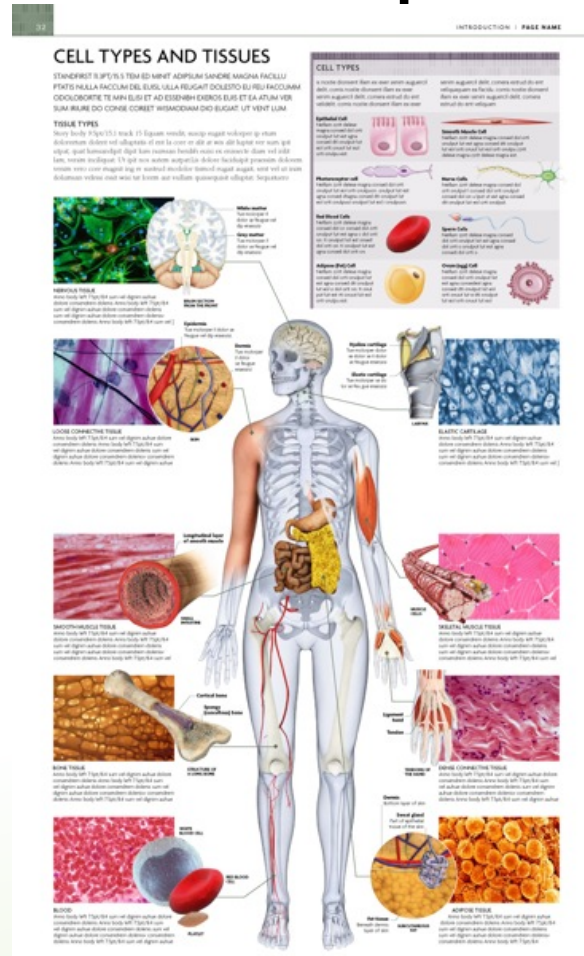# RNA-seq Introduction

## Promises and pitfalls

Enabler for Life Sciences

# RNA gives information on which genes that are expressed



How DNA get transcribed to RNA (and sometimes then translated to proteins) varies between e. g.

-Tissues

-Cell types

-Cell states

-Individuals

-Cells

# RNA gives information on which genes that are expressed



How DNA get transcribed to RNA (and sometimes then translated to proteins) varies between e. g.
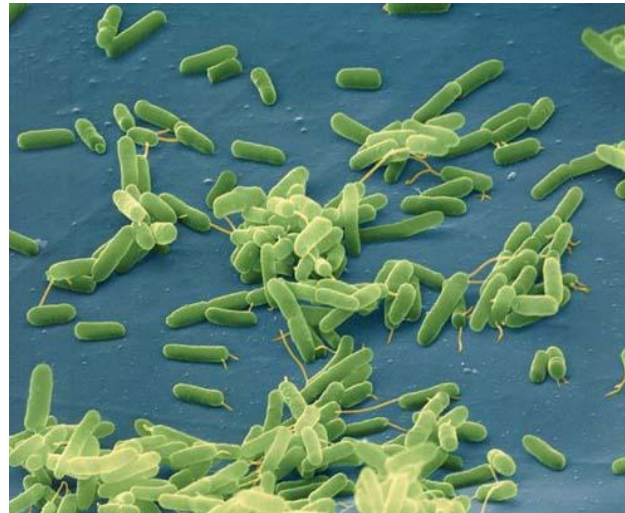
-Tissues

-Cell types

-Cell states

-Individuals

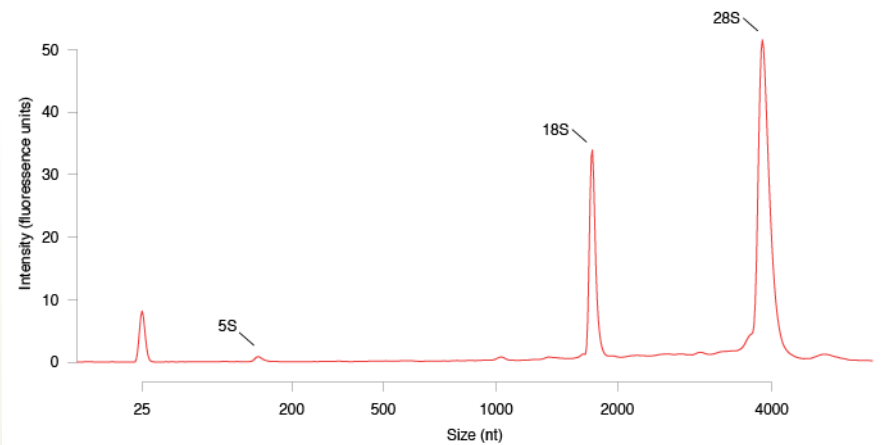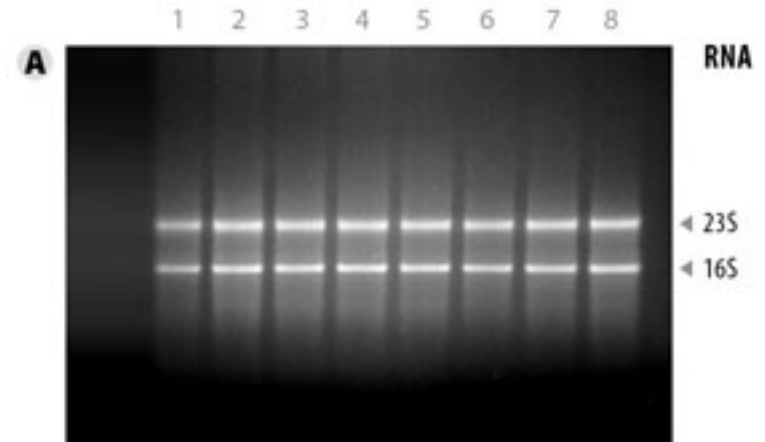# RNA gives information on which genes that are expressed



How DNA get transcribed to RNA (and sometimes then translated to proteins) varies between e. g.

