# RNAseq analysis
## *-it's complicated*

### March 2017

# RNA reads are not enough to identify functional RNAs



Defining functional DNA elements in the human genome
**Kellis M et al. PNAS 2014;111:6131-6138**

# Depending on the steps from sample to RNA seq will give different results

RNA->

enrichments ->

extraction of poly-A RNAs

conversion into ds-cDNA and shearing

amplification and adapter ligation

library ->

sequencing

single end (SET)

paired-end (PET)

reads ->

AAAAAAAA

AAAAAAA
AAAAA
AAAAA

TTTTT
AAAAA

A
T

AAAAAA
TTTTTT

AA
TT

AAA
TTT

A
T

AAAAAA
TTTTTT

TTTTT
AAAAA

AA
TT

AAA
TTT

A

A

---

PolyA          (mRNA)
RiboMinus    (- rRNA)
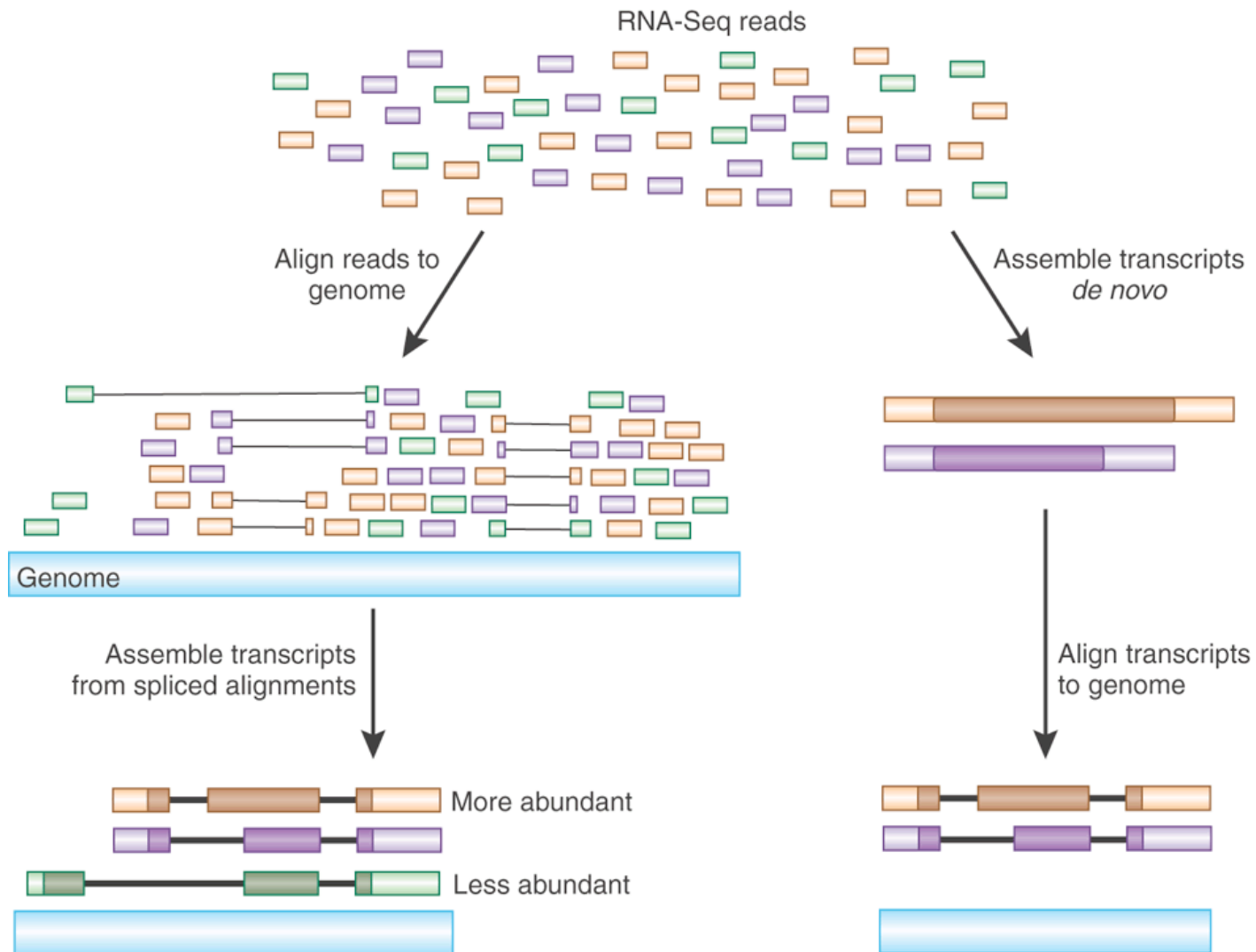Size  <50 nt    (miRNA )
…..

---

Size of fragment
Strand specific
5' end specific
3' end specific
…..

