

Introduction to
RNA-Seq applications and tools
26-27th September, 2018

Organised and delivered by Bioinformatics Core at WHG:

Helen Lockstone

Ben Wright

Eshita Sharma

Santiago Revale



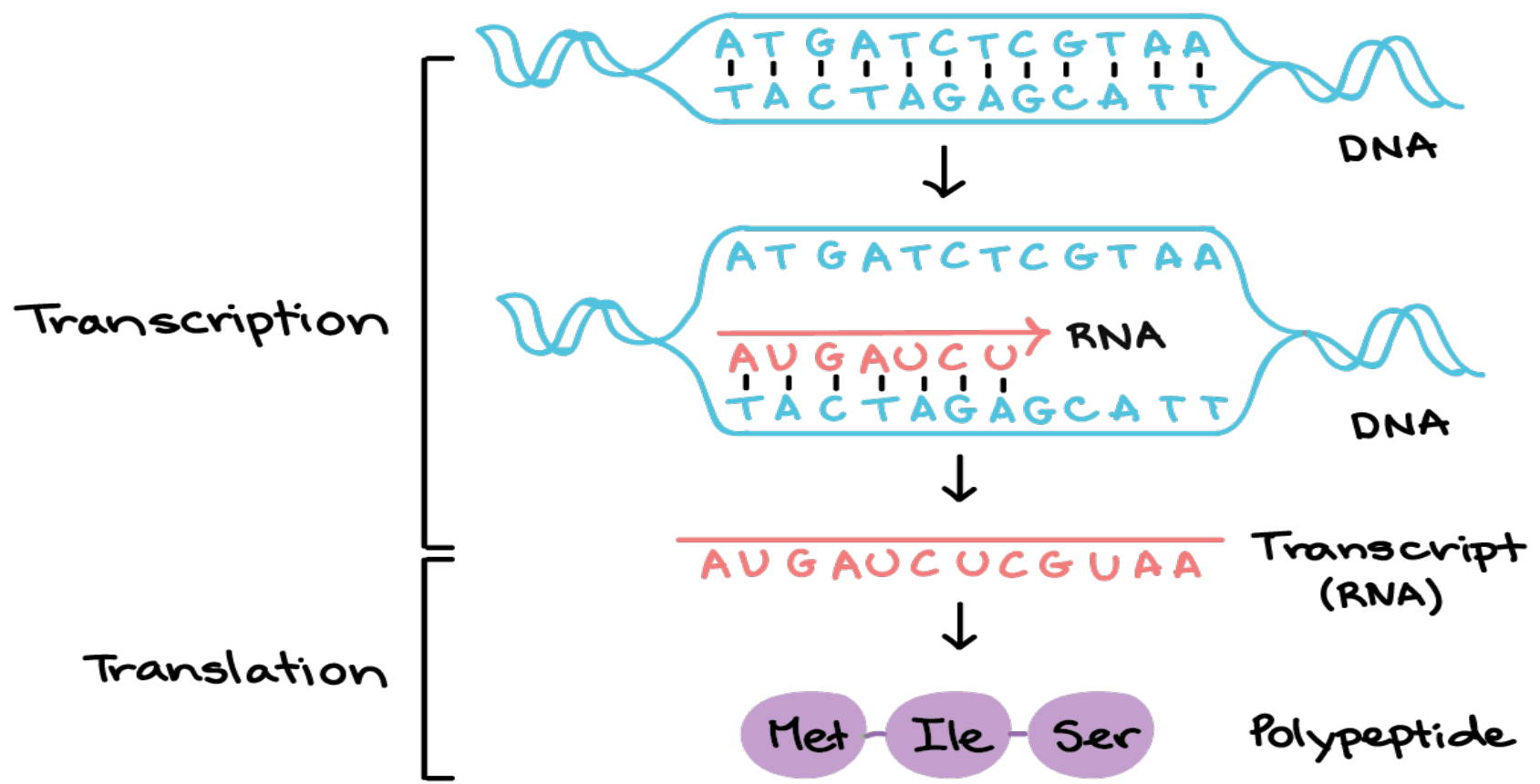
- **What it is?**
- **Scope of RNAseq**
- **Usual approaches for RNAseq library preparation?**
- **Considerations for RNAseq experiments**
- **General methods for RNAseq data analysis.**

Eshita Sharma,
eshita.sharma@well.ox.ac.uk

Research Associate in Functional Genomics, Bioinformatics Core

What it is?

RNA - Mid-point of the information cascade

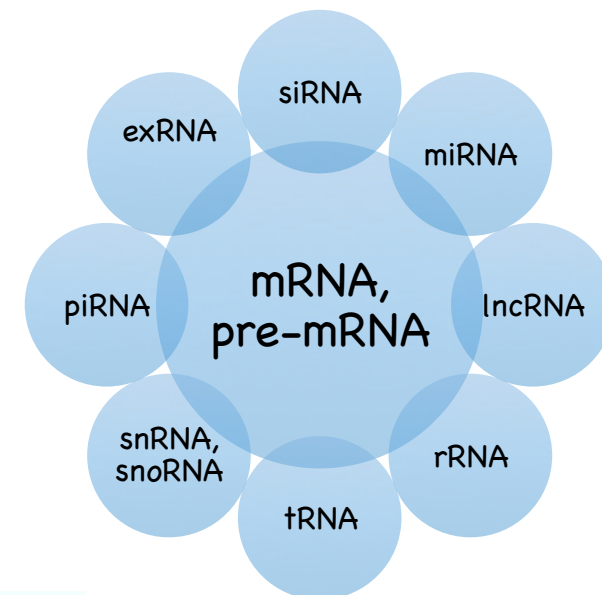
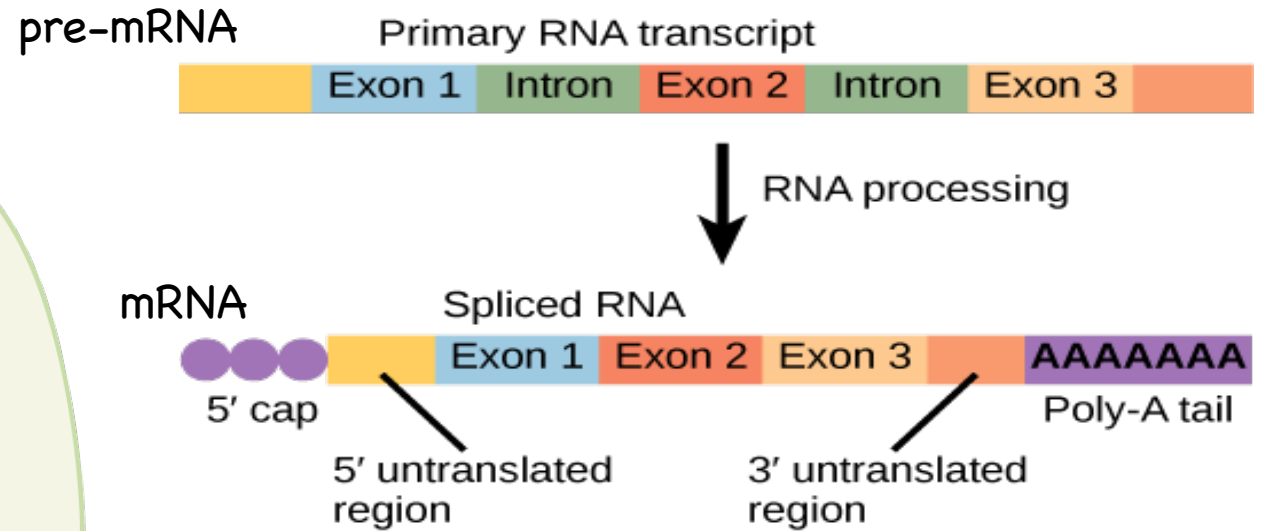
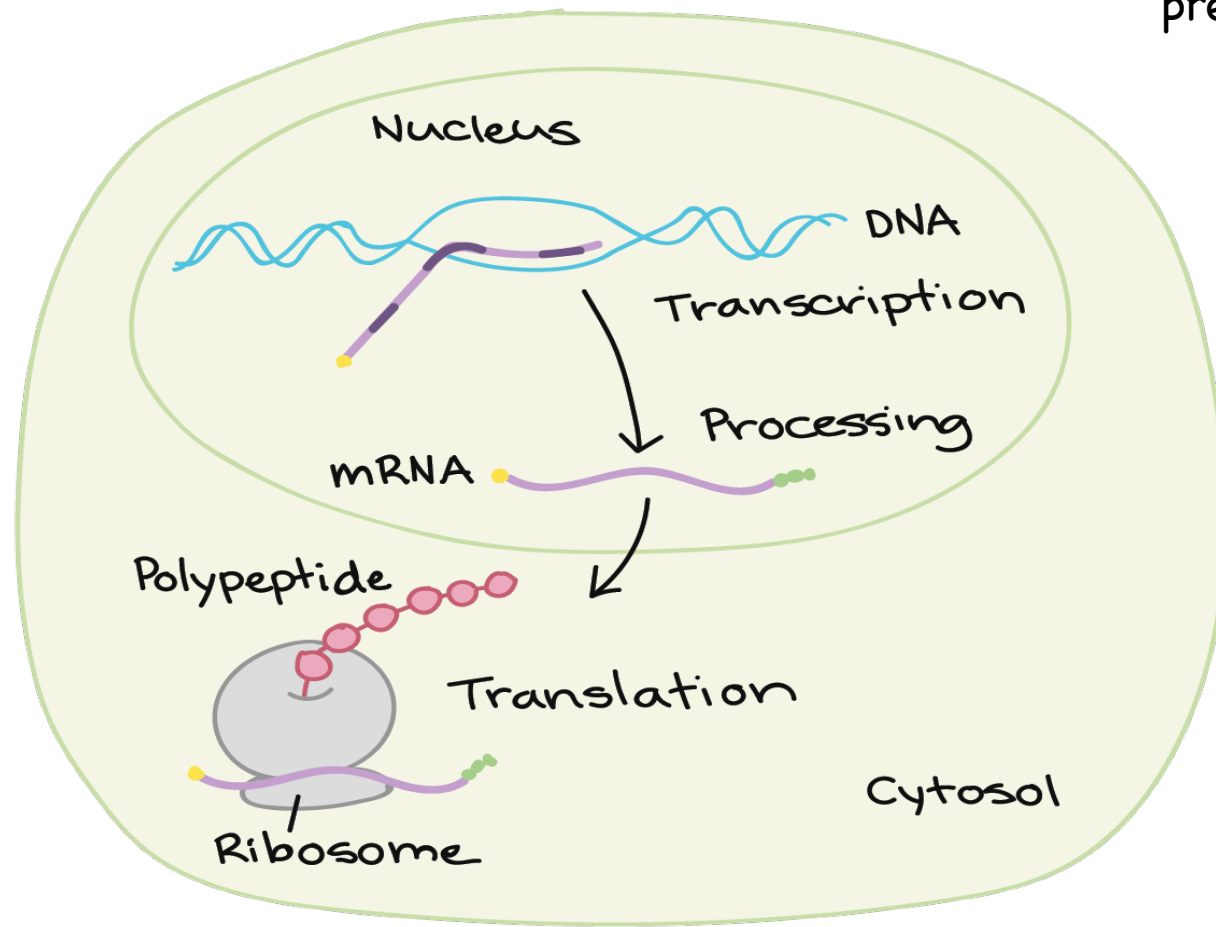


- Image credit Khan Academy Open Courses

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



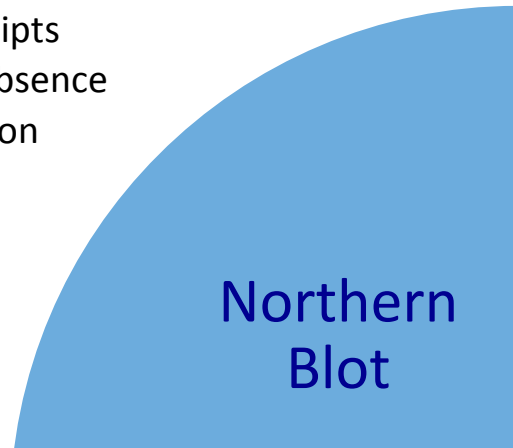
- Image credit Khan Academy Open Courses

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

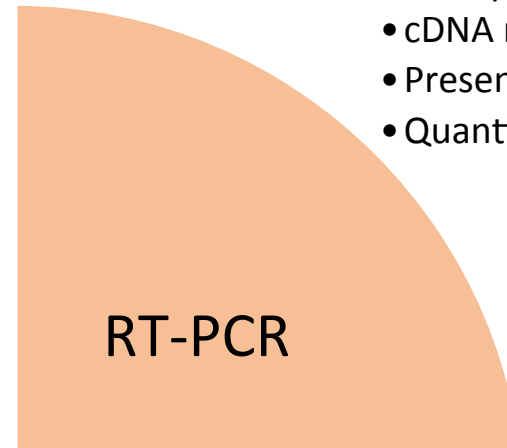
RNAseq is a method for Transcriptome profiling

Image of the transcribed genome at any point of time!

