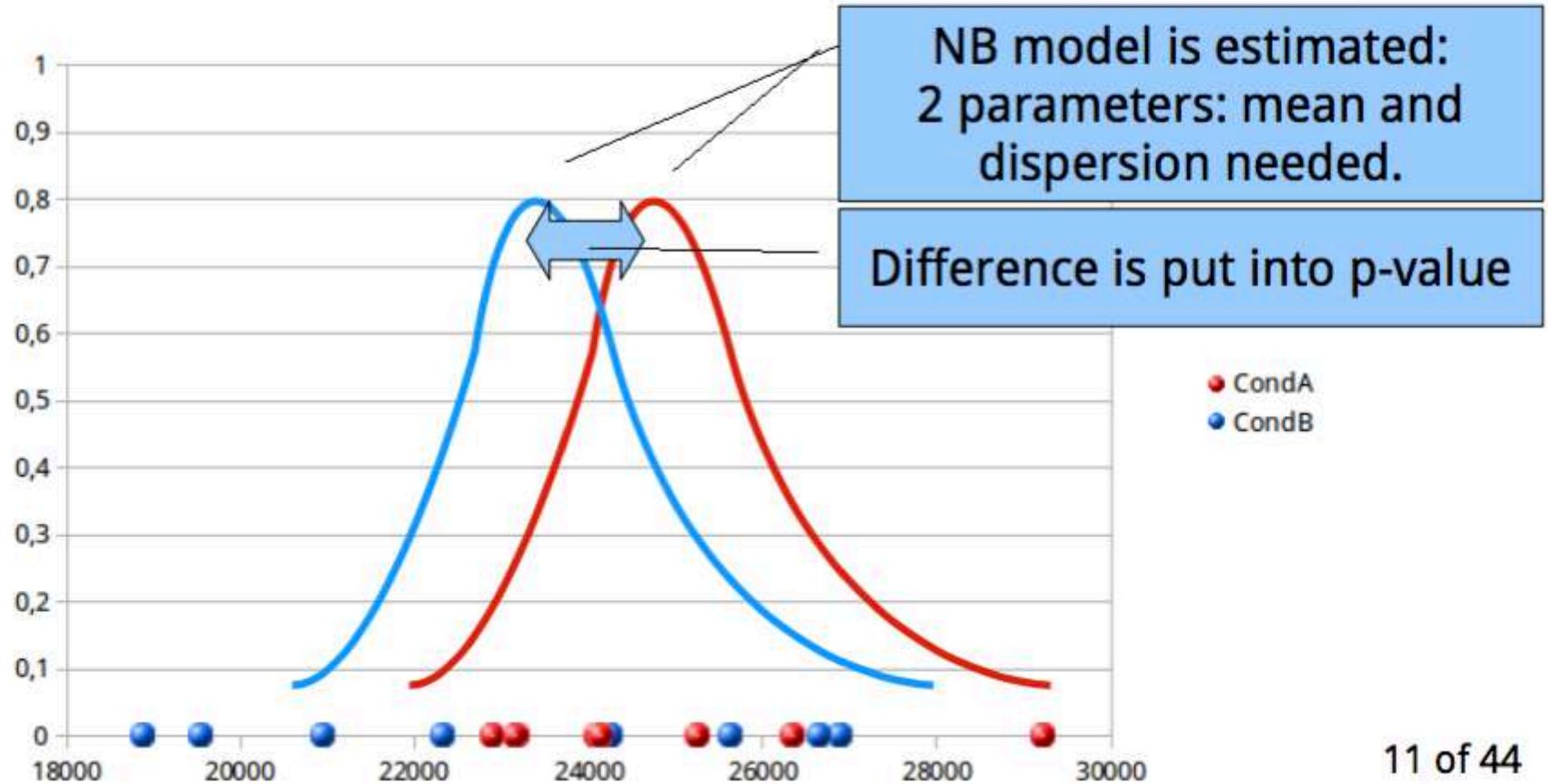
<b>Condition A</b>	sample1 23171	sample2 22903	sample3 29227	sample4 24072	sample5 23151	sample6 26336	sample7 25252	sample8 24122
<b>Condition B</b>	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629



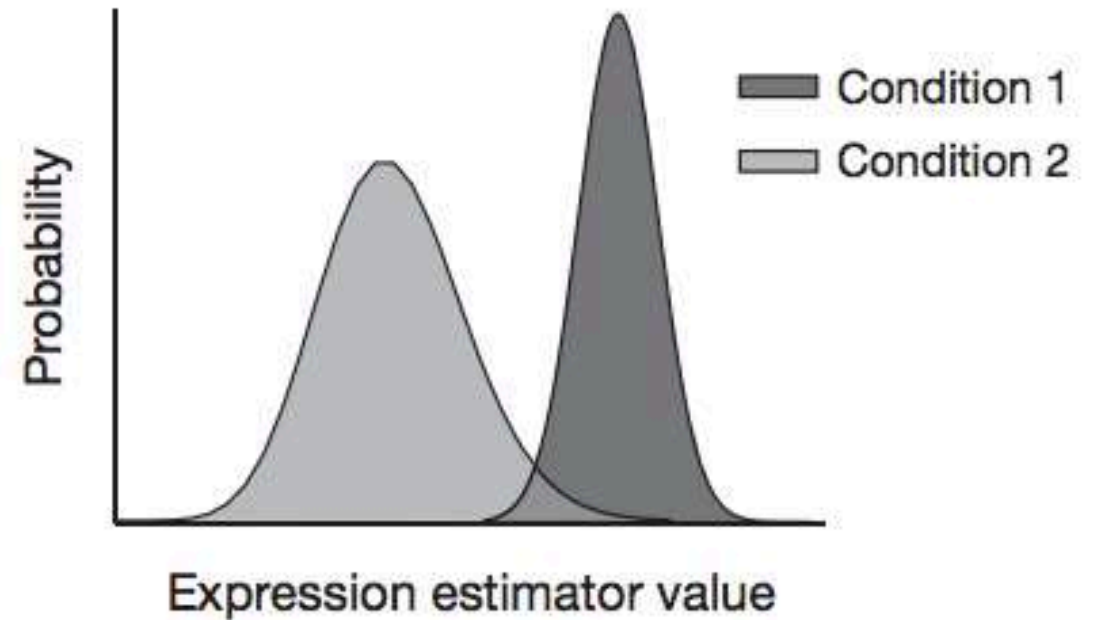
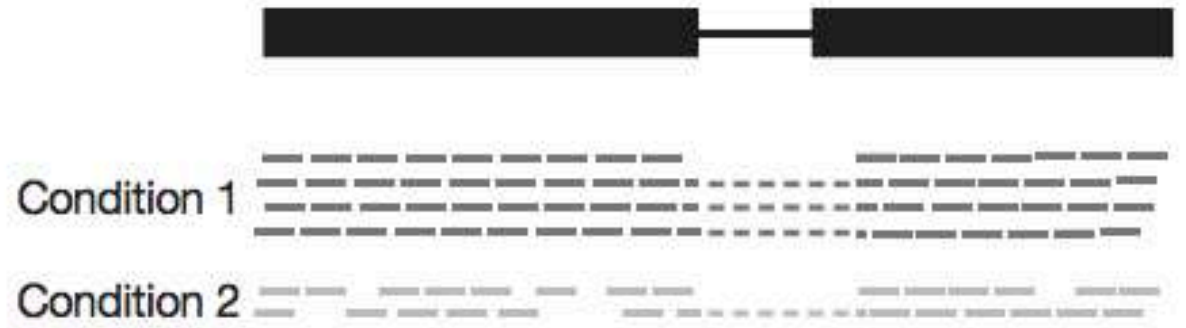
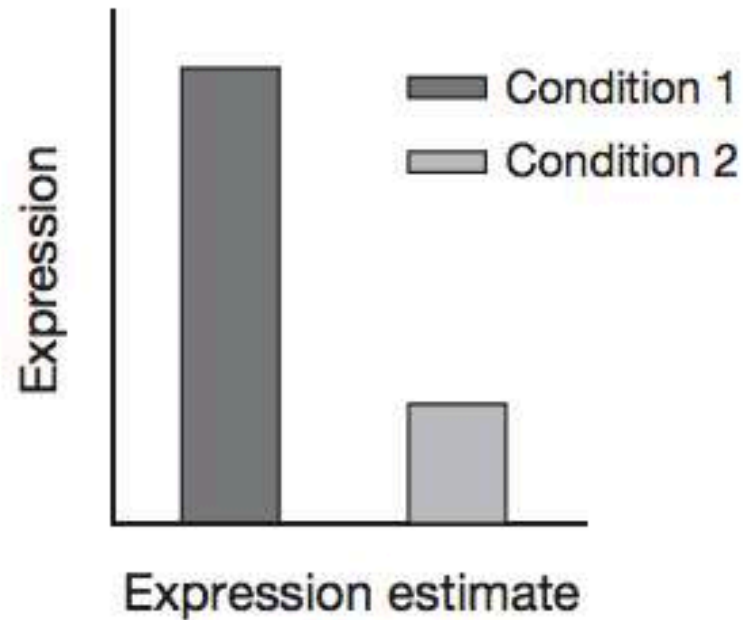
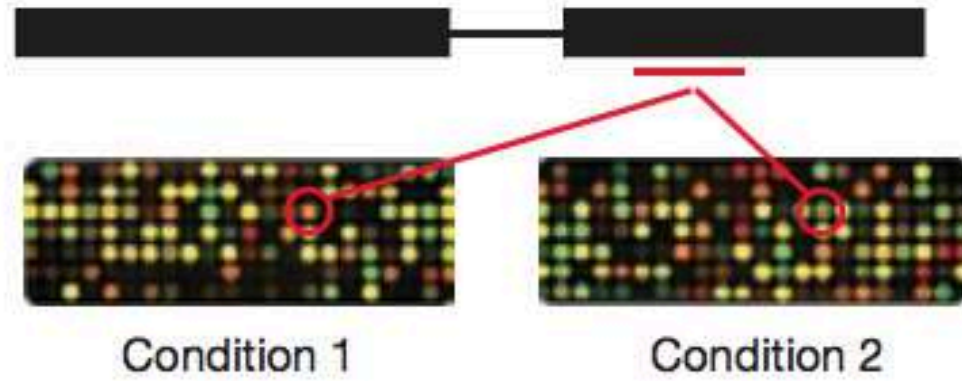
# Scenario

gene\_id CAF0006876

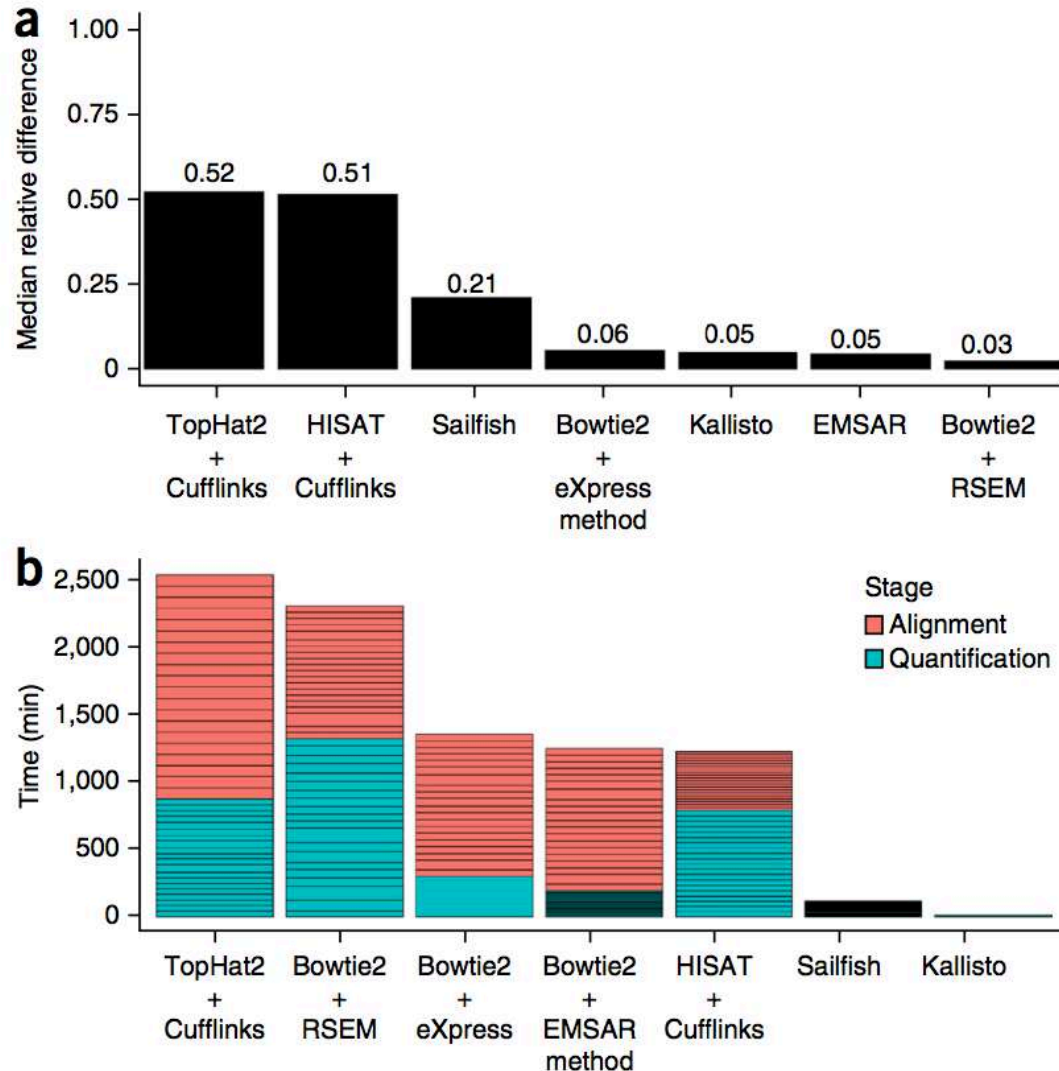
<b>Condition A</b>	sample1 23171	sample2 22903	sample3 29227	sample4 24072	sample5 23151	sample6 26336	sample7 25252	sample8 24122
<b>Condition B</b>	Sample9 19527	sample10 26898	sample11 18880	sample12 24237	sample13 26640	sample14 22315	sample15 20952	sample16 25629



# RNAseq vs Microarray



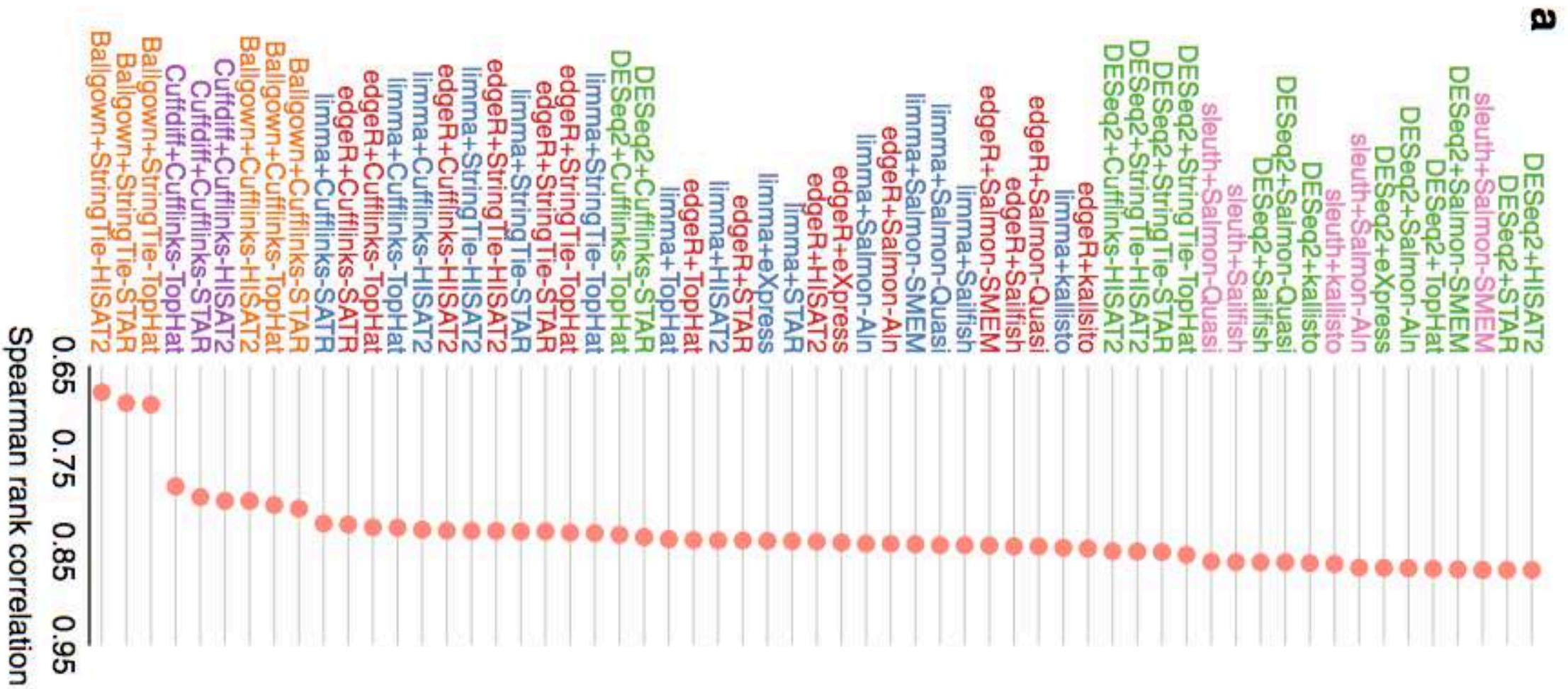
# Advances in quantification



We present kallisto, an RNA-seq quantification program that is two orders of magnitude faster than previous approaches and achieves similar accuracy. Kallisto pseudoaligns reads to a reference, producing a list of transcripts that are compatible with each read while avoiding alignment of individual bases. We use kallisto to analyze 30 million unaligned paired-end RNA-seq reads in <10 min on a standard laptop computer. This removes a major computational bottleneck in RNA-seq analysis.