---

Single end (1 read per fragment)
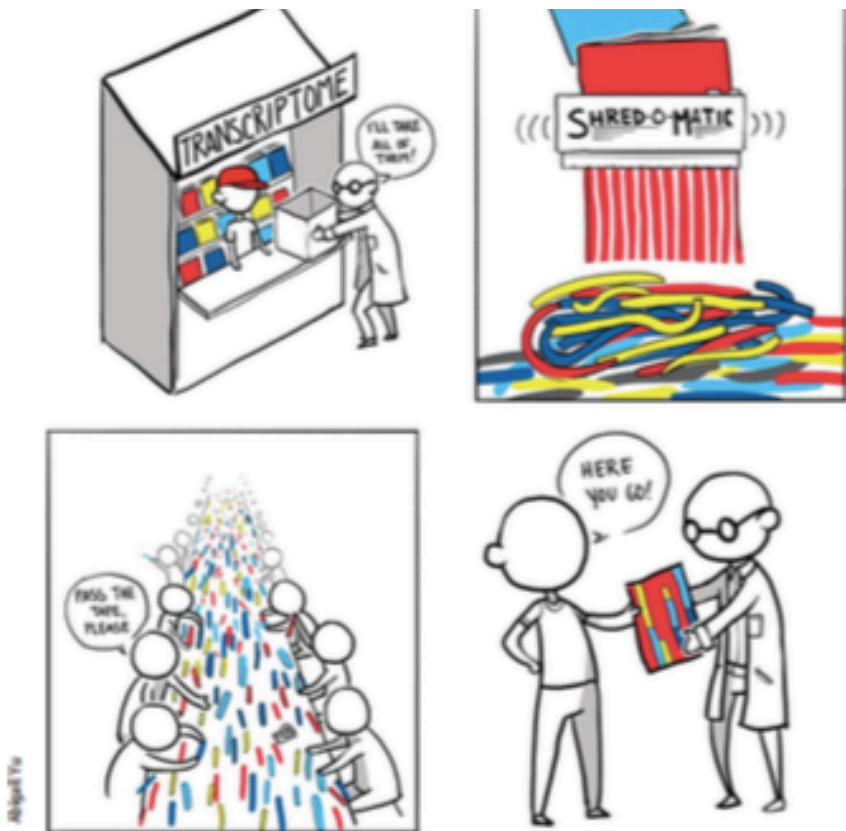Paired end (2 reads per fragment)

# Mapping (Pär Engström)

- Use RNA specific mapper

- Use a two-pass workflow

- STAR or HISAT

- For long (PacBio) reads, STAR, BLAT or GMAP can be used

# Gene and Isoform detection

# Long reads might be the way to go



"The way we do RNA-seq now is... you take the transcriptome, you **blow it up into pieces** and then you try to figure out **how they all go back together again**... If you think about it, it's kind of a **crazy way to do things.**"

Michael Snyder
Stanford University

Tal Nawy (2013) End-to-end RNA sequencing, *Nature Methods* 10: 1144–1145

**Figure 1** | Transcriptome reconstruction—akin to reassembling magazine articles after they have been through a paper shredder.

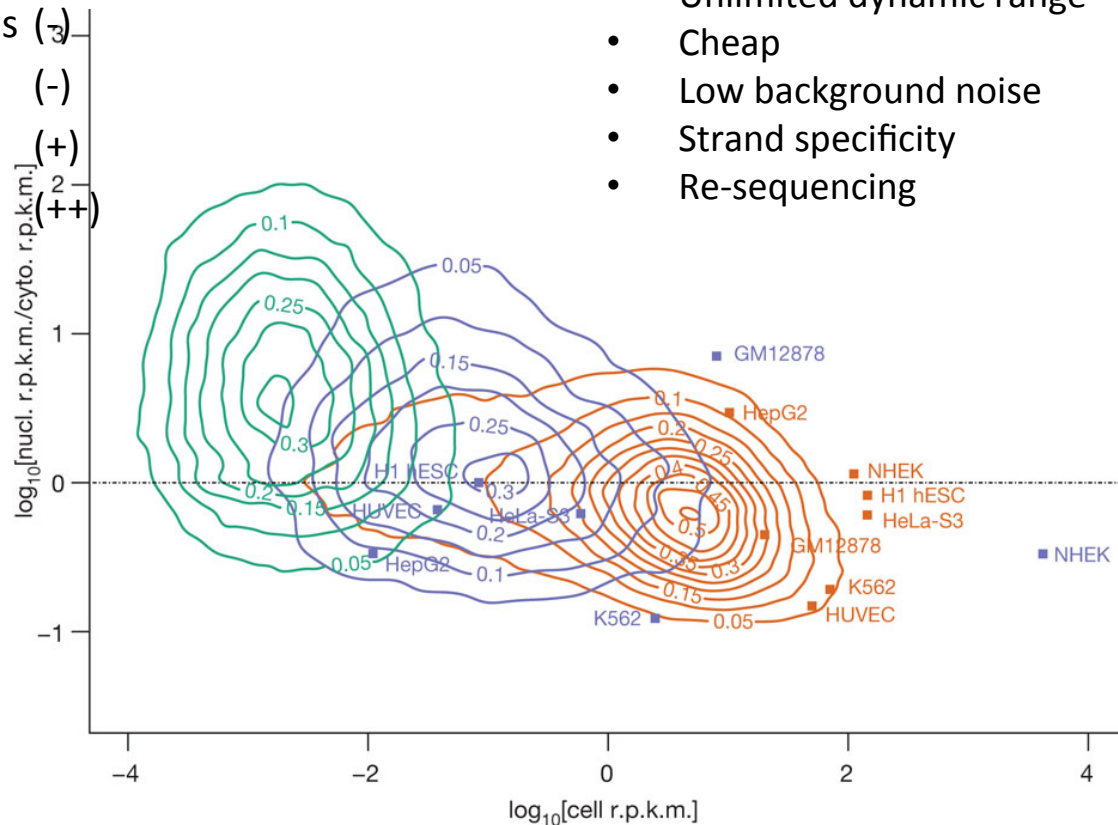Ian Korf (2013) Genomics: the state of the art in