-Tissues

-Cell types

-Cell states

-Individuals

# RNA flavors
# (pre sequencing era)

- House keeping RNAs
  - rRNAs, tRNAs, snoRNAs, snRNAs, SRP RNAs, catalytic RNAs (RNAse E)

- Protein coding RNAs
  - (1 coding gene ~ 1 mRNA)
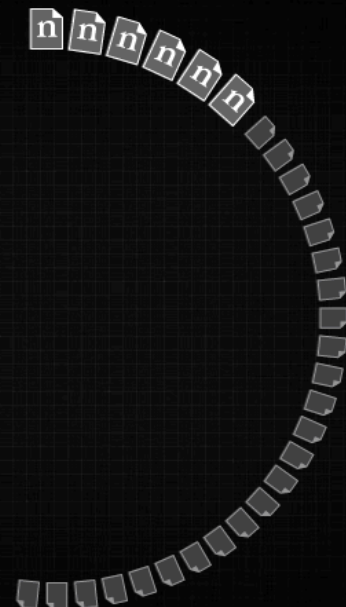
- Regulatory RNAs
  - Few rare examples

ENCODE, the Encyclopedia of DNA Elements, is a project funded by the National Human Genome Research Institute to identify all regions of transcription, transcription factor association, chromatin structure and histone modification in the human genome sequence.
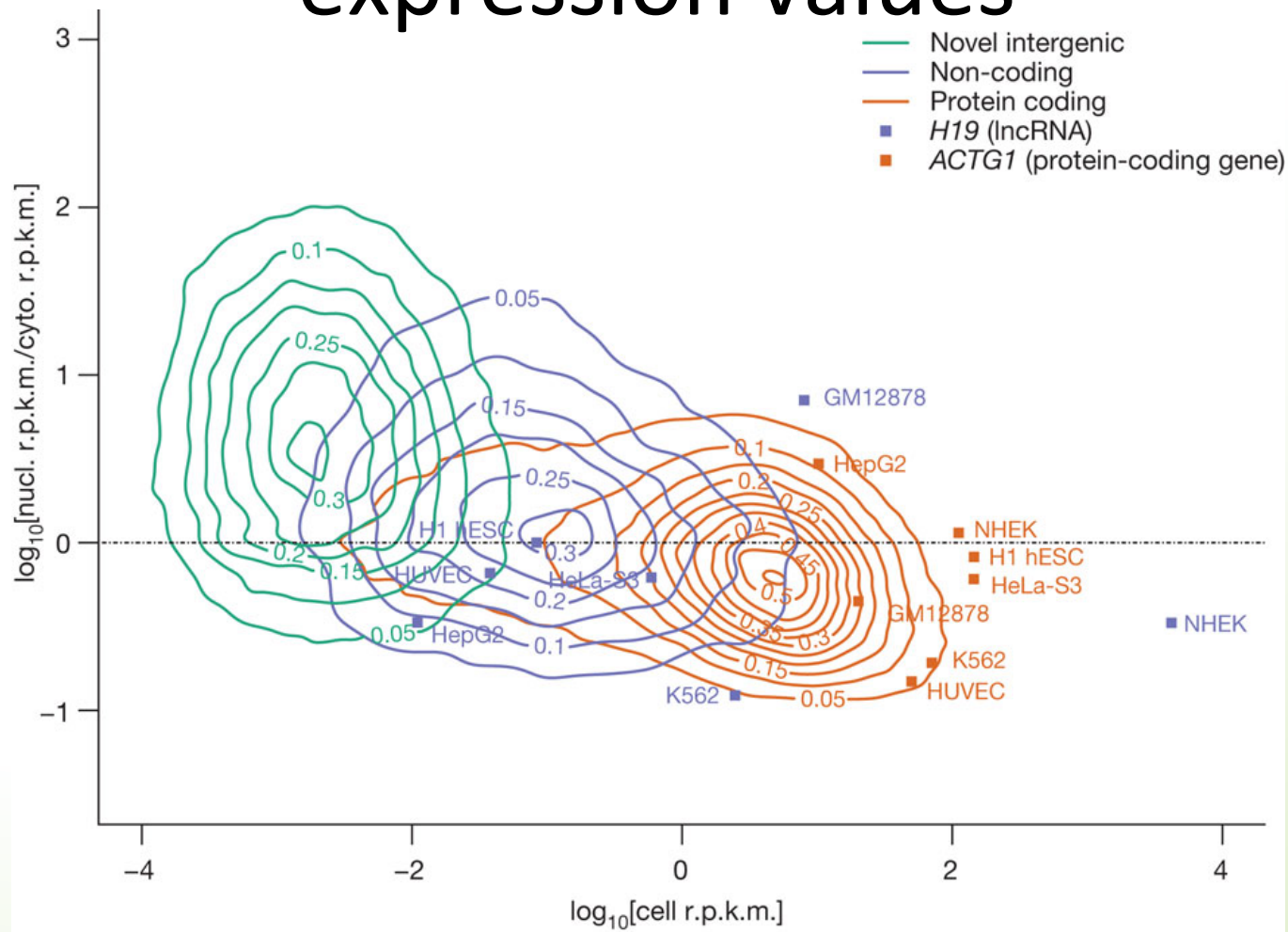
# ENCyclopedia Of Dna Elements
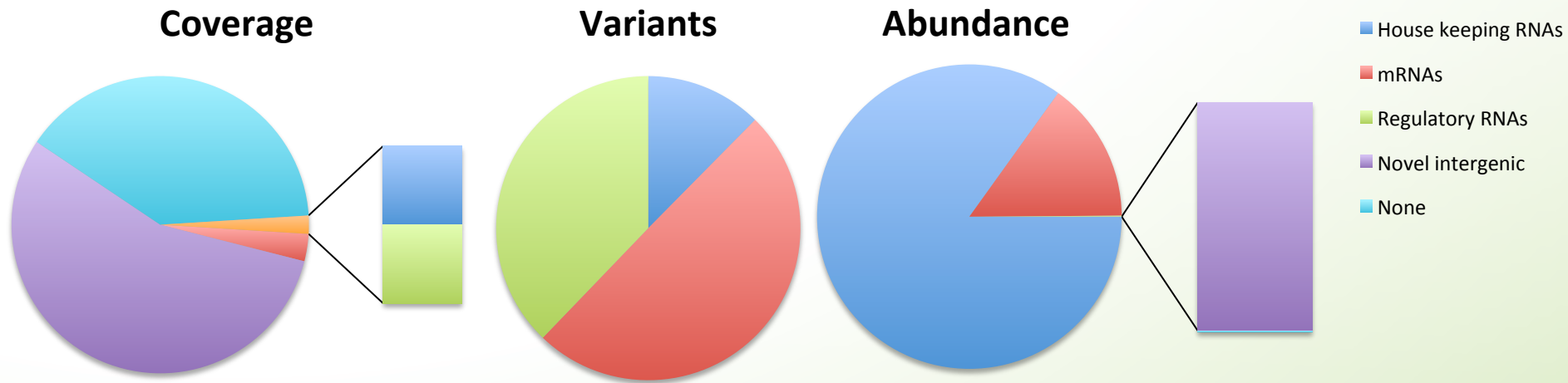
**ENCODE By the Numbers**

**147** cell types studied

**80%** functional portion of human genome

**20,687** protein-coding genes

**18,400** RNA genes

**1640** data sets

**30** papers published this week

**442** researchers

**$288 million** funding for pilot, technology, model organism, and current

Cumulatively, we observed a total of 62.1% and 74.7% of the human genome to be covered by either processed or primary transcripts, respectively, with no cell line showing more than 56.7% of the union of the expressed transcriptomes across all cell lines.

