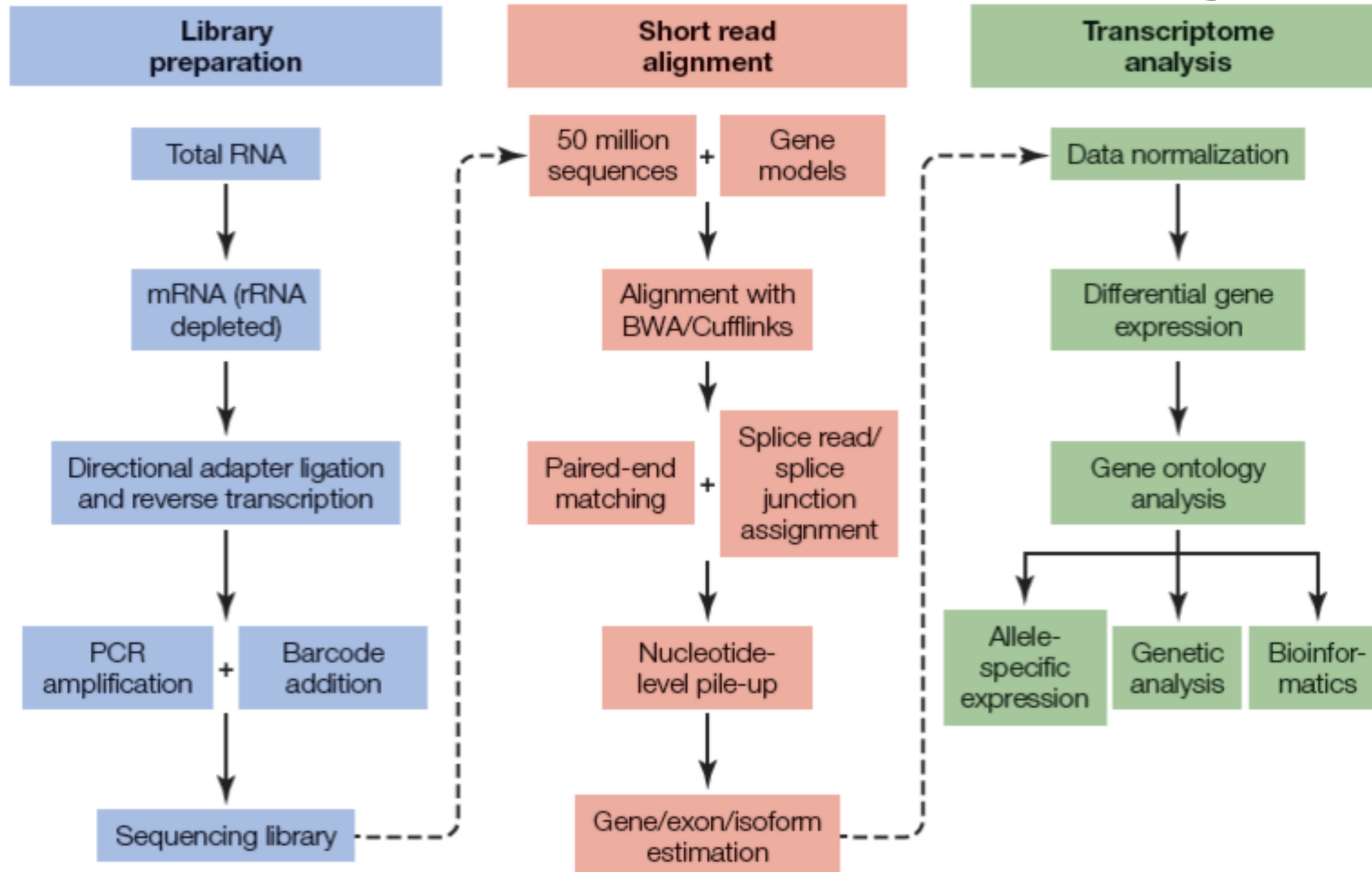# Introduction to RNAseq

## Experimental design and data analysis

Vitor Pavinato

# RNAseq in summary

# Pieces you'll need to put together

**Experimental Design**

Biological vs Technical Replicates

How many (sample size)?

Desired sequencing yield (coverage)

How many lanes?

Hoe many samples / lane?

Single-end vs Paired-end sequencing

Sequence length (75, 100, 300 bp)?