# Promises and pitfalls

## Long reads

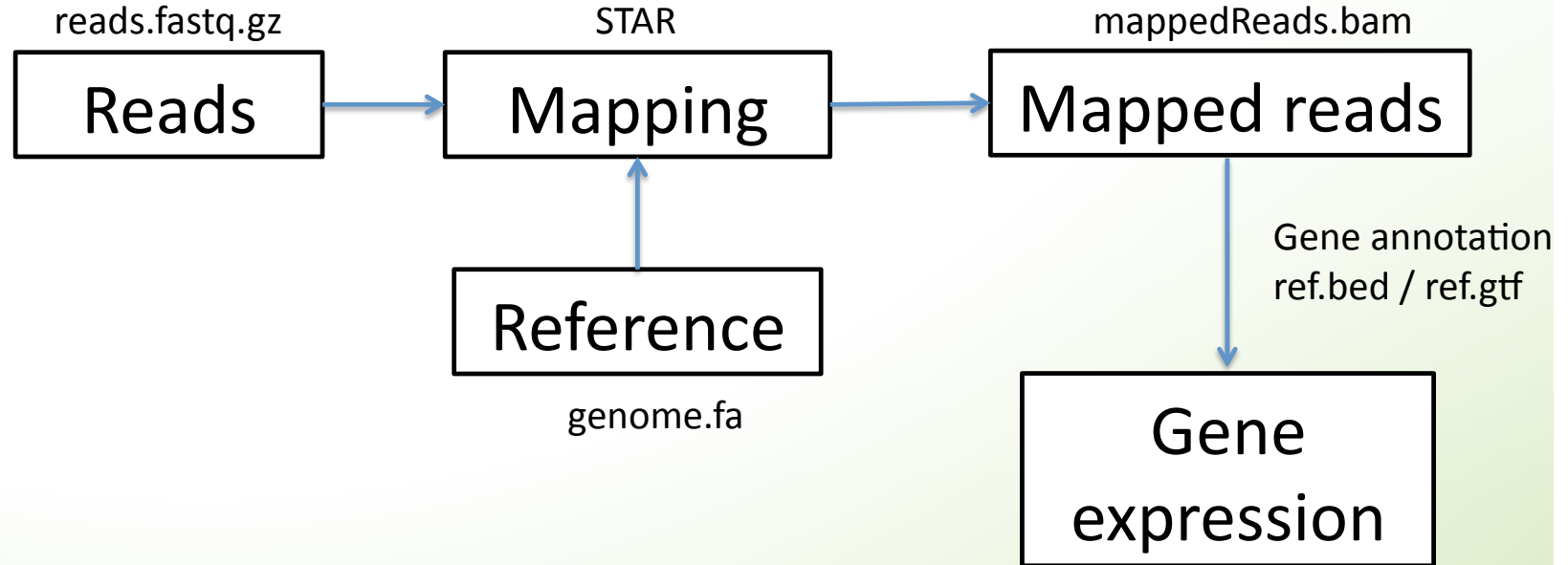- Low throughput                              (-)
- Complete transcripts                   (++)
- Not quantitative                          (-)
- Only highly expressed genes (-)
- Expensive                                     (-)
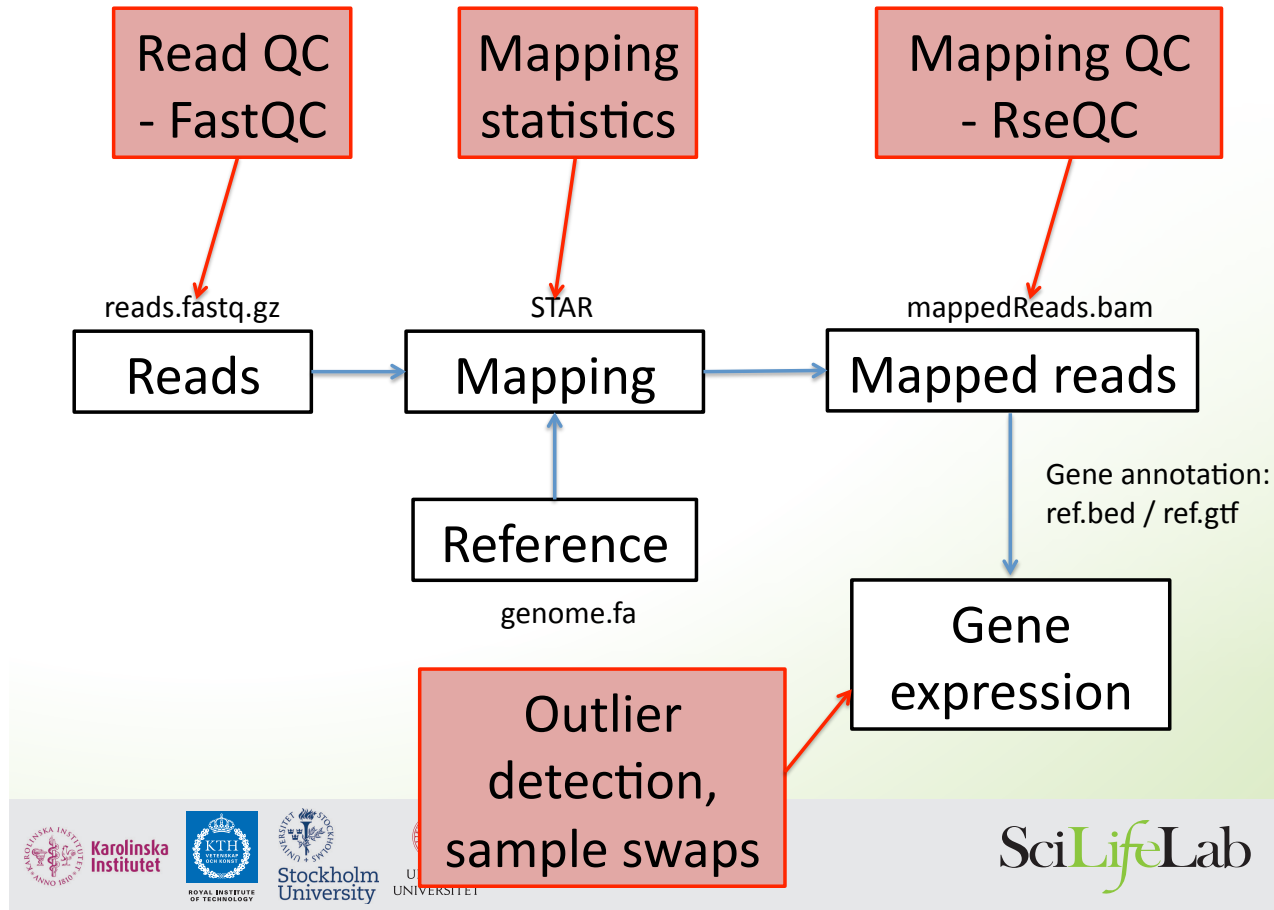- Low background noise                (+)
- Easy downstream analysis      (++)