- Single genes
- RNA transcripts
- Presence/Absence
- Quantification
- Length

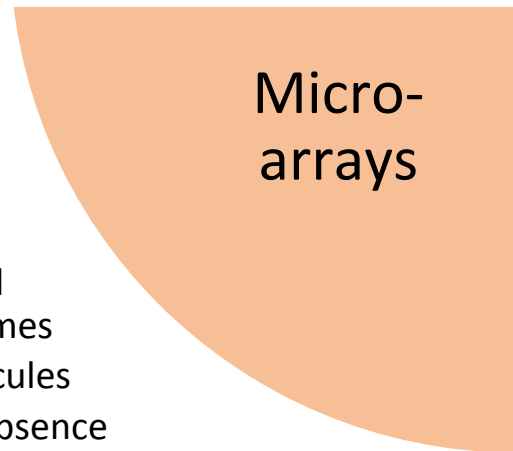


- Multiple genes
- cDNA molecules
- Presence/Absence
- Quantification

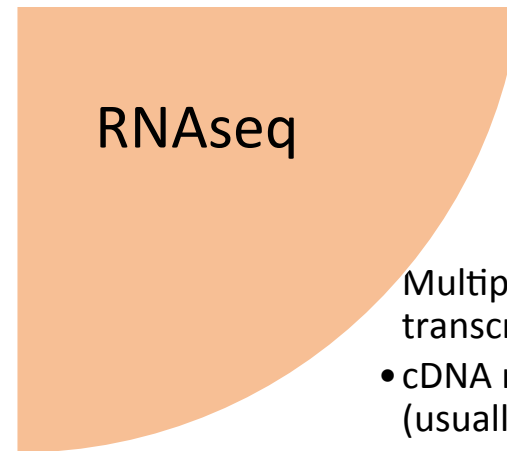


How do we take this image?

- Multiple full transcriptomes
- cDNA molecules
- Presence/Absence
- Quantification



- Multiple full transcriptomes
- cDNA molecules (usually)
- Presence/Absence
- Quantification
- Length/Splicing
- Sequence



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:

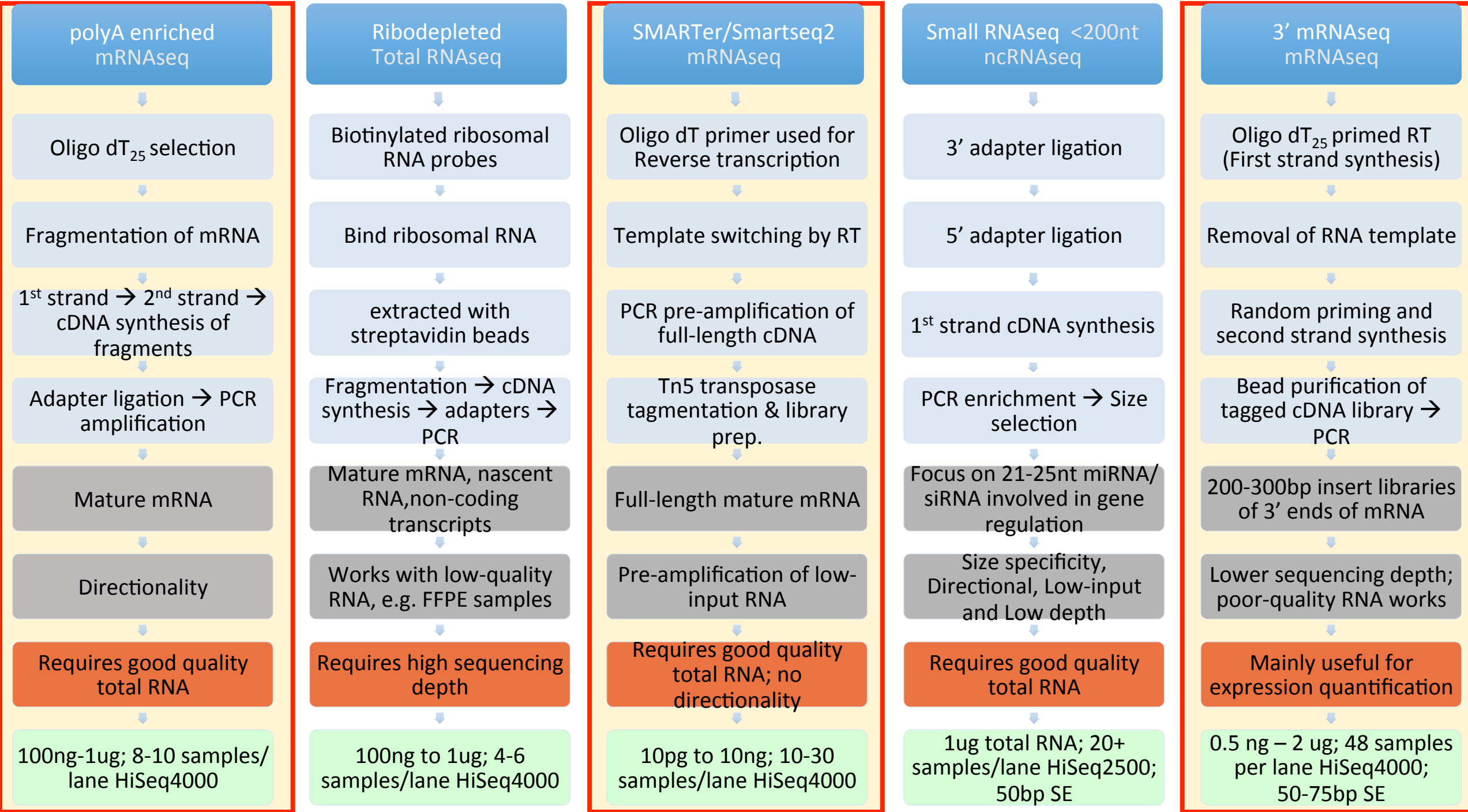
- ✓ total RNA,
- ✓ mRNA,
- ✓ small RNA (miRNA, tRNA, ribosomal profiling, etc.)

Can be used to:

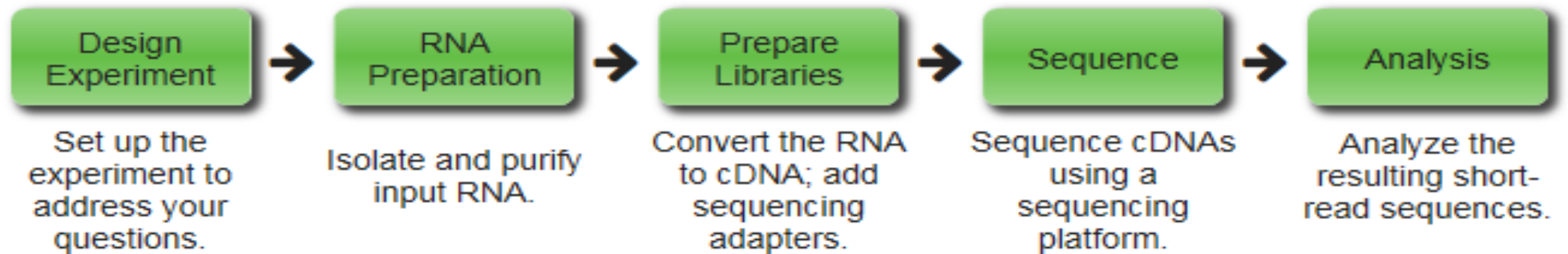
- ✓ determine exon/intron boundaries,
- ✓ verify or amend previously annotated 5' and 3' gene boundaries.

