

Introduction to  
RNA-Seq applications and tools  
22<sup>nd</sup> October, 2019

Organised and delivered by Bioinformatics Core at WHG:  
Eshita Sharma



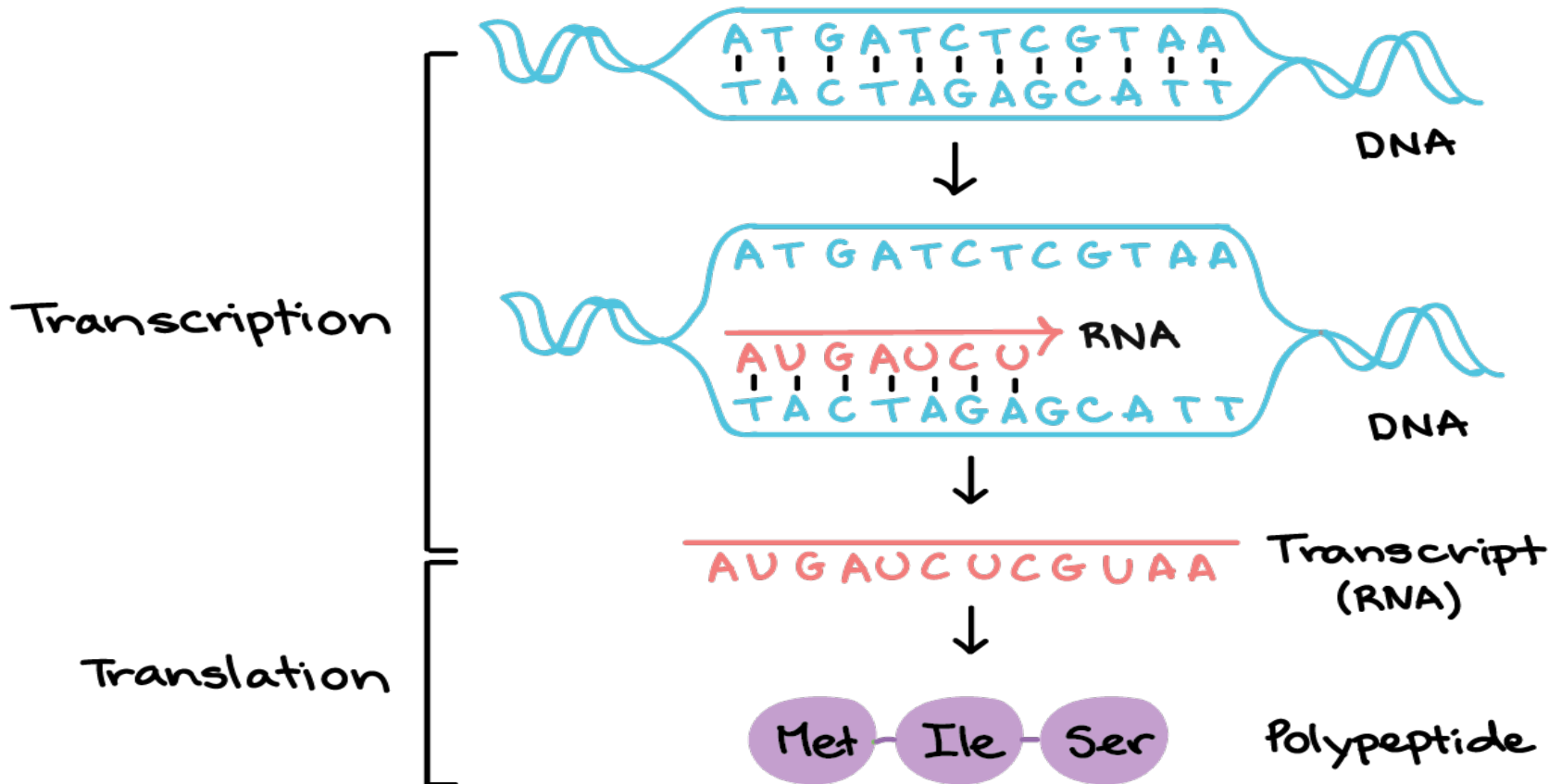
# What we do when we do RNAseq?

- **General Introduction to RNA world**
- **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](mailto:eshita.sharma@well.ox.ac.uk)

Research Associate in Functional Genomics, Bioinformatics Core

# RNA - Mid-point of the information cascade

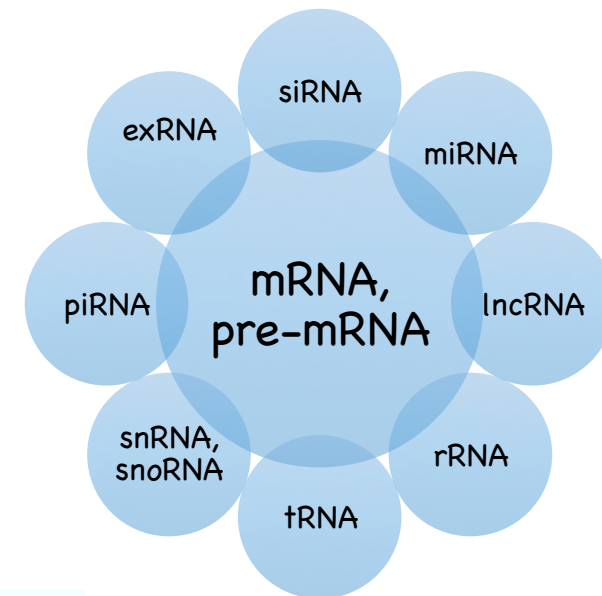
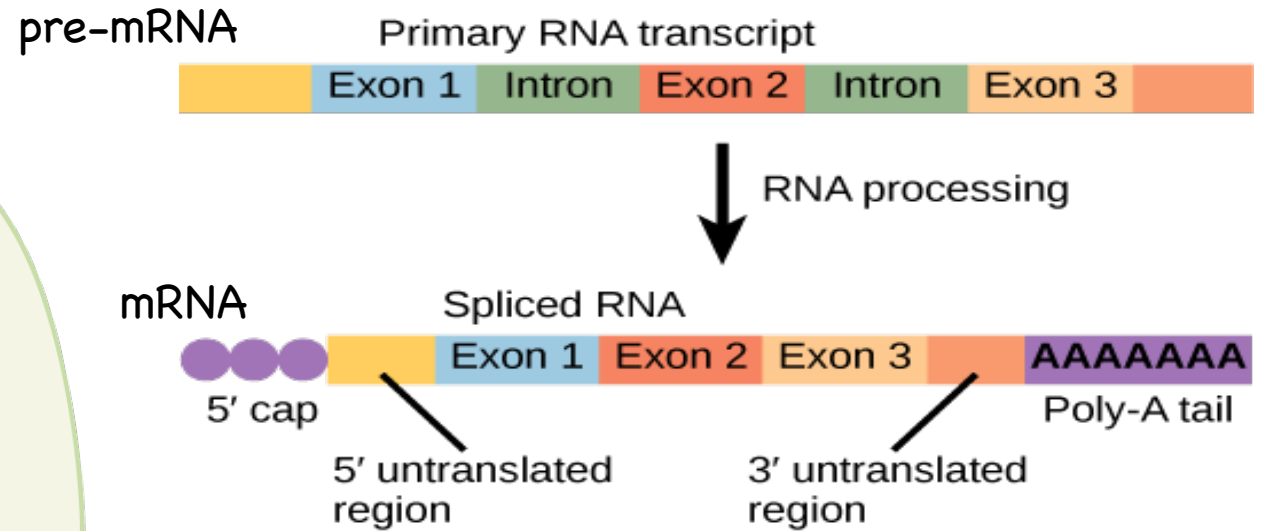
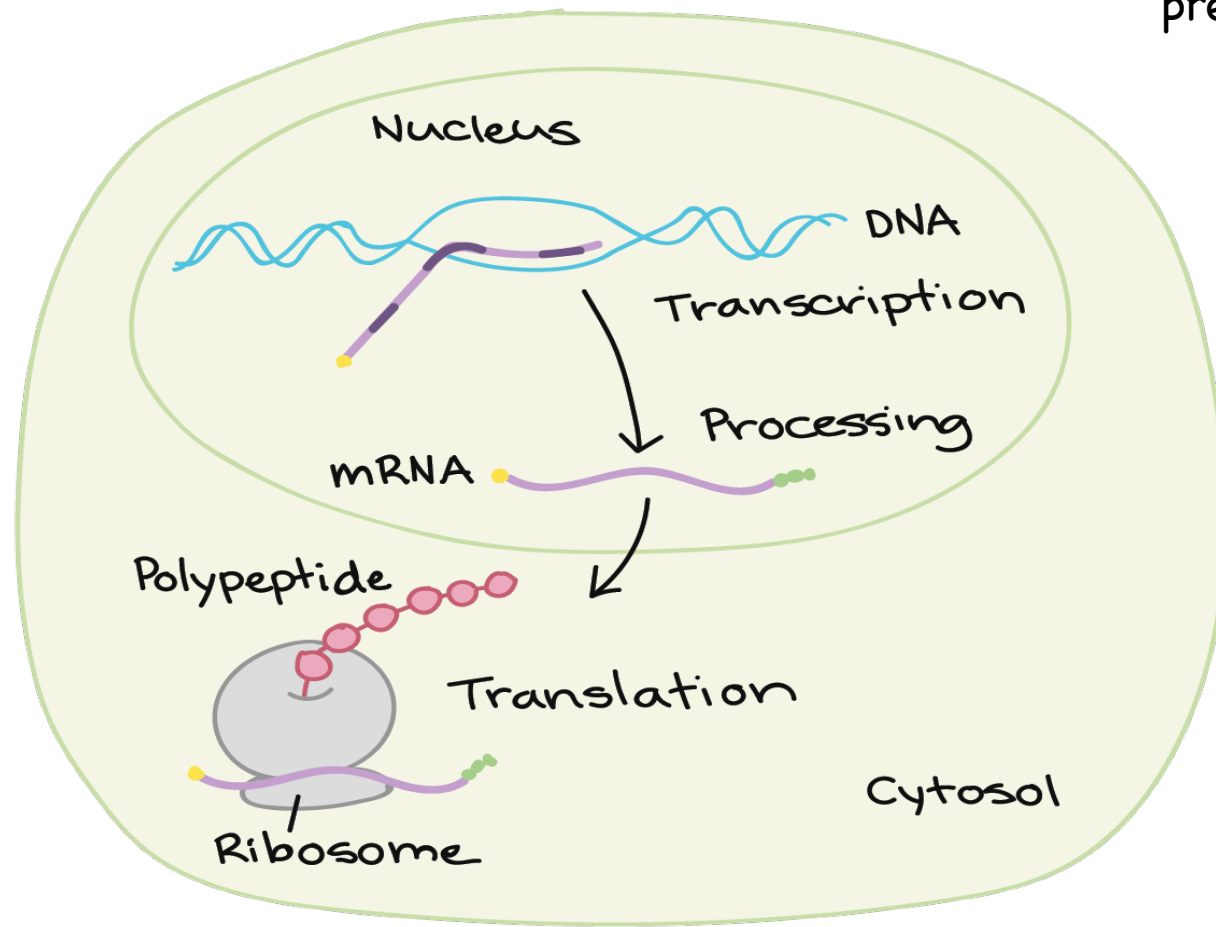


- Image credit Khan Academy Open Courses

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



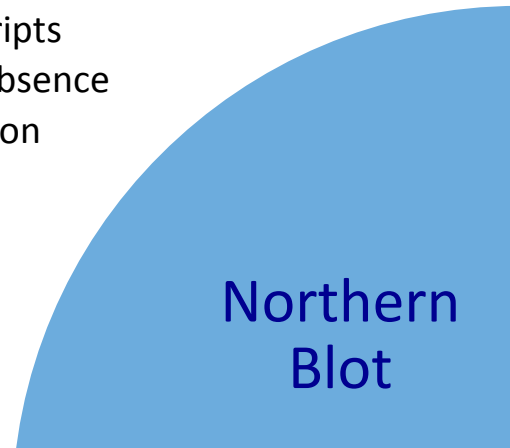
- Image credit Khan Academy Open Courses

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

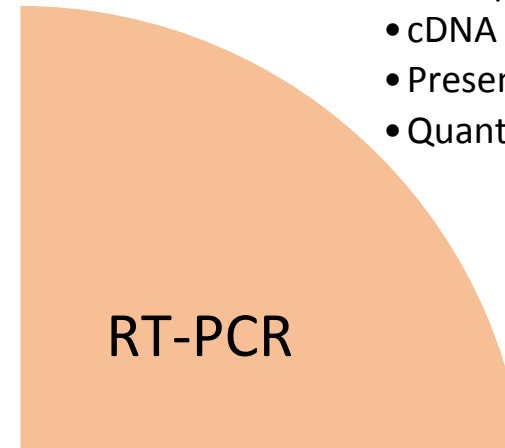
# RNAseq is a method for Transcriptome profiling

Image of the transcribed genome at any point of time!