## short reads

- High throughput                      (++)
- Quantitative                              (++)
- Fractions of transcripts            (-)
- Full dynamic range                   (+-)
- Unlimited dynamic range        (+)
- Cheap                                        (+)
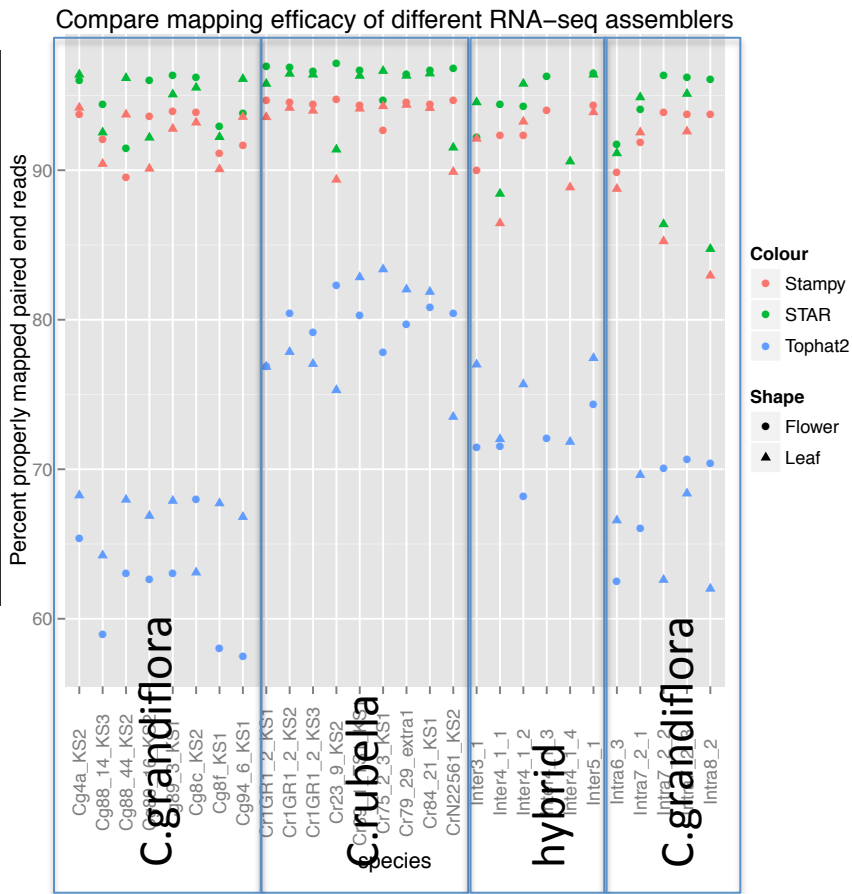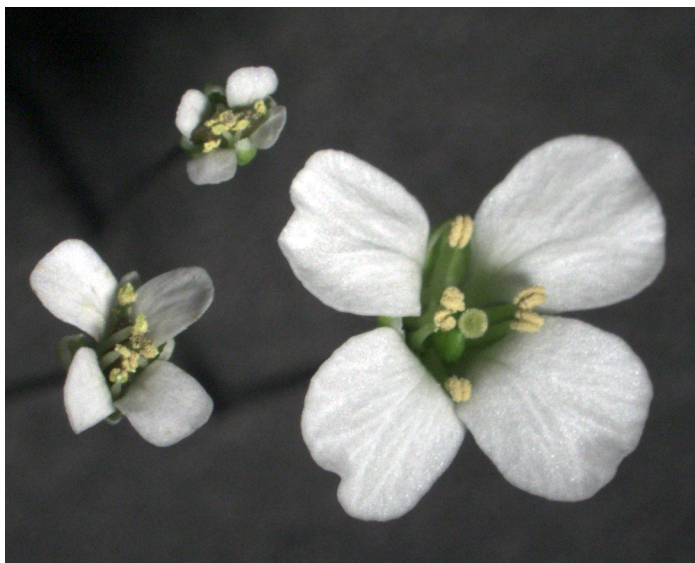- Low background noise              (+)
- Strand specificity                      (+)
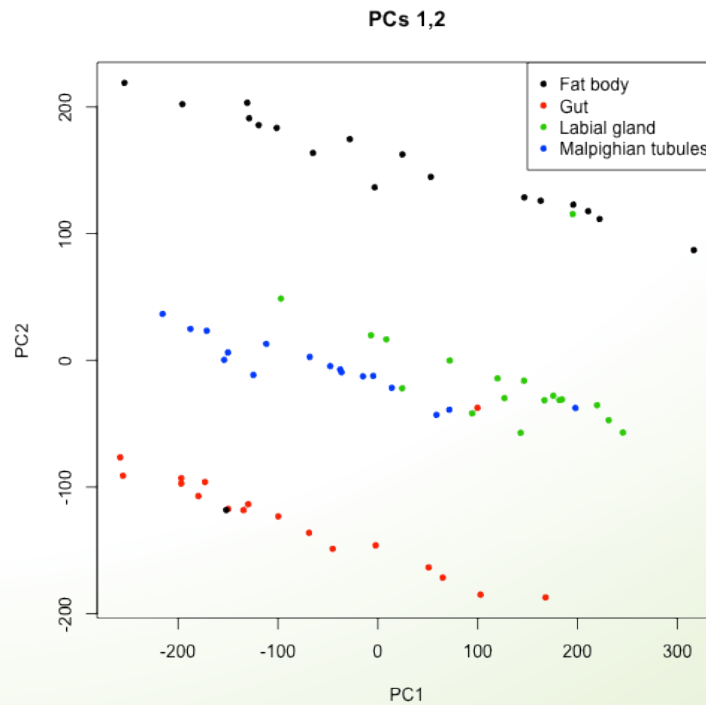- Re-sequencing                          (+)