- 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



Typical RNAseq experiment

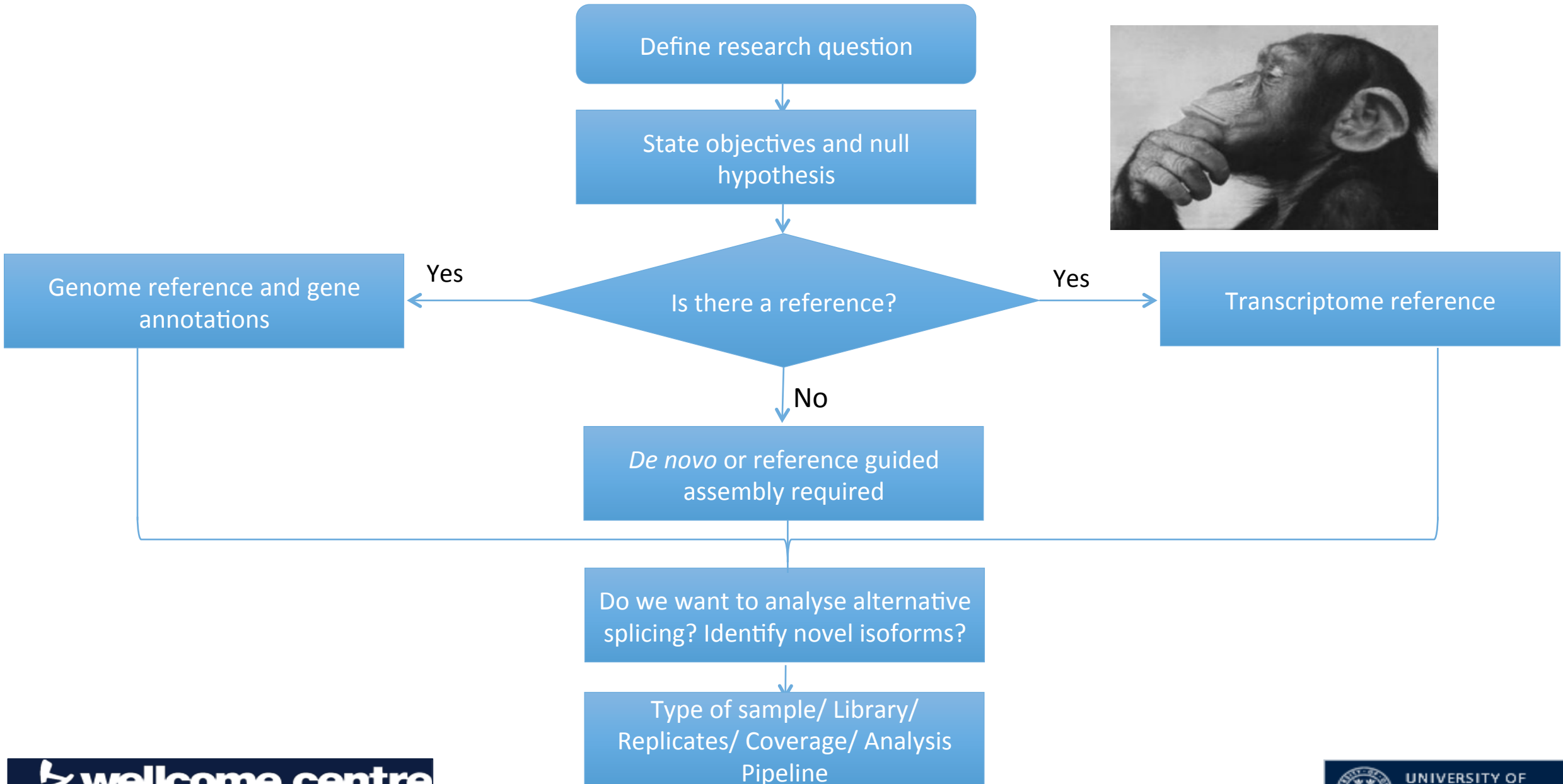


Key considerations for Experimental Design

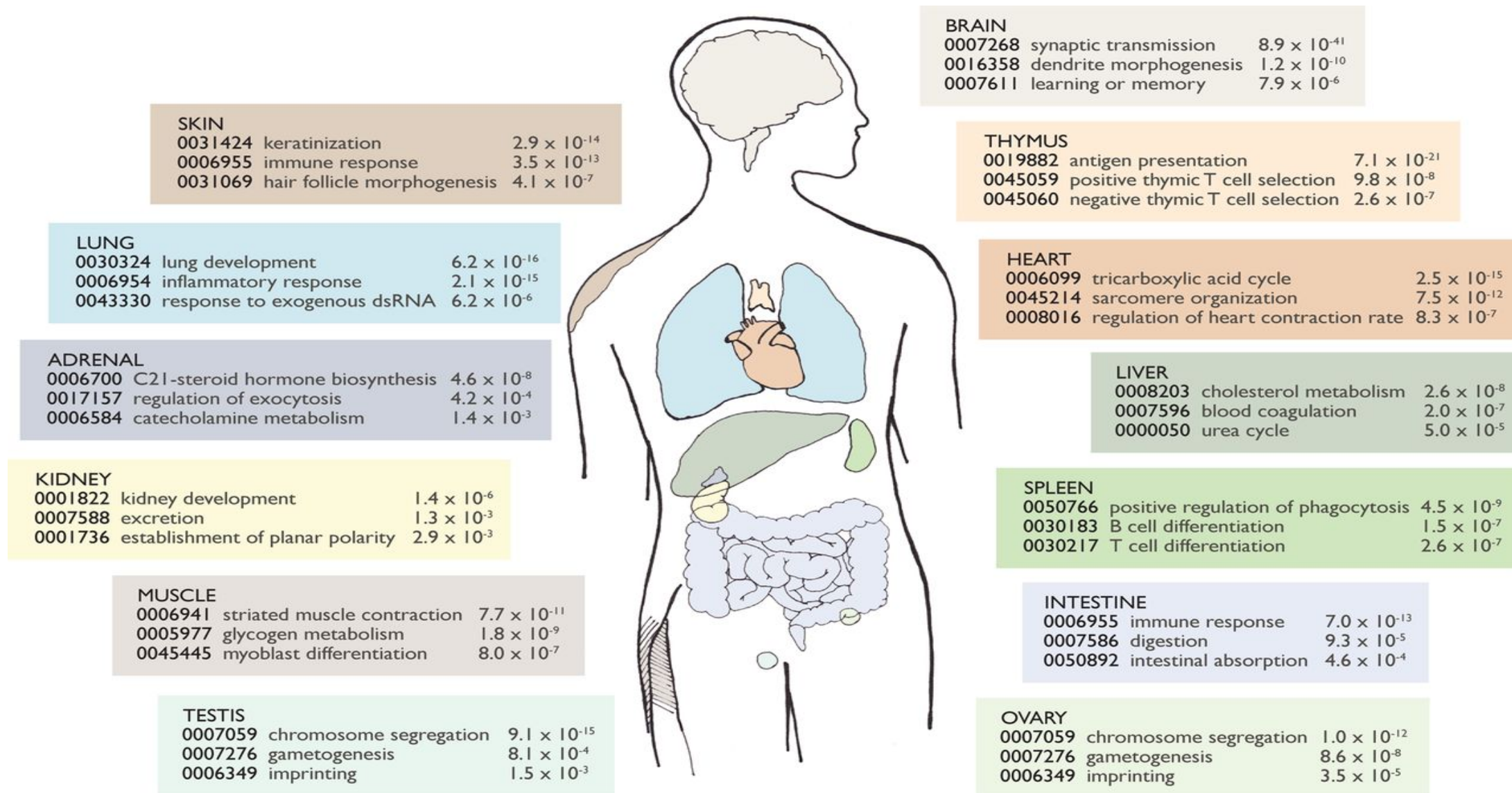
How is an experiment designed?

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

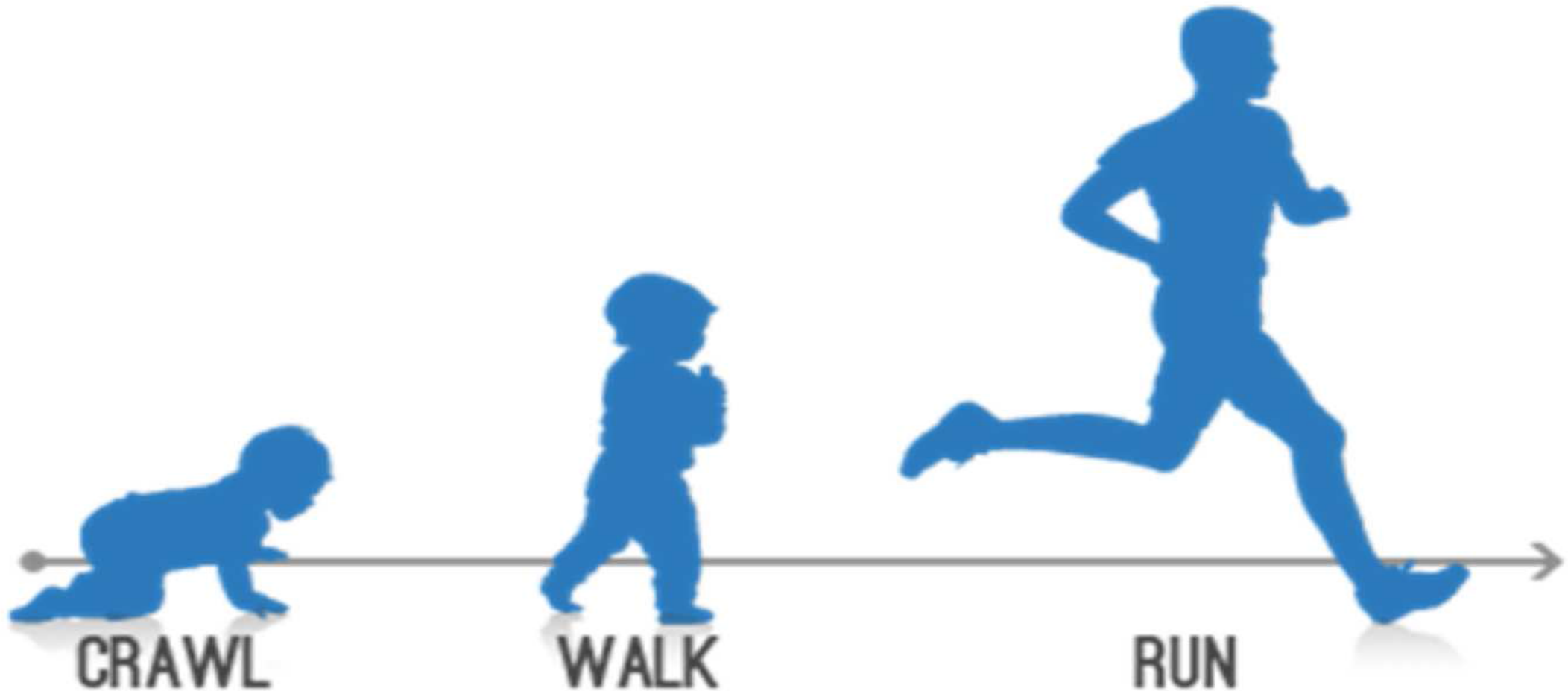
Before we begin...



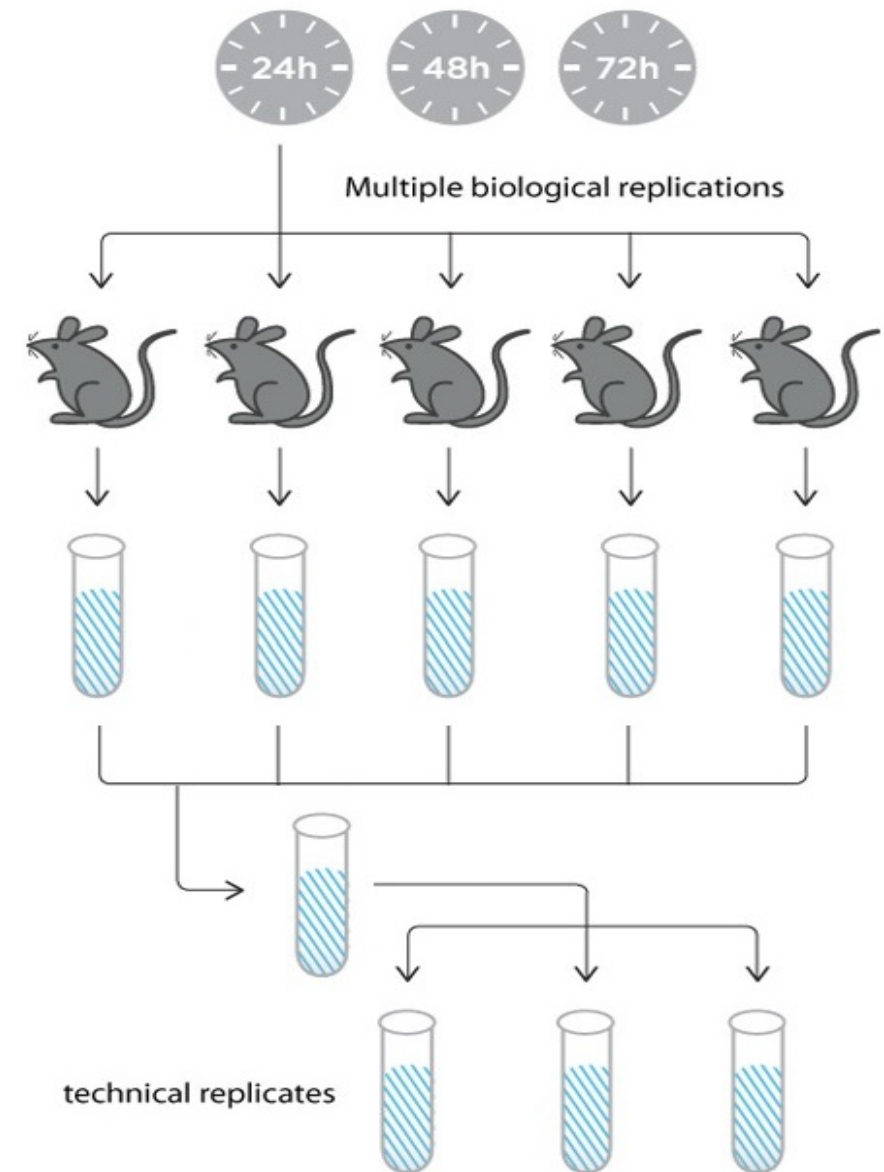
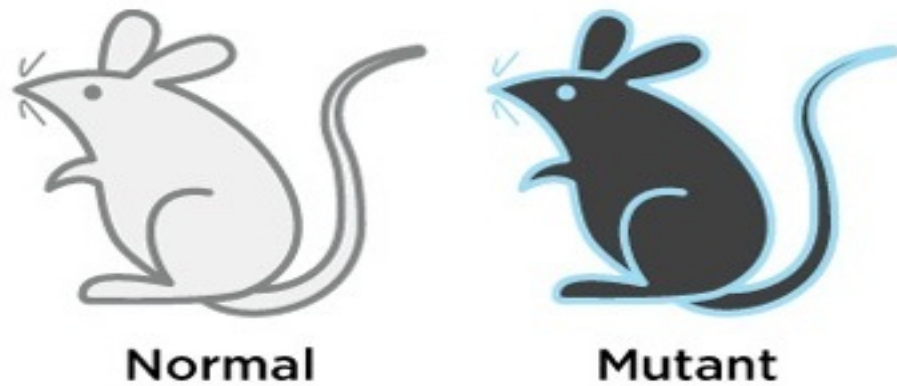
Sample selection