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

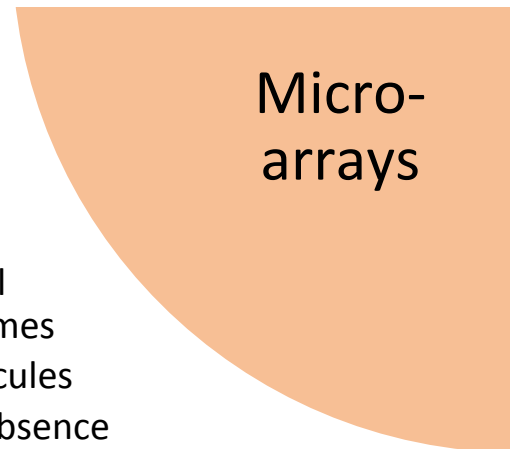


- Multiple genes
- cDNA molecules
- Presence/Absence
- Quantification

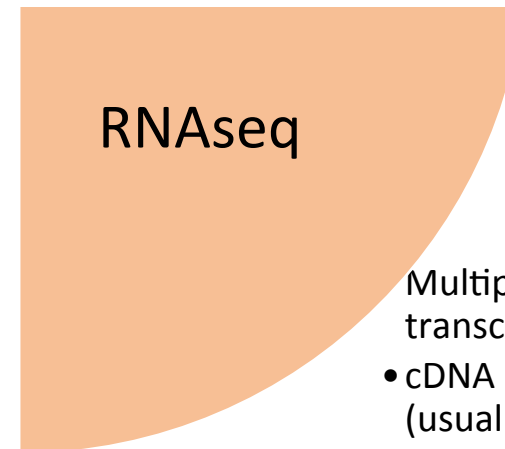


## How do we take this image?

- 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.)

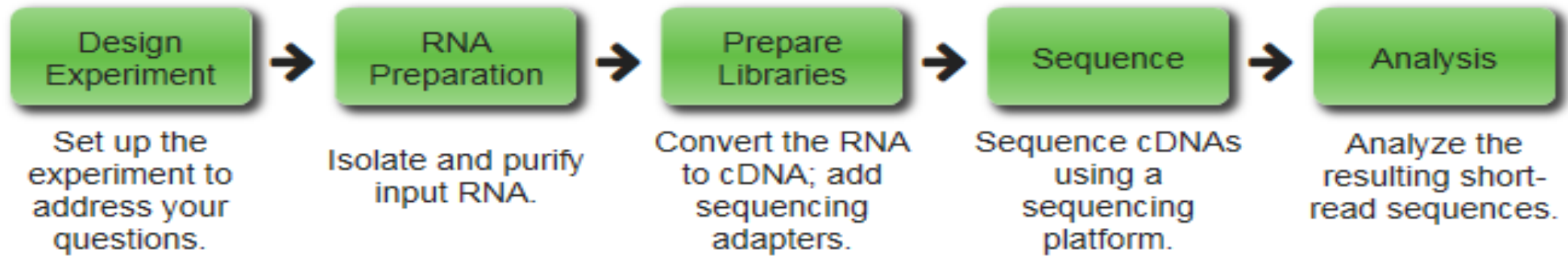
Can be used to:

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

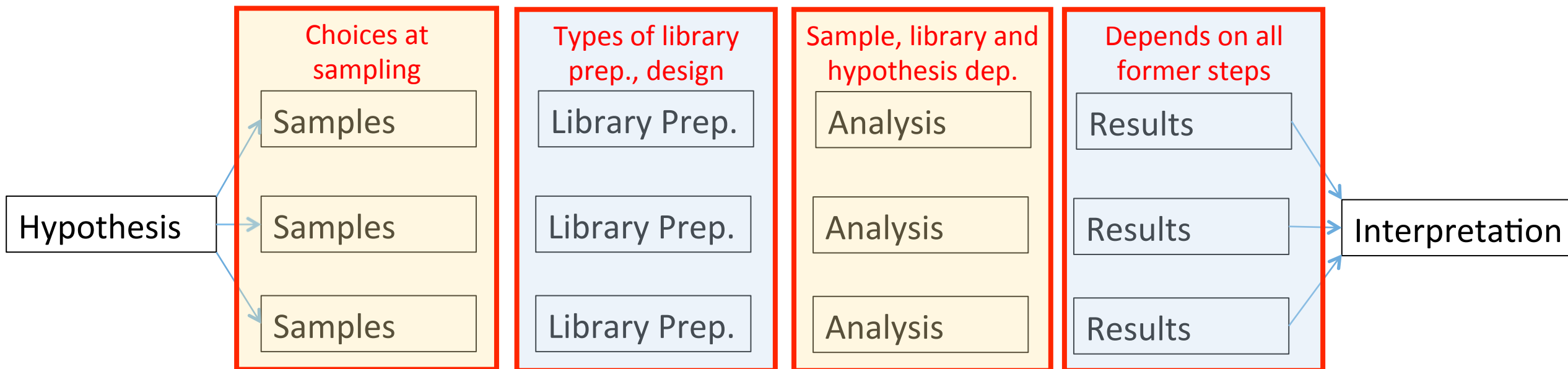
# Most experiments aim to -

- 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 gene/transcript in different conditions.

# Typical RNAseq experiment



**“All research is atypical”**



Choices in experiment setup

# Usual approaches to RNAseq library-prep

polyA enriched  
mRNAseq

Oligo dT<sub>25</sub> selection

Fragmentation of mRNA

1<sup>st</sup> strand → 2<sup>nd</sup> strand →  
cDNA synthesis of  
fragments

Adapter ligation → PCR  
amplification

Mature mRNA

Directionality

Requires good quality  
total RNA

100ng-1ug; 8-10 samples/  
lane HiSeq4000

Ribodepleted  
Total RNAseq

Biotinylated ribosomal  
RNA probes

Bind ribosomal RNA

extracted with  
streptavidin beads

Fragmentation → cDNA  
synthesis → adapters →  
PCR

Mature mRNA, nascent  
RNA, non-coding  
transcripts

Works with low-quality  
RNA, e.g. FFPE samples

Requires high sequencing  
depth

100ng to 1ug; 4-6  
samples/lane HiSeq4000

SMARTer/Smartseq2  
mRNAseq

Oligo dT primer used for  
Reverse transcription

Template switching by RT

PCR pre-amplification of  
full-length cDNA

Tn5 transposase  
tagmentation & library  
prep.

Full-length mature mRNA

Pre-amplification of low-  
input RNA

Requires good quality  
total RNA; no  
directionality

10pg to 10ng; 10-30  
samples/lane HiSeq4000

Small RNAseq <200nt  
ncRNAseq

3' adapter ligation

5' adapter ligation

1<sup>st</sup> strand cDNA synthesis

PCR enrichment → Size  
selection

Focus on 21-25nt miRNA/  
siRNA involved in gene  
regulation

Size specificity,  
Directional, Low-input  
and Low depth

Requires good quality  
total RNA

1ug total RNA; 20+  
samples/lane HiSeq2500;  
50bp SE

3' mRNAseq  
mRNAseq

Oligo dT<sub>25</sub> primed RT  
(First strand synthesis)

Removal of RNA template

Random priming and  
second strand synthesis

Bead purification of  
tagged cDNA library →  
PCR

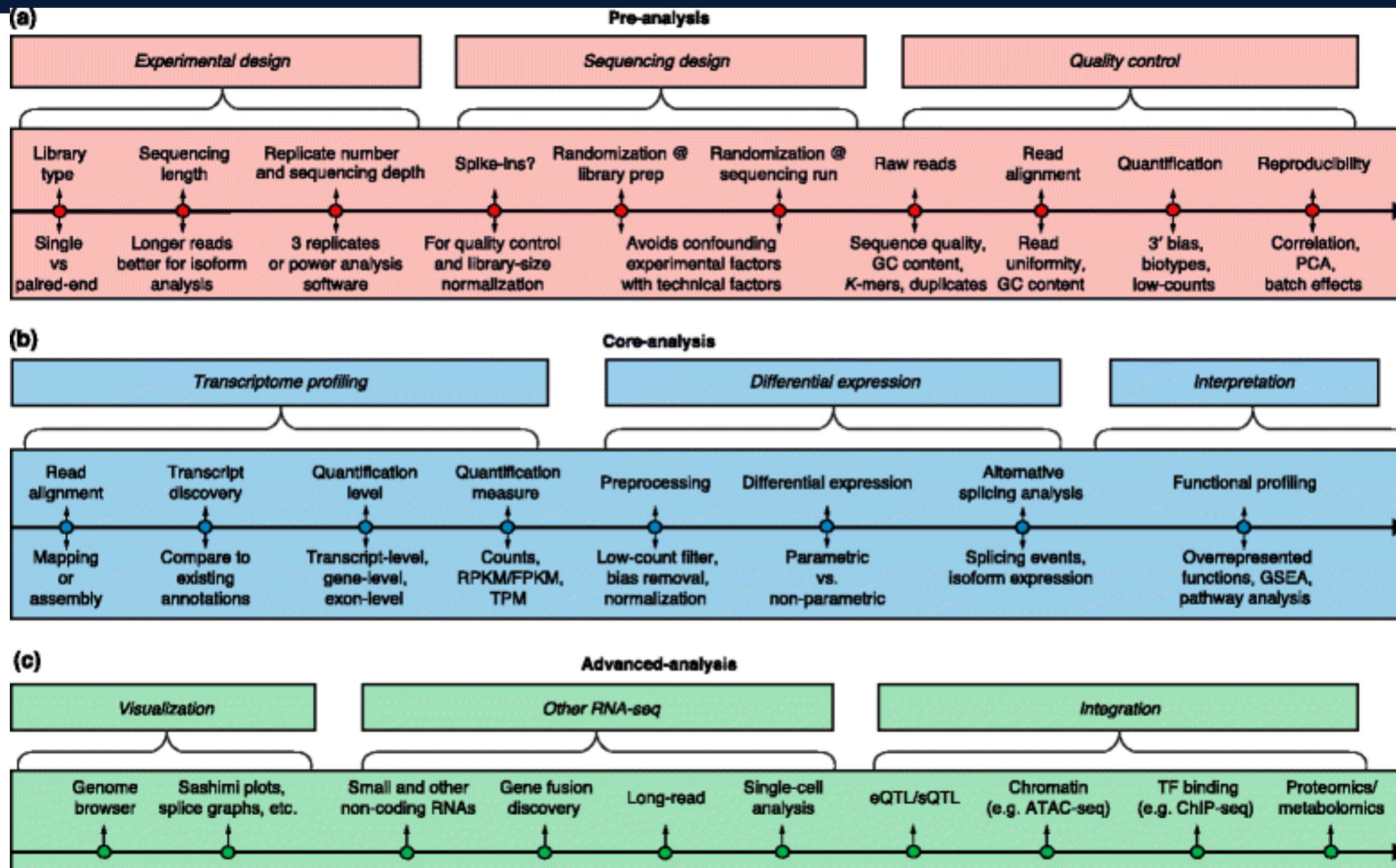
200-300bp insert libraries  
of 3' ends of mRNA