# Spearman rank correlation of DEG results to qPCR measured genes

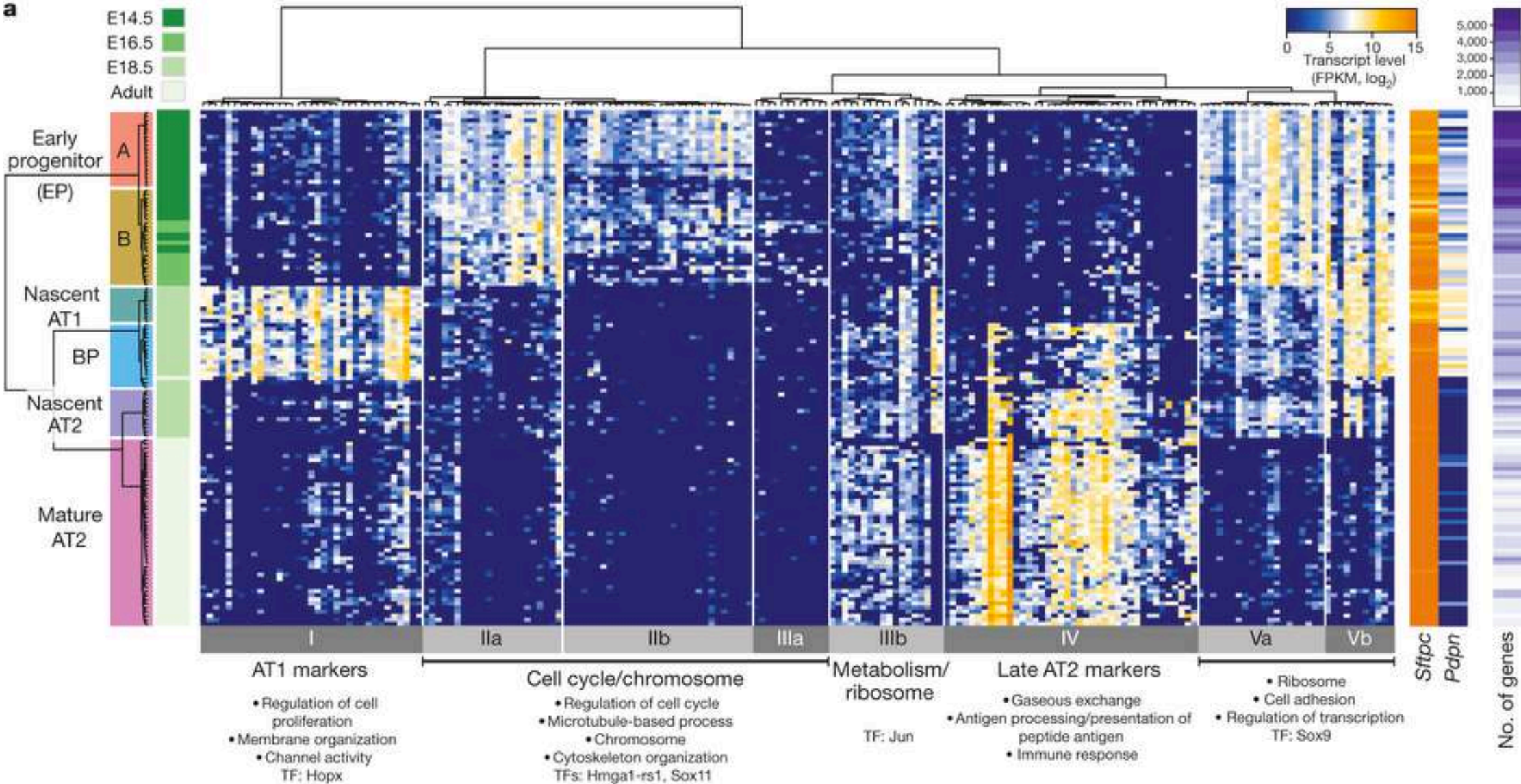




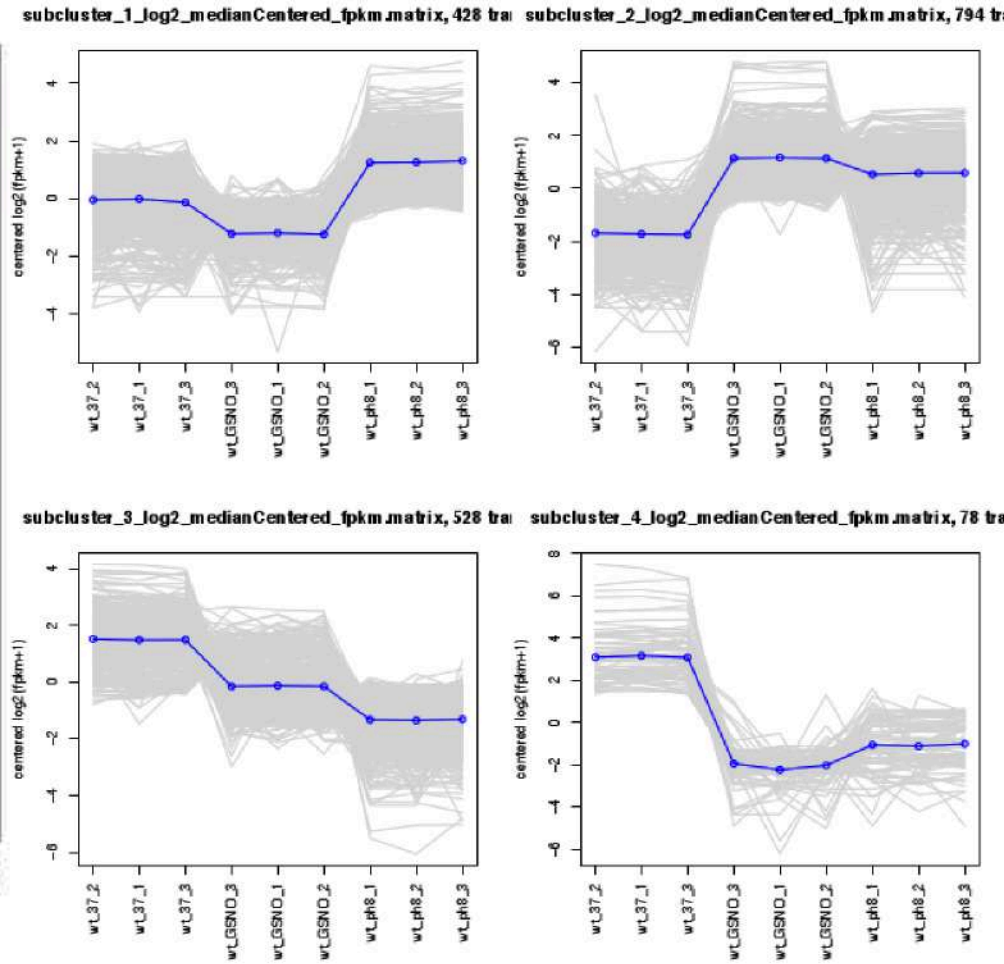
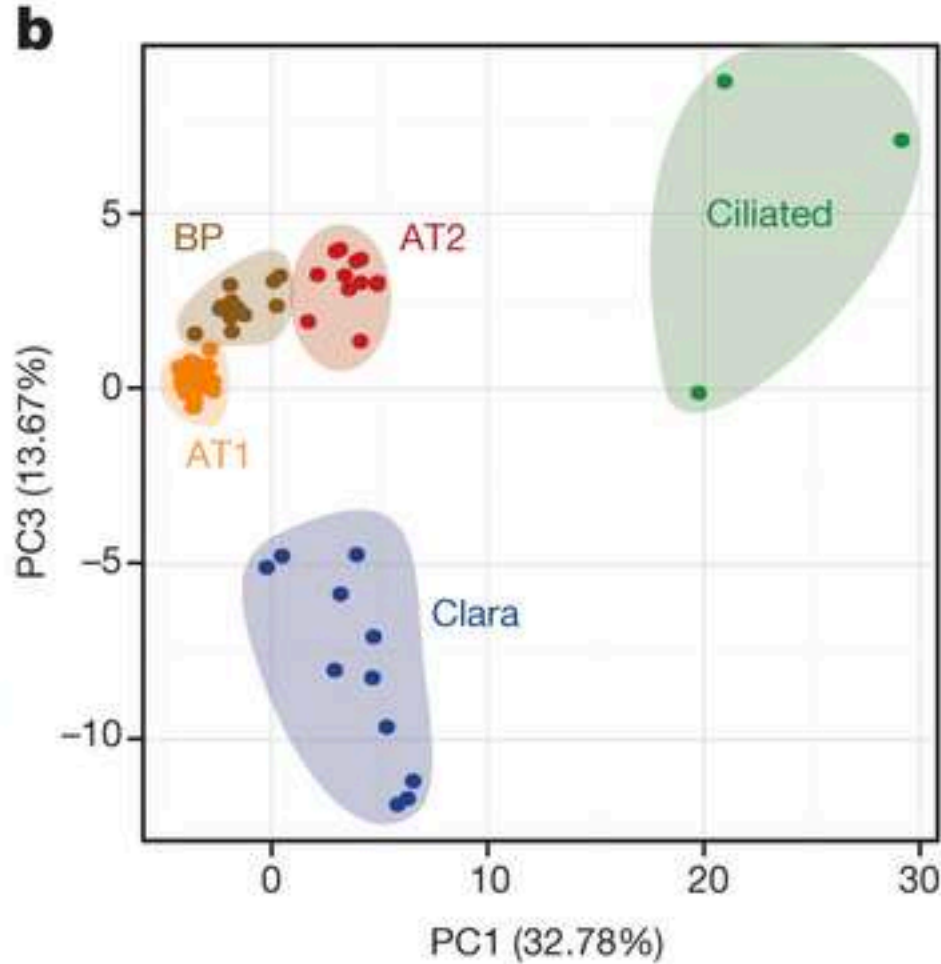
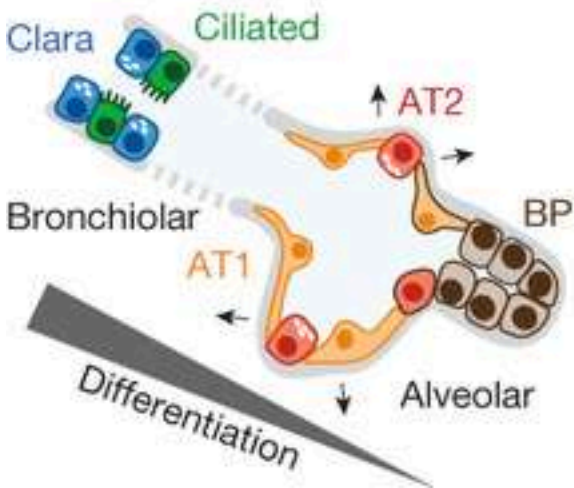
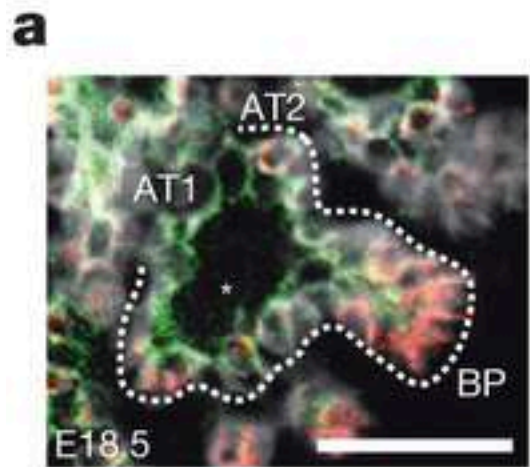
Further analyses

3.2 Once you have set of differentially expressed genes

# Summarization visualizing the expression data through heatmap ; Classification using Gene Ontology terms and metabolic annotations



# Clustering of the expression values and principal component analysis to reduce the variables.

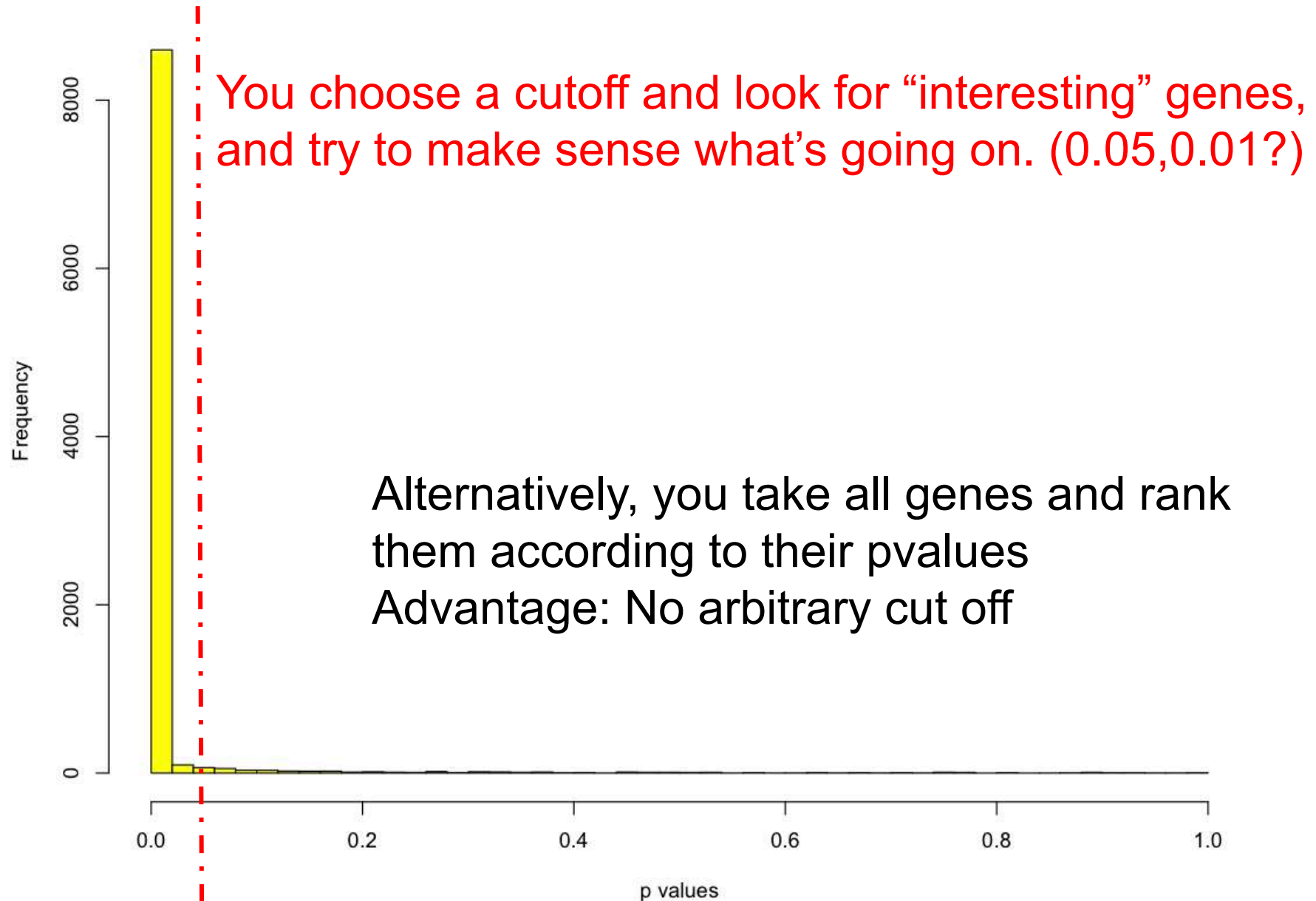




# Gene Ontology Enrichment analysis

<b>GO.ID</b>	<b>Term</b>	<b>Annotated Terms in list</b>		<b>Expected</b>	<b>p values</b>
GO:0044281	small molecule metabolic process	481	150	54.18	< 1e-30
GO:0017144	drug metabolic process	155	72	17.46	3.20E-29
GO:0055114	oxidation-reduction process	308	103	34.7	4.30E-27
GO:0009126	purine nucleoside monophosphate metaboli.	79	47	8.9	2.10E-25
GO:0009167	purine ribonucleoside monophosphate meta	79	47	8.9	2.10E-25
GO:0072521	purine-containing compound metabolic pro..	129	61	14.53	2.30E-25
GO:0006163	purine nucleotide metabolic process	122	59	13.74	3.40E-25
GO:0009150	purine ribonucleotide metabolic process	119	58	13.41	5.30E-25
GO:0007218	neuropeptide signaling pathway	108	55	12.17	5.80E-25
GO:0019693	ribose phosphate metabolic process	138	62	15.55	3.00E-24
GO:0009161	ribonucleoside monophosphate metabolic p.	87	48	9.8	6.50E-24
GO:0009259	ribonucleotide metabolic process	129	59	14.53	1.30E-23
GO:0009117	nucleotide metabolic process	178	70	20.05	4.20E-23
GO:0006082	organic acid metabolic process	246	84	27.71	9.50E-23
GO:0019752	carboxylic acid metabolic process	232	81	26.13	1.20E-22
GO:0006753	nucleoside phosphate metabolic process	181	70	20.39	1.30E-22
GO:0009123	nucleoside monophosphate metabolic proce	97	49	10.93	4.00E-22
GO:0043436	oxoacid metabolic process	242	82	27.26	5.80E-22
GO:0055086	nucleobase-containing small molecule met..	204	74	22.98	7.10E-22
GO:0006091	generation of precursor metabolites and ...	111	52	12.5	1.70E-21