# Different kind of RNAs have different expression values

# What defines RNA depends on how you look at it



Coverage · Variants · Abundance

Legend: House keeping RNAs · mRNAs · Regulatory RNAs · Novel intergenic · None

# Defining functional DNA elements in the human genome

- Statement
  - A priori, we should not expect the transcriptome to consist exclusively of functional RNAs.
- Why is that
  - Zero tolerance for errant transcripts would come at high cost in the proofreading machinery needed to perfectly gate RNA polymerase and splicing activities, or to instantly eliminate spurious transcripts.
  - In general, sequences encoding RNAs transcribed by noisy transcriptional machinery are expected to be less constrained, which is consistent with data shown here for very low abundance RNA

- Consequence
  - Thus, one should have high confidence that the subset of the genome with large signals for RNA or chromatin signatures coupled with strong conservation is functional and will be supported by appropriate genetic tests.
  - In contrast, the larger proportion of genome with reproducible but low biochemical signal strength and less evolutionary conservation is challenging to parse between specific functions and biological noise.

# This is of course not without an debate

# Biochemical evidence not enough to identify functional RNAs



**Genetic evidence?** (generates phenotype)

**Biochemical evidence** (ENCODE, by level of activity) low    medium    high

**Evolutionary evidence** (mammalian conservation)

Protein-coding

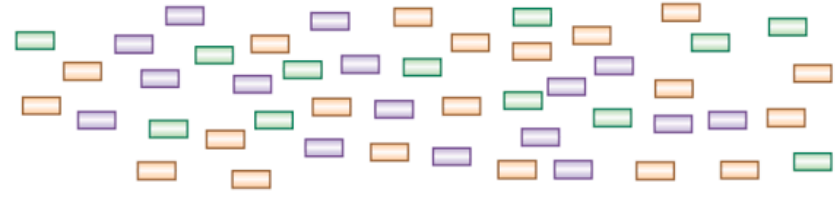Whole genome

# One gene many different mRNAs
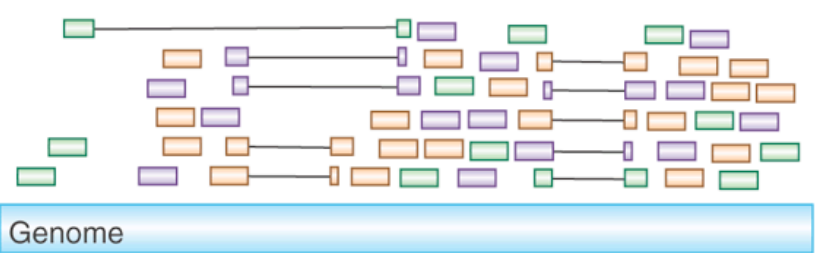
- RNA seq course

# The RNA seq course

- From RNA seq to reads
- Mapping reads programs
- Transcriptome reconstruction using reference
- Transcriptome reconstruction without reference
- QC analysis
- sRNA analysis
- Differential expression analysis
  - mRNAs
  - miRNAs
- Genome annotation using RNA and other sources
- Differential expression using multi variate analysis
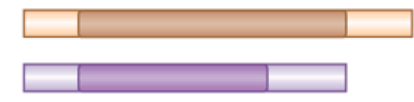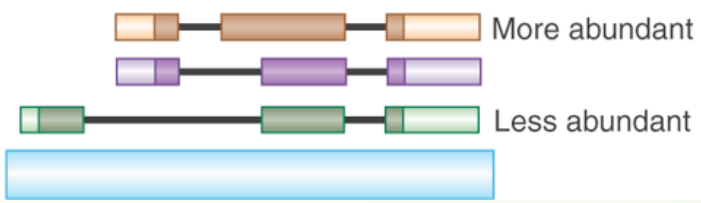- RNA long read analysis
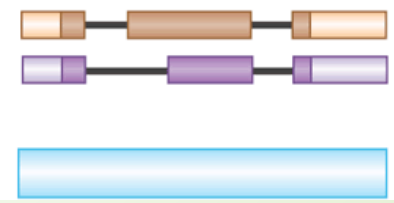
From RNA to short reads