Lower sequencing depth;  
poor-quality RNA works

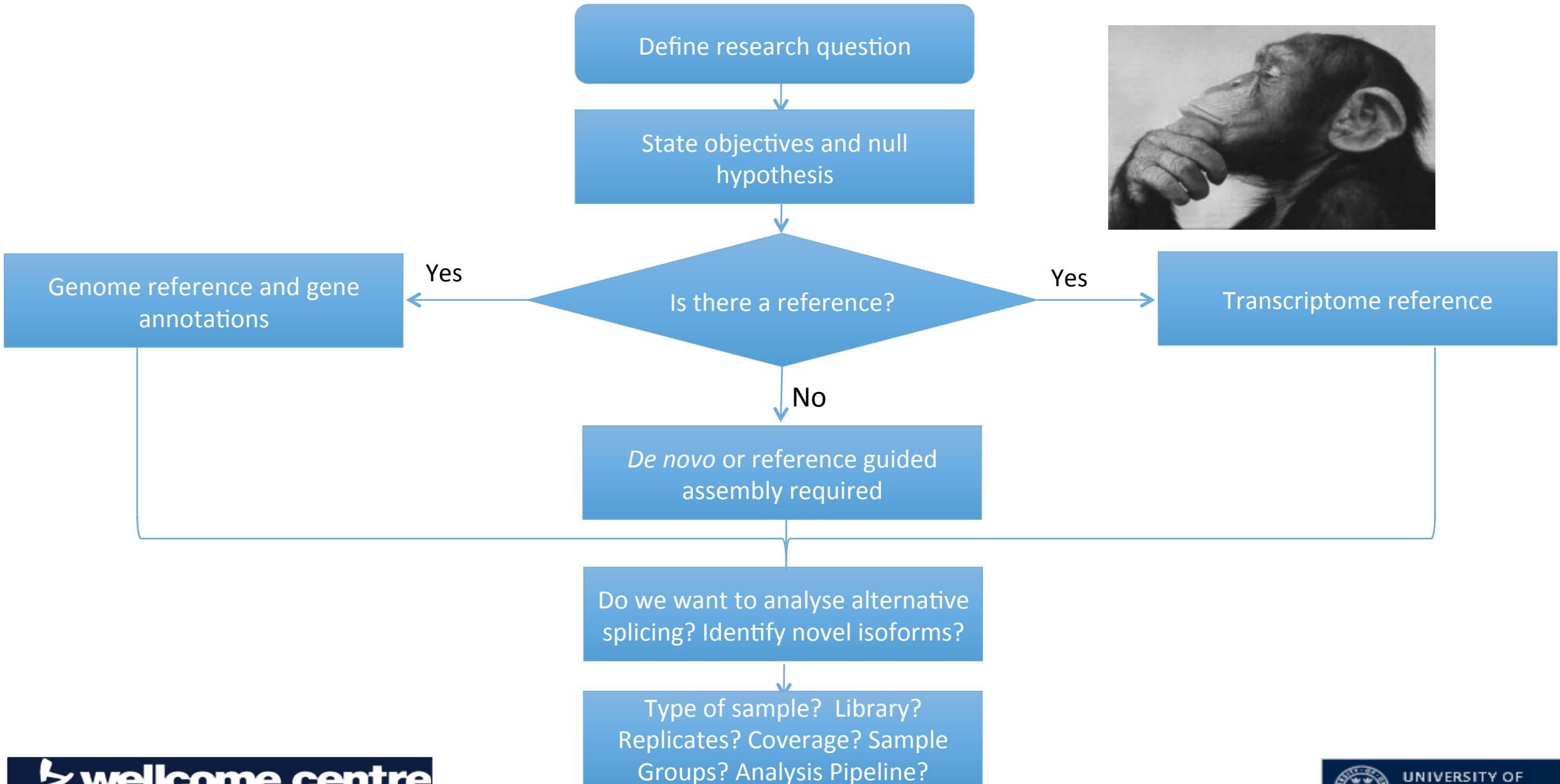
Mainly useful for  
expression quantification

0.5 ng – 2 ug; 48 samples  
per lane HiSeq4000;  
50-75bp SE

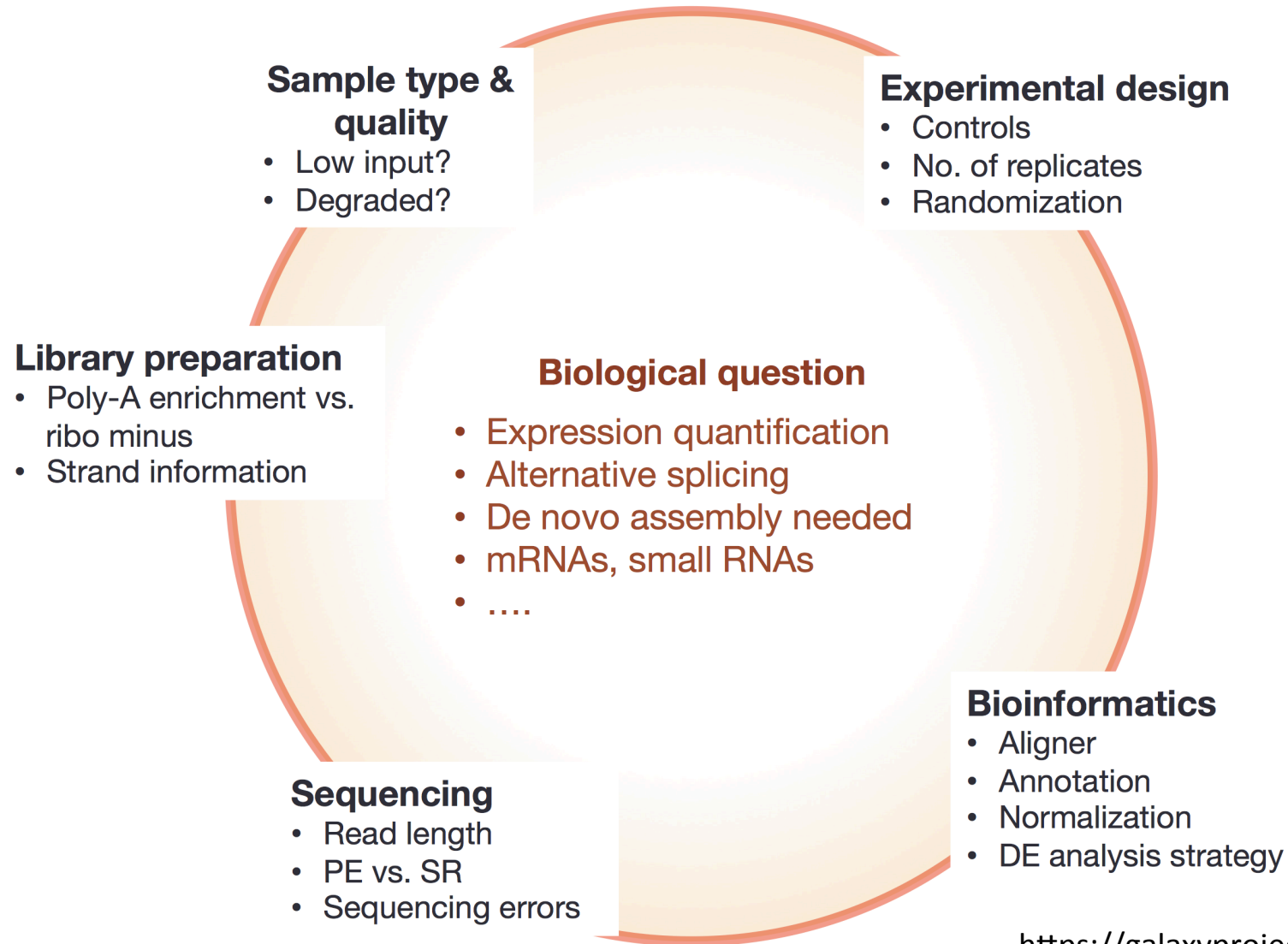
# Generic roadmap of RNAseq analysis



# What do we know? What are we looking for?



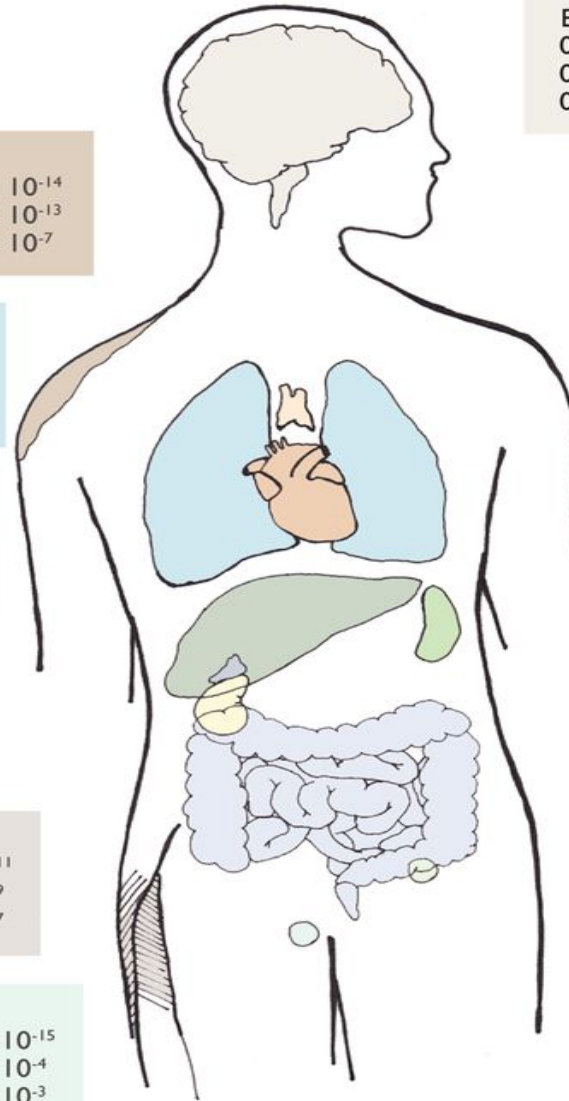
# Everything is connected...