# Experimental design

**Technical vs biological replicates**

Biological replicates are necessary because there is a large variance between samples in the same condition.

Technical replicates are necessary if there were factors during sample and library preparation and sequencing that may increase the variance among replicates

# Examples of biological and technical replicates

Often you will have a fixed budget that constrains how many arrays can be processed. So your first task is to determine what levels of replication you can afford, and how they will impact statistical power.

Technical Replication:

  - RNA preparation (eg. from adjacent biopsies)
  - cDNA synthesis (pooling minimizes outlier effects)
  - library preparation
  - sequencing lane or array hybridization (usually a minimal effect)


Biological Replication:

Fixed effects:
  - sex
  - treatment (drug, growth regimen, tissue)
  - time of sampling (repeated measures in some cases)
  - genotype (IF specifically chosen and resampled)

Random effects -  individual from a population
  - field plot

# Contrast of interest

At the same time, you need to be aware of the contrasts you wish to make since by tweaking the design you may gain a lot in terms of what you can infer.

Suppose you want to compare B cells and T cells from Healthy controls and COVID-19 patients, and you have the funds to generate 24 RNASeq profiles
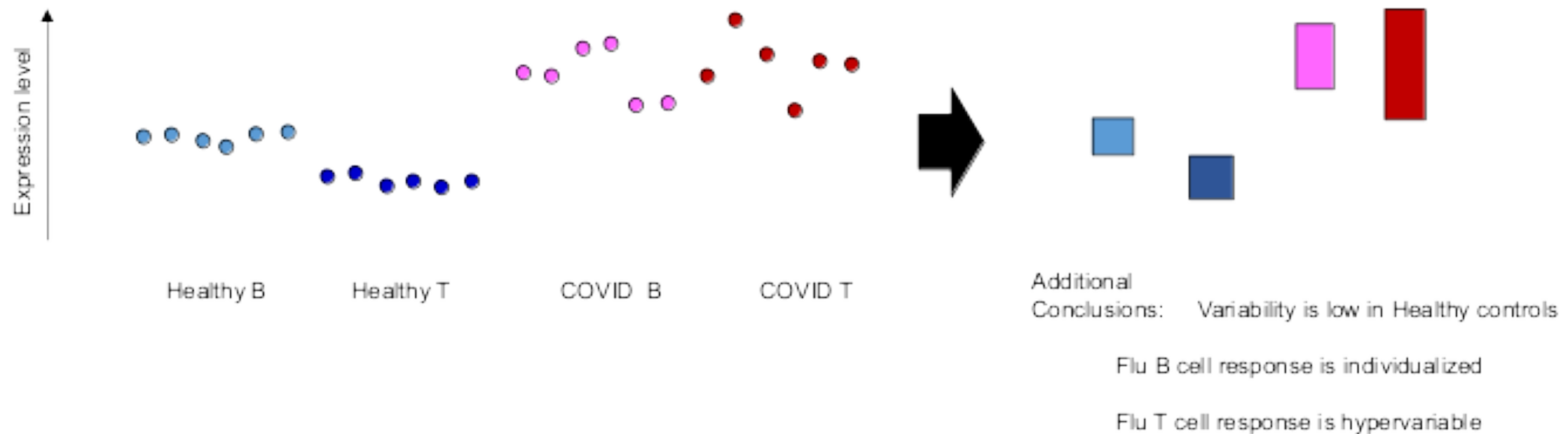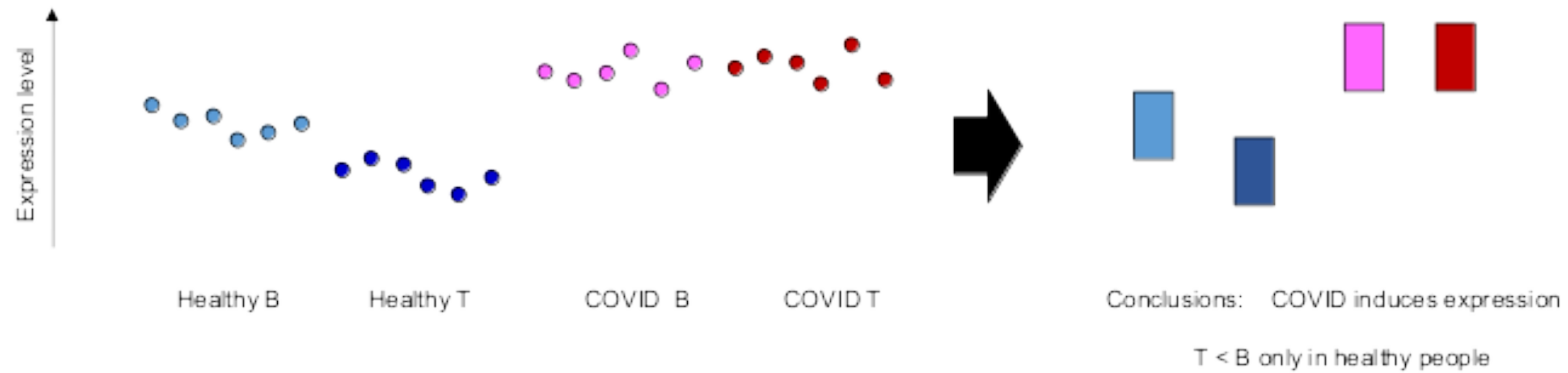
What is the best design?