Align reads to genome

Assemble transcripts de novo

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant
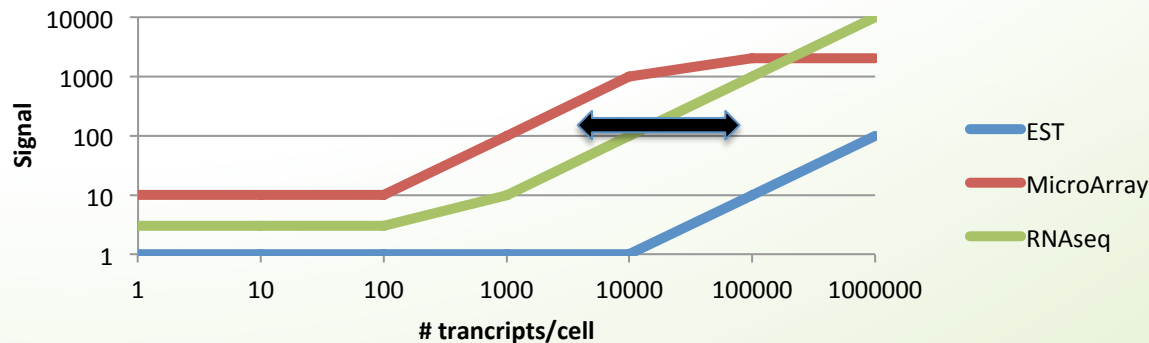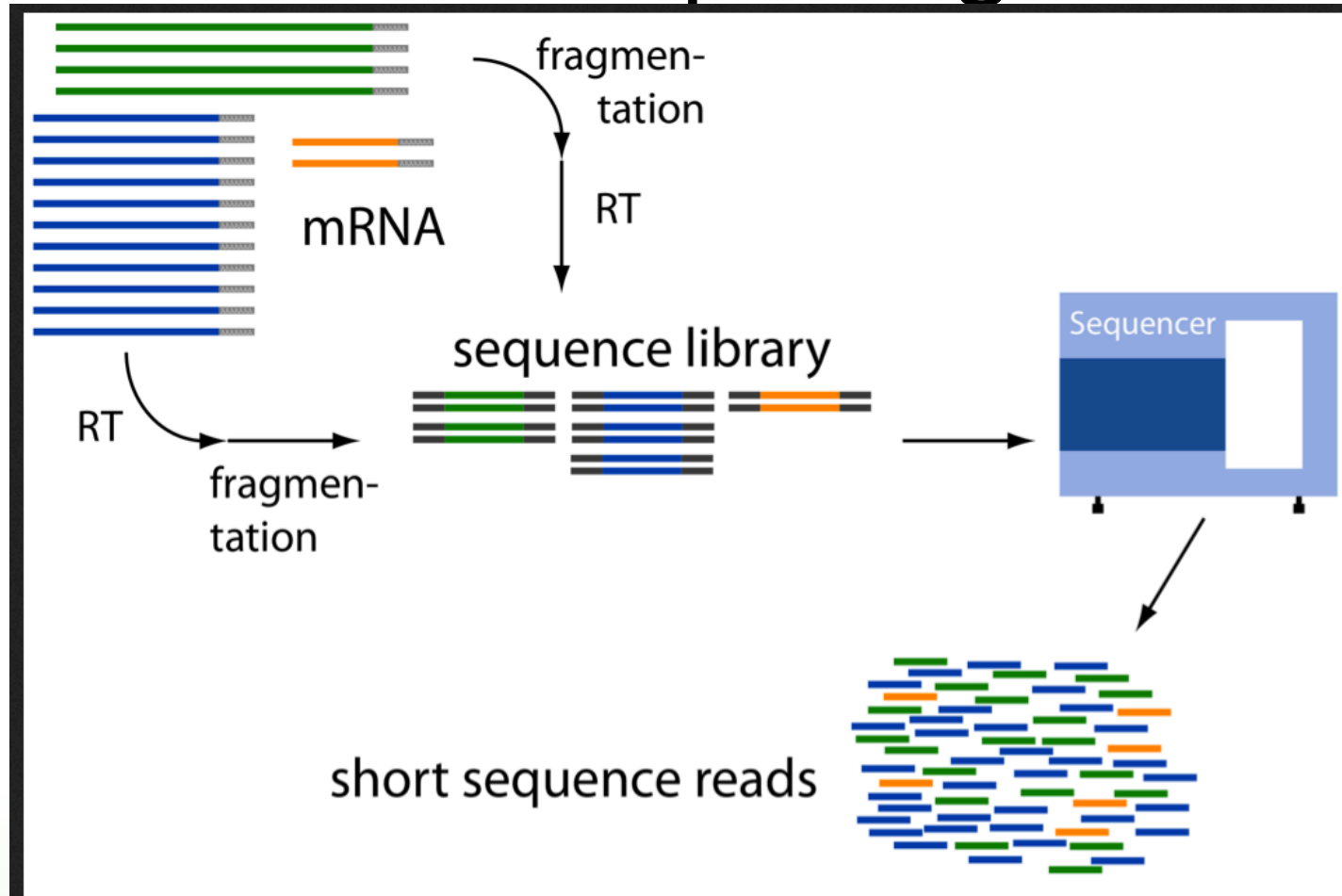
Less abundant

# Promises and pitfalls

## Long reads

- Low throughput (-)
- Complete transcripts (+)
- Only highly expressed genes (--)
- Expensive (-)
- Low background noise (+)
- Easy downstream analysis (+)

## Micro Arrays

- High throughput (+)
- Only known sequences (-)
- Limited dynamic range (-)
- Cheap (+)
- High background noise (-)
- Not strand specific (-)
- Well established downstream methods (+)

## RNAseq

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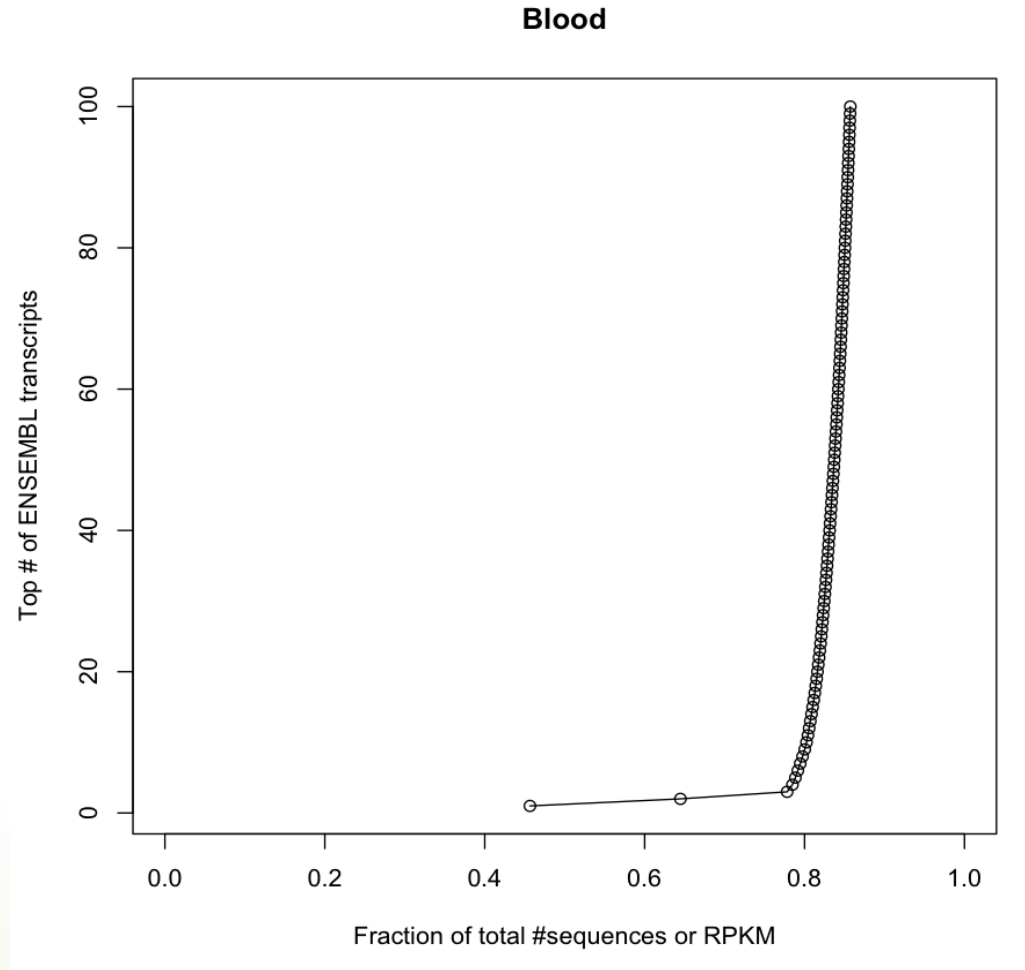- Strand specificity (+)
- Re-sequencing (+)