[https://galaxyproject.org/tutorials/rb\\_rnaseq/#transcript-quantification](https://galaxyproject.org/tutorials/rb_rnaseq/#transcript-quantification)

# Common Decision Points

# Sample selection



SKIN		
0031424	keratinization	$2.9 \times 10^{-14}$
0006955	immune response	$3.5 \times 10^{-13}$
0031069	hair follicle morphogenesis	$4.1 \times 10^{-7}$

LUNG		
0030324	lung development	$6.2 \times 10^{-16}$
0006954	inflammatory response	$2.1 \times 10^{-15}$
0043330	response to exogenous dsRNA	$6.2 \times 10^{-6}$

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

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

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

TESTIS		
0007059	chromosome segregation	$9.1 \times 10^{-15}$
0007276	gametogenesis	$8.1 \times 10^{-4}$
0006349	imprinting	$1.5 \times 10^{-3}$

BRAIN		
0007268	synaptic transmission	$8.9 \times 10^{-41}$
0016358	dendrite morphogenesis	$1.2 \times 10^{-10}$
0007611	learning or memory	$7.9 \times 10^{-6}$

THYMUS		
0019882	antigen presentation	$7.1 \times 10^{-21}$
0045059	positive thymic T cell selection	$9.8 \times 10^{-8}$
0045060	negative thymic T cell selection	$2.6 \times 10^{-7}$

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

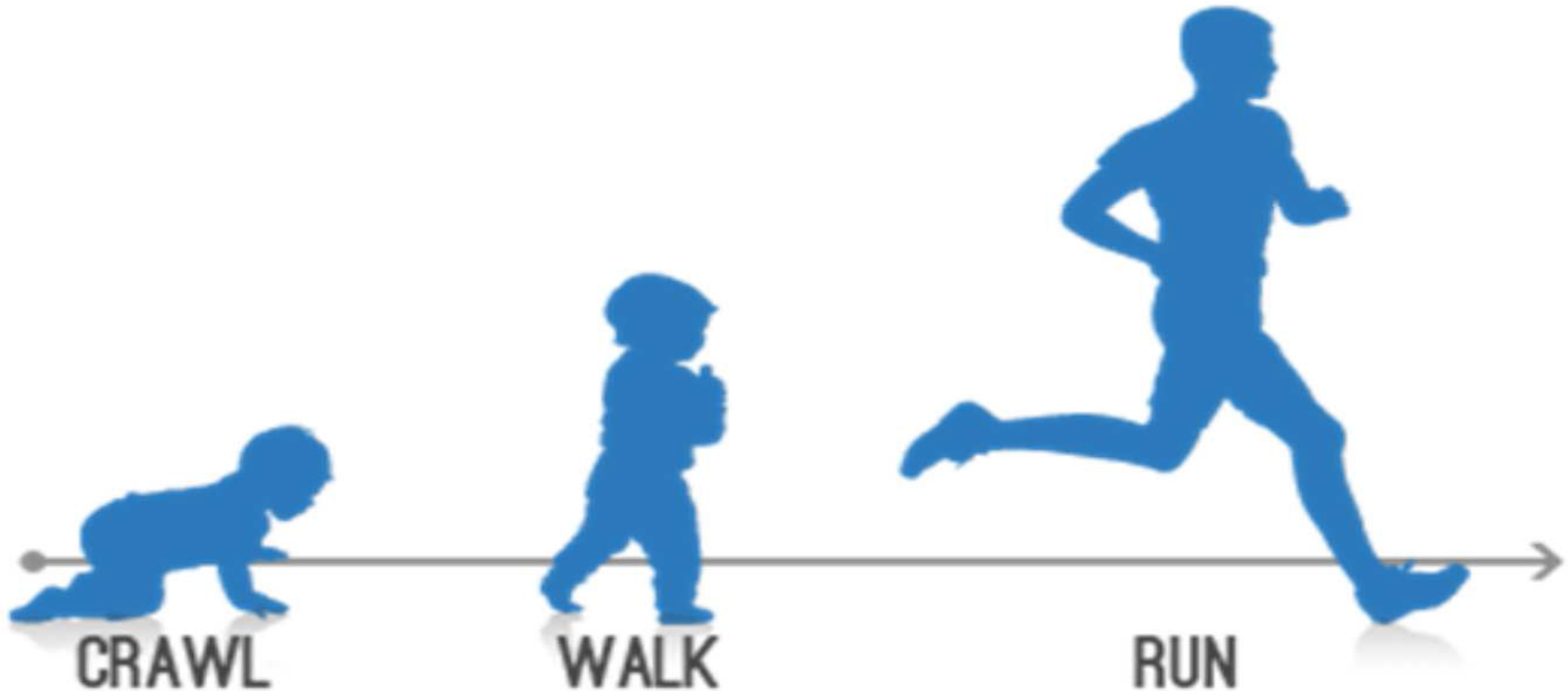
LIVER		
0008203	cholesterol metabolism	$2.6 \times 10^{-8}$
0007596	blood coagulation	$2.0 \times 10^{-7}$
0000050	urea cycle	$5.0 \times 10^{-5}$

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

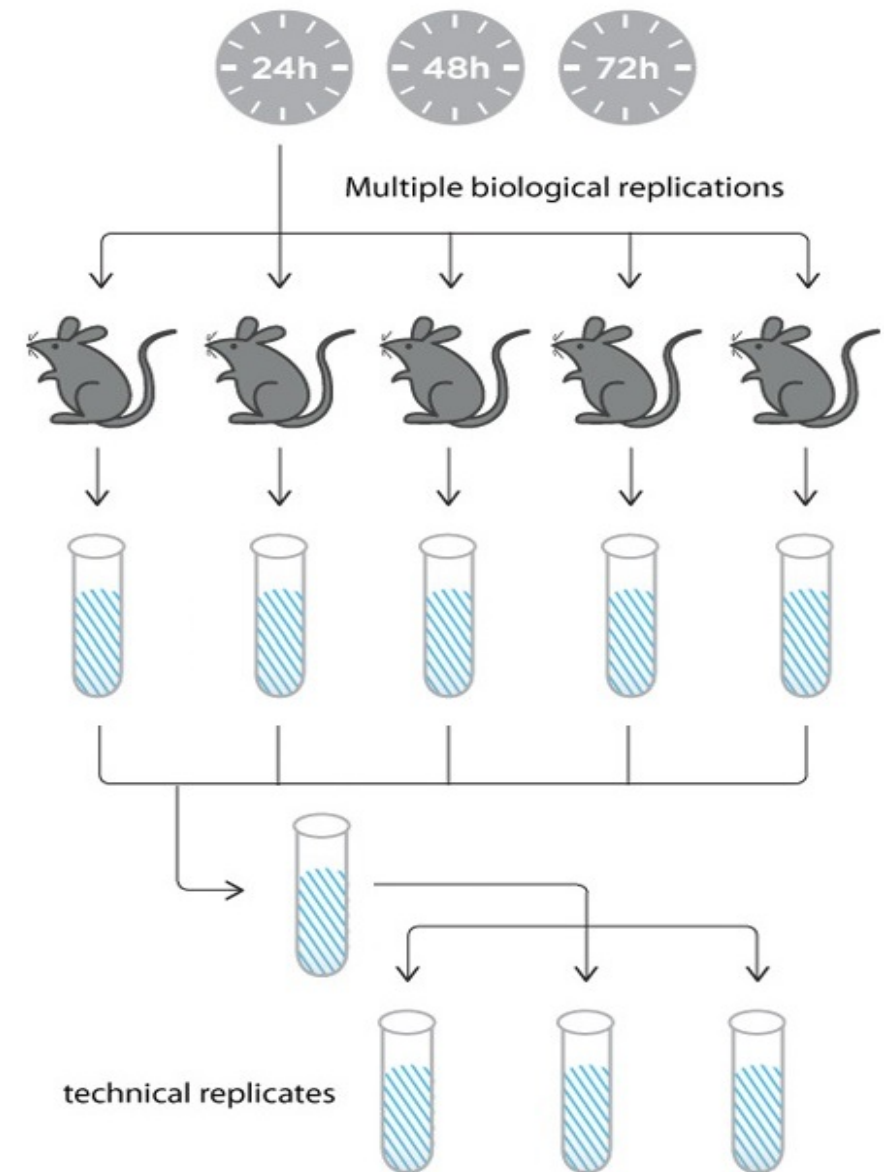
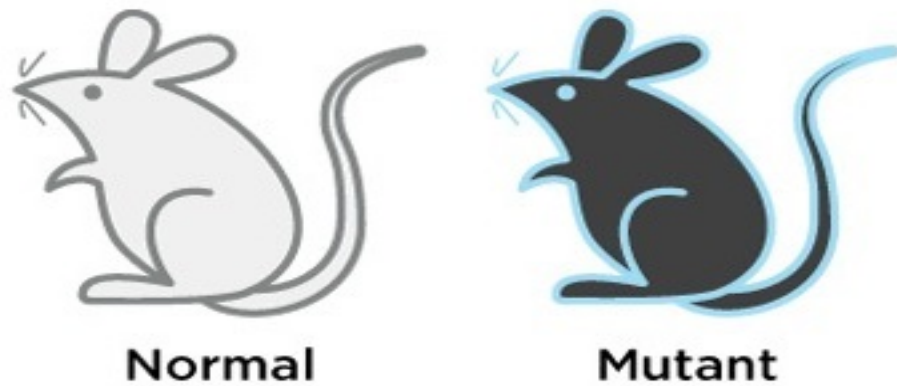
INTESTINE		
0006955	immune response	$7.0 \times 10^{-13}$
0007586	digestion	$9.3 \times 10^{-5}$
0050892	intestinal absorption	$4.6 \times 10^{-4}$

OVARY		
0007059	chromosome segregation	$1.0 \times 10^{-12}$
0007276	gametogenesis	$8.6 \times 10^{-8}$
0006349	imprinting	$3.5 \times 10^{-5}$