# Now, setting a cut-off?

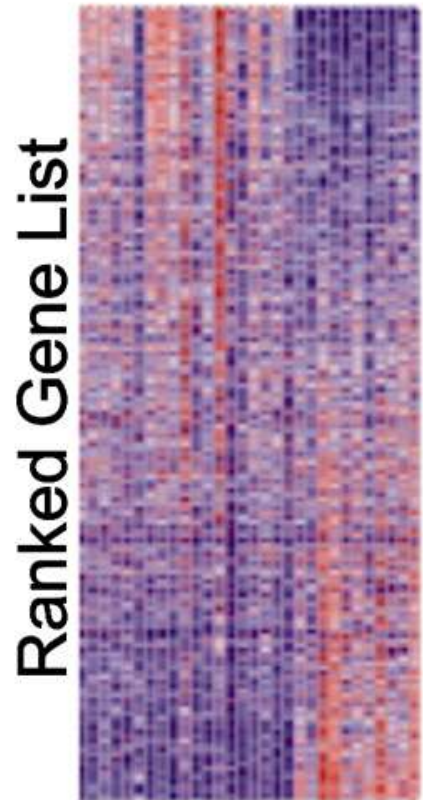




# GSEA (Gene Set Enrichment Analysis) methods (cut-off free approach)

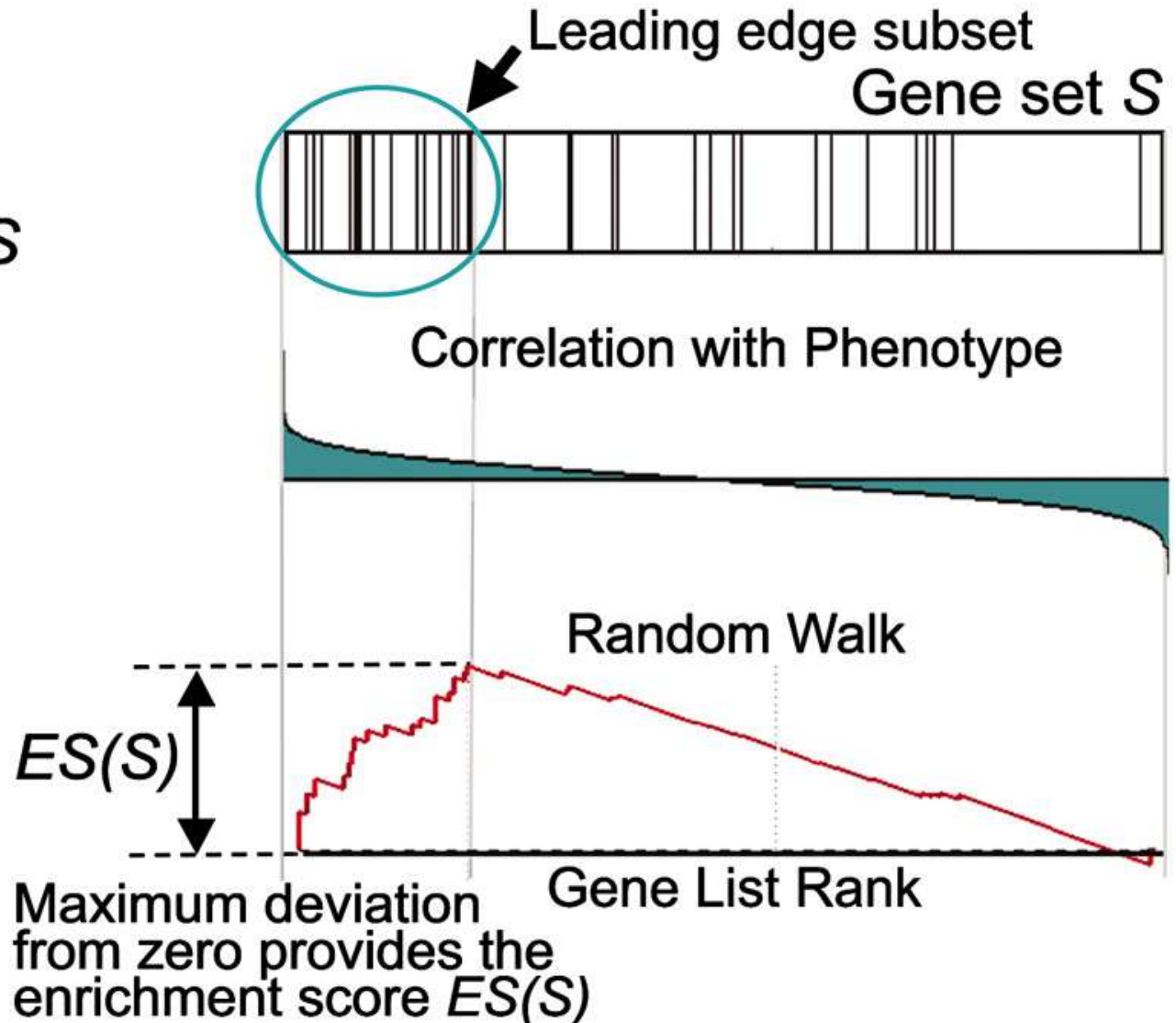
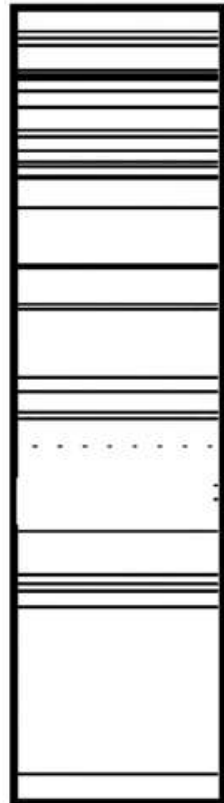
**A** Phenotype  
Classes

A B

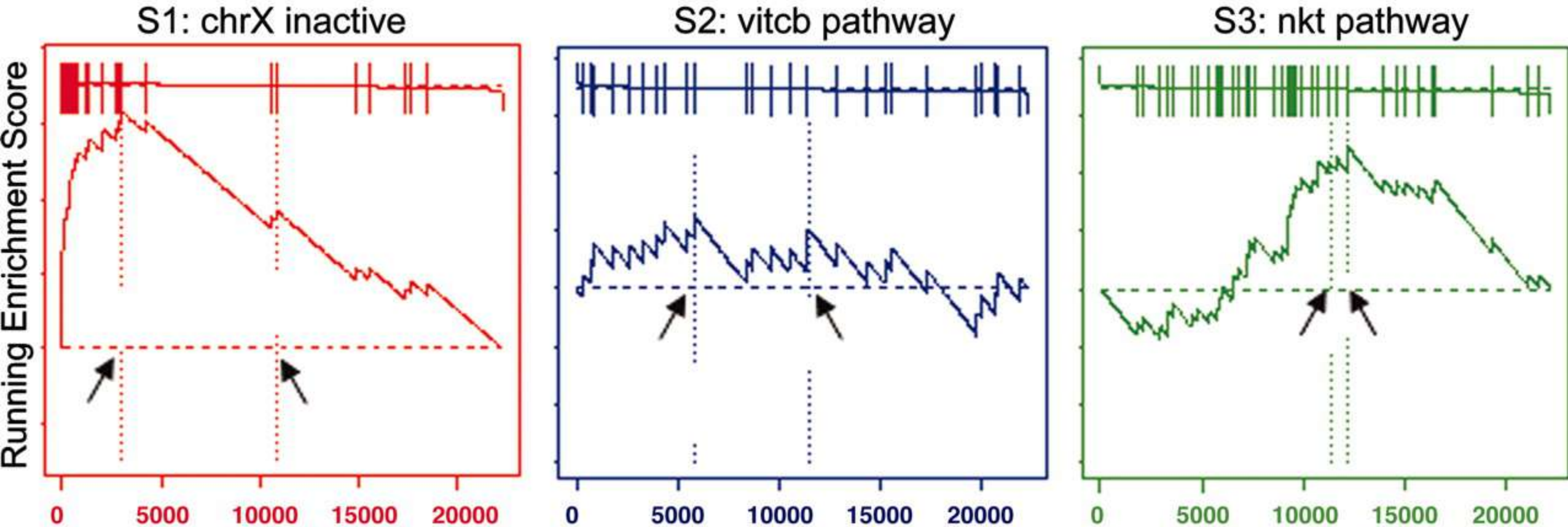


**B**

Gene set S

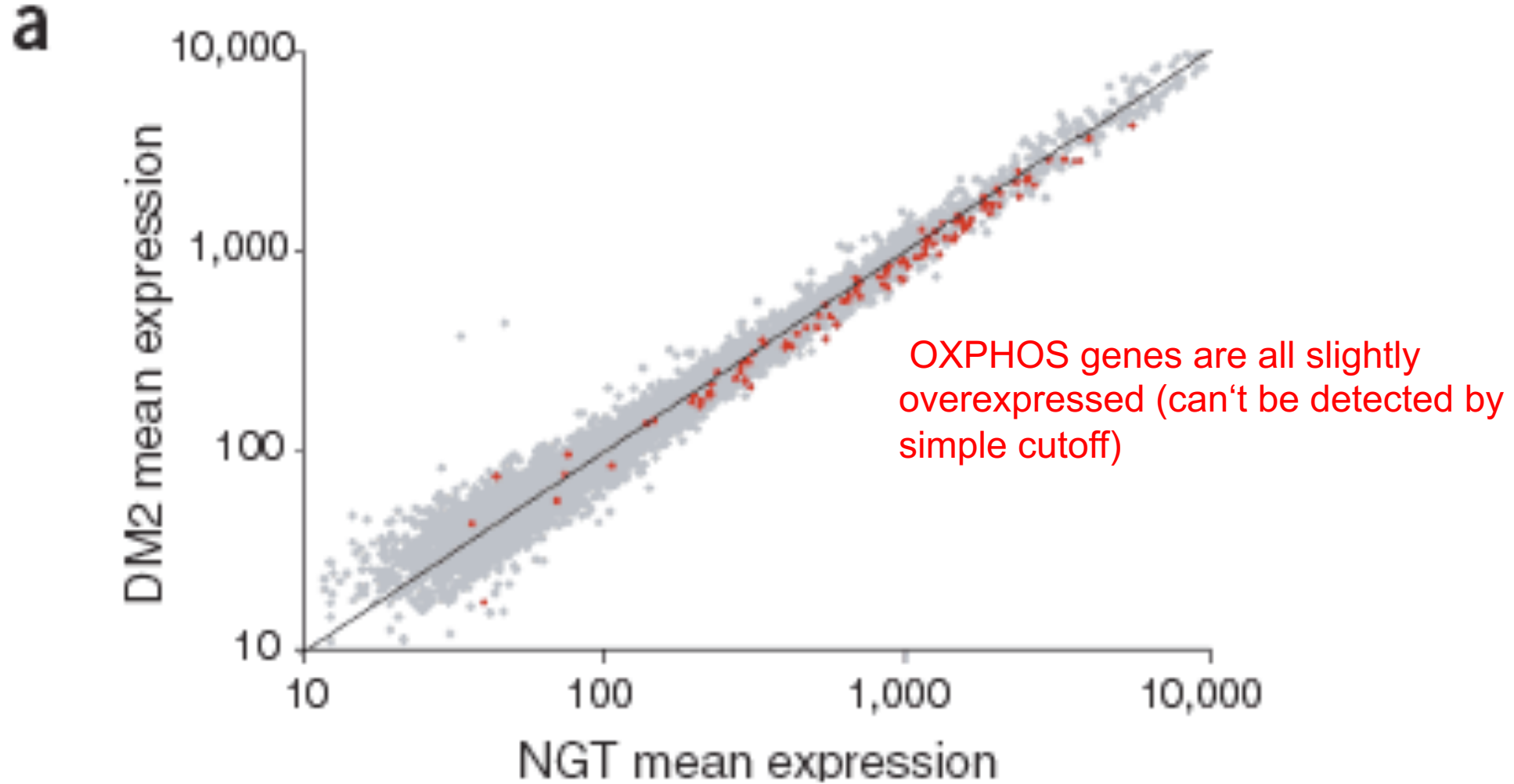


# GSEA (Gene Set Enrichment Analysis) methods (cut-off free approach)

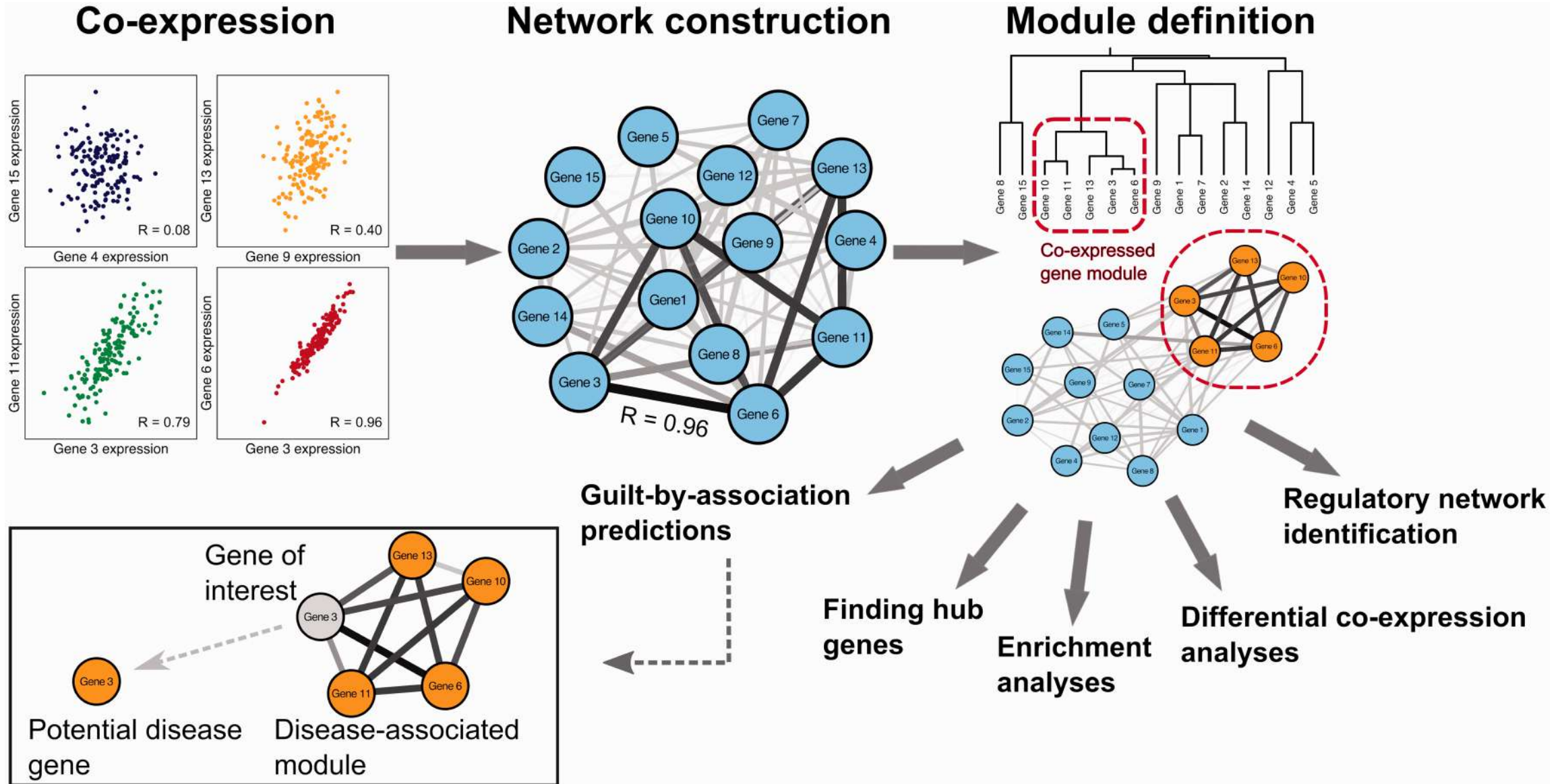


S1 is significantly enriched in females as expected, S2 is randomly distributed and scores poorly, and S3 is not enriched at the top of the list but is nonrandom, so it scores well.

# GSEA (Gene Set Enrichment Analysis) methods (cut-off free approach)



# Gene co-expression network





# Gene co-expression network

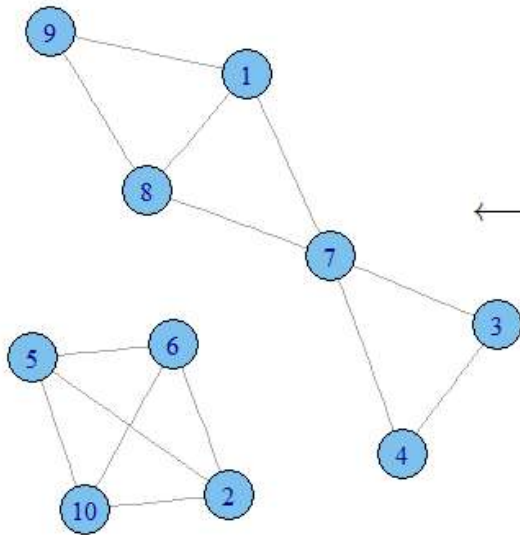
	$S_1$	$S_2$	$S_3$
$G_1$	43.26	40.89	5.05
$G_2$	166.6	41.87	136.65
$G_3$	12.53	39.55	42.09
$G_4$	28.77	191.92	236.56
$G_5$	114.7	79.7	99.76
$G_6$	119.1	80.57	114.59
$G_7$	118.9	156.69	186.95
$G_8$	3.76	2.48	136.78
$G_9$	32.73	11.99	118.8
$G_{10}$	17.46	56.11	21.41

Gene expression values

$|r(G_i, G_j)|$   
 Pearson correlation

	$G_1$	$G_2$	$G_3$	$G_4$	$G_5$	$G_6$	$G_7$	$G_8$	$G_9$	$G_{10}$
$G_1$	1.00	0.23	0.61	0.71	0.03	0.35	<b>0.86</b>	<b>1.00</b>	<b>0.97</b>	0.37
$G_2$	0.23	1.00	0.63	0.52	<b>0.98</b>	<b>0.99</b>	0.29	0.30	0.46	<b>0.99</b>
$G_3$	0.61	0.63	1.00	<b>0.99</b>	0.77	0.53	<b>0.93</b>	0.56	0.41	0.51
$G_4$	0.71	0.52	<b>0.99</b>	1.00	0.69	0.41	<b>0.97</b>	0.66	0.52	0.40
$G_5$	0.03	<b>0.98</b>	0.77	0.69	1.00	<b>0.95</b>	0.48	0.09	0.27	<b>0.94</b>
$G_6$	0.35	<b>0.99</b>	0.53	0.41	<b>0.95</b>	1.00	0.17	0.41	0.57	<b>1.00</b>
$G_7$	0.86	0.29	<b>0.93</b>	<b>0.97</b>	0.48	0.17	1.00	<b>0.83</b>	0.72	0.16
$G_8$	<b>1.00</b>	0.30	0.56	0.66	0.09	0.41	0.83	1.00	<b>0.98</b>	0.42
$G_9$	<b>0.97</b>	0.46	0.41	0.52	0.27	0.57	0.72	<b>0.98</b>	1.00	0.58
$G_{10}$	0.37	<b>0.99</b>	0.51	0.40	<b>0.94</b>	<b>1.00</b>	0.16	0.42	0.58	1.00