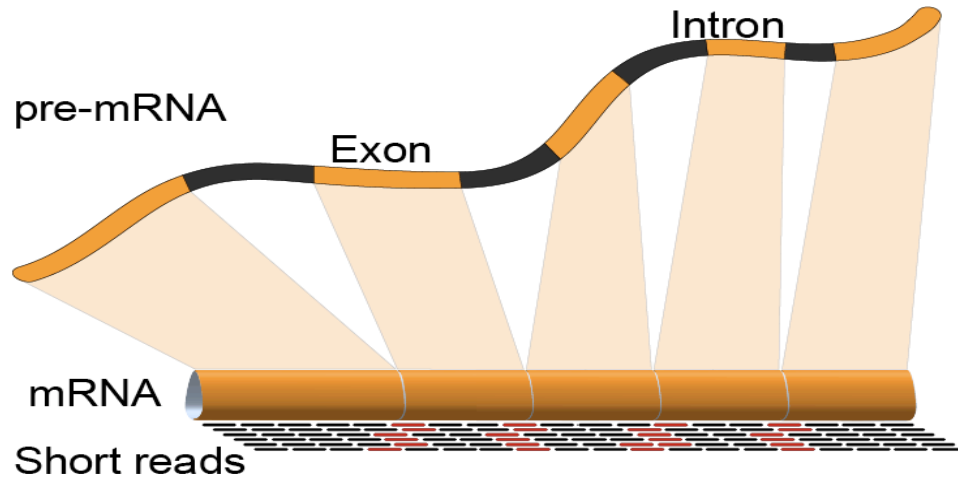
Time of sampling



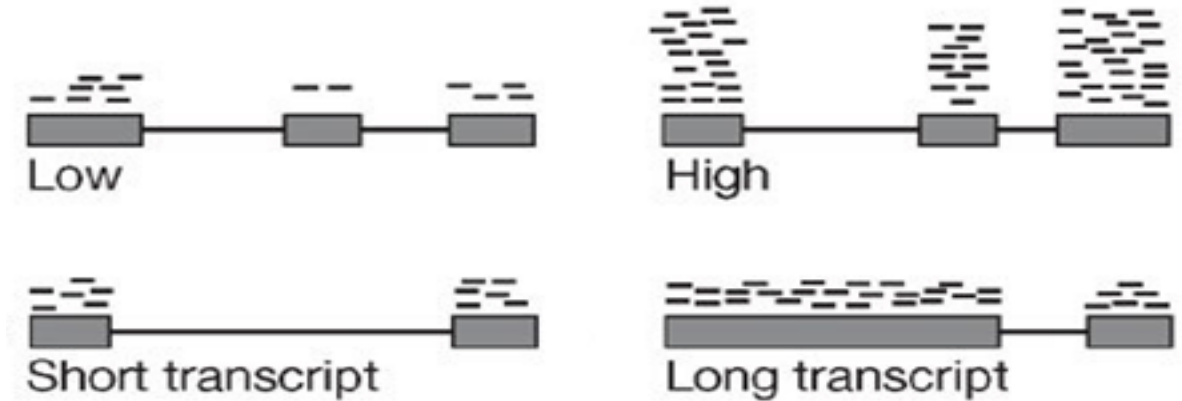
Replicates (technical / biological)



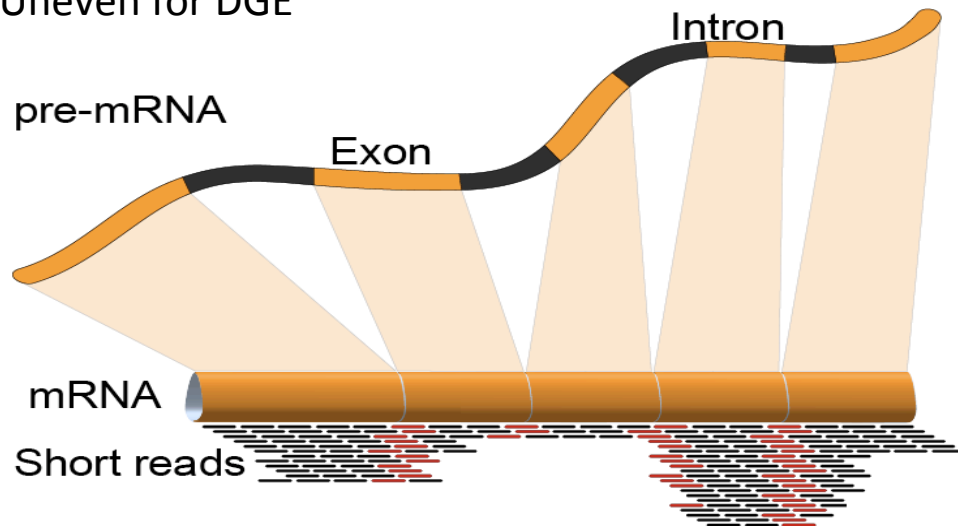
- Even for annotation



- Target transcript properties (low abundant vs high abundant transcripts)

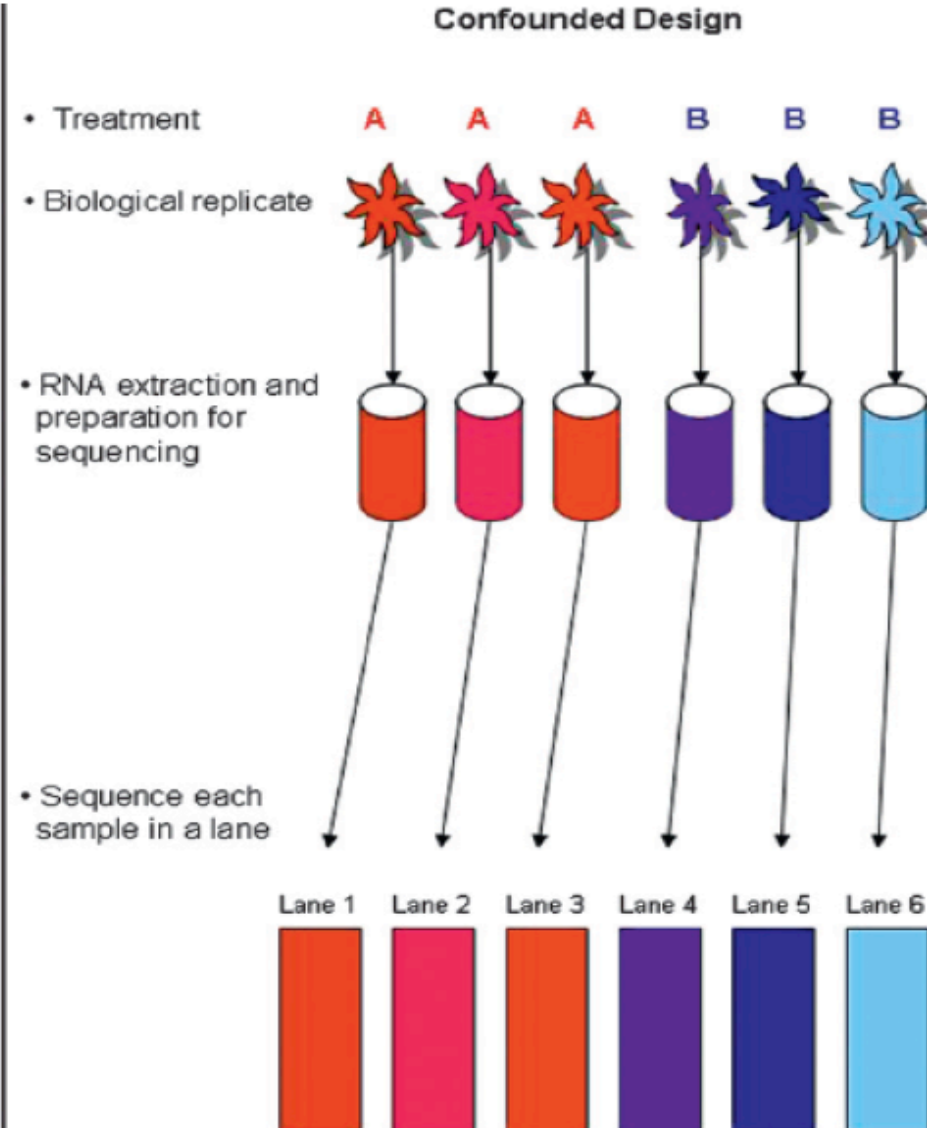
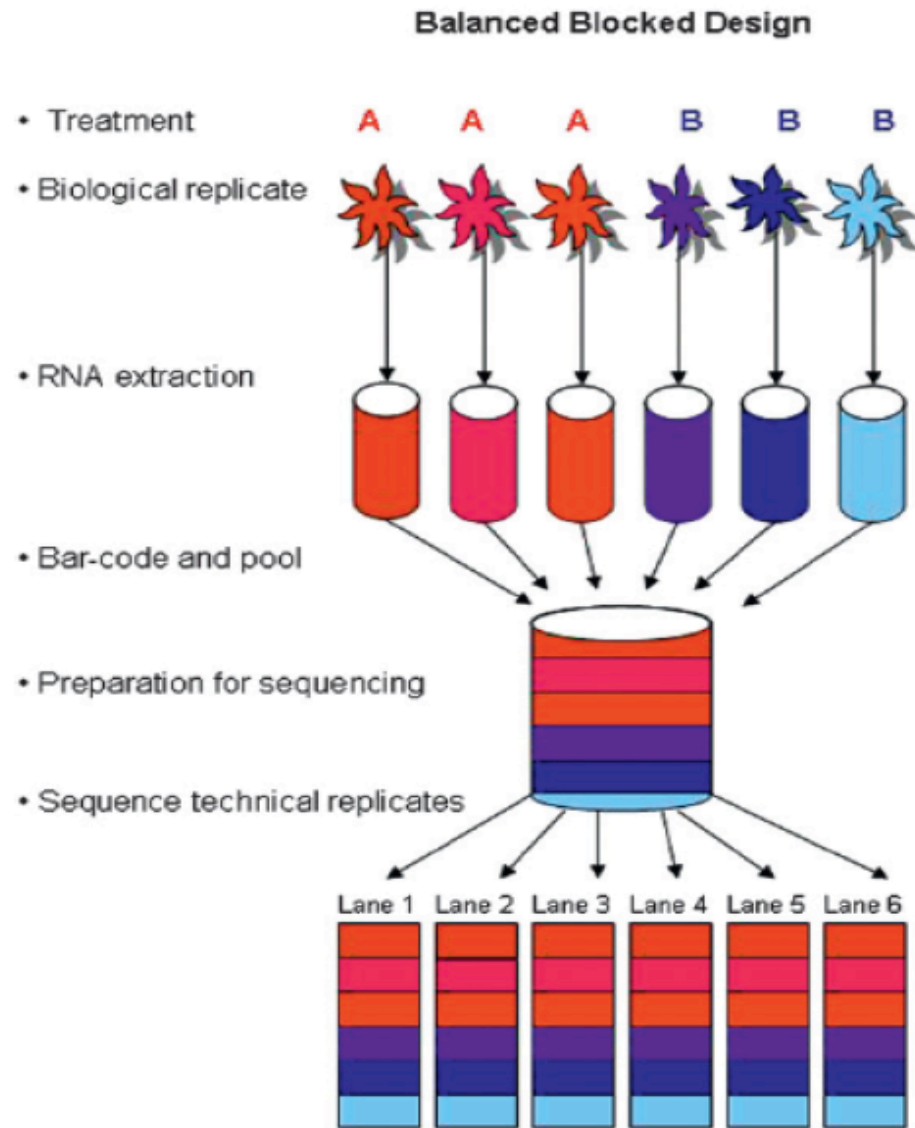