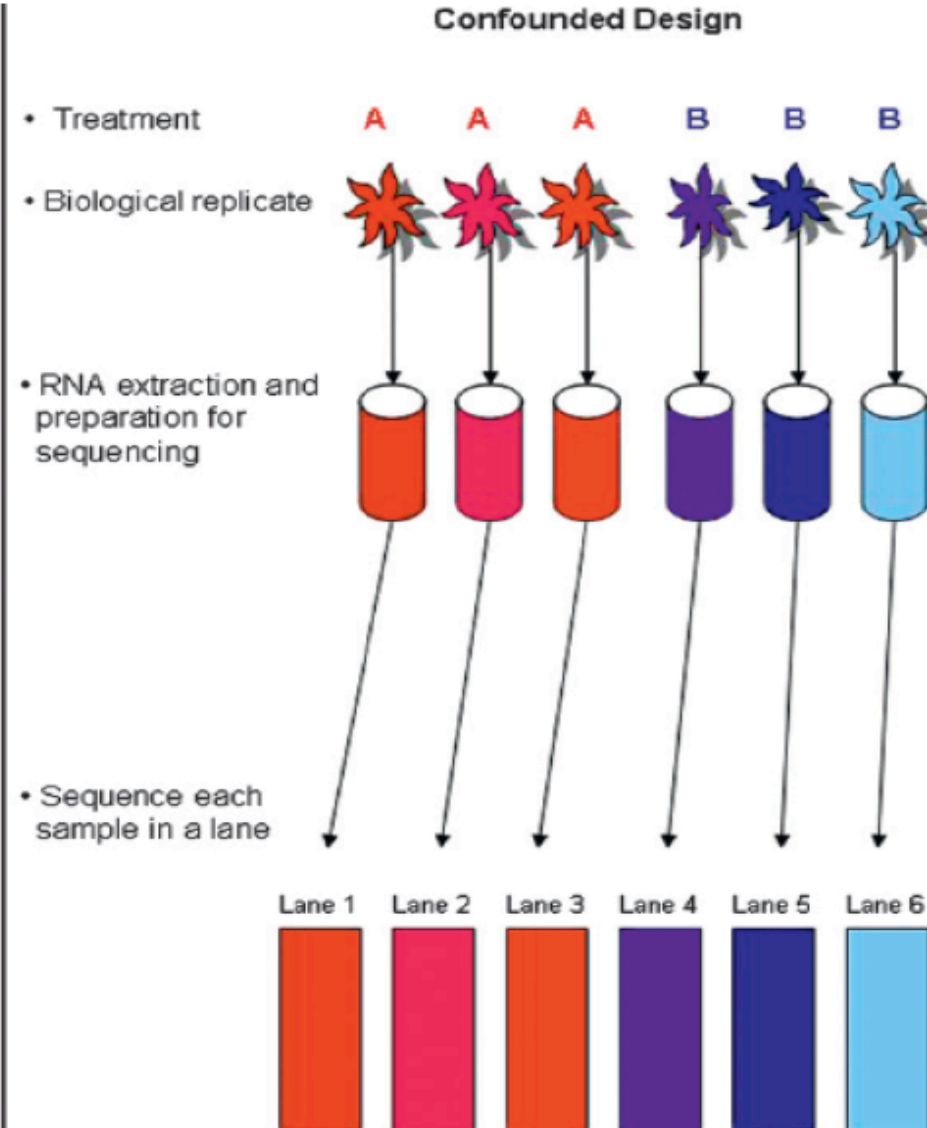
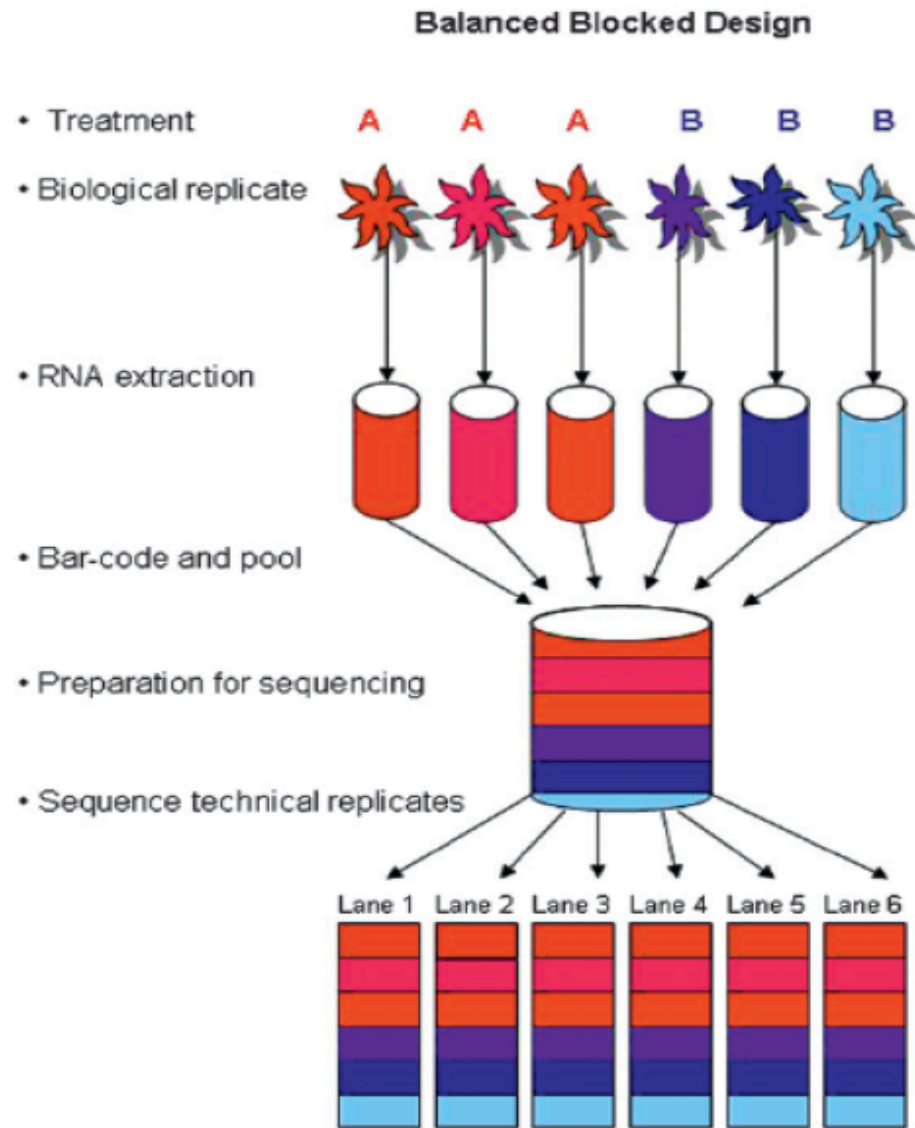
# Time of sampling



# Replicates (technical / biological)



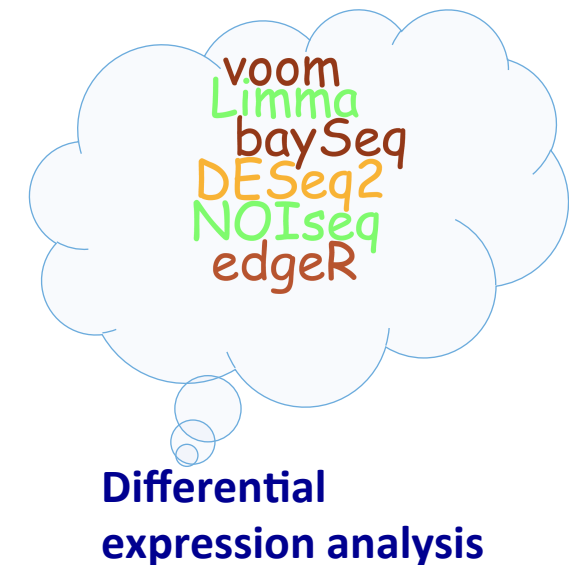
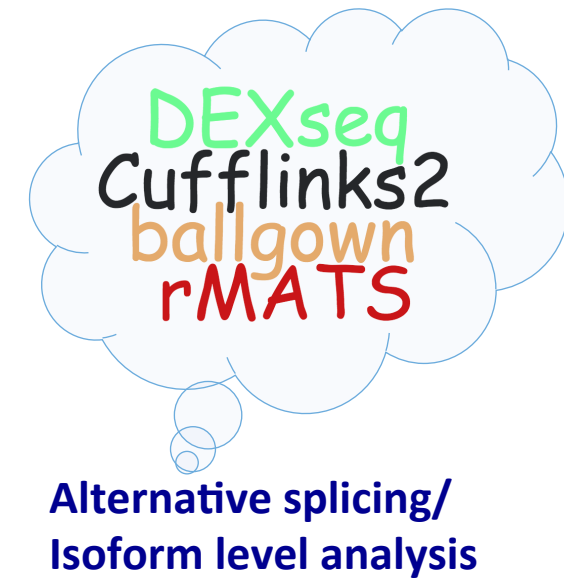
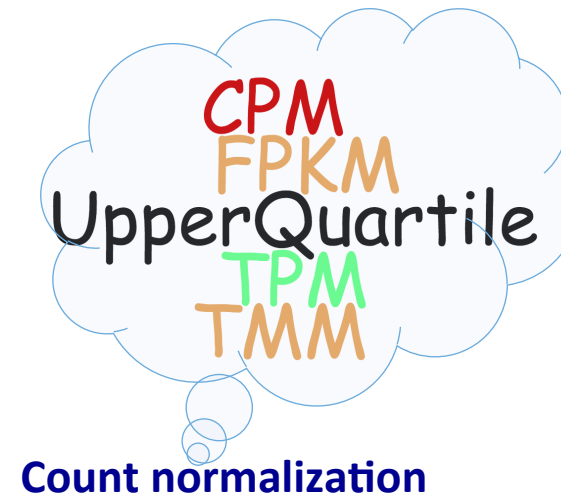
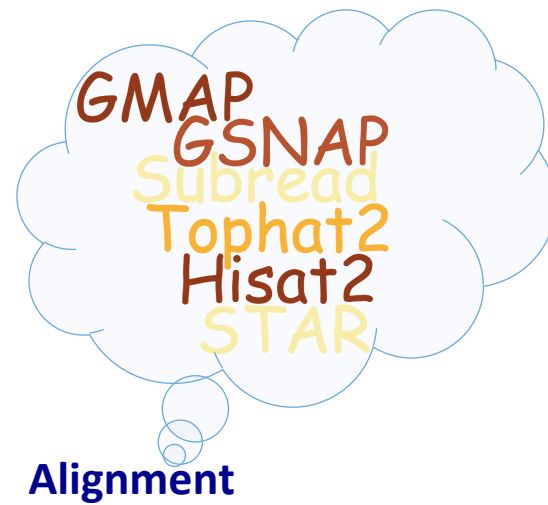
# Experiment Setup: Randomization and Blocking





# Subjectivity of the analysis

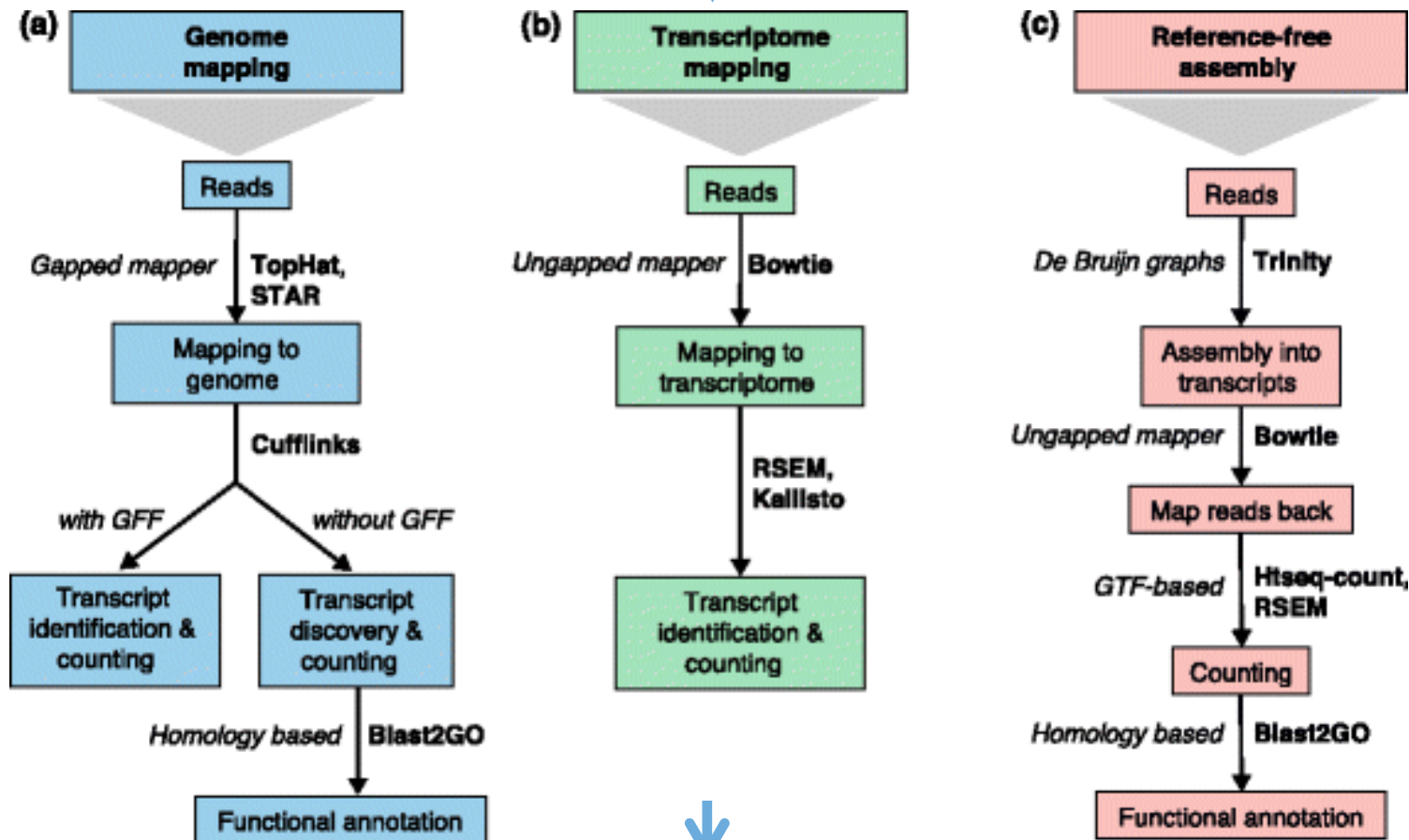
- ✓ Multitude of algorithms and pipelines available.
- ✓ Most approaches correct, but have to be tailored to the needs of the investigators in order to better capture the desired effect.