Similarity (Co-expression) score



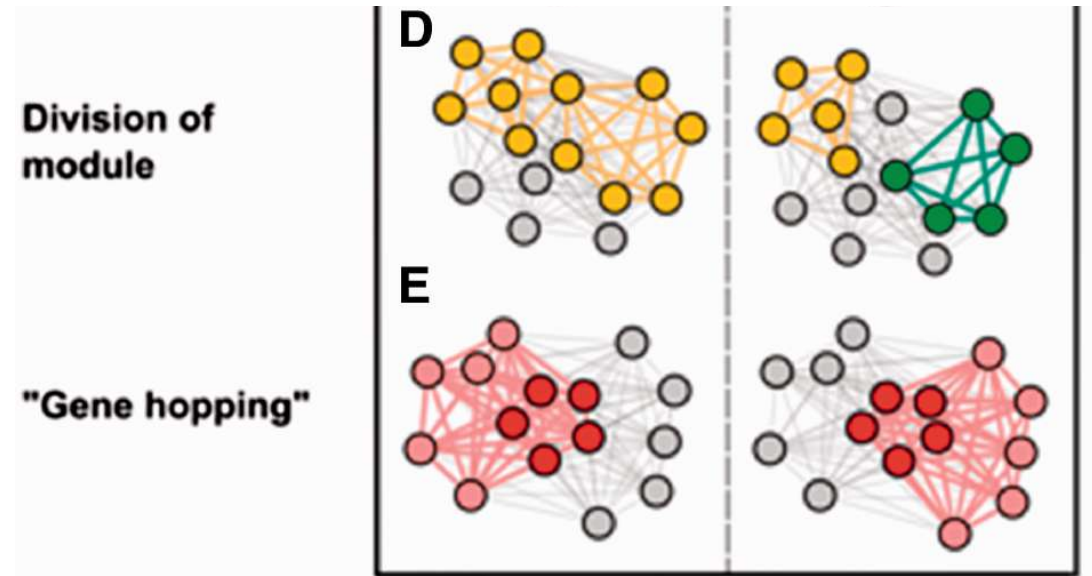
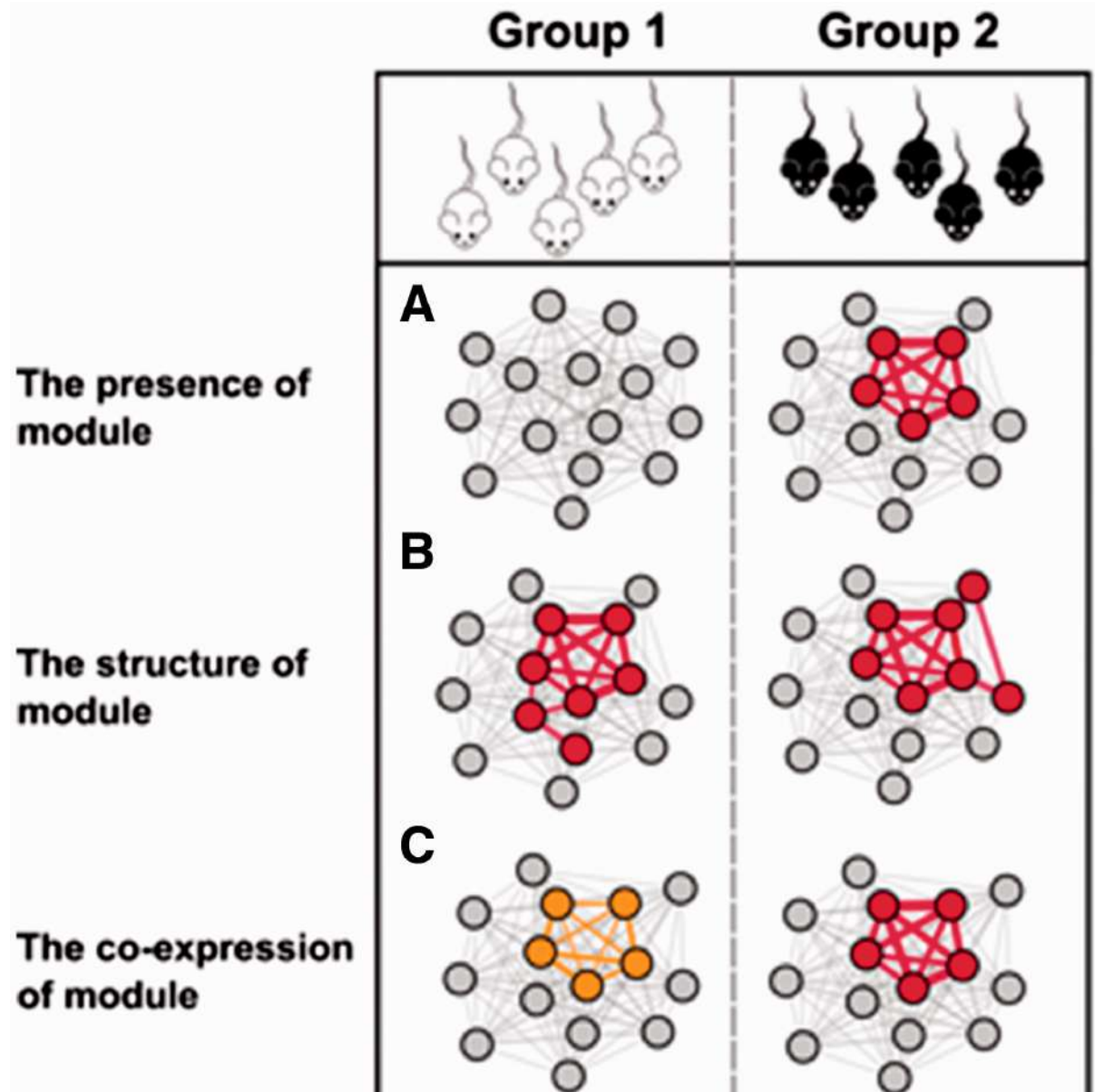
	$G_1$	$G_2$	$G_3$	$G_4$	$G_5$	$G_6$	$G_7$	$G_8$	$G_9$	$G_{10}$
$G_1$	0	0	0	0	0	0	1	1	1	0
$G_2$	0	0	0	0	1	1	0	0	0	1
$G_3$	0	0	0	1	0	0	1	0	0	0
$G_4$	0	0	1	0	0	0	1	0	0	0
$G_5$	0	1	0	0	0	1	0	0	0	1
$G_6$	0	1	0	0	1	0	0	0	0	1
$G_7$	1	0	1	1	0	0	0	1	0	0
$G_8$	1	0	0	0	0	0	1	0	1	0
$G_9$	1	0	0	0	0	0	0	1	0	0
$G_{10}$	0	1	0	0	1	1	0	0	0	0

Network adjacency matrix

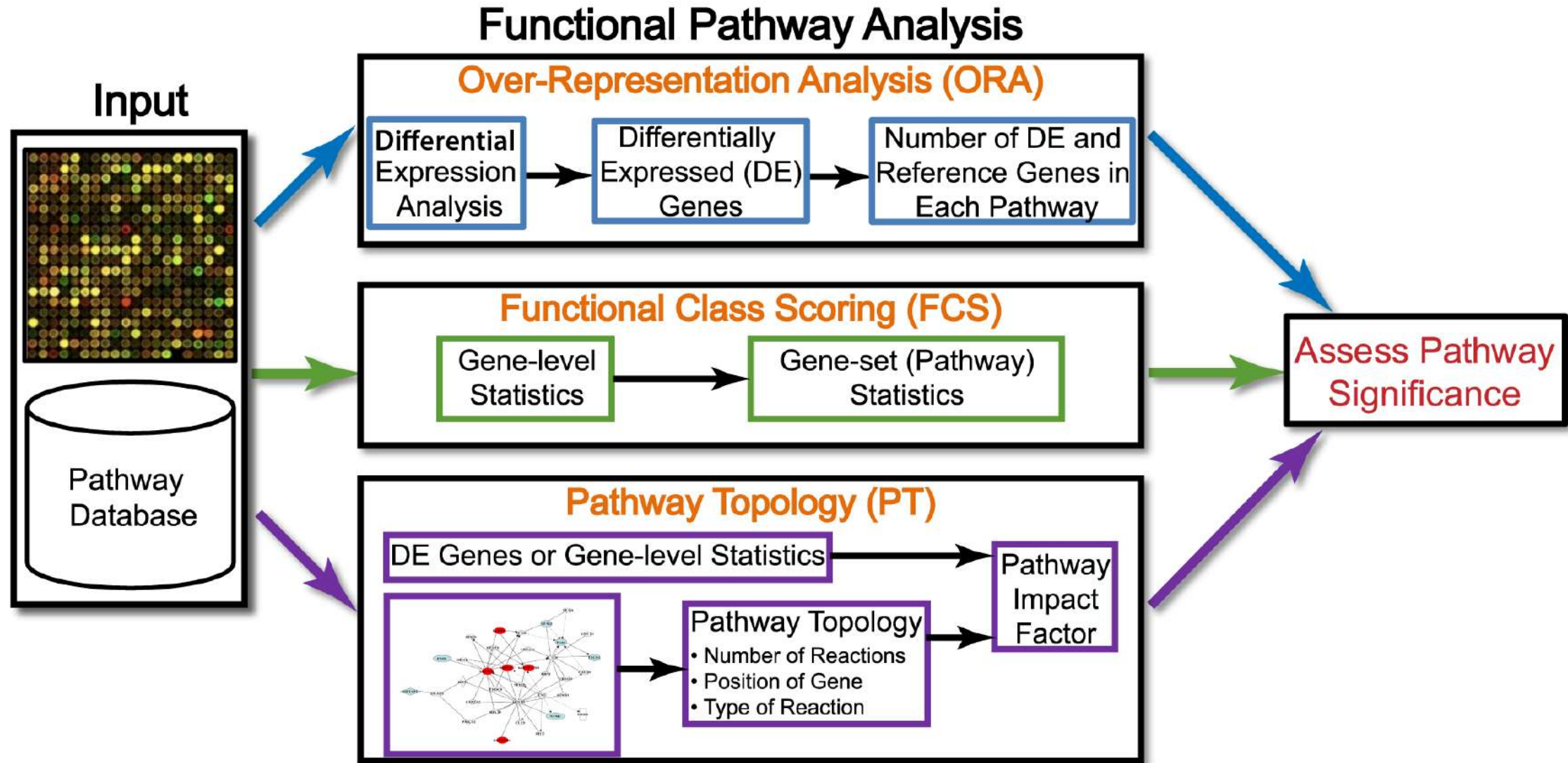
$|r(G_i, G_j)| \geq 0.8$   
 Significance threshold



# Gene co-expression network



# Overview of existing pathway analysis methods using gene expression data as an example (only applicable to model species)

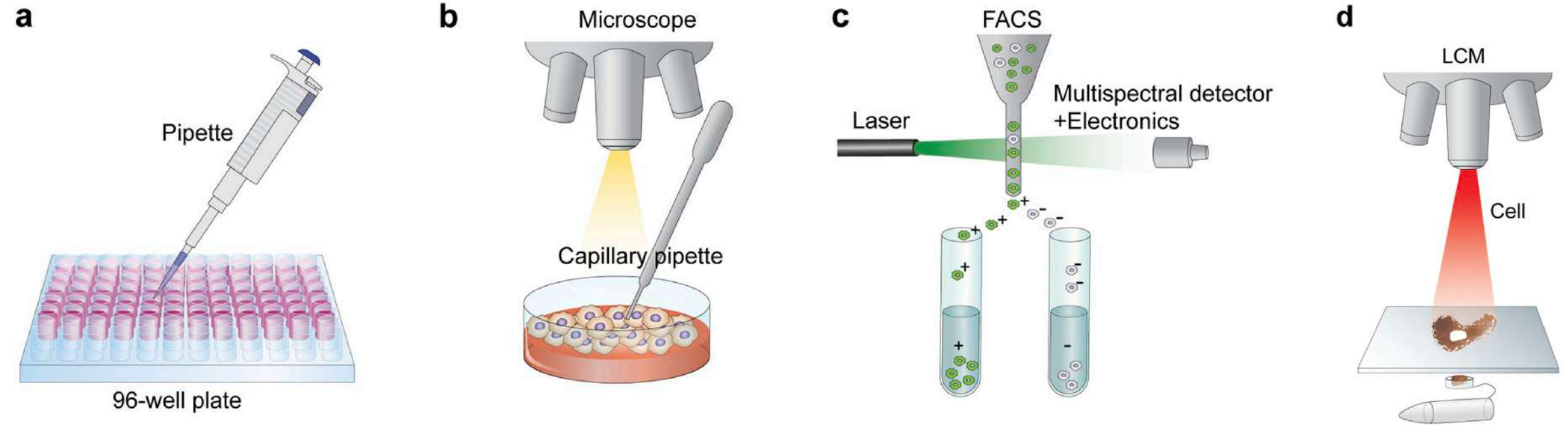


Break

Further advances

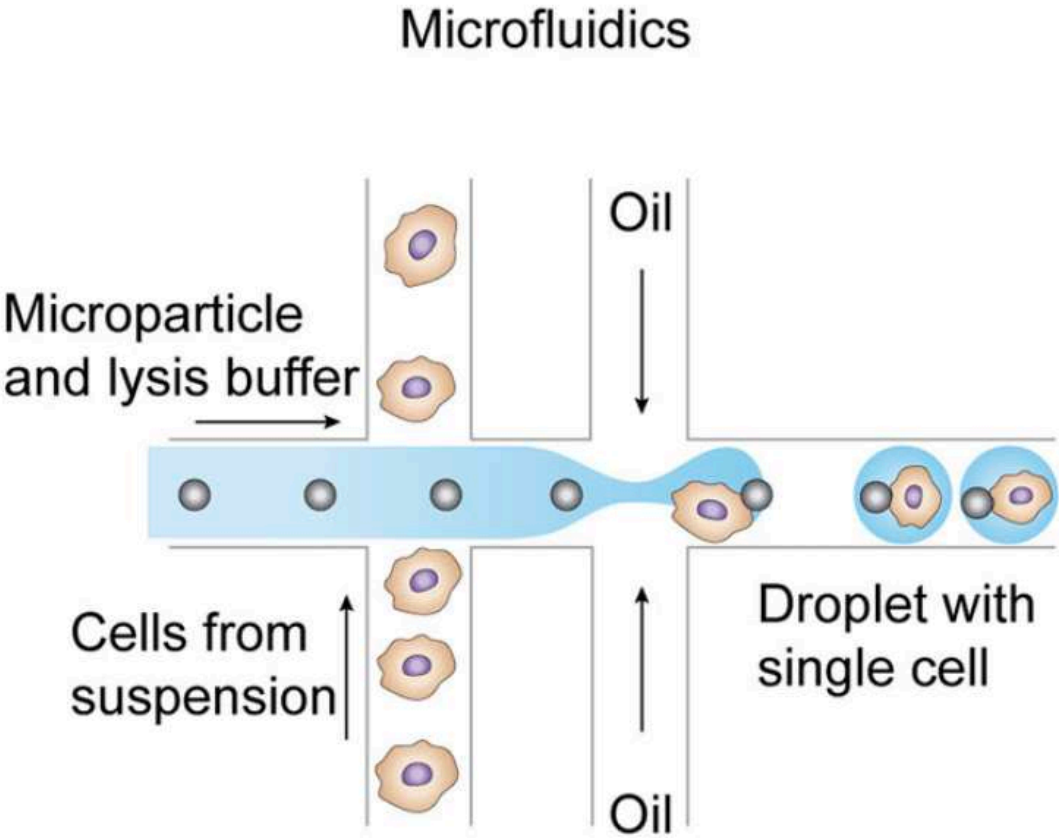
Single cell RNAseq (ScRNAseq)

# Evolution of single-cell isolation

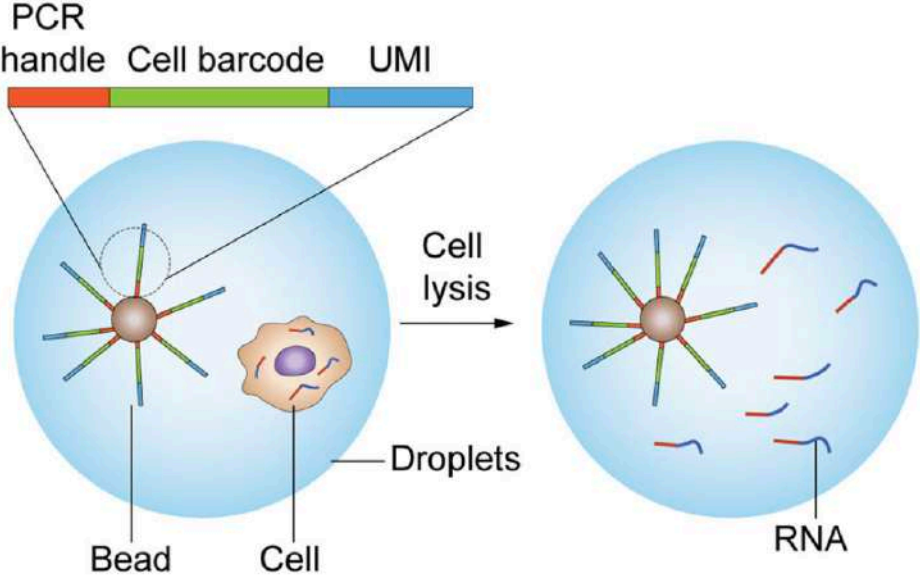




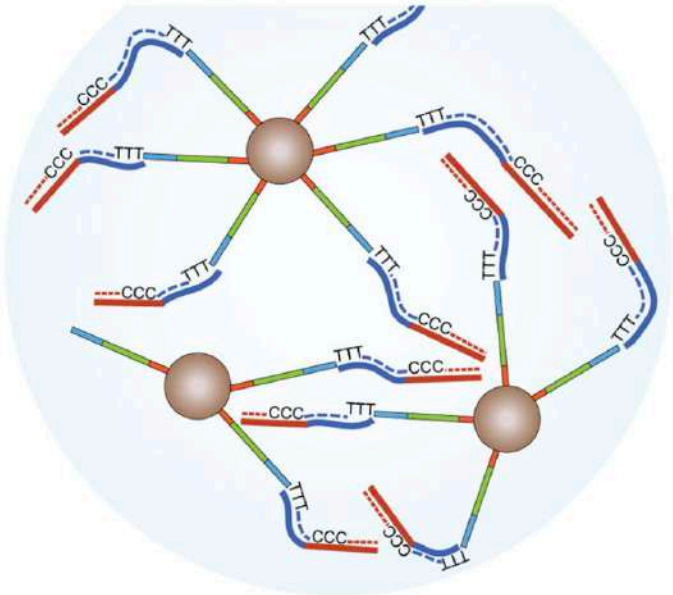
# Microfluidic isolation in reagent- filled droplets



## Structure of the barcode primer bead

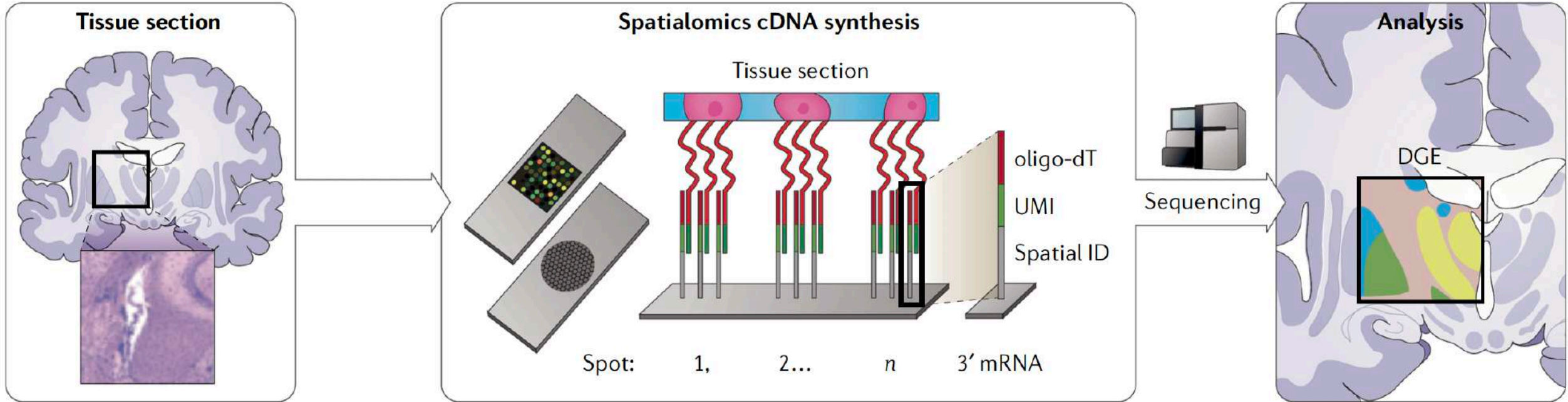


## Reverse transcription with template switching



# Spatialomics

1. Spatial encoding requires a frozen tissue section to be applied to oligo- arrayed microarray slides or to 'pucks' of densely packed oligo- coated beads.