- Uneven for DGE



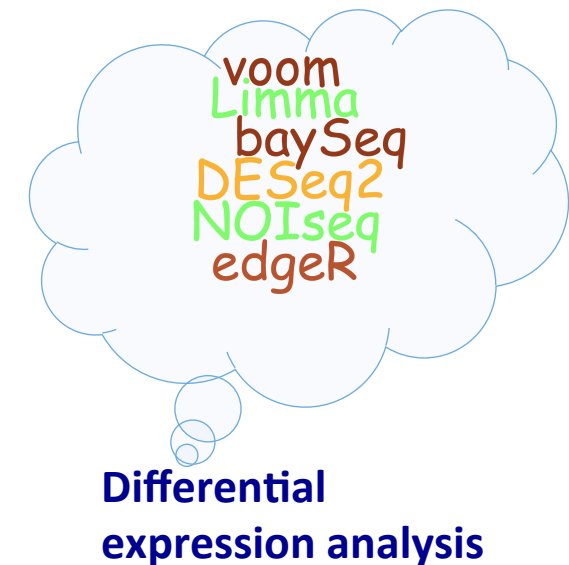
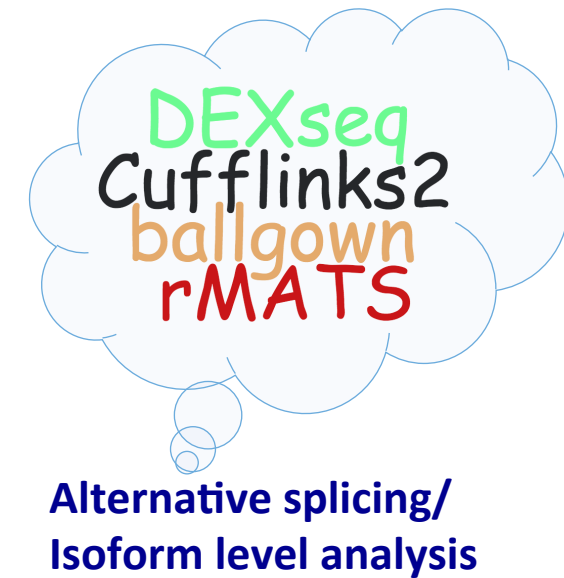
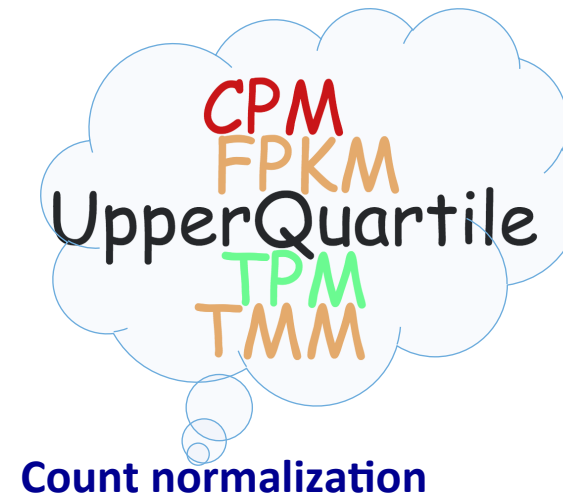
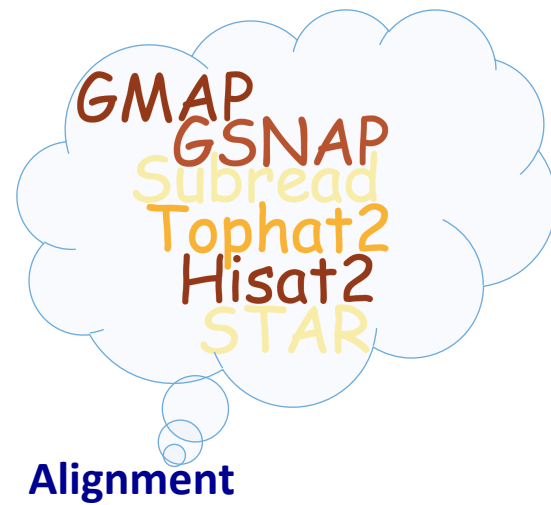
- Allele might not be detected (not in the genome/not being expressed)
- Estimate expression of each allele

Randomization and Blocking



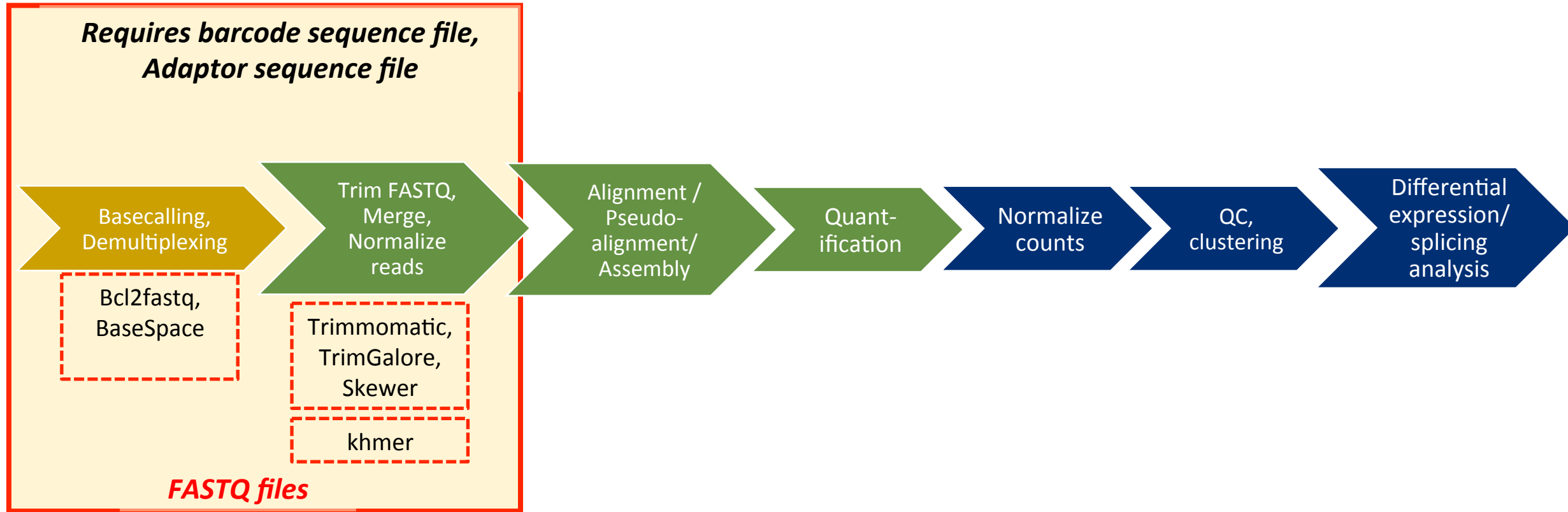
Subjectivity of the analysis

- ✓ 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.



Data Analysis

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

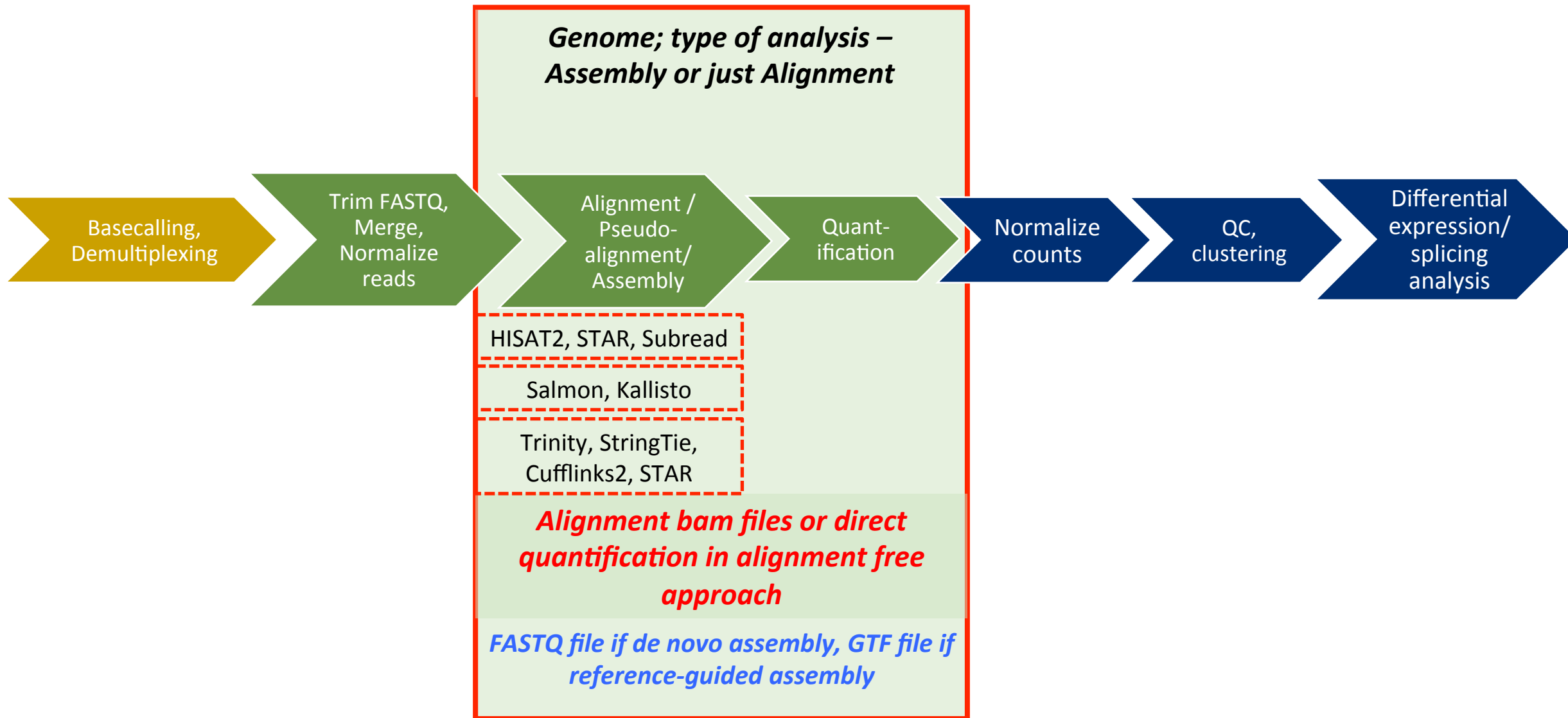


Read processing

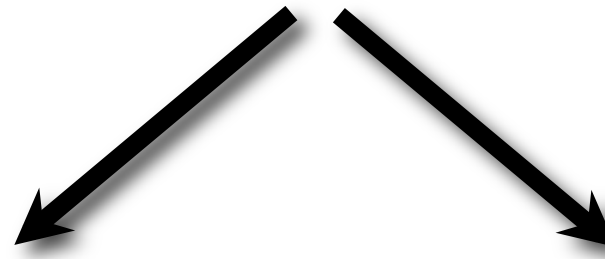
For de novo assembly

Software	De-multiplexing	Adaptor Trimming	Quality Filtering/ Trimming	K-mer Filtering	K-mer Normalization
FASTX-Toolkit	✓	✓	✓		
Goby	✓	✓			
khmer				✓	✓
NGS_backbone		✓	✓		
Stacks	✓	✓	✓	✓	✓
trimmomatic		✓	✓		
biopieces	✓	✓	✓		

Read processing, alignment/assembly+alignment, quantification



When to use each?