# RNA-seq analysis workflow

# Do a lot of QC

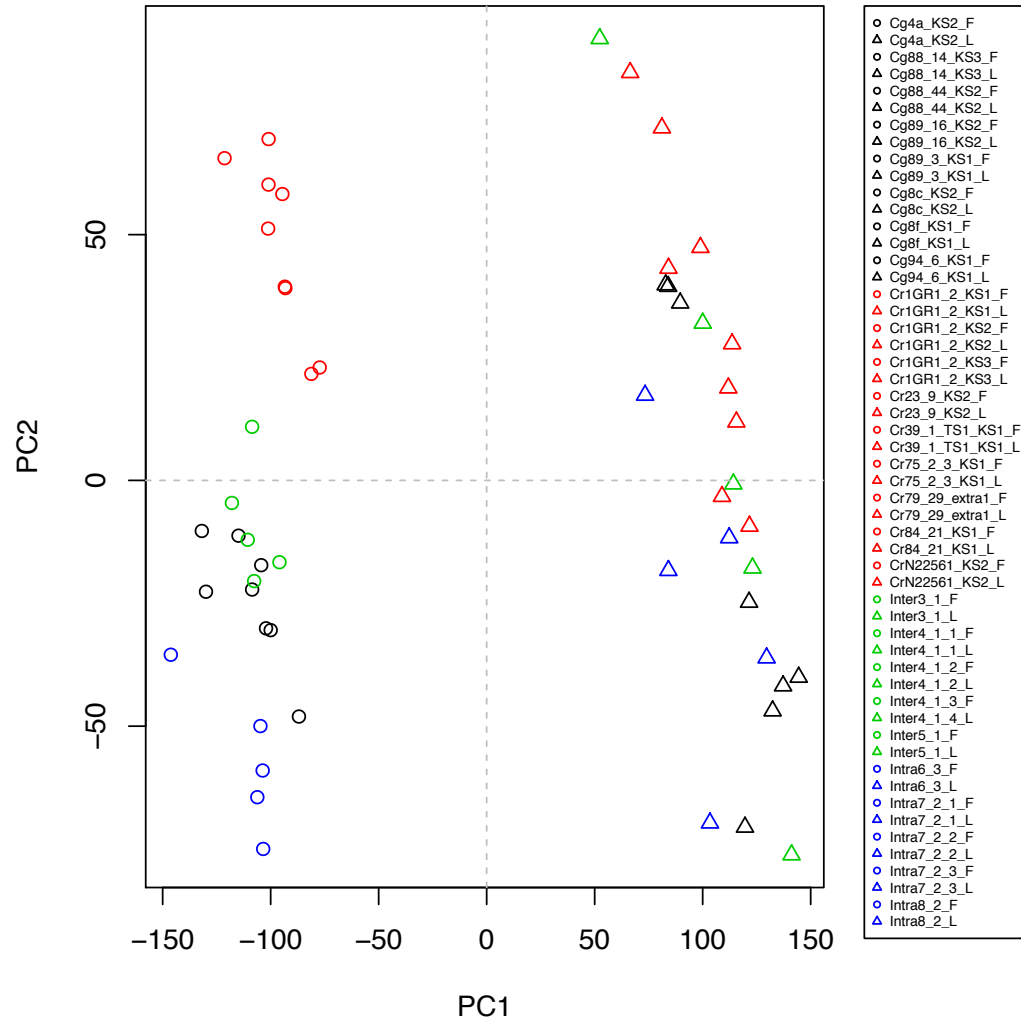# More variation when using top hat 2 with default settings than when using STAR or Stampy with default setting