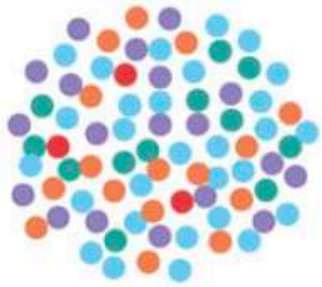
2. The mRNA diffuses to the slide surface and hybridizes to oligo- dT cDNA synthesis primers that encode UMIs and spatial barcodes. It is then reverse transcribed to produce cDNA, which is pooled for library preparation and sequencing.

3. Computational analysis of the spatialomics data maps sequence reads back to their spatial coordinates after DGE analysis and allows differential spatial expression to be visualized.

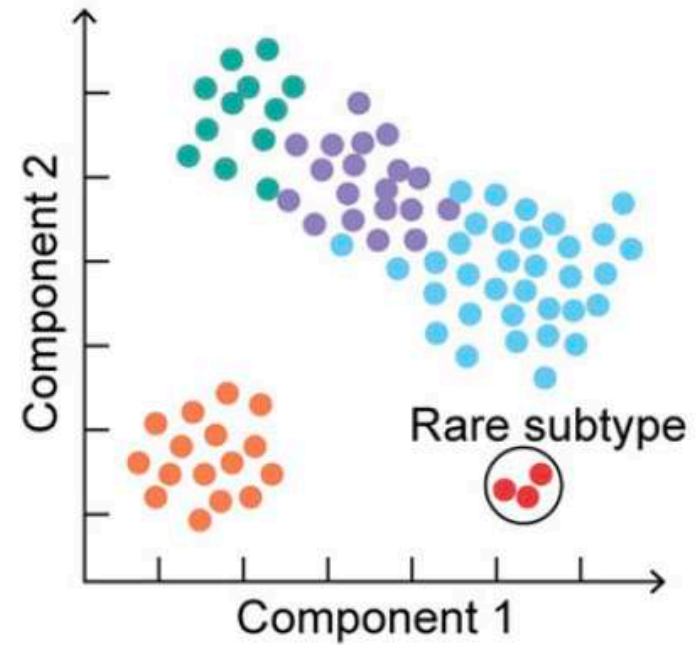
# Applications of scRNAseq computational approaches

## 1. Cell type identification

Heterogeneous tissue or tumor



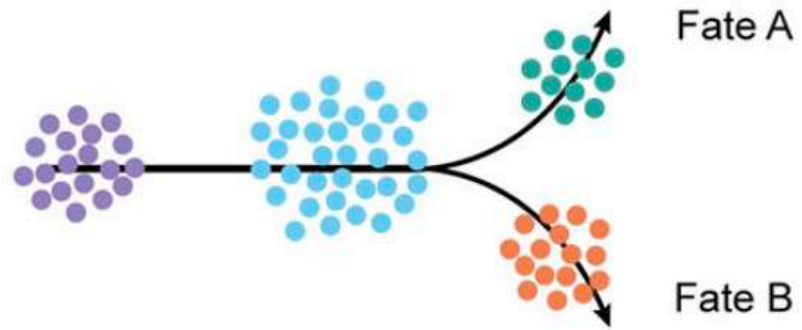
Dimensionality reduction  
(e.g. PCA)



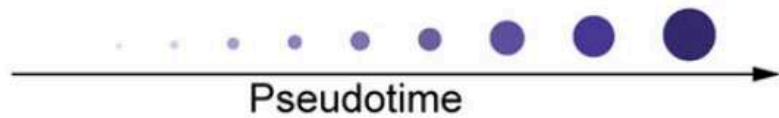
# Applications of scRNAseq computational approaches

## 2. Cell hierarchy reconstruction

Cell differentiation, or response to stimulus



Linear

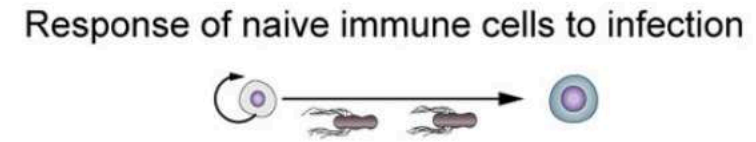
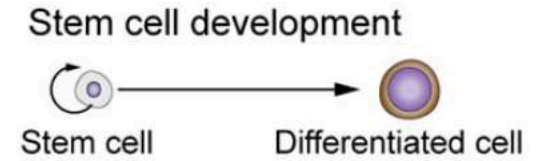
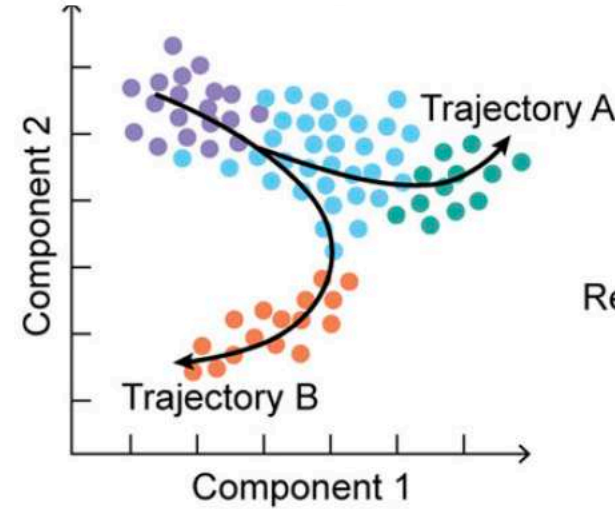
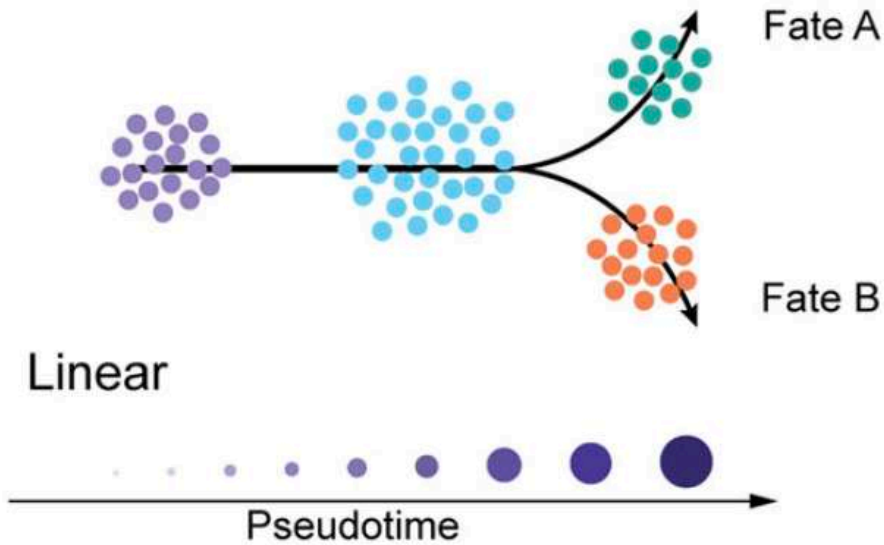


➔  
Trajectory analysis  
pipeline (Monocle)

# Applications of scRNAseq computational approaches

## 2. Cell hierarchy reconstruction

Cell differentiation, or response to stimulus

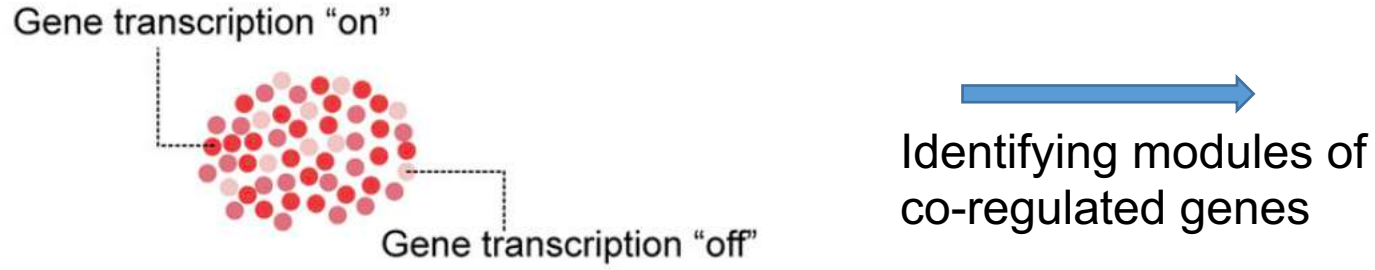


➔  
Trajectory analysis  
pipeline (Monocle)



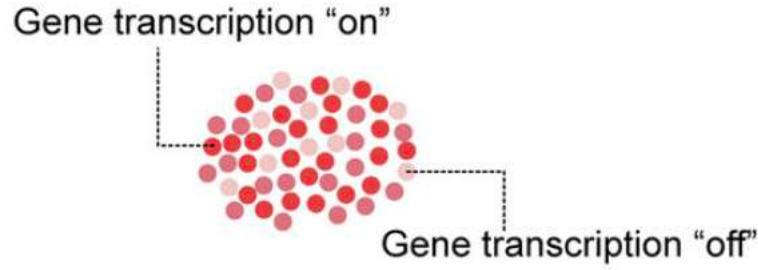
# Applications of scRNAseq computational approaches