- 6 controls and 6 patients, each donating both a B and a T cell sample
- 12 controls and 12 patients, each donating either a B or a T cell sample
- 3 controls and 3 patients, each donating a B and a T cell sample, processed twice
- 3 controls and 3 patients, each donating 2 B and 2 T cell samples, on separate days
- same as above, but only men or only women
- 12 controls and 12 patients, each donating either a B or a T cell sample, but pooling two visits

Main effects can only be contrasted if you have biological replicates:
        reducing the number of individuals may allow you to address intra-individual variability

Interaction effects allow you to ask questions like whether B cells and T cells differ more between healthy volunteers or patients

Two Hypothetical Sets of Results Illustrating Design Principles

https://cqs-vumc.shinyapps.io/rnaseqsamplesizeweb/

# Pieces you'll need to put together

**How to process your data**

Is there a reference genome that I can use?

If not, what can I do?

If I have a reference genome and transcriptome,
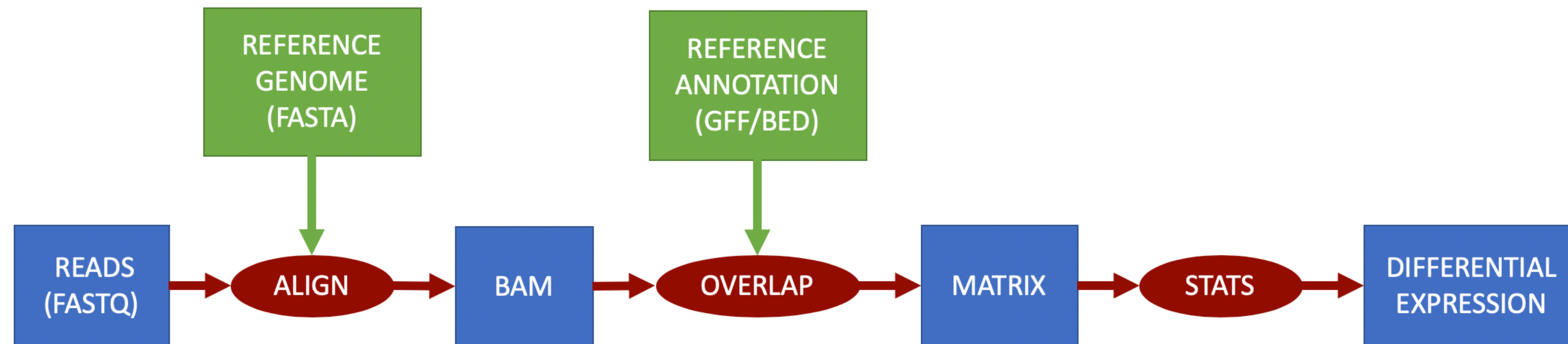Which approach is the best?

# DEG analysis

```
   name        control      shock        fold_change    pvalue
Gene A         100          200              2.0         0.000035
Gene B          80           60              0.75        0.234
Gene C         120          180              1.5         0.013
...
```

**How do I get this?**

# Methods to quantify expression

**Figure 4**

An illustration of spliced alignment of RNA sequencing (RNA-seq) fragments to a genome (*a*) and direct alignment to a transcriptome (*b*). Reads are designated by thick solid lines, while dashed arcs represent the pairing relationship between paired-end reads. This illustration depicts alignment to a single four-exon gene consisting of three distinct transcripts. In the spliced alignment (*a*), the left read of the rightmost pair is a junction-spanning alignment to the red–green exon boundary. In the direct alignment to the transcriptome (*b*), one observes how the same alignment (e.g., the alignment to the blue exon) is repeated for each transcript.