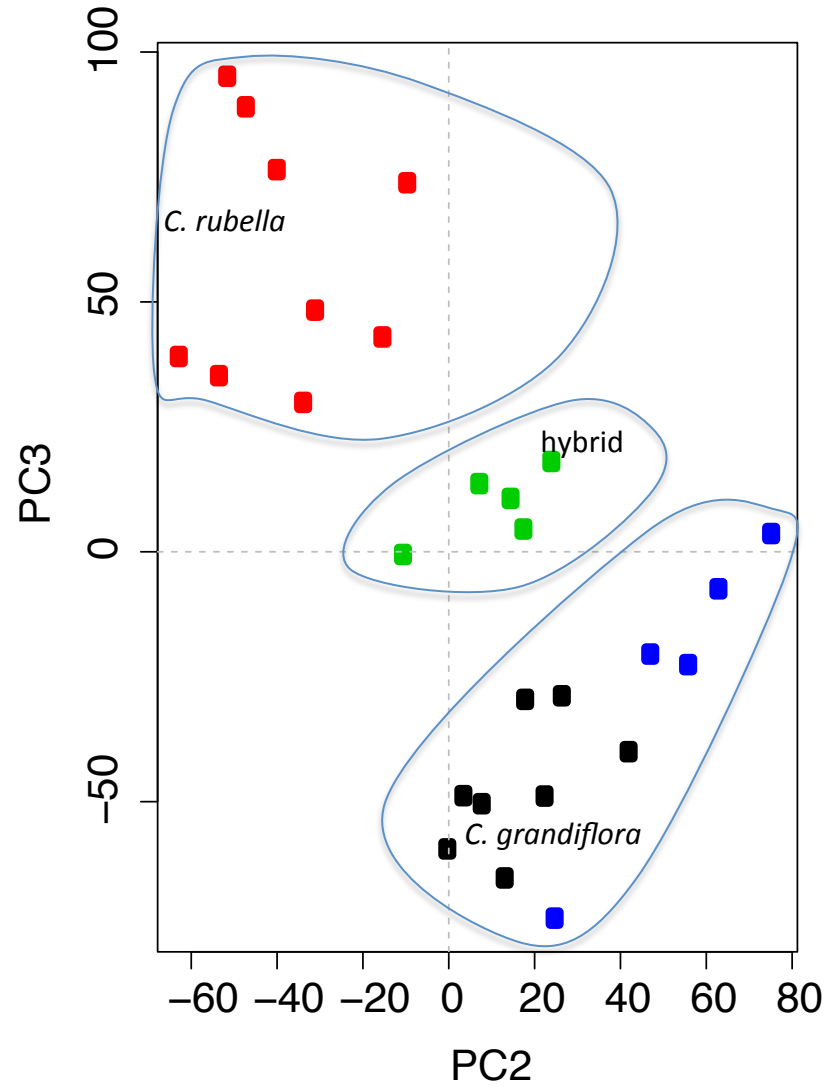Compare mapping efficacy of different RNA−seq assemblers

# RNA QC

## PCA analysis detected potential sample swaps

# Principal component 1 separates samples from flowers and leaves

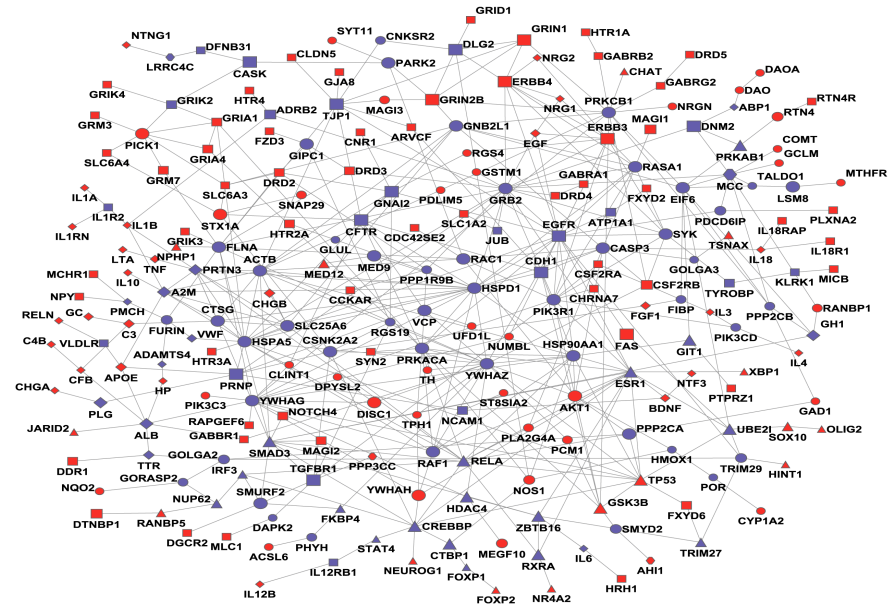# Principal component 2 and 3 separates the different species

# Differential expression analysis
# Mikael Huss

The identification of genes (or other types of genomic features, such as transcripts or exons) that are expressed in significantly different quantities in distinct groups of samples, be it biological conditions (drug-treated vs. controls), diseased vs. healthy individuals, different tissues, different stages of development, or something else.