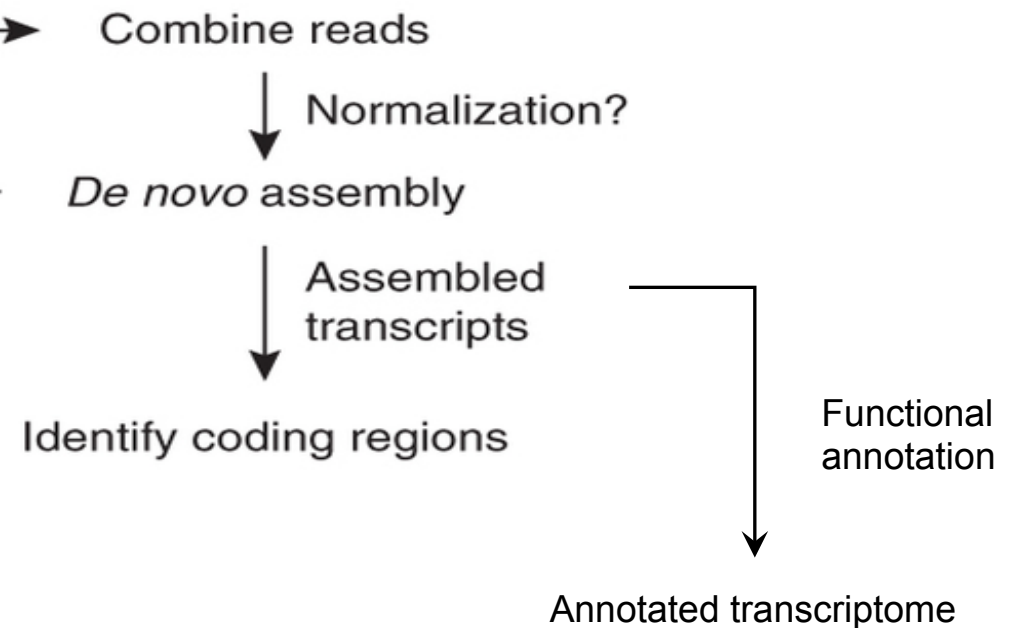
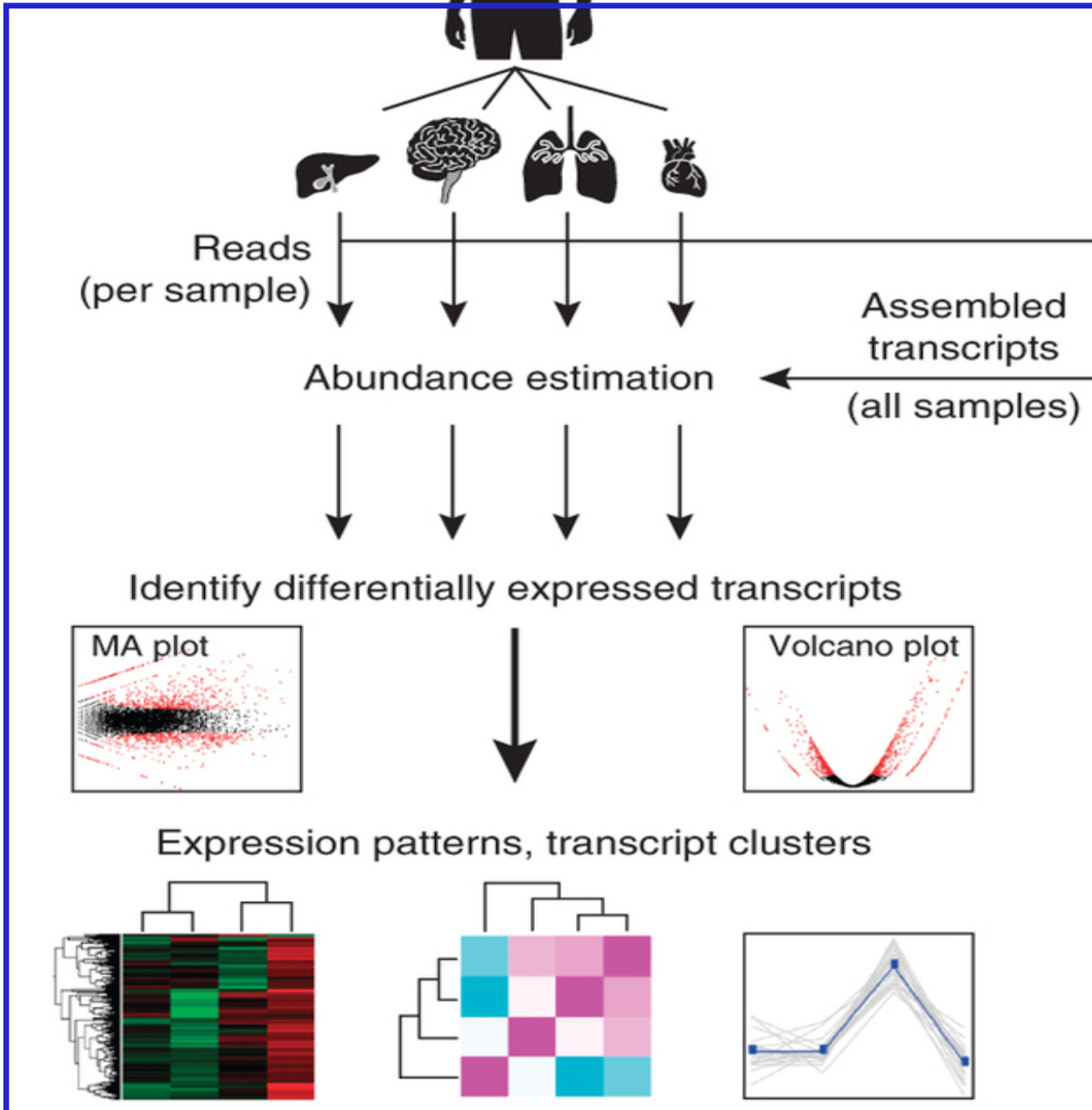
de novo

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

reference

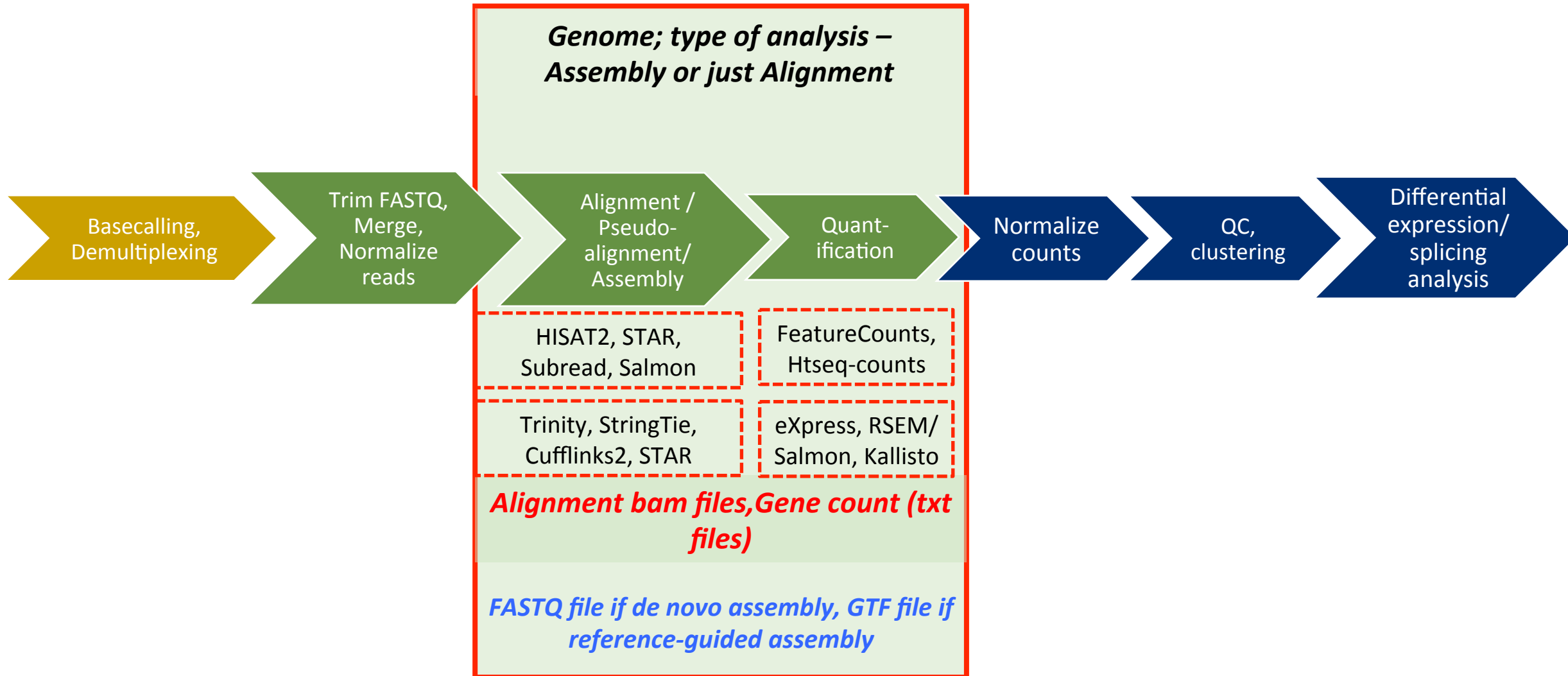
(do know the transcriptome)
(main goal is to **quantify** NOT to
discover)