# Data Processing and Analysis

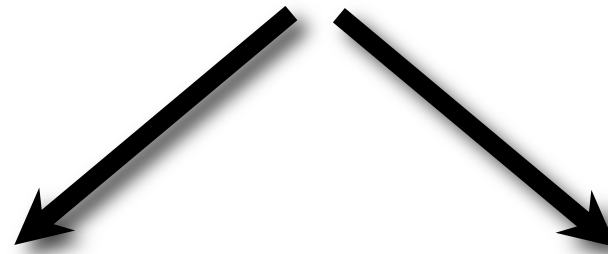
# Workflows for mapping/assembly, transcript identification and quantification

Demultiplex, filter, and trim sequencing reads



Perform statistical analysis to identify differential expression/splicing

## When to use each?

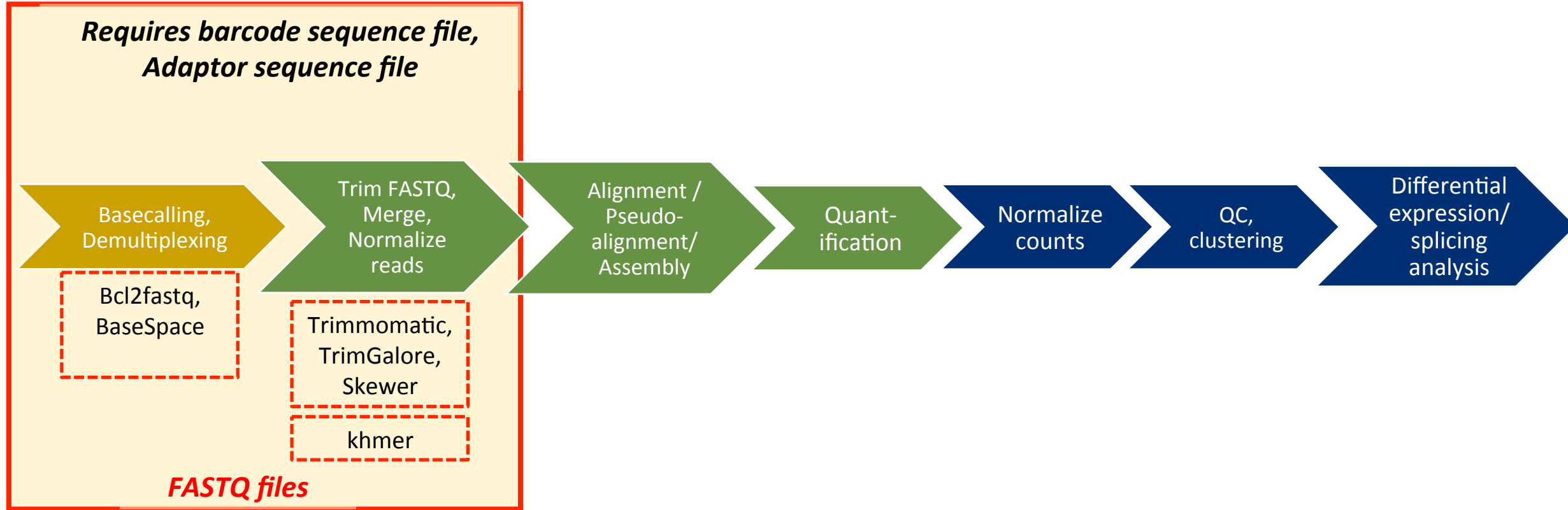


### *de novo assembly*

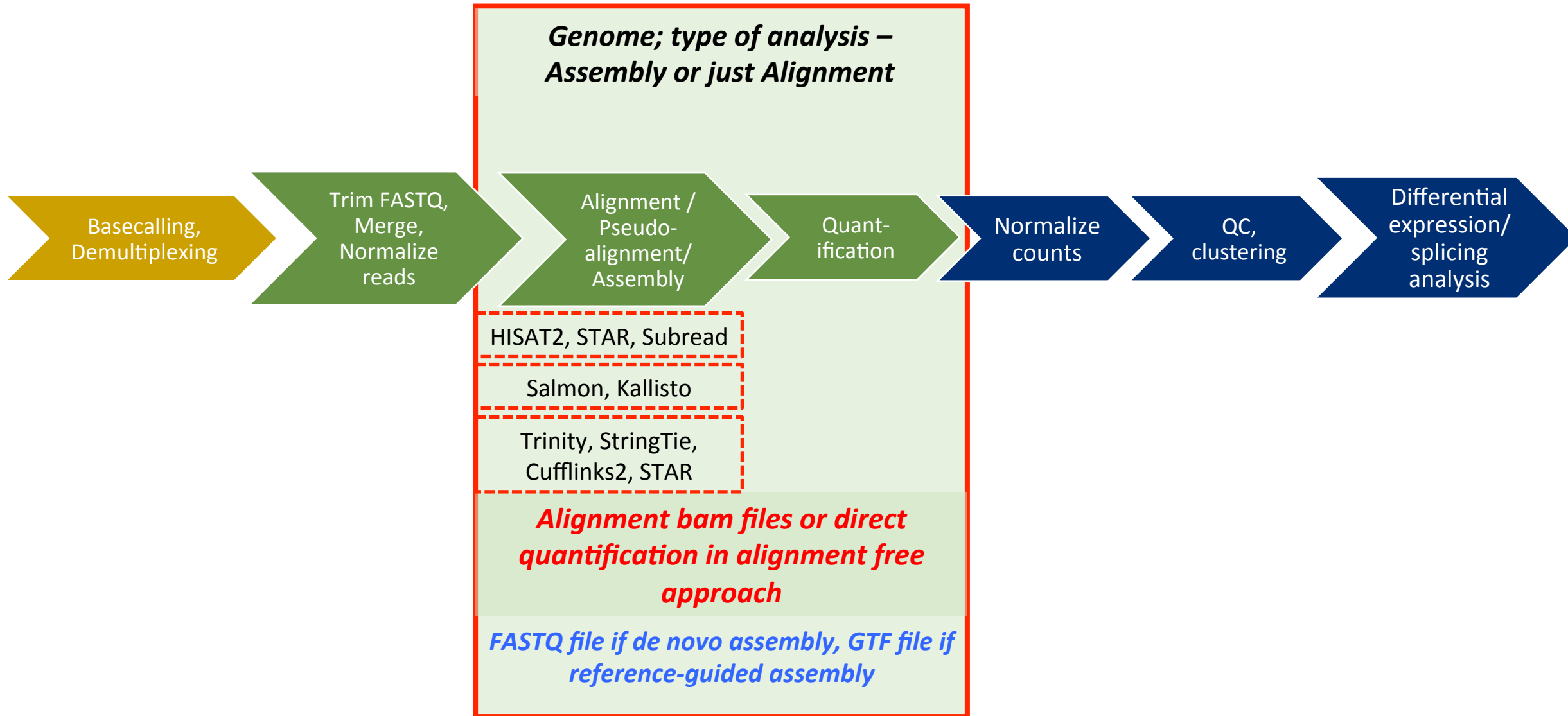
(do **not** know the transcriptome)  
(main goal is to **discover** NOT to quantify)

### reference

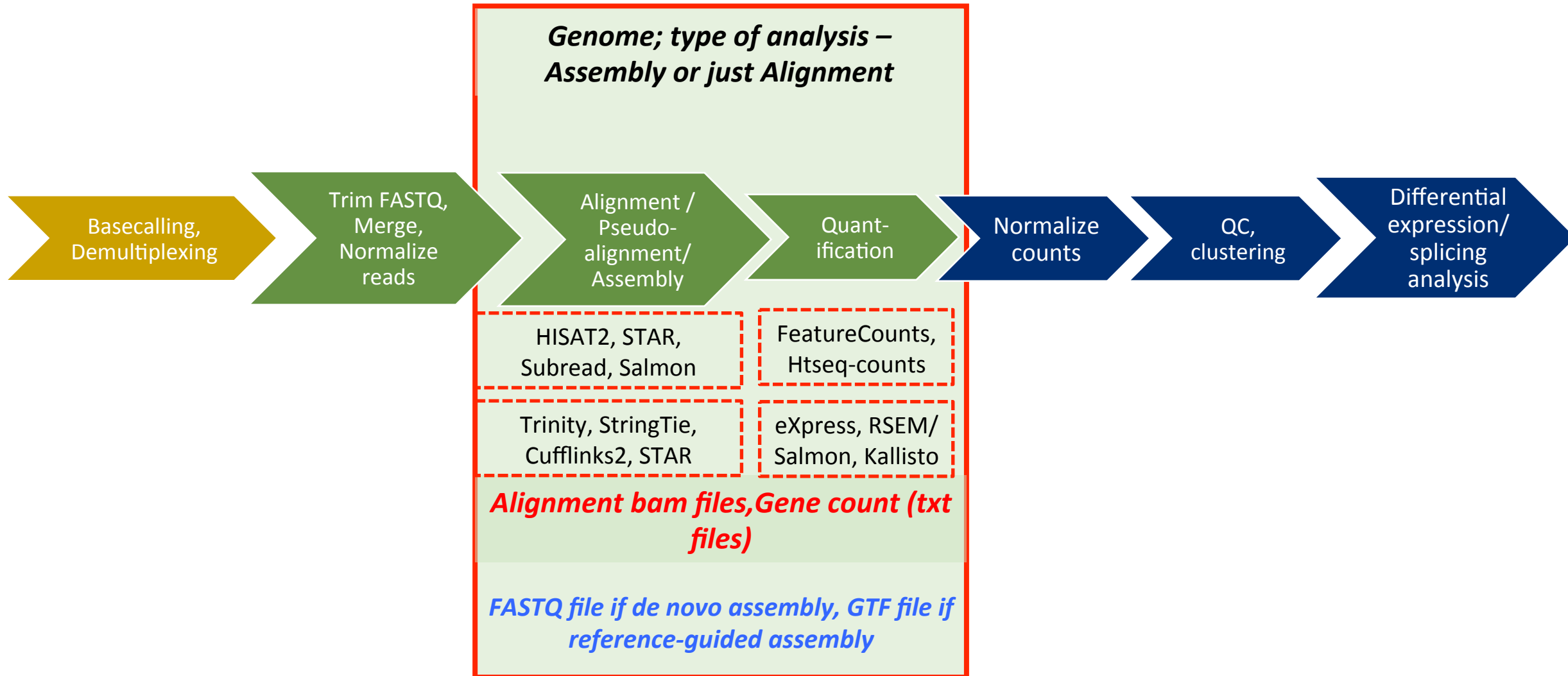
(do know the transcriptome)  
(main goal is to **quantify** NOT to discover)