## 3. Inferring regulatory networks

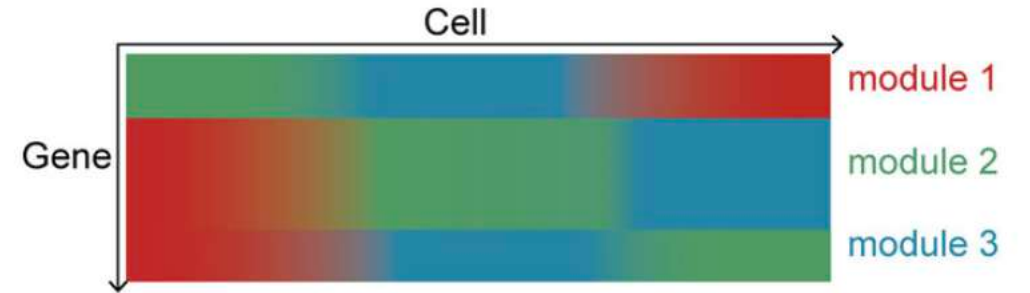


# Applications of scRNAseq computational approaches

## 3. Inferring regulatory networks

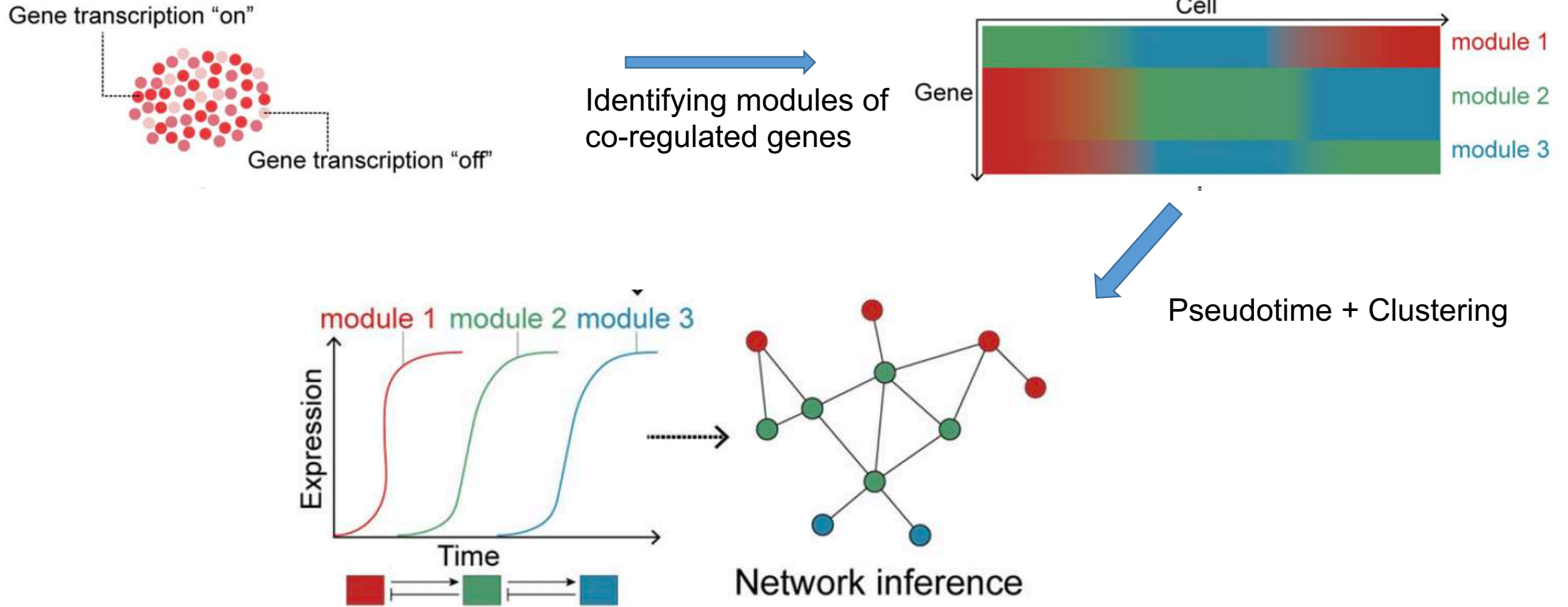


Identifying modules of co-regulated genes



# Applications of scRNAseq computational approaches

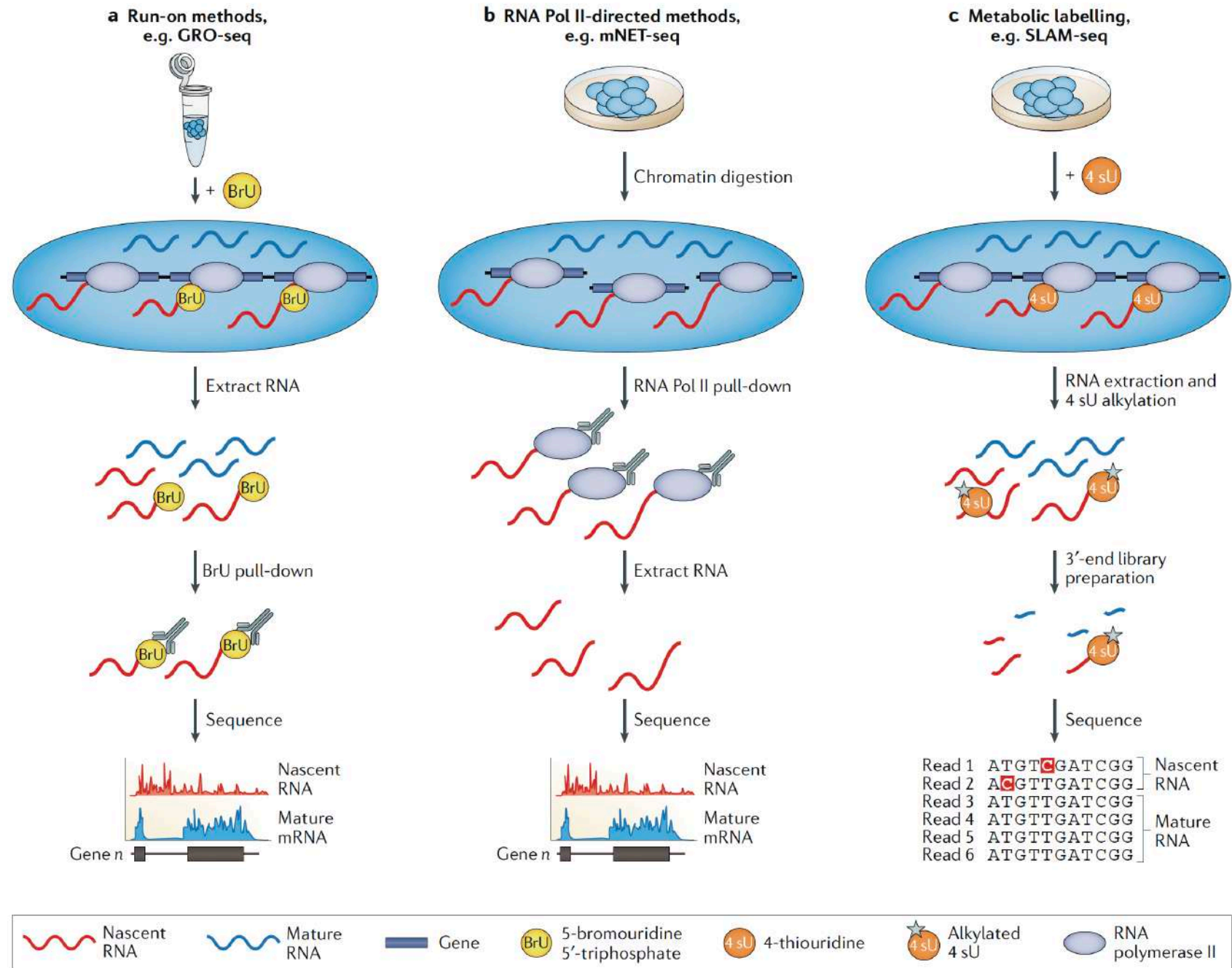
## 3. Inferring regulatory networks



Other approaches

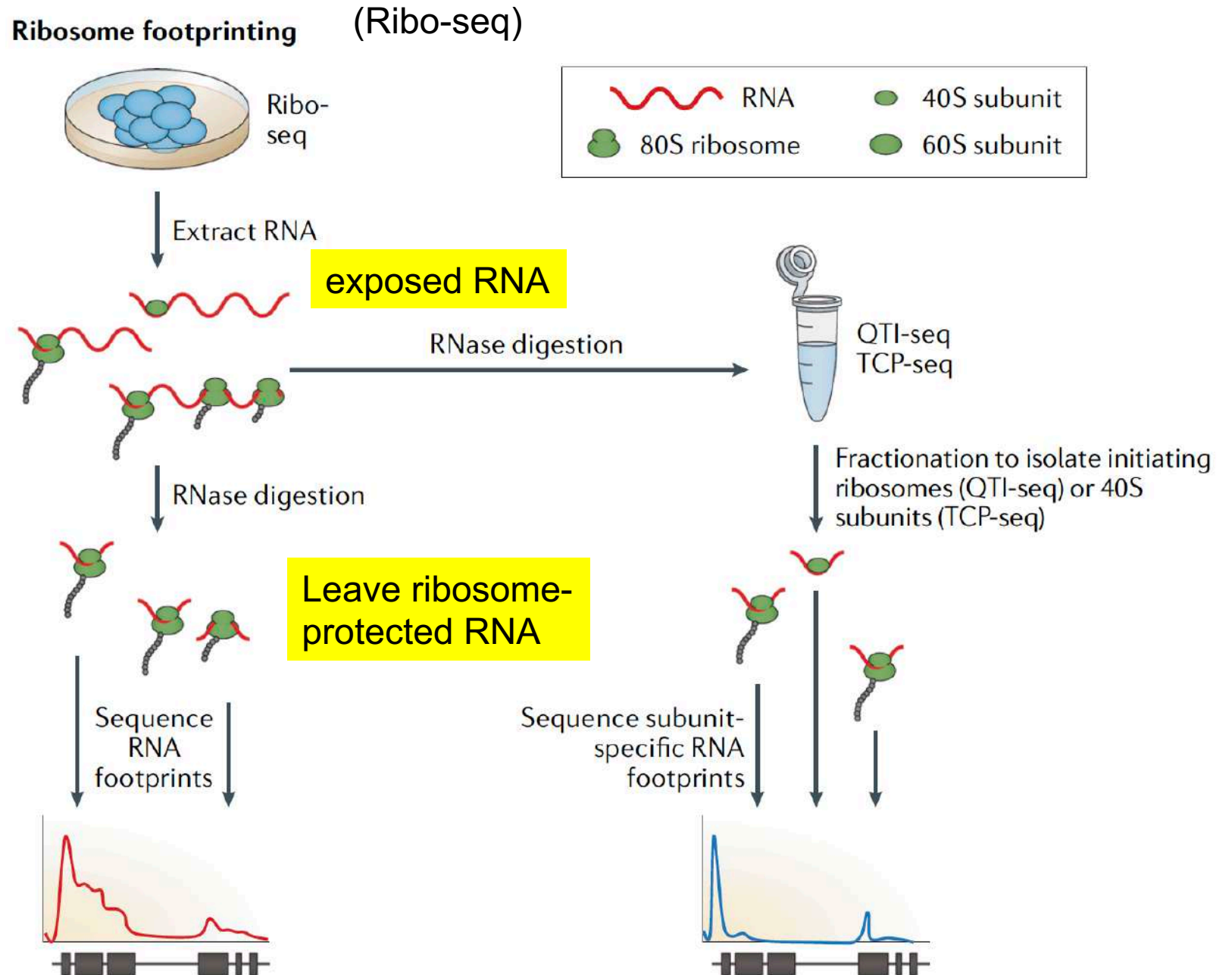
# nascent RNA

Essentially enrich newly transcribed RNAs in a cell and compare to control (mature RNA)





# translatome

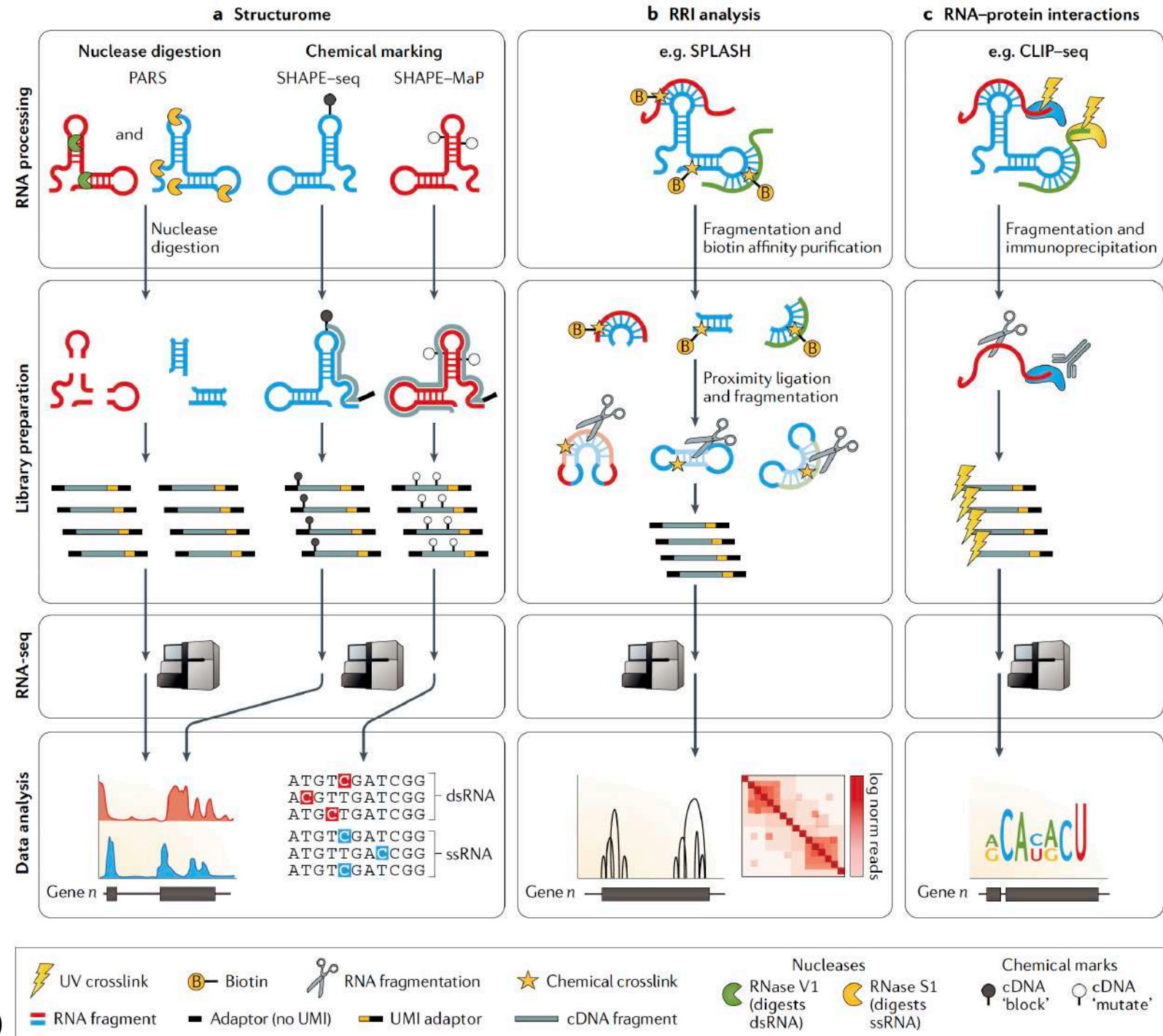


- RNA-sequencing from ribosomally bound RNA
- mRNA ribosome density correlates with the protein synthesis level

# RNA-RNA interaction

## RNA-protein interaction

- A) Probe structured (ddRNA) or unstructured (ssRNA) RNA in transcriptome level
- B) Crosslinking interacting RNA with biotinylated psoralen
- C) Crosslinking immunoprecipitation of RNA followed by sequencing



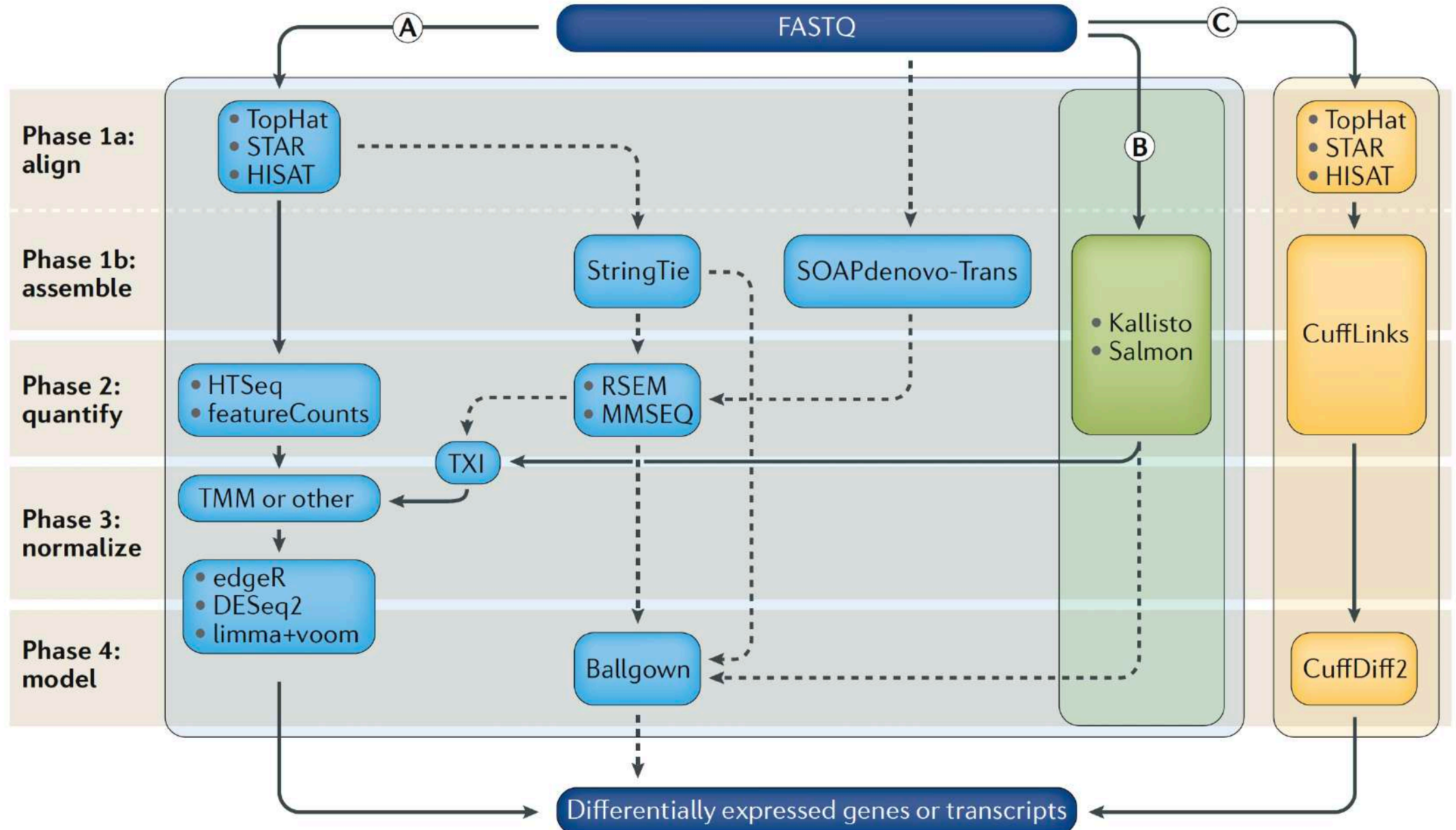
# Summary (I)



# Overview of technologies

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

# RNAseq analysis workflow for differential expression (generalized)



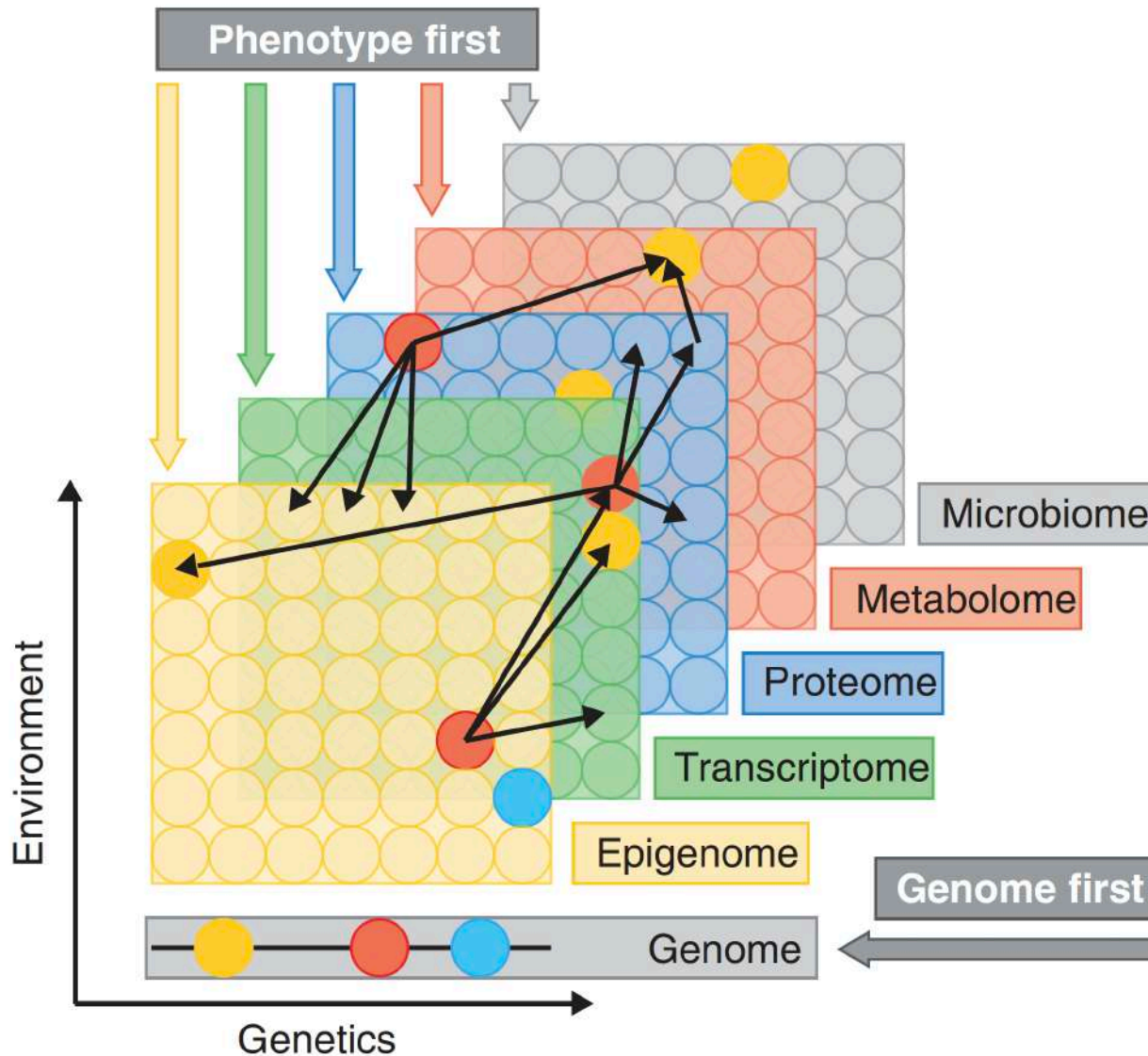


## Summary / Our experiences

- Experimental design is key to correctly address your biological question
- Always use replicates (at least 5)
- Avoid *de novo* transcriptome assembly if you can
  
- DEseq2 are easy to use and have been standardised
- Cuffdiff2 are theoretically better but for some reasons are worse (since we used mostly 2-3 replicates)
  
- Still many challenges ahead
  
- **Question: What will be integrated/obsoleted within 5? years with the arrival of long read sequencing**