Protocol for *de novo* RNAseq

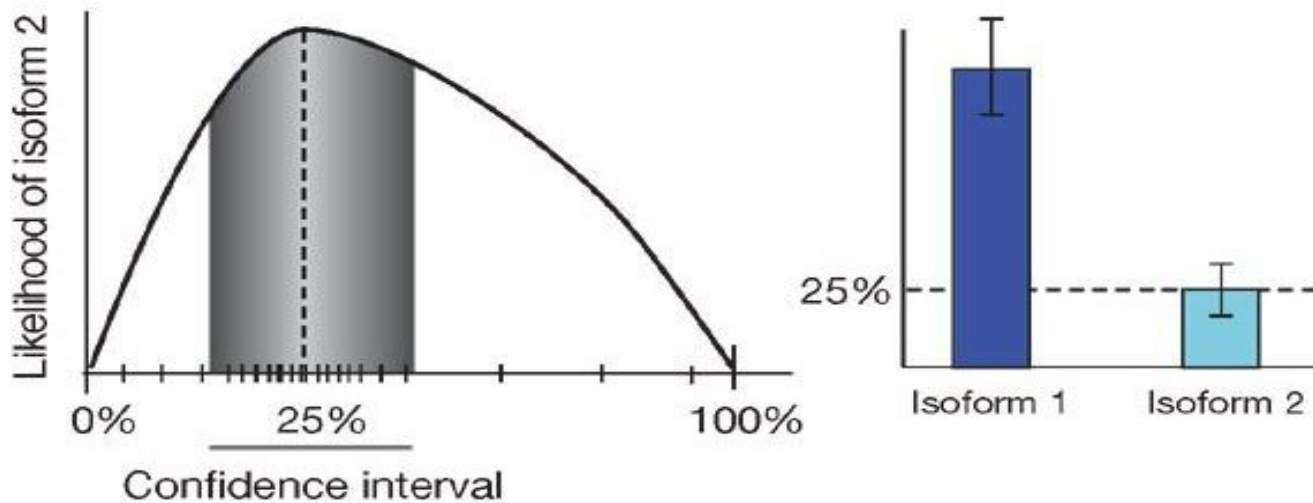
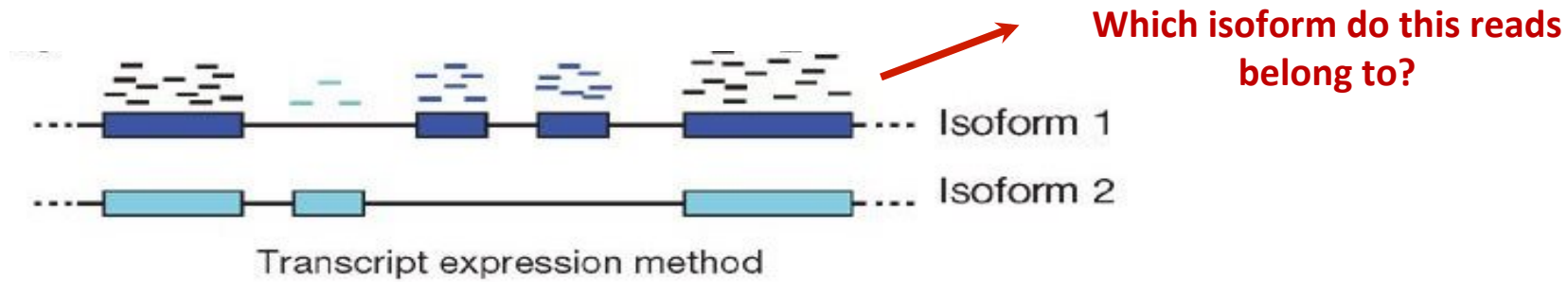


Protocol for reference-based RNAseq

Read processing, alignment/assembly+alignment, quantification

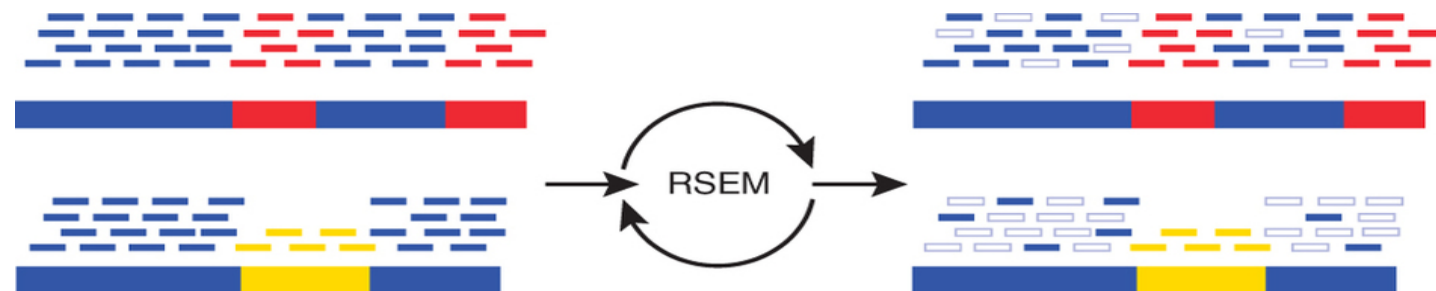


Expression quantification

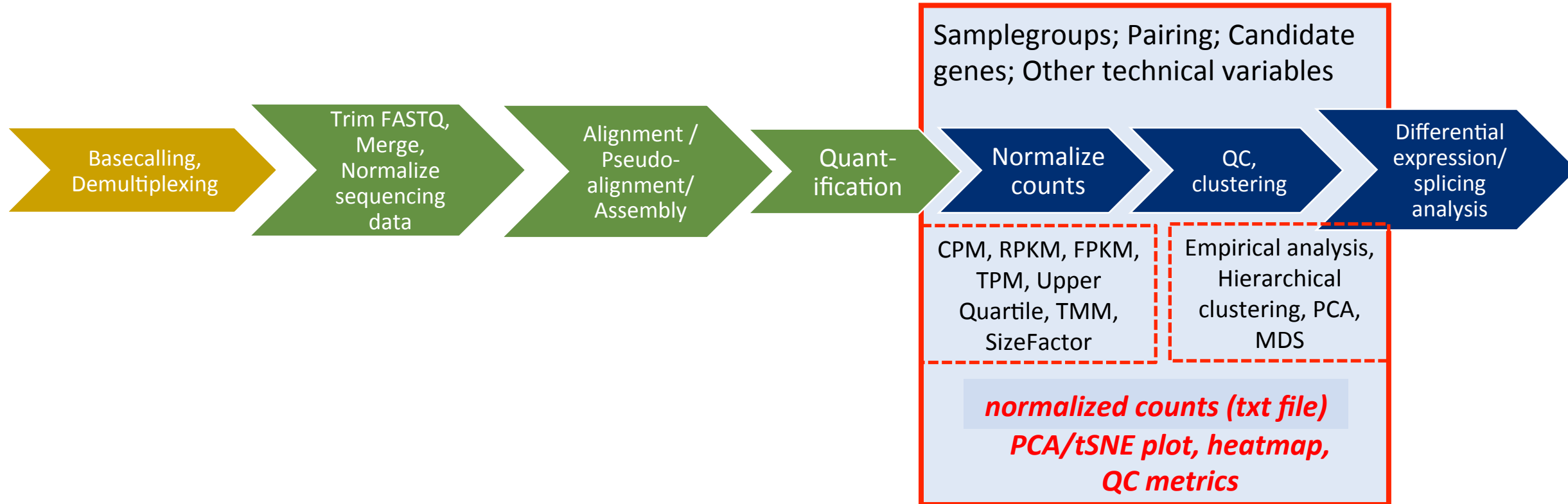


Genome guided: Cufflinks2, StringTie

Transcriptome guided: RSEM, eXpress, Salmon, Kallisto

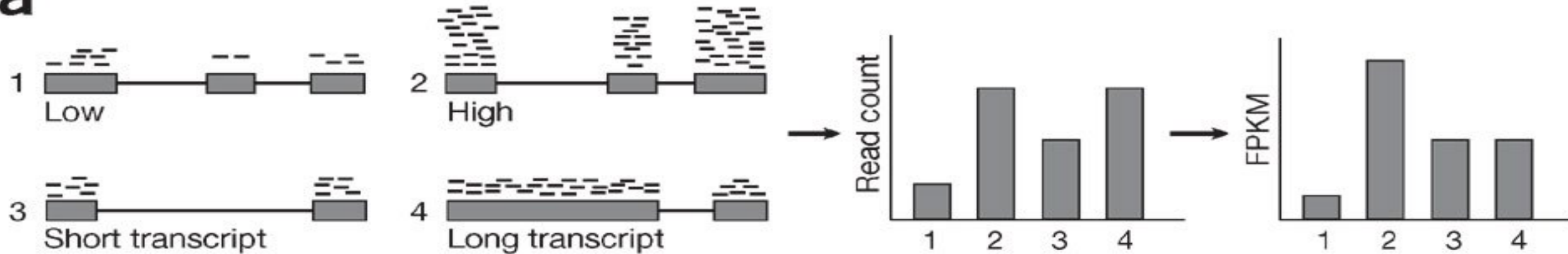


Normalization, Expression quantification



Count normalization

a



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.

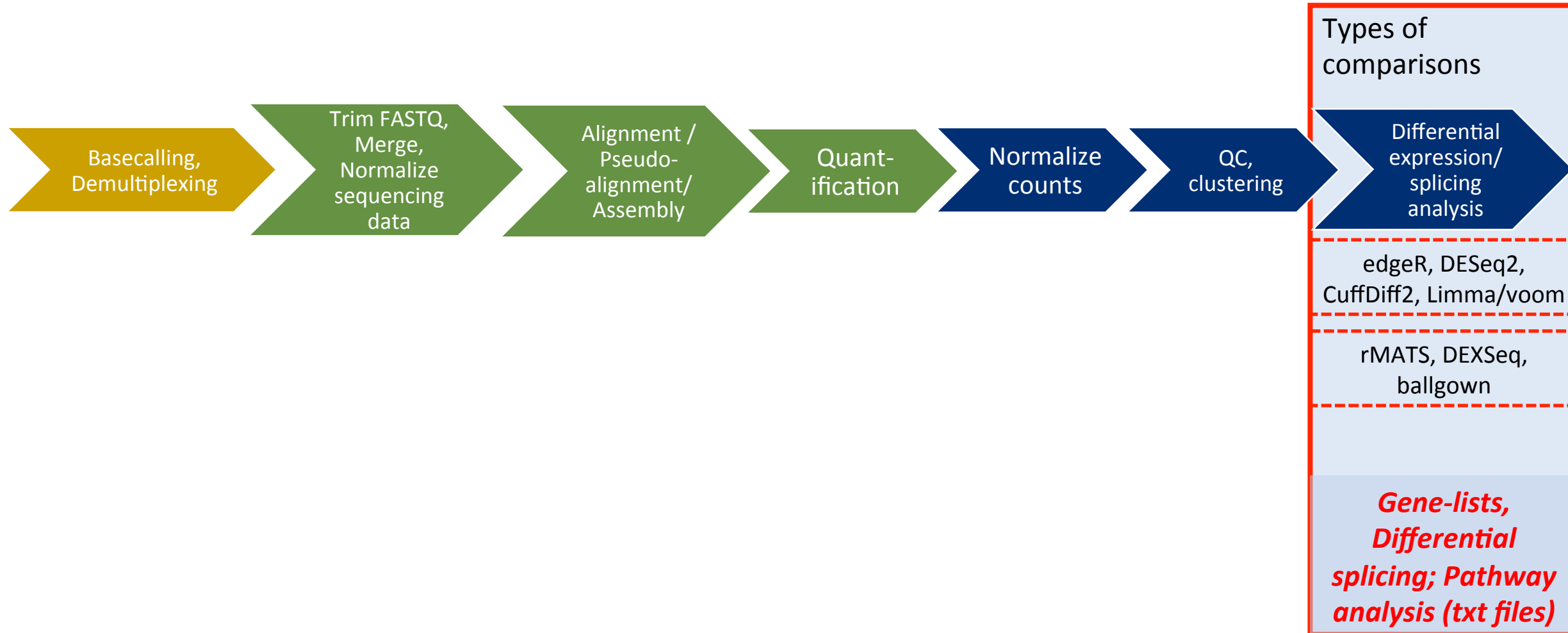
”Gene counts” should be corrected in order to minimize these biases:
normalization.