# How are RNA-seq data generated?



**Sampling process**

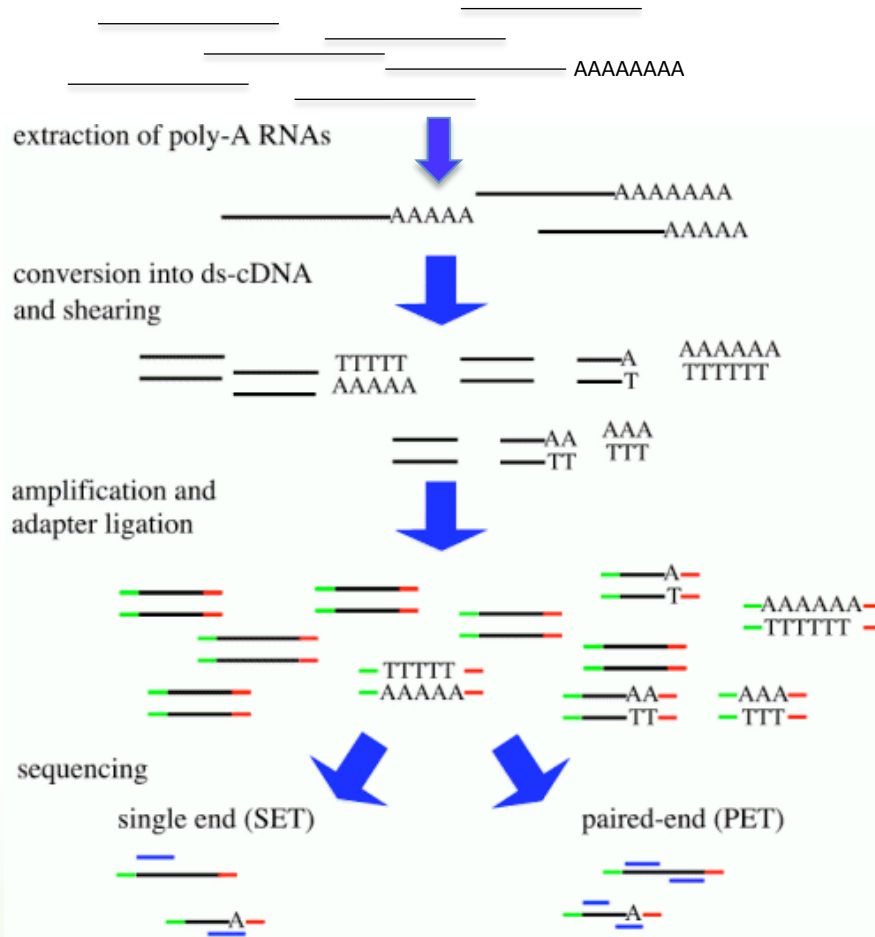# RNA seq reads correspond directly to abundance of RNAs in the sample

# RNA to reads



RNA->

enrichments ->
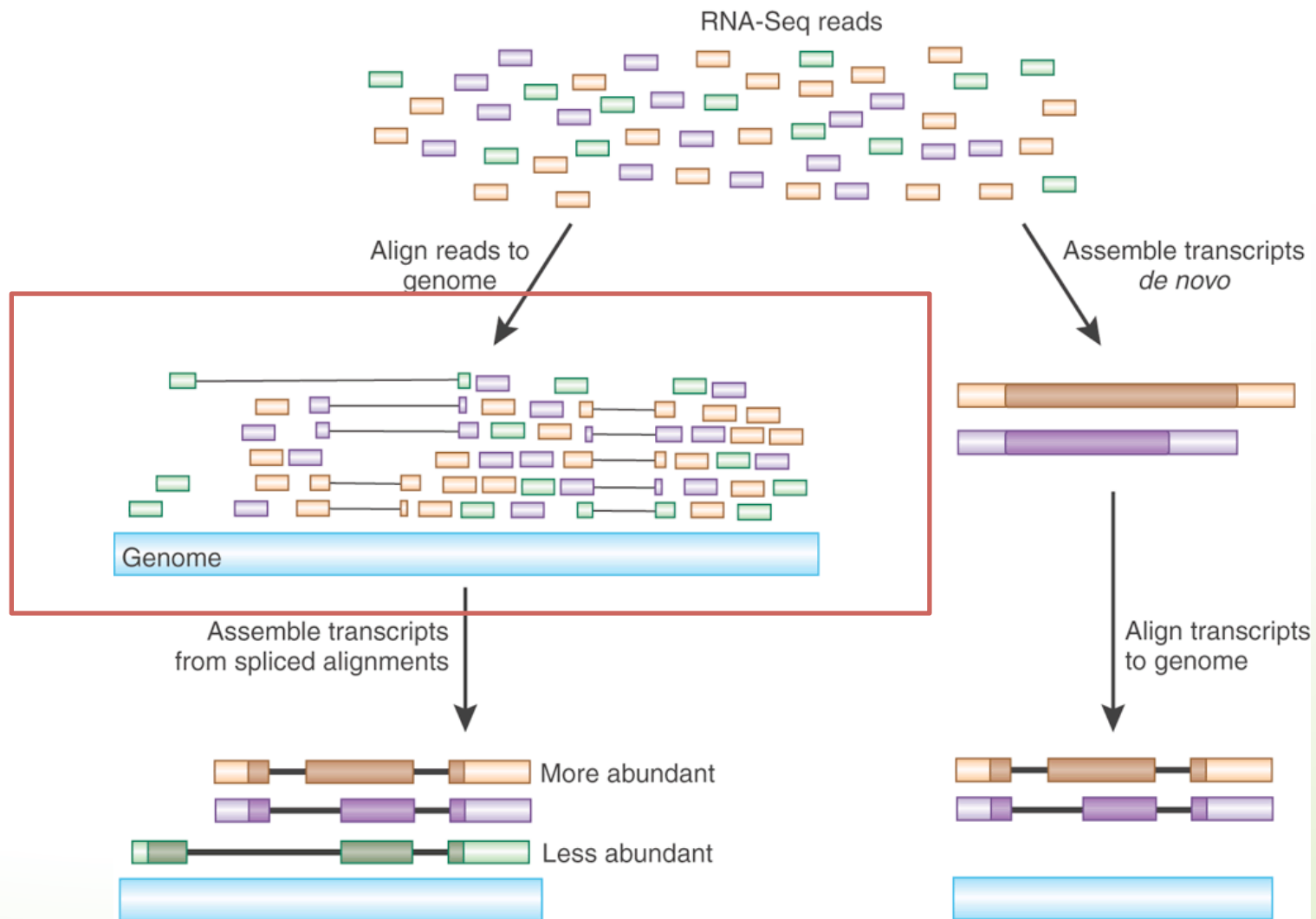
extraction of poly-A RNAs

