# RNA-seq Introduction
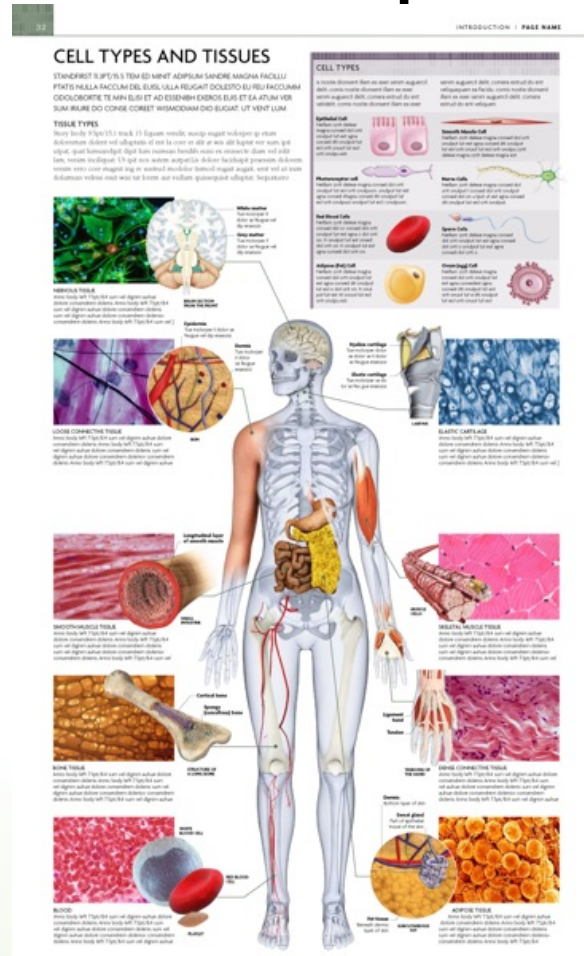
## Promises and pitfalls

# 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.
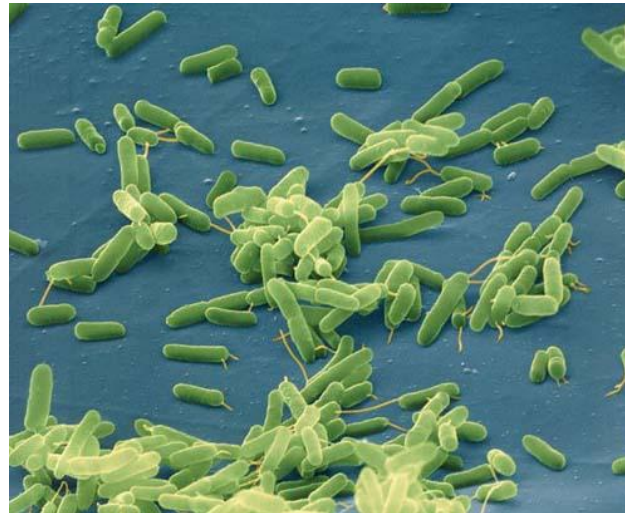
-Tissues

-Cell types

-Cell states

-Individuals

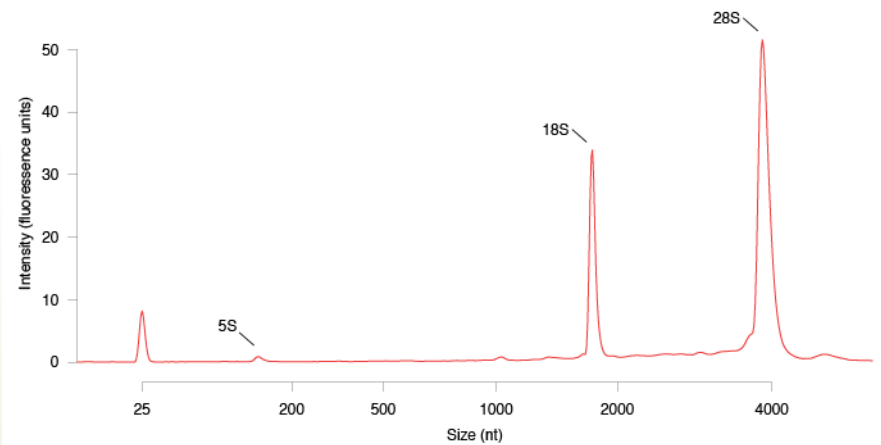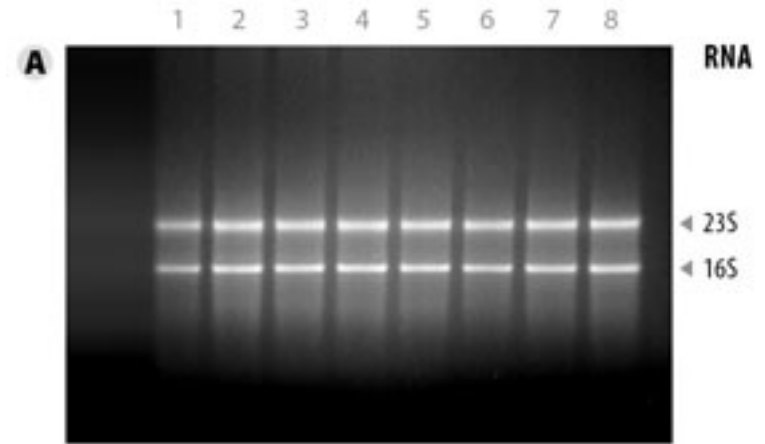# RNA gives information on which genes that are expressed