Which leads us to... a multi-omic perspective

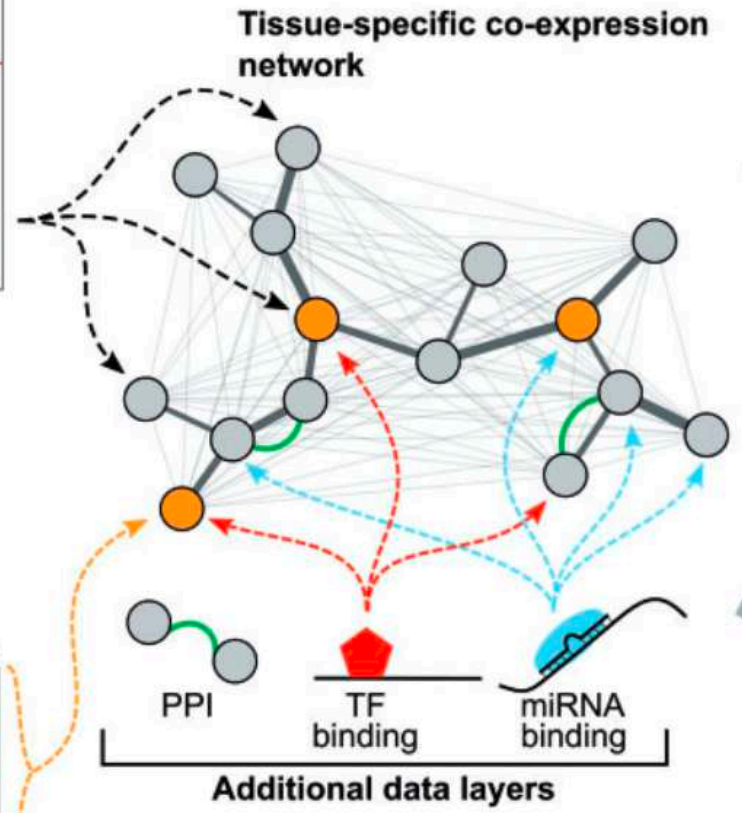
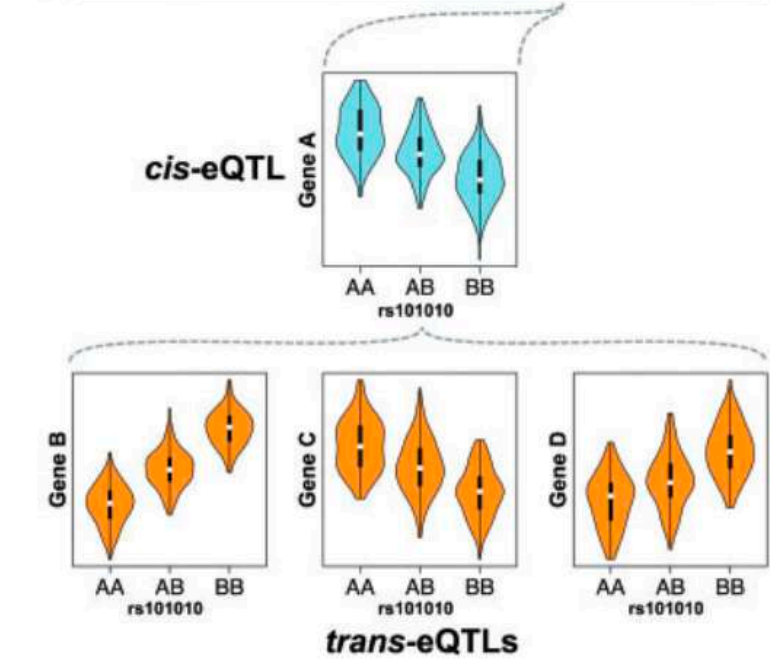
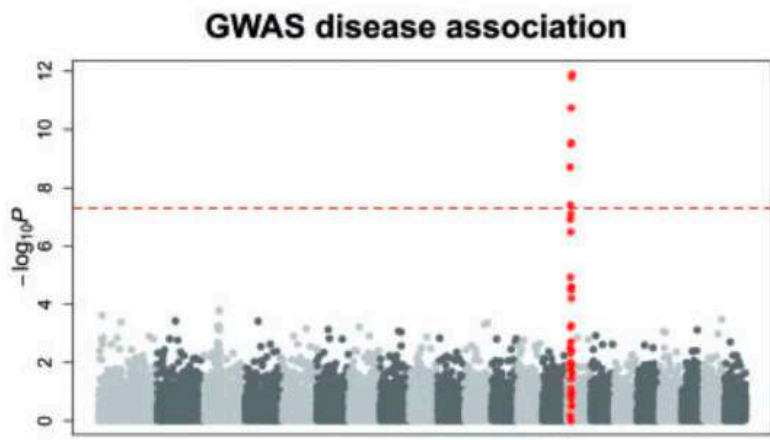
# Multiple omics data types



- Genome first or Phenotype first or environment first?
- Genome first -> GWAS
- “Locus-centered integration of additional omics layers can help to identify causal single nucleotide polymorphisms(SNPs) and genes at GWAS loci and then to examine how these perturb pathways leading to disease”

# Integrating multi-omics to network

Various additional data can then be used to enrich and extract biological relevant information from the network

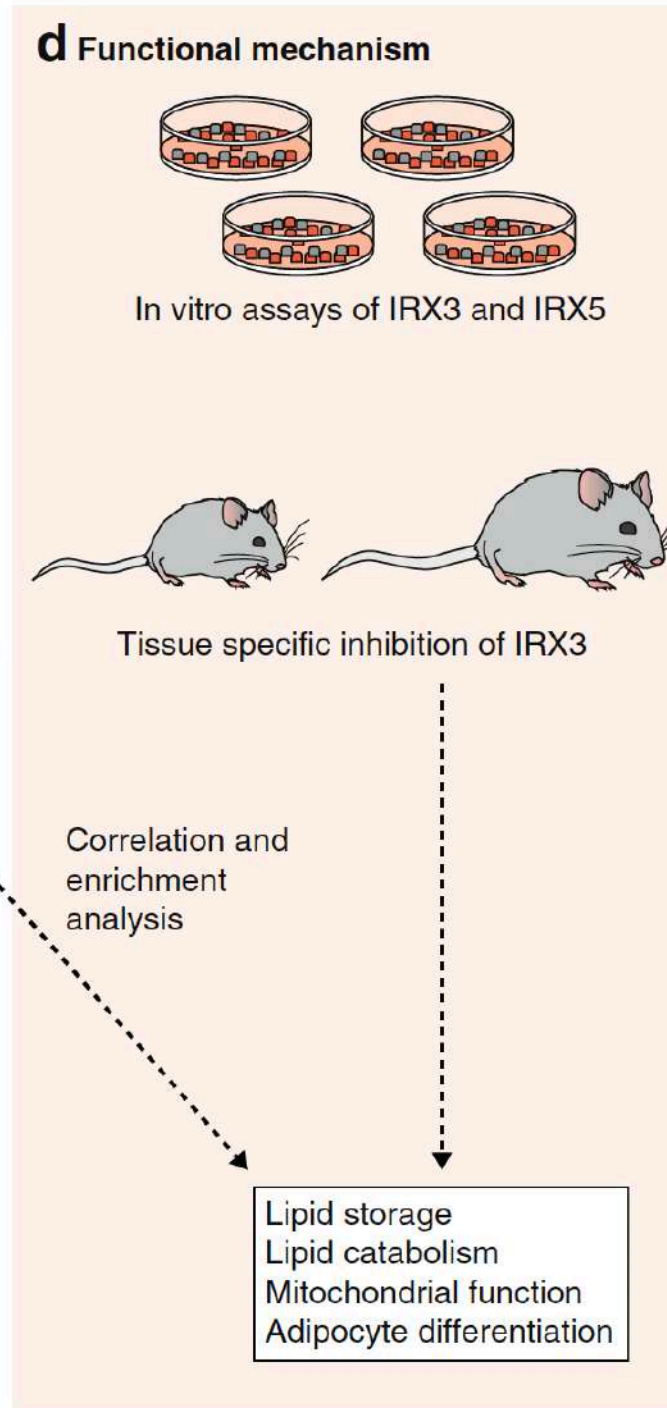
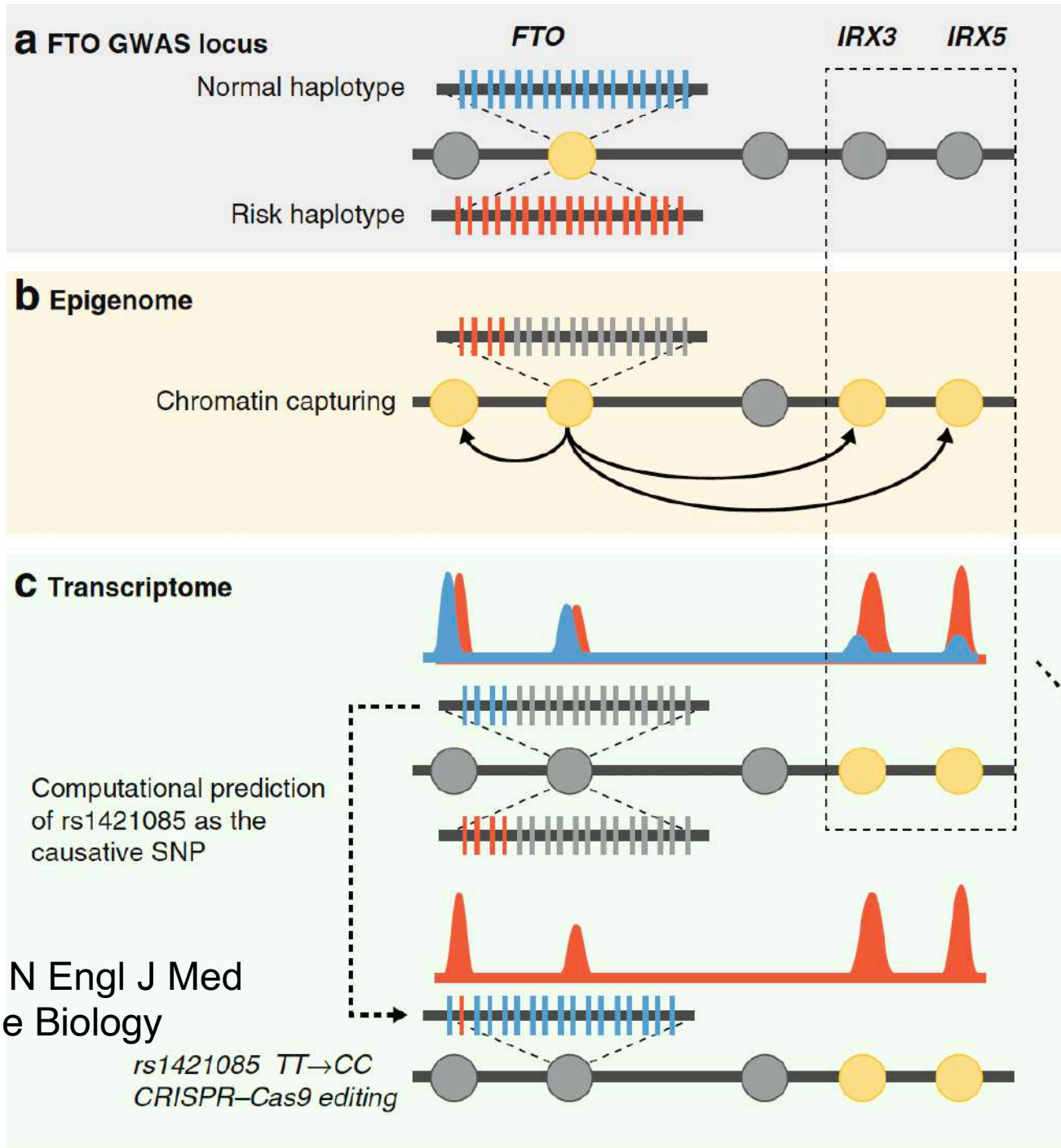


- Network analyses**
- Identification of modules
  - Identification of hub genes
  - Guilt-by-association predictions
  - Regulatory network construction

- Module interpretation**
- TFBS enrichment
  - miRNA binding site enrichment
  - GWAS enrichment
  - eQTL enrichment
  - OMIM enrichment
  - Phenotype Ontology enrichment
  - Gene Ontology enrichment
  - Pathway enrichment
  - ...

- Supporting data**
- Differential expression
  - Differential co-expression
  - Differential methylation
  - ...

# Example: FTO GWAS locus



Claussnitzer *et al* (2015) *N Engl J Med*  
Hain *et al* (2017) *Genome Biology*

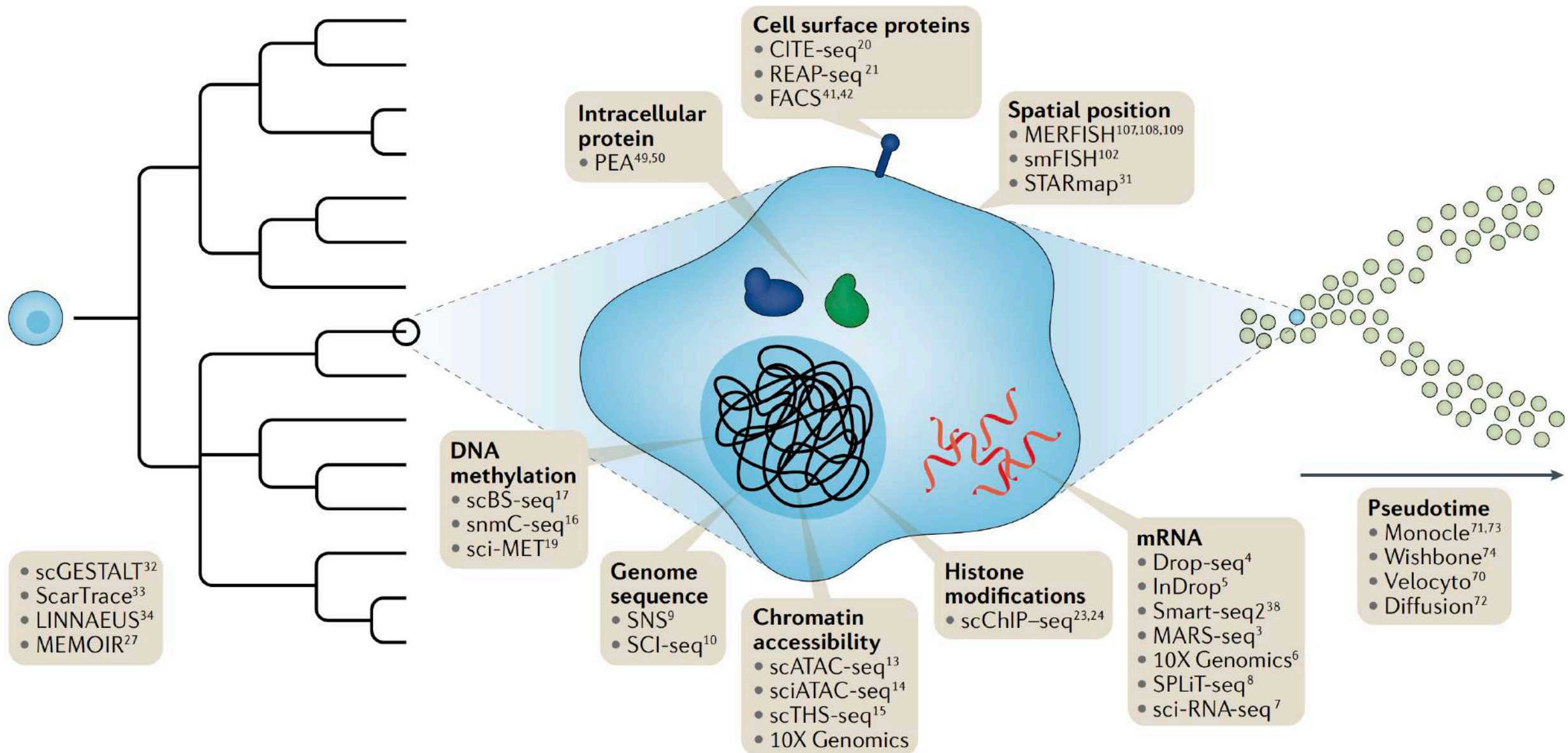


# Overview of current methods for single cell data integration

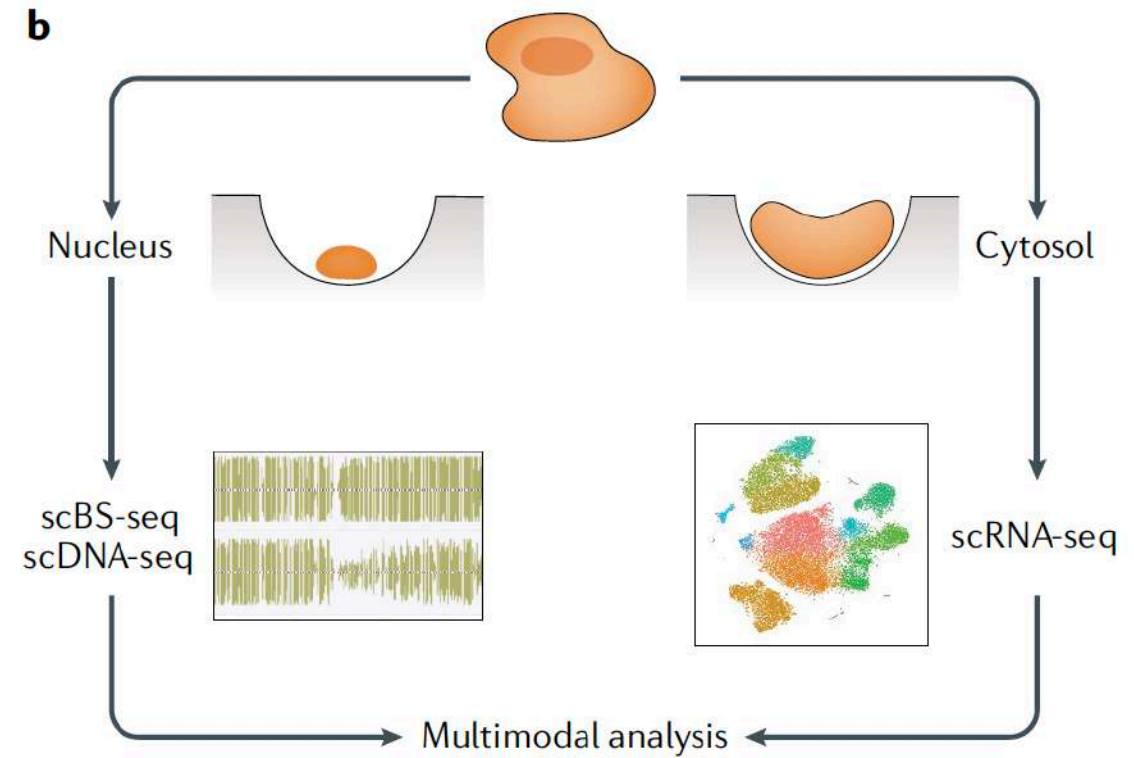
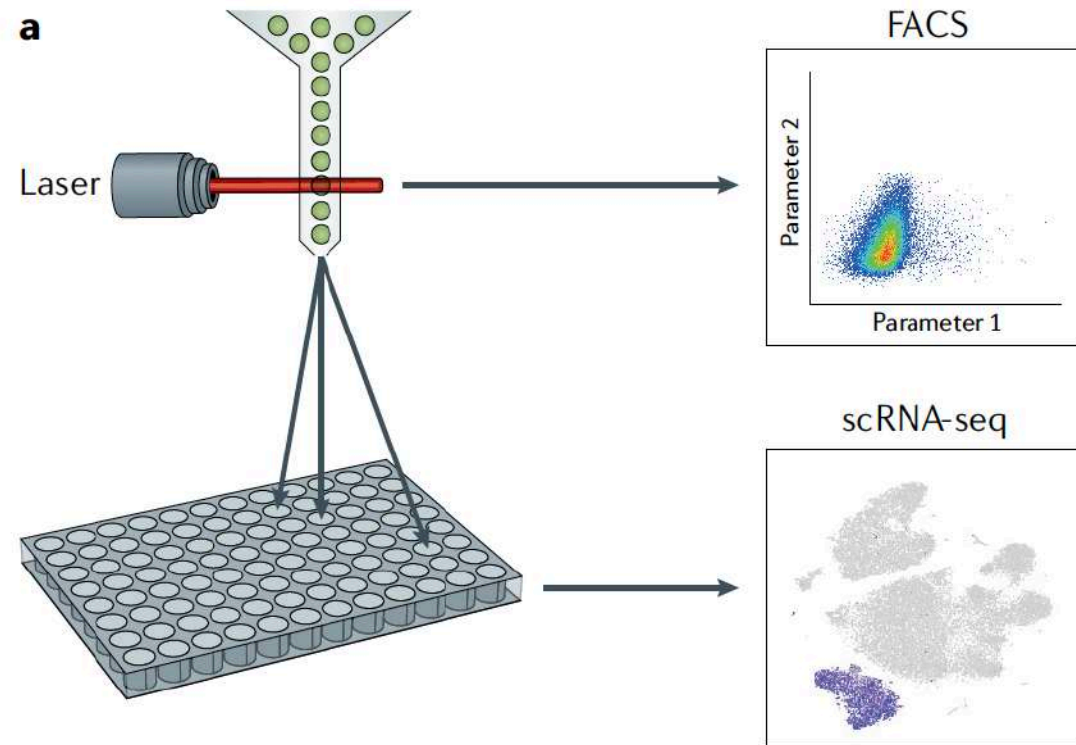
Lineage