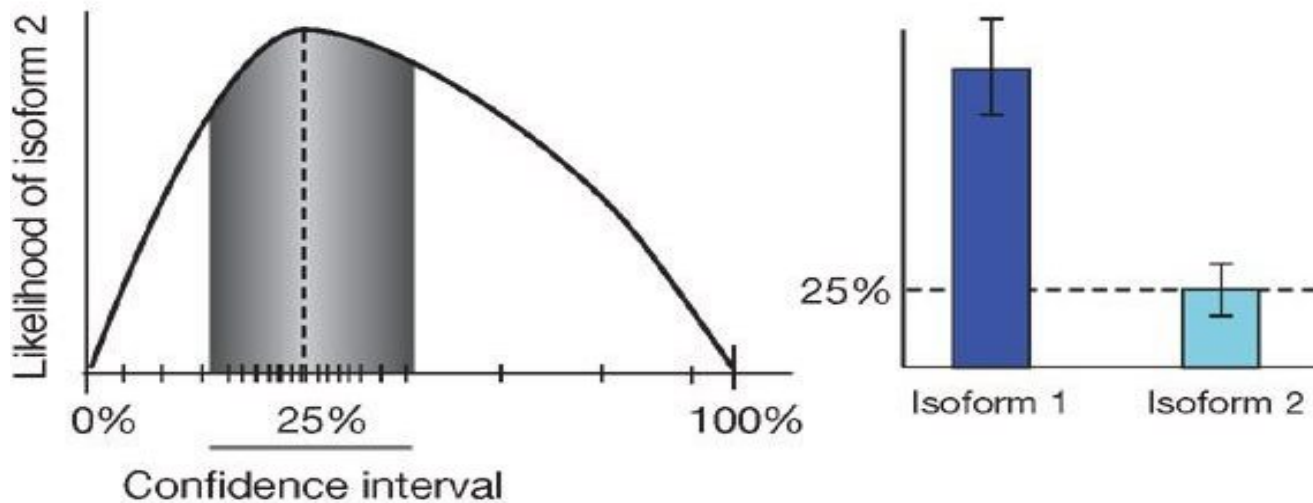
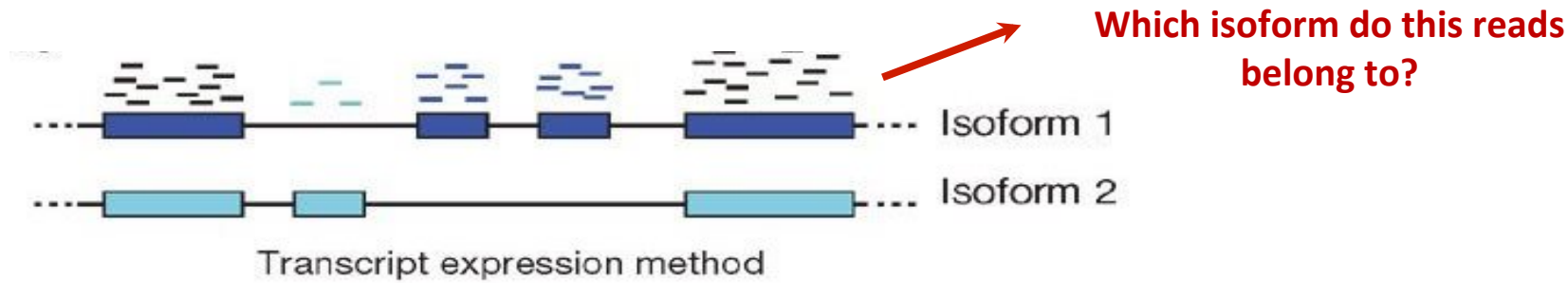
# Alignment/assembly+alignment, quantification



# Read processing, alignment/assembly+alignment, quantification

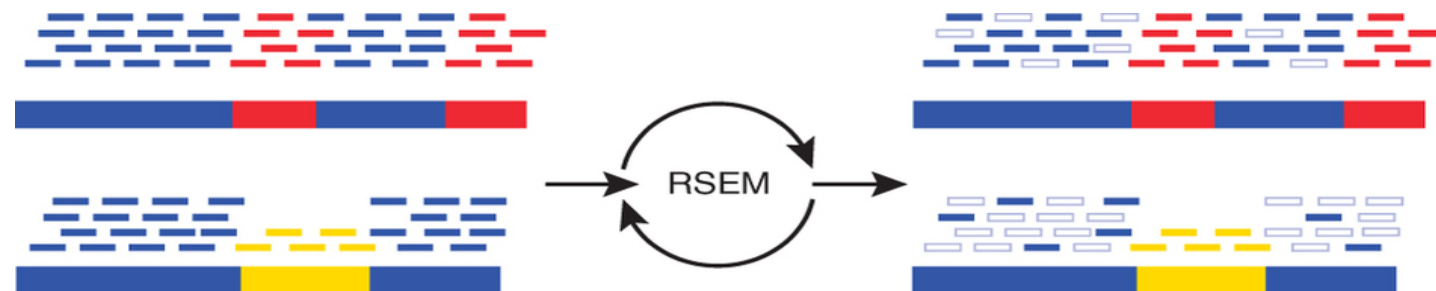


# Expression quantification

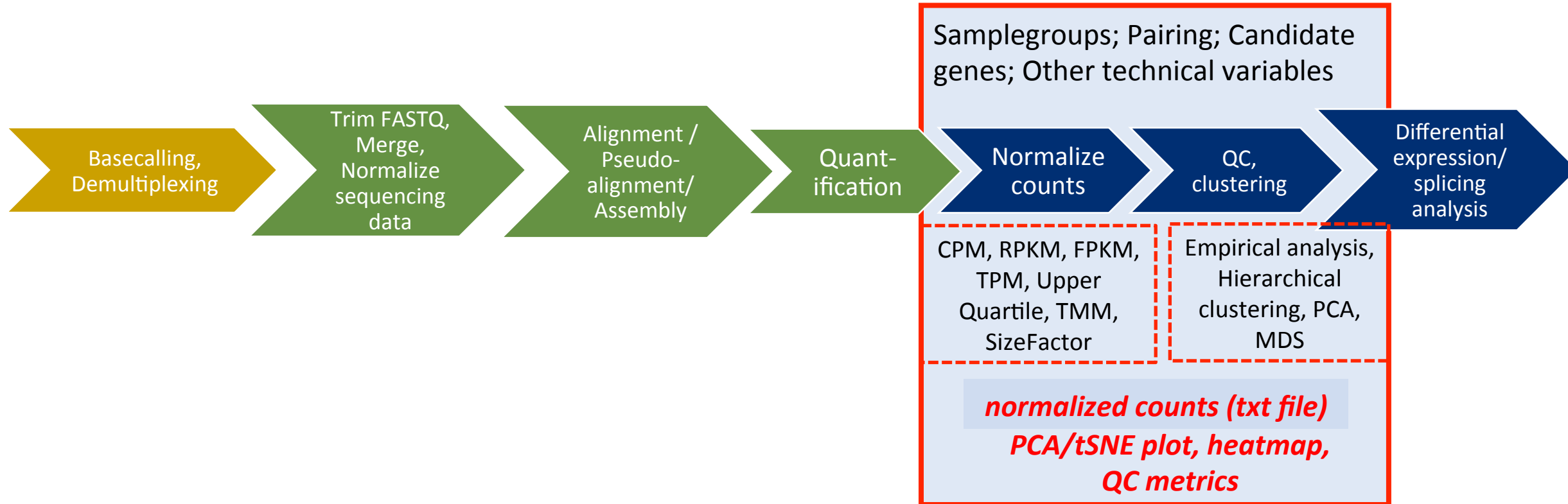


**Genome guided: Cufflinks2, StringTie**

**Transcriptome guided: RSEM, eXpress, Salmon, Kallisto**

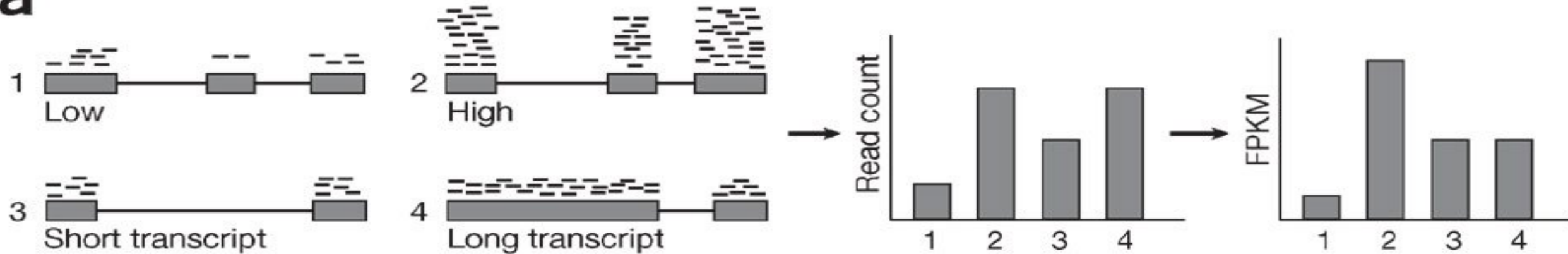


# Normalization, Expression quantification



# 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.

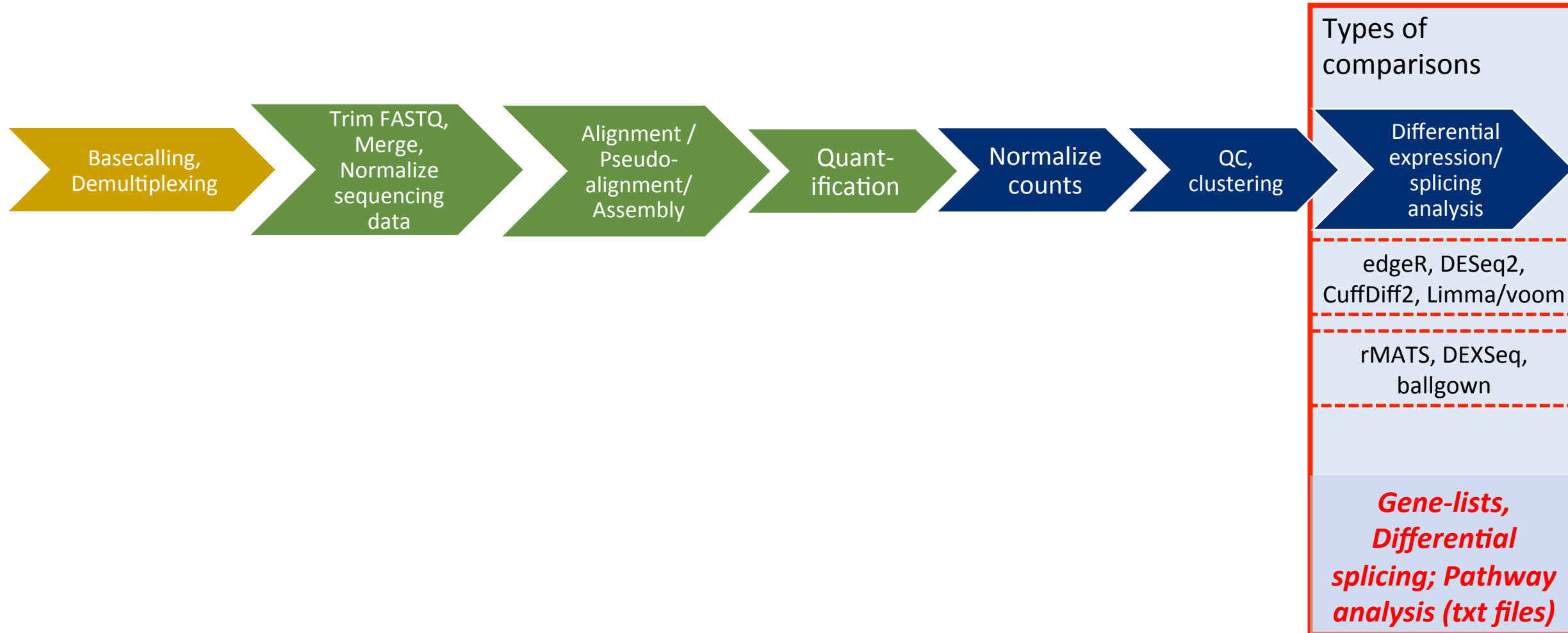
”Gene counts” should be corrected in order to minimize these biases:  
**normalization.**