conversion into ds-cDNA and shearing

amplification and adapter ligation

library ->

sequencing

single end (SET)    paired-end (PET)

reads ->

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

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

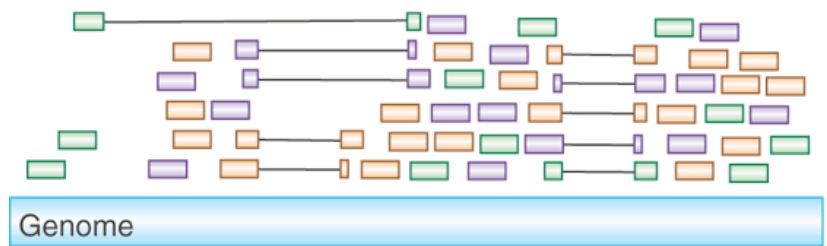Single end (1 read per fragment)
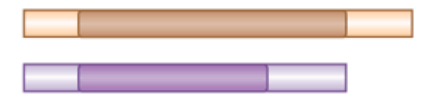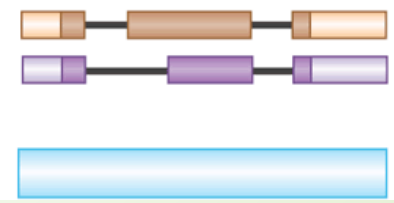Paired end (2 reads per fragment)

RNA-Seq reads
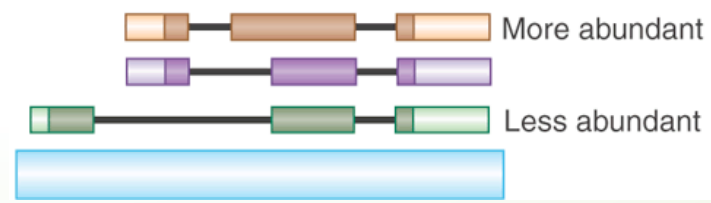
Align reads to genome

Assemble transcripts *de novo*

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome
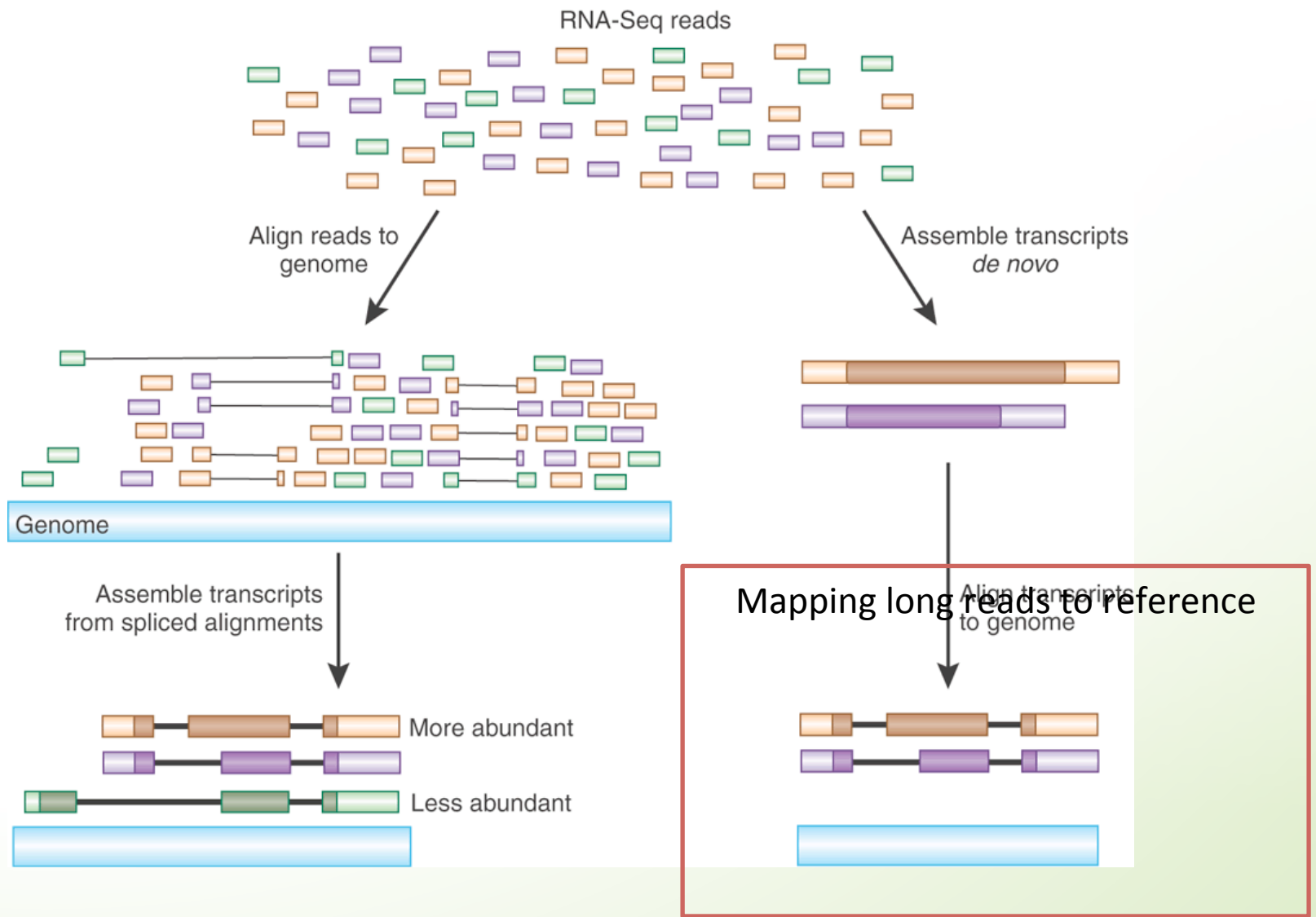
More abundant