State

Trajectory

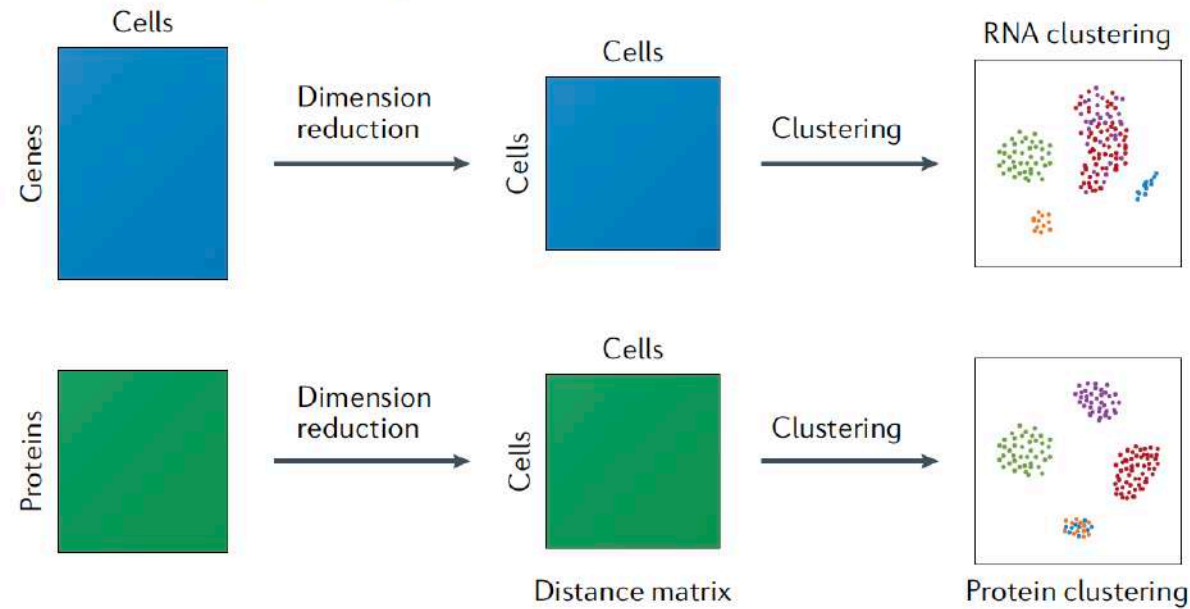


# Example of experimental methods for performing single-cell multimodal measurements

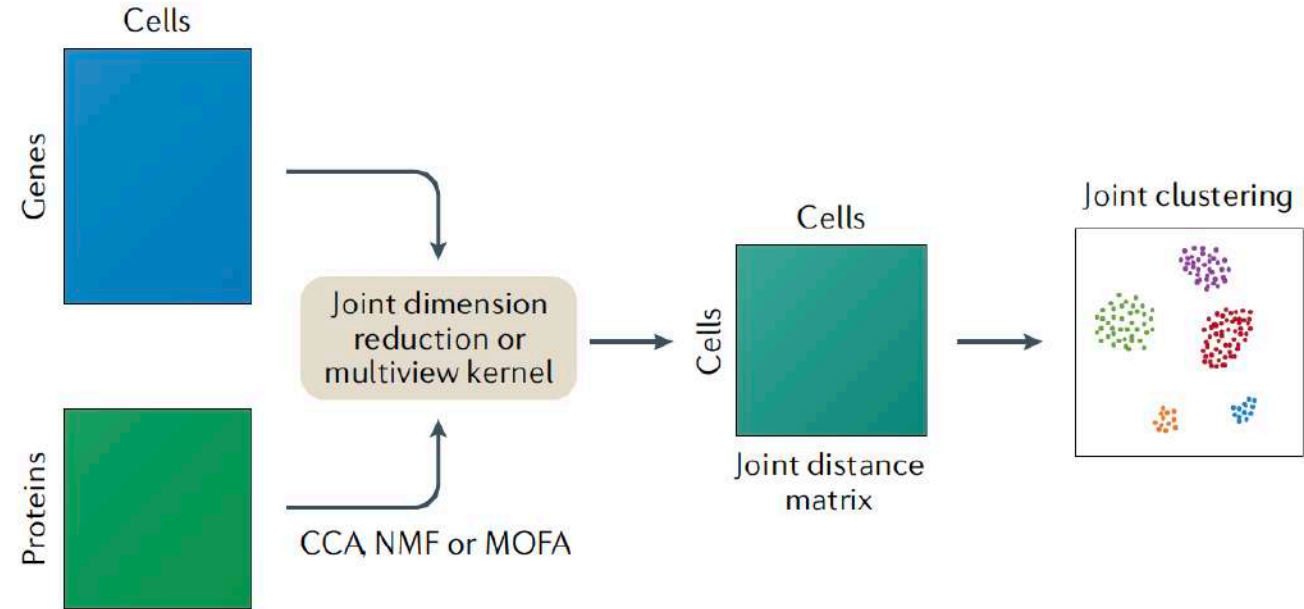


# Multi-modal data can lead to better power at identifying cell states

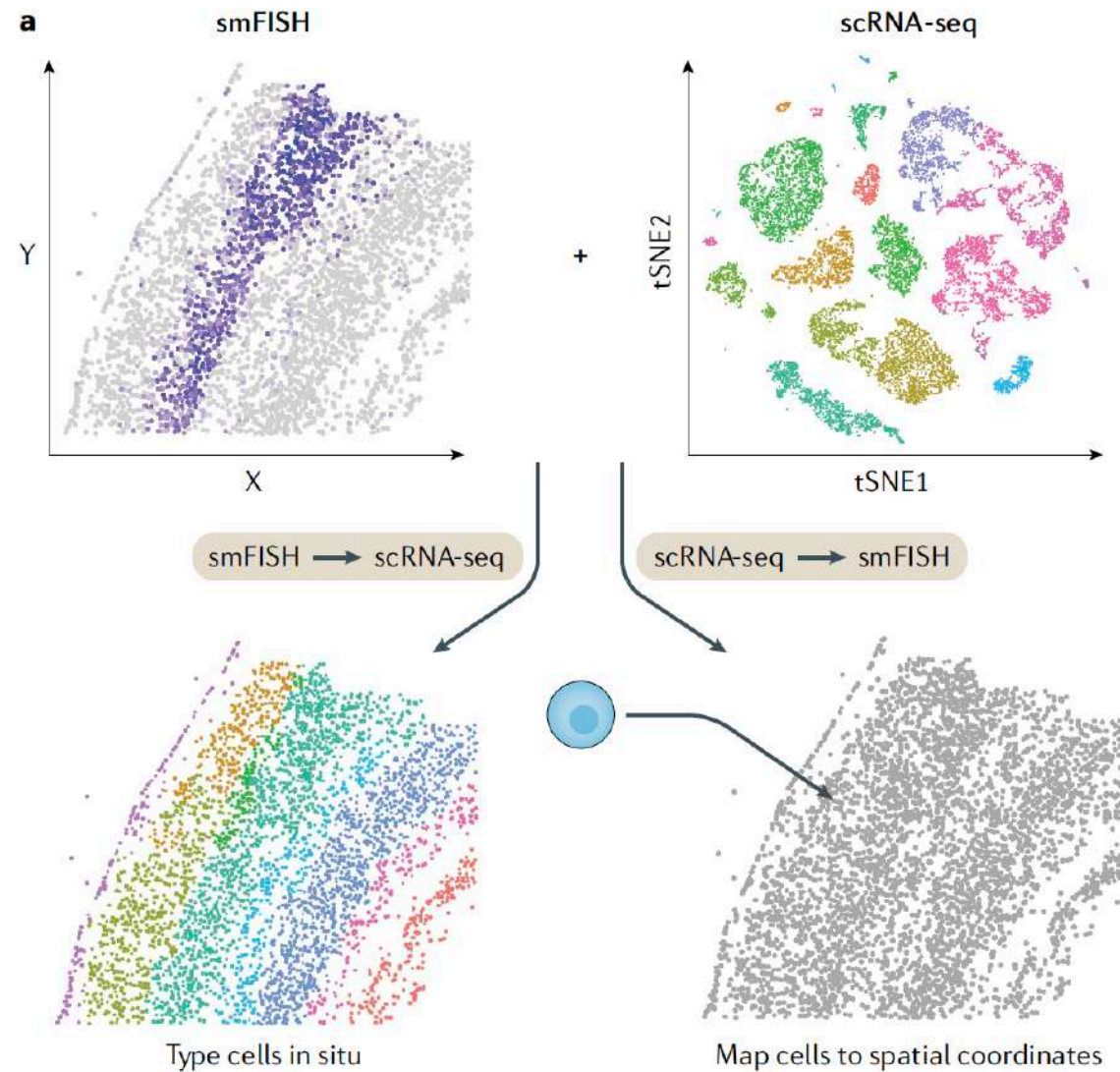
## a Separate analysis of multiple modalities



## b Joint analysis of multiple modalities

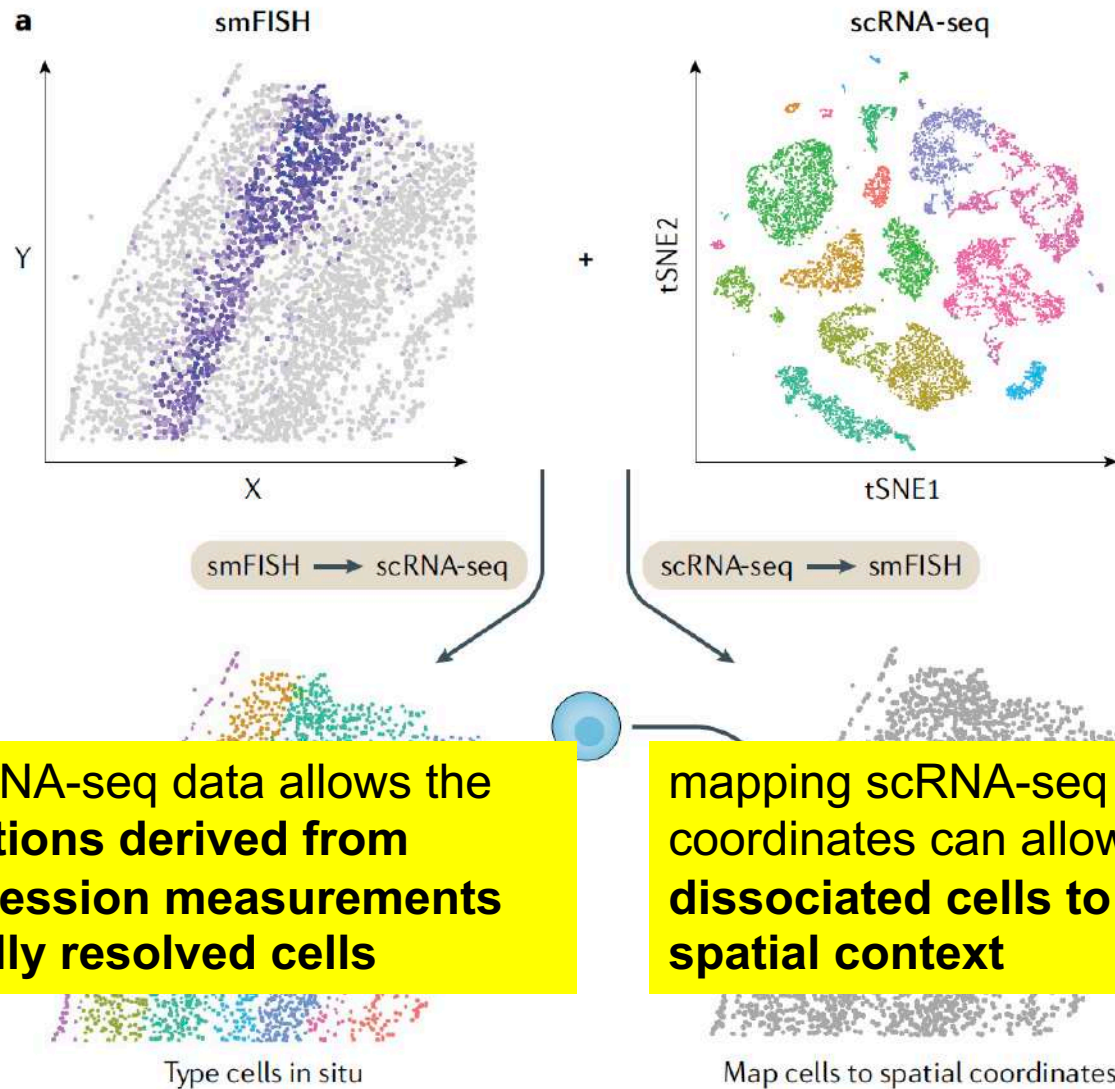


# Spatial omics + scRNA-seq





# Spatial omics + scRNA-seq



Mapping smFISH cells onto scRNA-seq data allows the **transfer of cell-type classifications derived from transcriptome-wide gene expression measurements to be transferred to the spatially resolved cells**

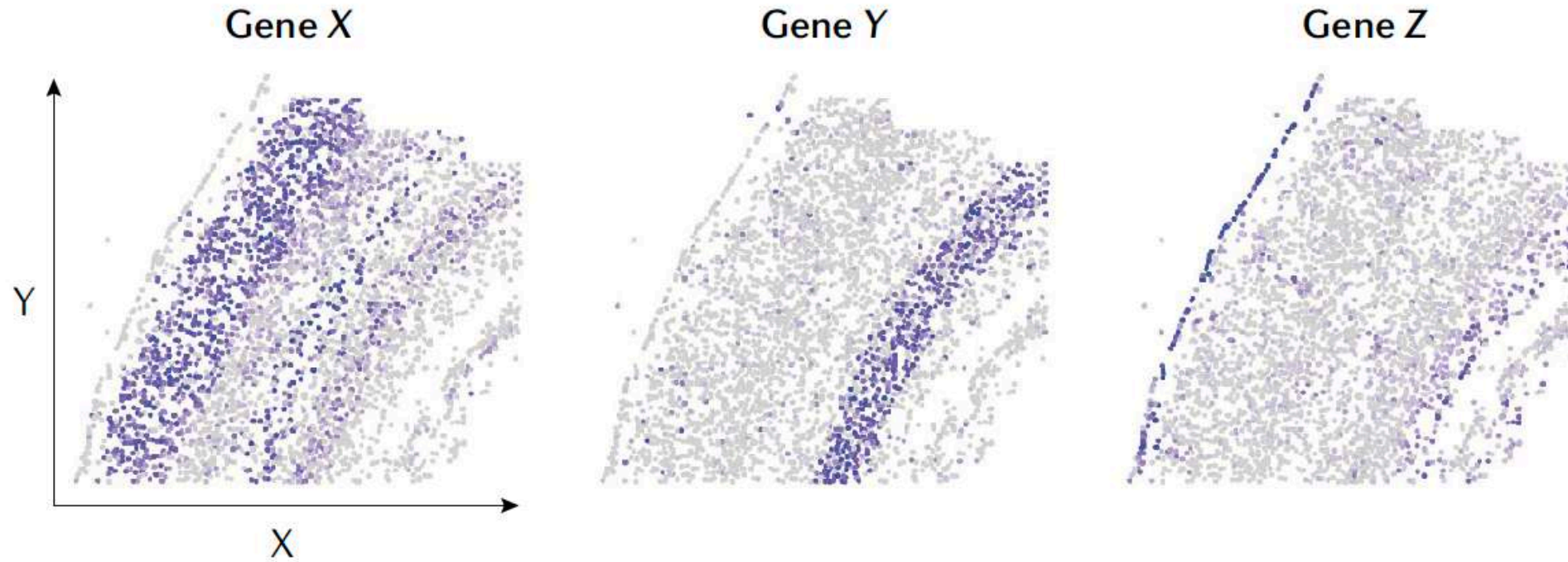
mapping scRNA-seq data onto smFISH-profiled spatial coordinates can allow **scRNA-seq data from dissociated cells to be placed back into their spatial context**



# Spatial omics + scRNA-seq



## c Novel spatial patterns of gene expression



By mapping scRNA-seq-profiled cells onto spatially resolved coordinates through the integration with smFISH data, **spatial patterns of gene expression can be predicted for any gene measured in the scRNA-seq data set**. Through these predictions, novel spatial patterns of gene expression may be identified through the analysis of genes that were not profiled by smFISH

# Summary (II) and Conclusion

## **Potential**

- Single cell/nucleus RNAseq + spatial information + long read sequencing + direct RNA sequencing?!
- It is an exciting time to be in

## **Challenges**

- Data type gets extremely complicated
- Integrating different sources of data are powerful

# References

<https://www.notion.so/References-papers-links-in-start-learning-genomics-b7e57b28e9194bb29a02f483e0b894ad>