Typically **univariate** analysis (one gene at a time) – even though we know that genes are not independent

# Decision tree for software selection (2016)

Differentially expressed **exons** => *DEXSeq*  
                                                    *Sleuth*

Differentially expressed **isoforms** => *BitSeq, ~~Cuffdiff~~ or ebSeq*

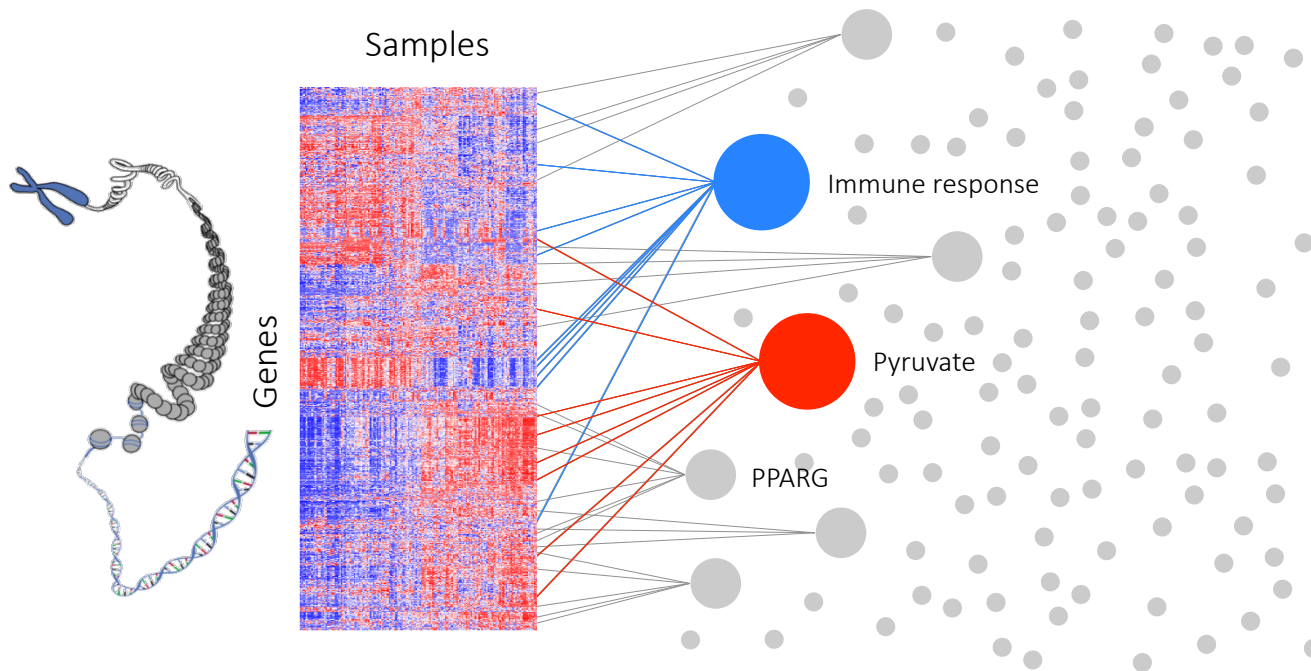Differentially expressed genes => **Select type of experimental design**

    Complex design (more than one varying factor) => *DESeq, edgeR, limma*, *Sleuth*

    Simple comparison of groups => **How many biological replicates?**

        More than about 5 biological replicates per group => ~~SAMSeq~~

        Less than 5 biological replicates per group => *DESeq, edgeR, limma*
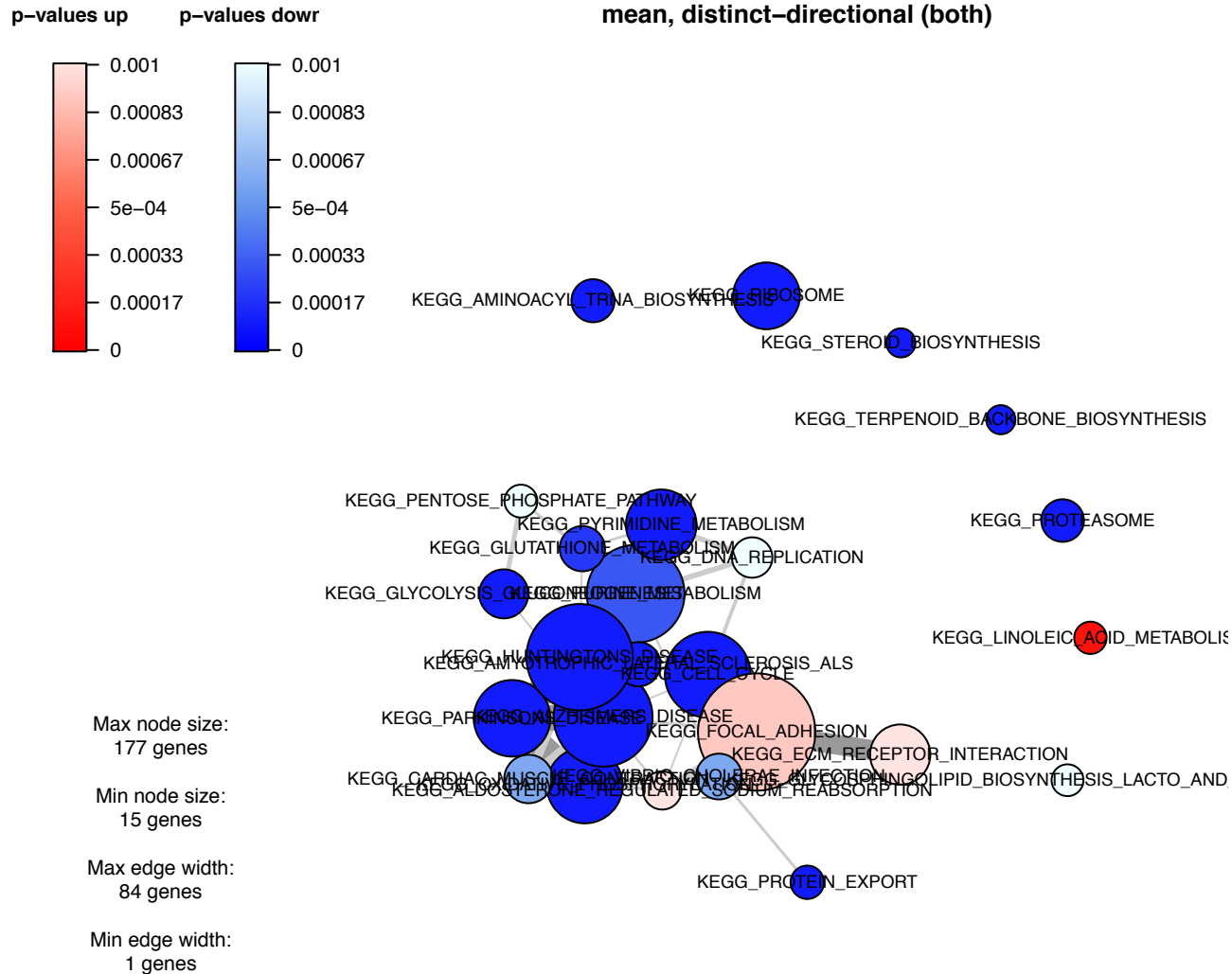
?