Less abundant

Transcriptome assembly using reference

RNA-Seq reads

Align reads to genome

Transcriptome assembly without reference

*de novo*

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant

Less abundant

RNA-Seq reads

Align reads to genome

Assemble transcripts de novo

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant

Less abundant

Mapping long reads to reference
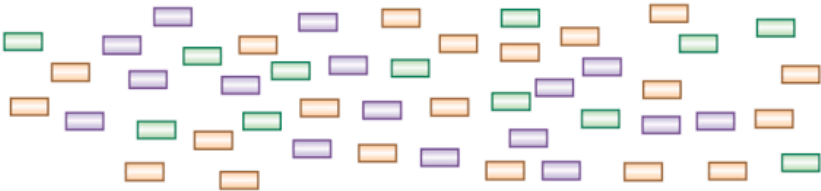
# Quality control
## -samples might not be what you think they are

- Experiments go wrong
  - 30 samples with 5 steps from samples to reads has 150 potential steps for errors
  - Error rate 1/100 with 5 steps suggest that one of every 20 samples the reads does not represent the sample

- Mixing samples
  - 30 samples with 5 steps from samples to reads has ~24M potential mix ups of samples
  - Error rate 1/ 100 with 5 steps suggest that one of every 20 sample is mislabeled

- Combine the two steps and approximately one of every 10 samples are wrong

# RNA QC



Read quality

RNA-Seq reads

Align reads to genome

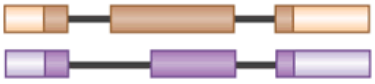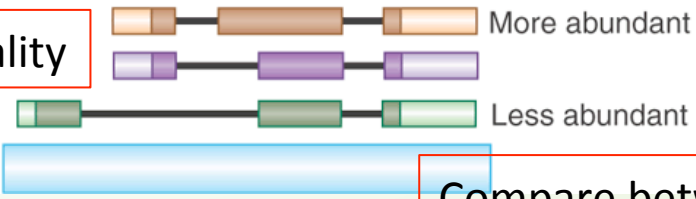Assemble transcripts *de novo*

Mapping statistics

Genome

Assemble transcripts from spliced alignments
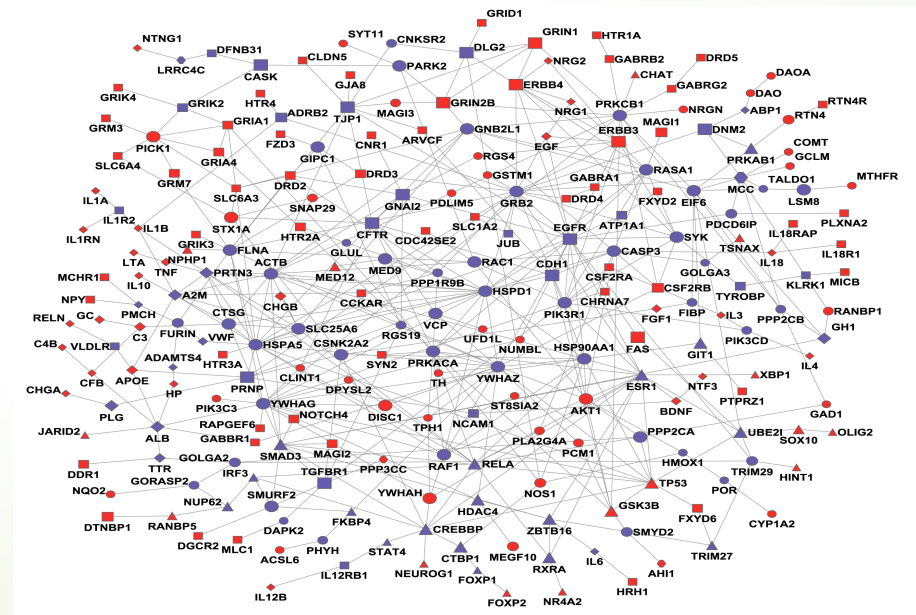
Align transcripts to genome