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

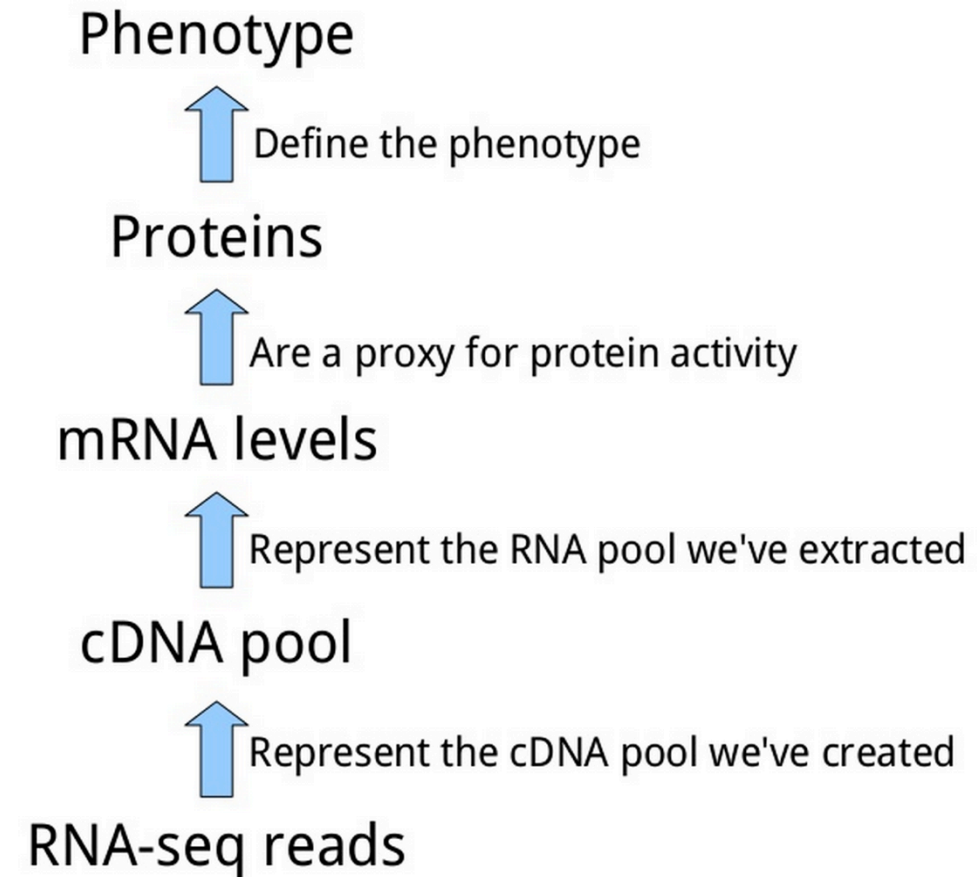
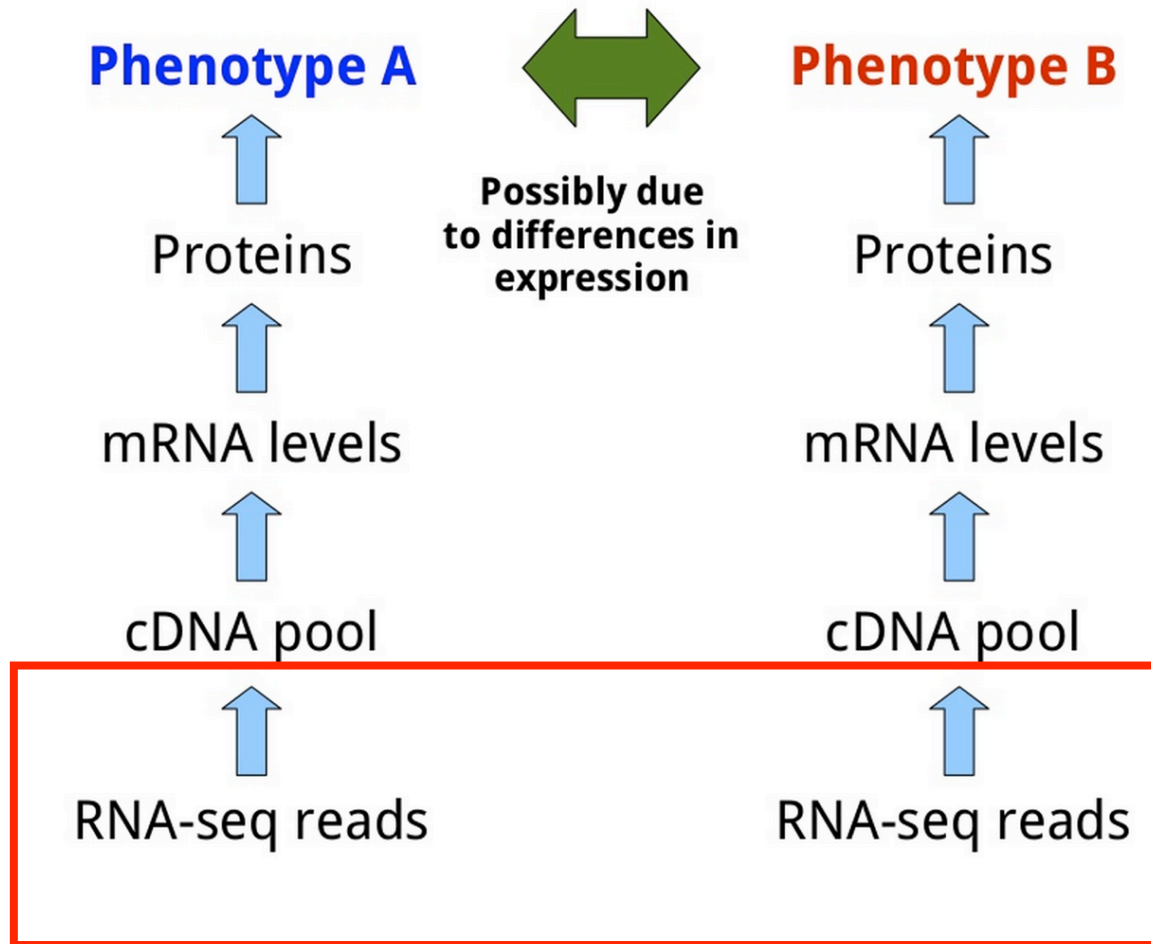
# Count normalization

	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

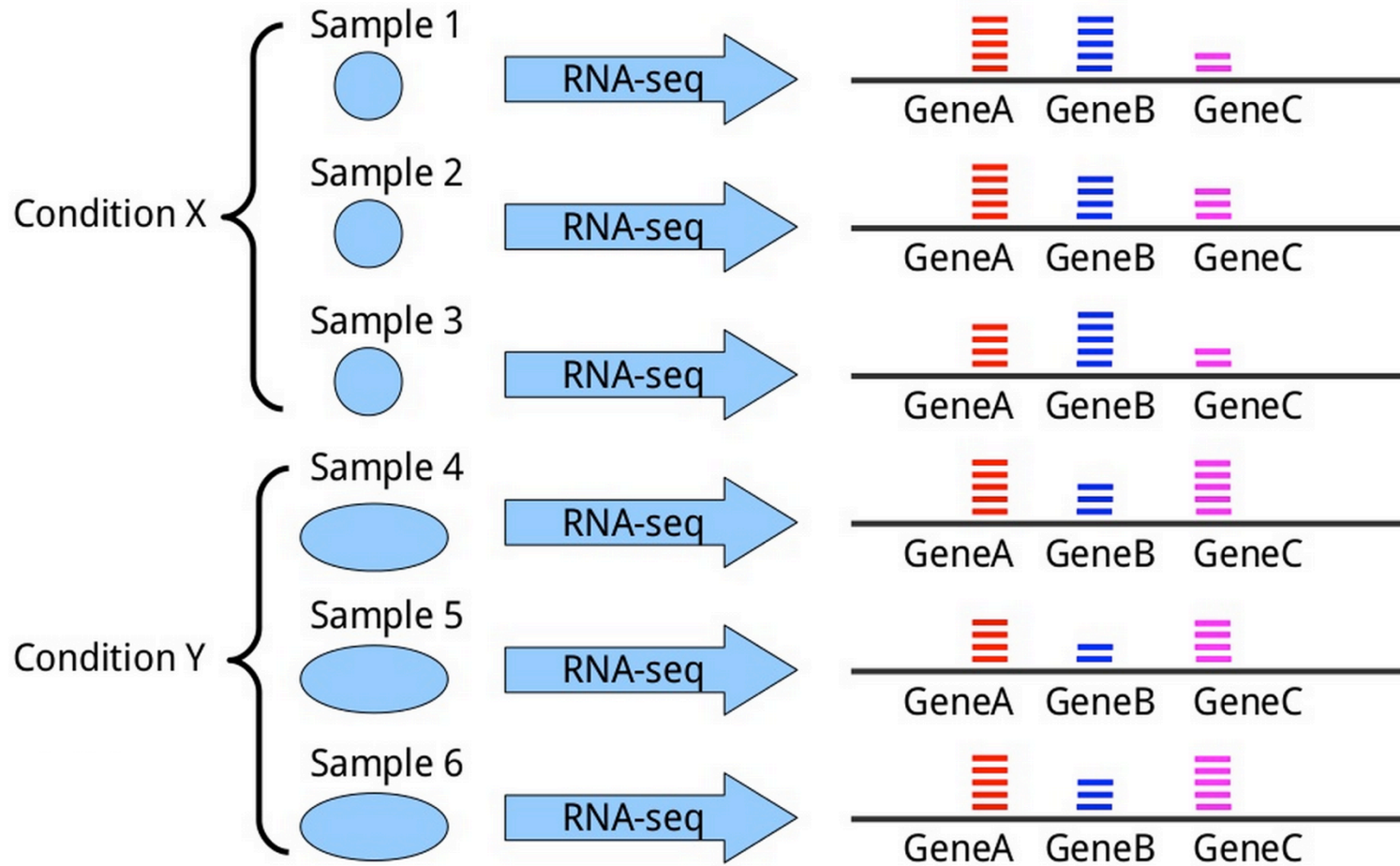


# Our assumptions and comparison



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

# Statistical testing for Differential expression



Read in raw count data

Remove genes with low-expression (<10 reads per sample across group)

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

Unsupervised clustering to identify technical effects and biological effects

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

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

Identify DE genes with pre-defined statistical criteria (~FDR < 0.05)

# How many replicates? What depth? Fold-change cut-off?

Statistical power to detect differential expression varies with effect size, sequencing depth and number of replicates

	Replicates per group		
	3	5	10
Effect size (fold change)			
1.25	17 %	25 %	44 %
1.5	43 %	64 %	91 %
2	87 %	98 %	100 %
Sequencing depth (millions of reads)			
3	19 %	29 %	52 %
10	33 %	51 %	80 %
15	38 %	57 %	85 %

***Probability of Detection of Differential Expression at 5% significance***

# What we do when we do RNAseq?

- **General introduction to RNA world**
- **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](http://en.wikipedia.org/wiki/List_of_RNA-Seq_bioinformatics_tools)
- RNA-seq analysis is as easy as 1-2-3 with limma, Glimma and edgeR.  
<https://f1000research.com/articles/5-1408/>

Thank You!



Eshita.sharma@well.ox.ac.uk