# Gene-set analysis (GSA)



Samples

Genes

Immune response

Pyruvate

PPARG

GO-terms
Pathways
Chromosomal locations
Transcription factors
Histone modifications
Diseases
etc...
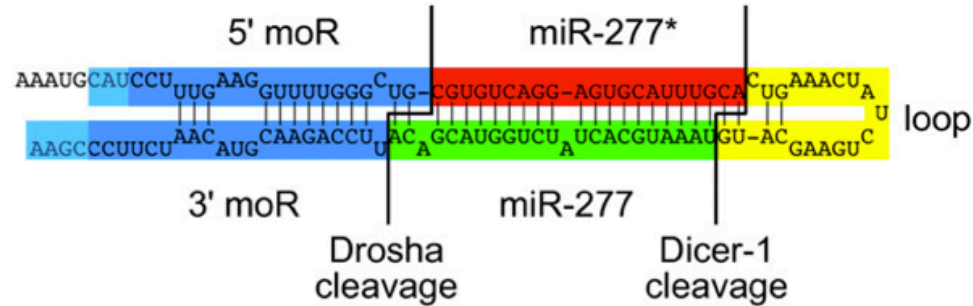
Gene-level data    → Gene-set analysis →    Gene-set data (results)

We will focus on transcriptomics and differential expression analysis
However, GSA can in principle be used on all types of genome-wide data.

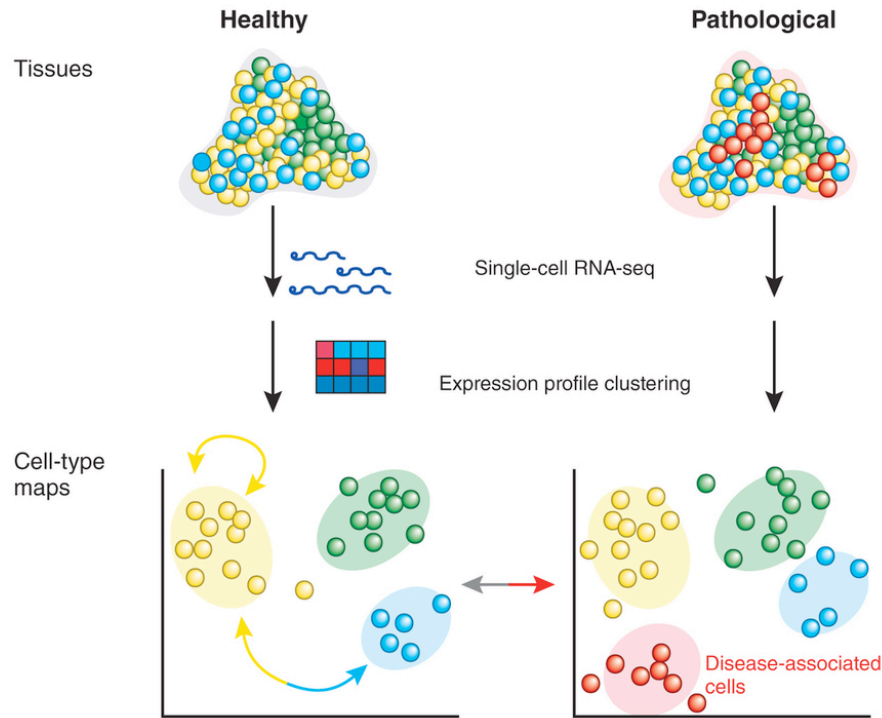# Analysis regarding Type II Diabetes

# Expression of genes on pathway

# miRNA seq analysis



(Berezikov et al. Genome Research, 2011.)

# Single cell sequencing



(Sandberg, Nature Methods 2014)

# Exercises

- Mapping
  - STAR
  - HISAT2
- Tutorial for reference guided assembly
  - Cufflinks
  - Stringtie
- Tutorial for de novo assembly
  - Trinity
- Visualise mapped reads and assembled transcripts on reference
  - IGV
- RNA quality controll
  - Tutorial for RNA seq Quality Control
- Differential expression analysis
  - DEseq2
  - Calisto and Sleuth
  - multi variate analysis in SIMCA

- small RNA analysis
  - miRNA analysis

- **Introductory**
  - Introduction to the RNA seq data provided
  - Short introduction to R
  - Short introduction to IGV
- **Beta labs**
  - Single cell RNA PCA and clustering
  - Gene set analysis
- **UPPMAX**
  - sbatch script example

# Need help??

- We are here for you. Apply for help.