# What are gene isoforms



## Alternative splicing

# How the data looks like

Example FASTQ file



FASTQ scoring

## Example FASTA file (Genome)



## Example GFF file (Annotation)

# Classification-based

# Transcripts or genes

## Features

**Reference-based**

**Genes (IDs)**

**Classification-based**

**Transcripts (isoforms)**

# Pipeline: from reads to candidate genes



Fig. 4 Bioinformatics tools commonly used in RNA-seq data analysis. These tools are primarily used in the four main processes of RNA-seq data analysis, including quality control, read alignment and transcript assembly, expression quantification and differential expression analysis

# Pieces you'll need to put together

## Normalization



Observed gene expression = Biological source + Technical source + Stochastic

# Approaches to Normalization

Mean or Median transform, simply centers the distribution
- Something like this is essential to control for overall distributional effects (eg RNA concentration)

Variance transforms, such as standardization or inter-quartile range
- Depends on whether you think the overall distributions should have similar variance

Quantile normalization
- Transforms the ranks to the average expression value for each rank

Gene-level model fitting
- Remove technical or biological effects before model fitting on the residuals

Supervised normalization
- Optimally estimate the biological effect while fitting technical factors across the entire experiment

# DEG analysis

| name | control | shock |
|------|---------|-------|
| Gene A | 100 | 200 |
| Gene B | 80 | 60 |
| Gene C | 120 | 180 |
| ... | | |

| fold_change | pvalue |
|-------------|--------|
| 2.0 | 0.000035 |
| 0.75 | 0.234 |
| 1.5 | 0.013 |

**How do we get this?**

**How about this part?**

# Pieces you'll need to put together

**Statistical analysis and hypothesis testing**

Two-conditions:
T-test

$$t = \frac{\bar{X}_1 - \bar{X}_2}{s_p \sqrt{2/n}} \qquad s_p = \sqrt{\frac{s_{X_1}^2 + s_{X_2}^2}{2}}$$

Between

Within

# Any questions?

## Multiple testing and FDR?

FDR-controlling procedures are designed to control the FDR, which is the expected proportion of "discoveries" (rejected null hypotheses) that are false (incorrect rejections of the null). Equivalently, the FDR is the expected ratio of the number of false positive classifications (false discoveries) to the total number of positive classifications (rejections of the null). The total number of rejections of the null include both the number of false positives (FP) and true positives (TP).

https://en.wikipedia.org/wiki/False_discovery_rate

| Primary category | Tool name | Notes |
|---|---|---|
| Splice-aware read alignment | GEM | Filtration-based approach to approximate string matching for alignment |
| | GSNAP | Based on seed and extend alignment algorithm aware of complex variants |
| | MapSplice | Based on Burrows-Wheeler Transform (BWT) algorithm |
| | RUM | Integrates alignment tools Blat and Bowtie to increase accuracy |
| | STAR | Based on seed searching in an uncompressed suffix arrays followed by seed clustering and stitching procedure; fast but memory-intensive |
| | TopHat | Uses Bowtie, based on BWT, to align reads; resolves spliced reads using exons by split read mapping |
| Transcript assembly and quantification | Cufflinks | Assembles transcripts to reference annotations or de novo and quantifies abundance |
| | FluxCapacitor | Quantifies transcripts using reference annotations |
| | iReckon | Models novel isoforms and estimates their abundance |
| Differential expression (DE) | BaySeq | Count-based approach using empirical Bayesian method to estimate posterior likelihoods |
| | Cuffdiff2 | Isoform-based approach based on beta negative binomial distribution |
| | DESeq | Exon-based approach using the negative binomial model |
| | DEGSeq | Isoform-based approach using the Poisson model |
| | EdgeR | Count-based approach using empirical Bayes method based on the negative binomial model |
| | MISO | Isoform-based model using Bayes factors to estimate posterior probabilities |
| Other tools | HCP | Normalizes expression data by inferring known and hidden factors with prior knowledge |
| | PEER | Normalizes expression data by inferring known and hidden factors using a probabilistic estimation based on the Bayesian framework |
| | Matrix eQTL | Fast eQTL detection tool that uses linear models (linear regression or ANOVA) |