Transcript quality
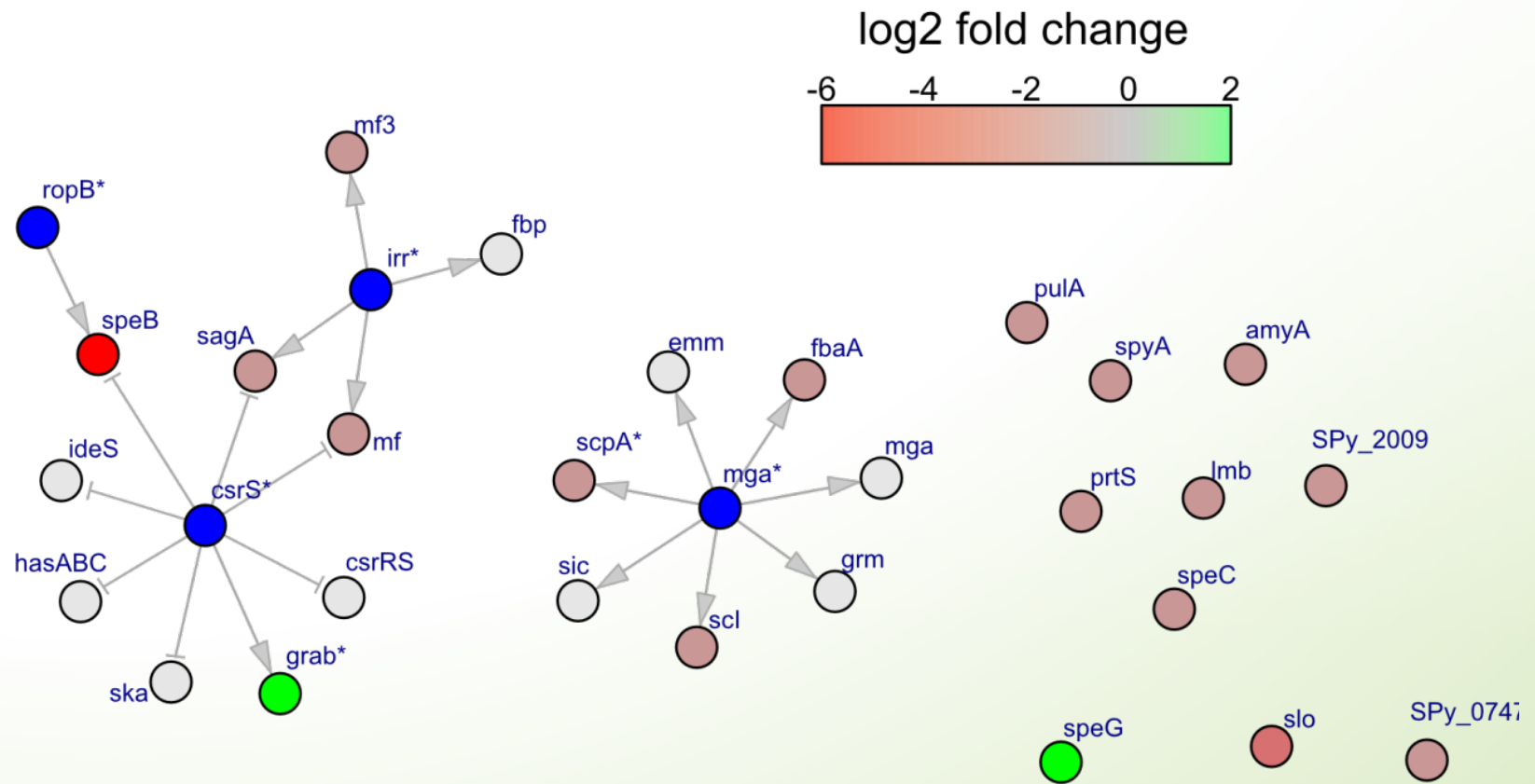
More abundant

Less abundant

Compare between samples

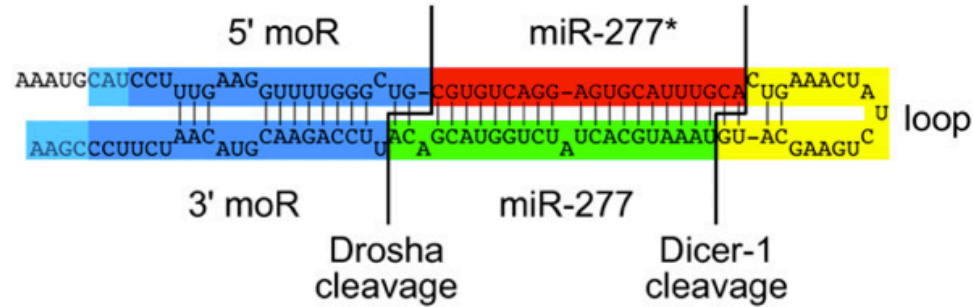# Differential expression analysis using univariate analysis

Typically **univariate** analysis (one gene at a time) – even though we know that genes are not independent
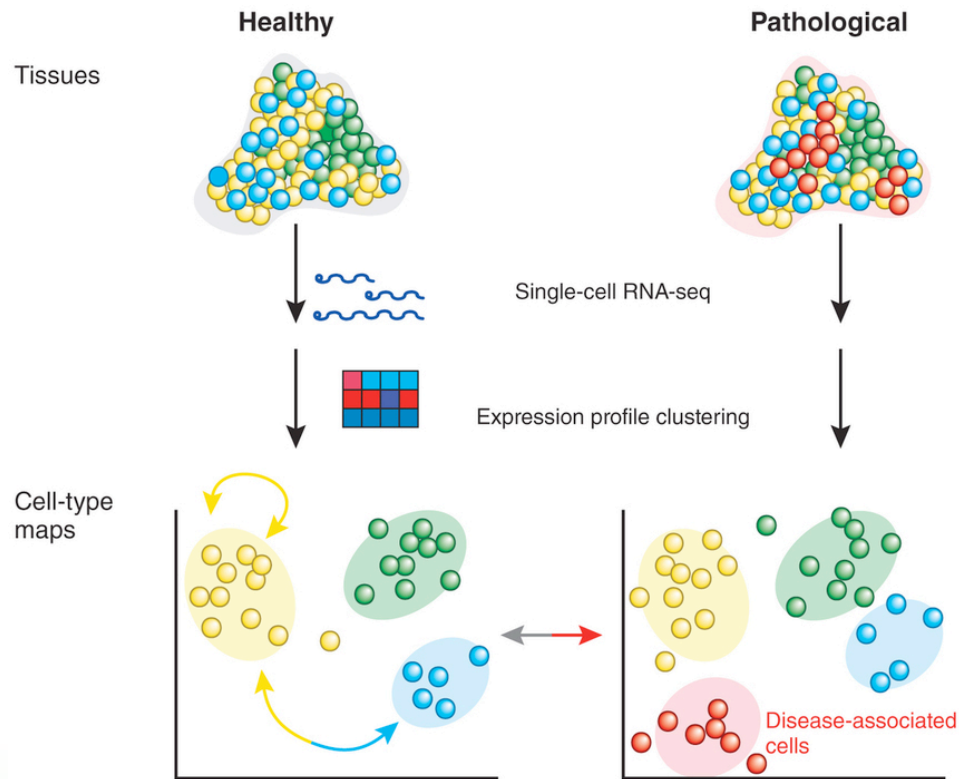
# Gene set analysis and data integration

# microRNA analysis (Jakub)



(Berezikov et al. Genome Research, 2011.)

# Single cell RNA-seq analysis



(Sandberg, Nature Methods 2014)

# Long reads

Short reads

# Long reads



Short reads

Long reads