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

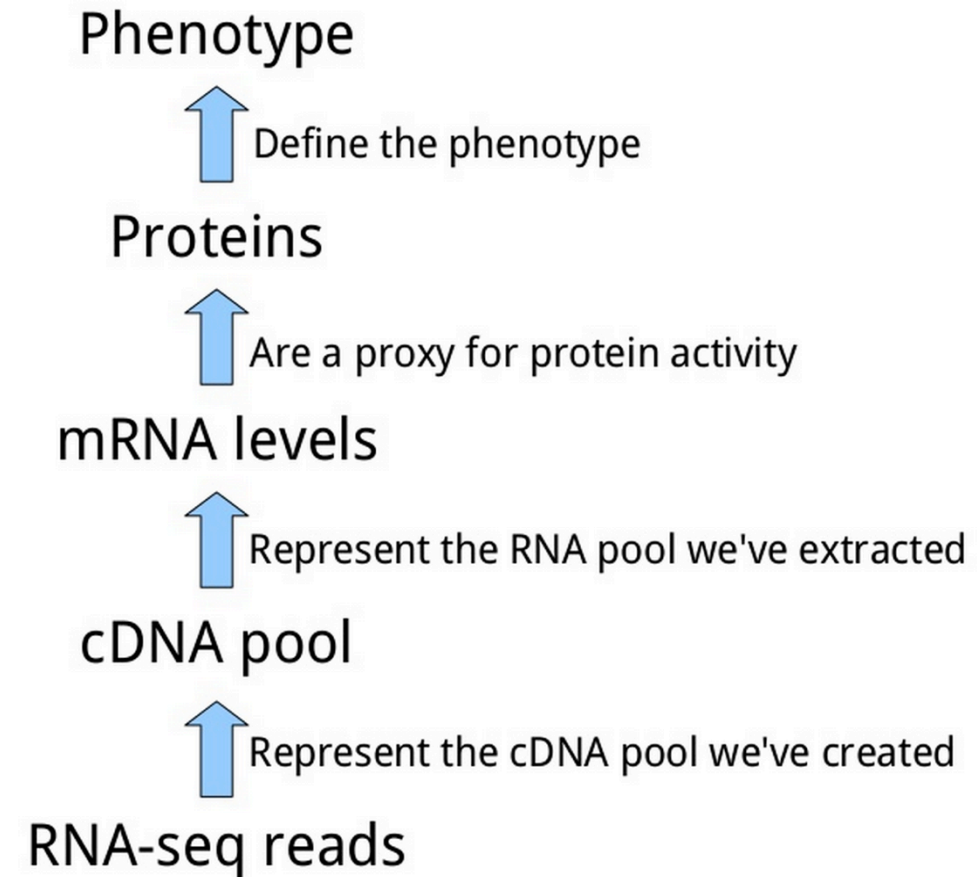
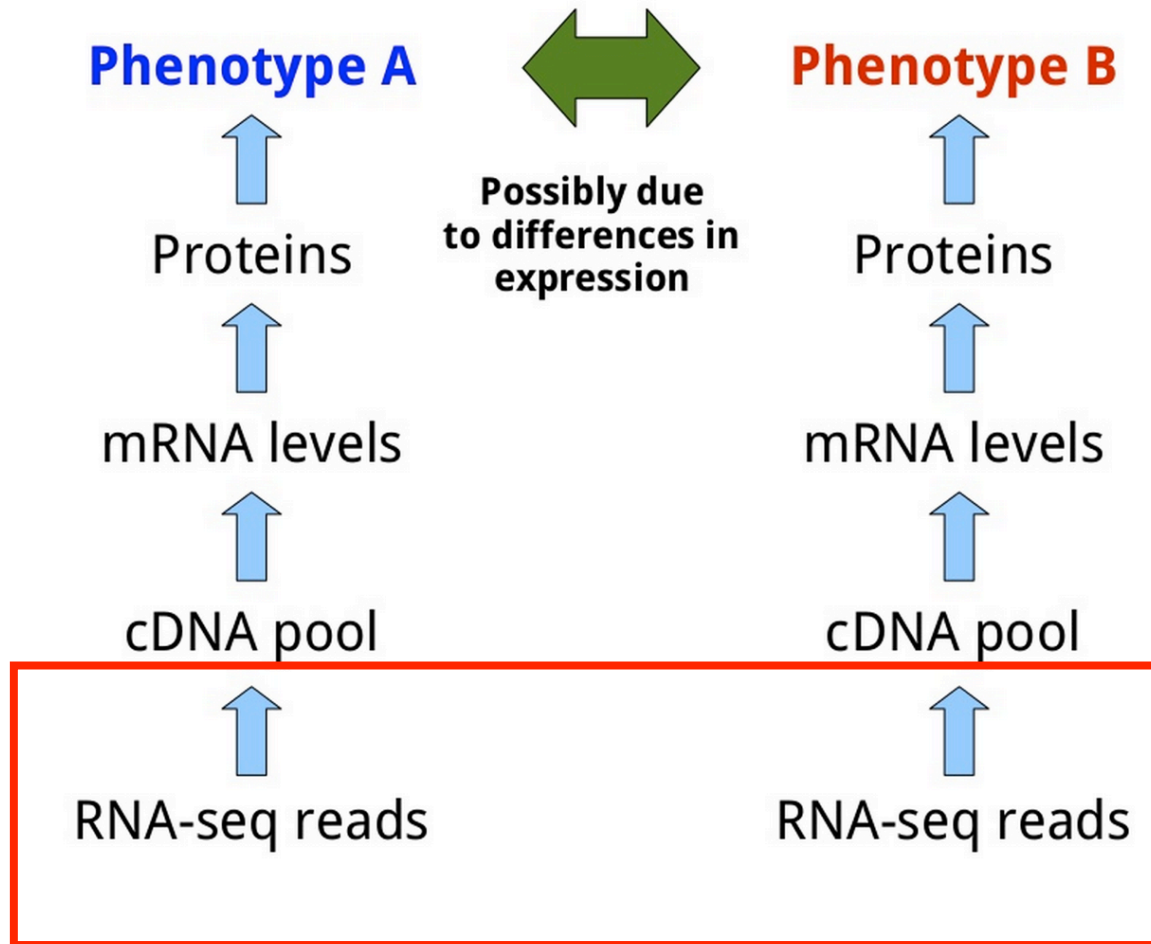
Count normalization

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

Differential expression analysis

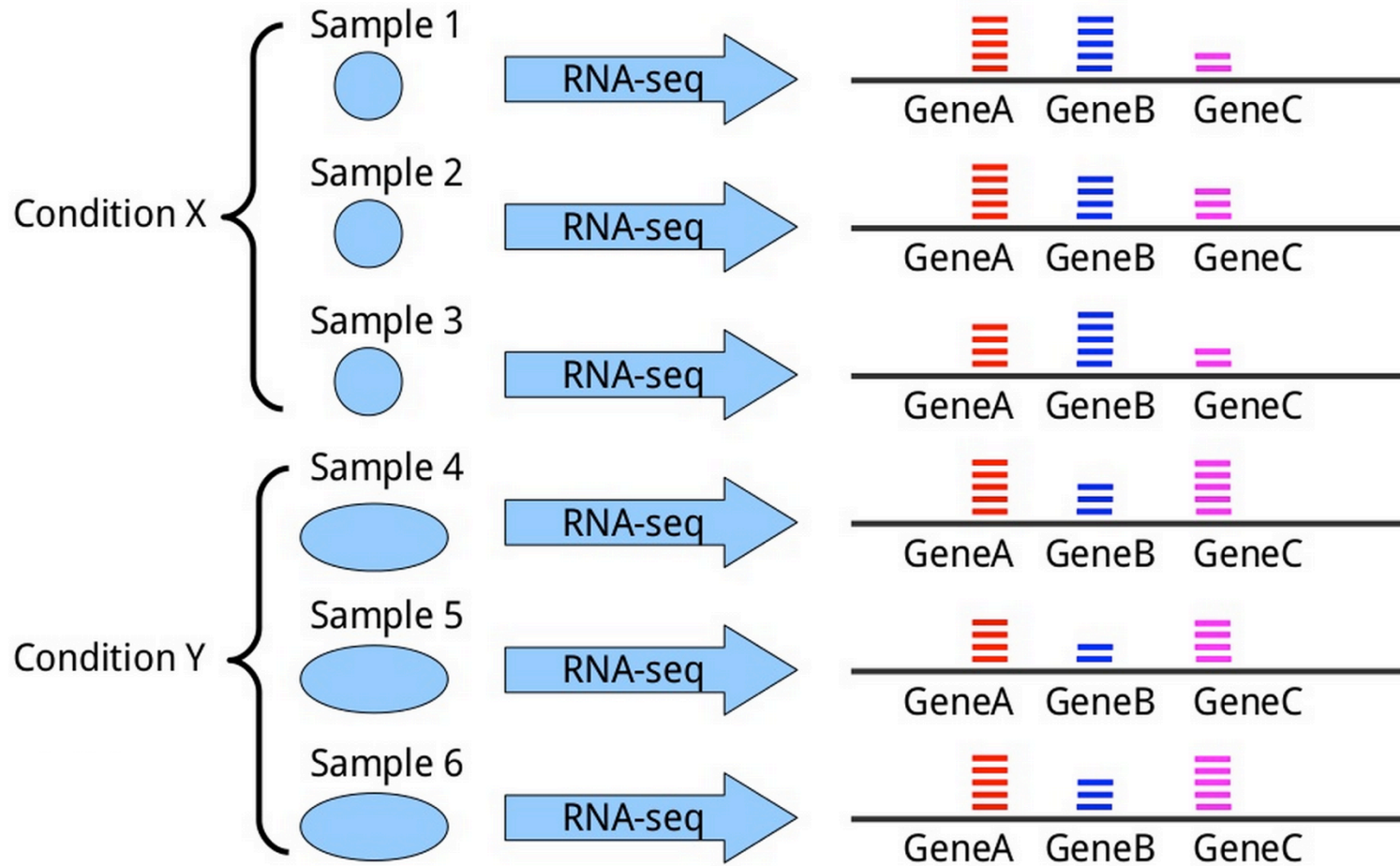


Our assumptions and comparison

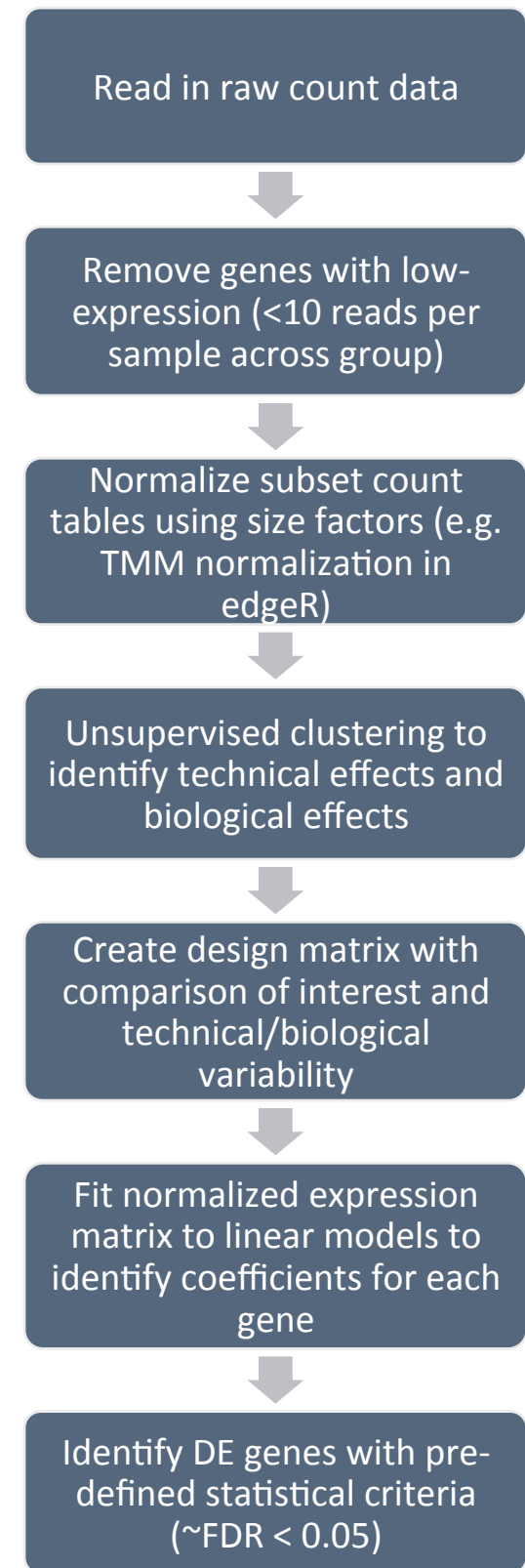


<http://www.slideshare.net/joachimjacob/1rna-seqpart1working-tothegoal?related=2>

Statistical testing for Differential expression



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



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.**

- 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