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)

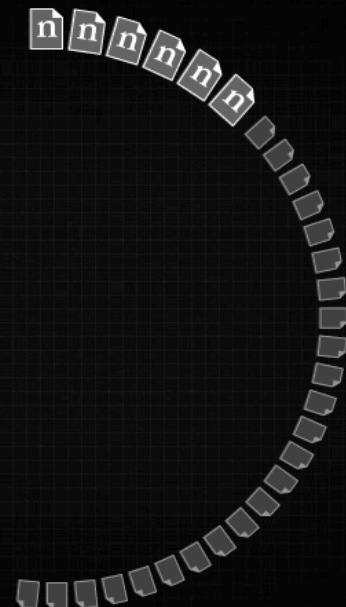- ## Regulatory RNAs
  - Few rare examples

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.
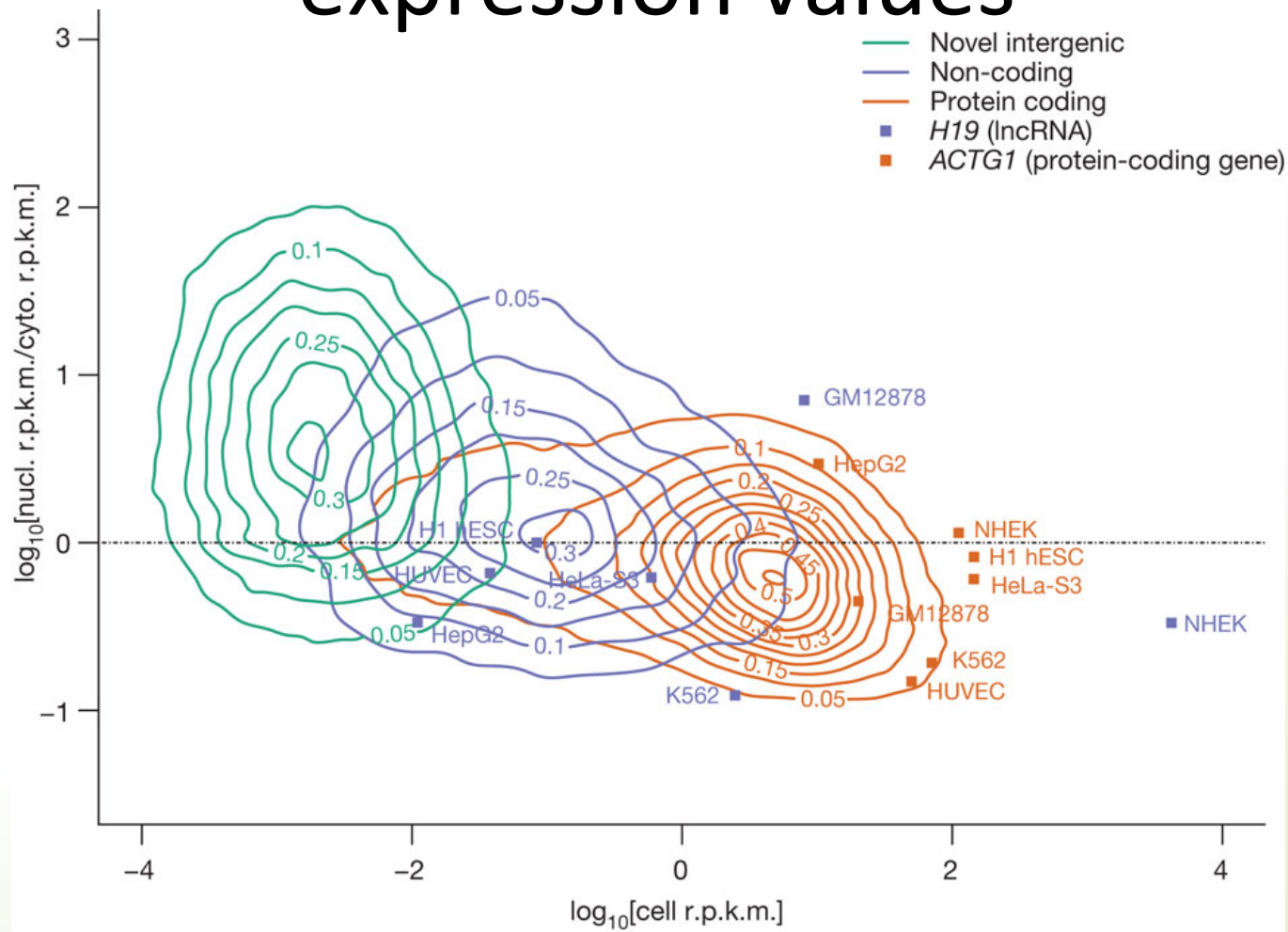
# ENCyclopedia Of Dna Elements
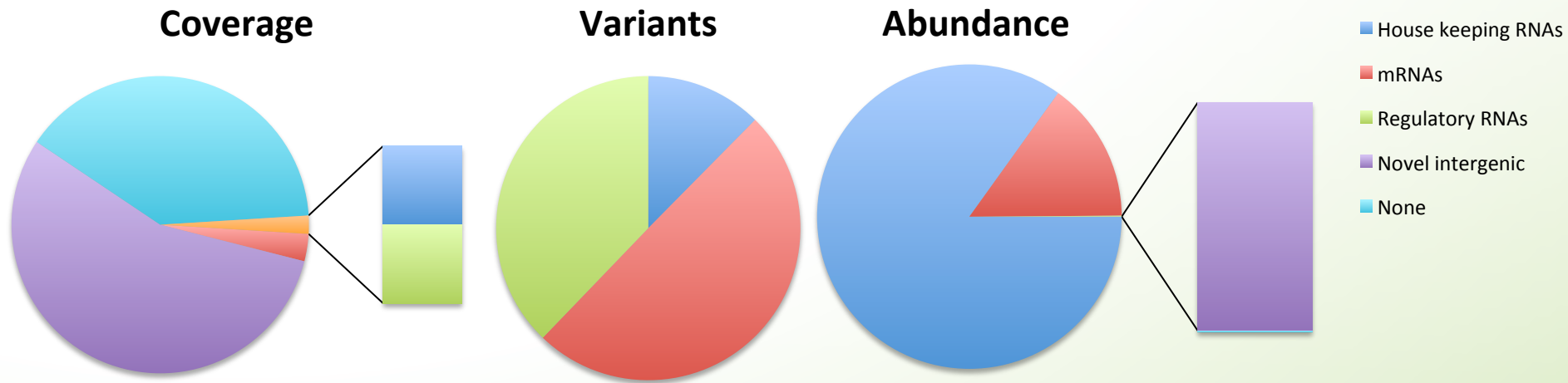
## ENCODE By the Numbers

**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



Landscape of transcription in human cells, S Djebali *et al.* **Nature 2012**

# 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

# One gene many different mRNAs

- RNA seq course

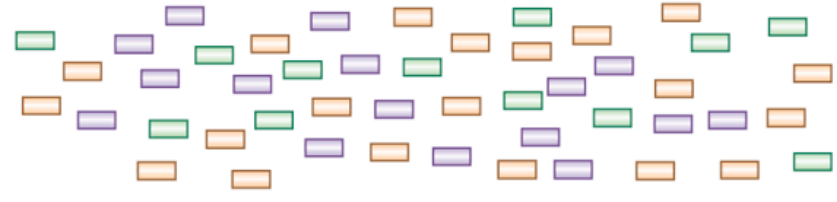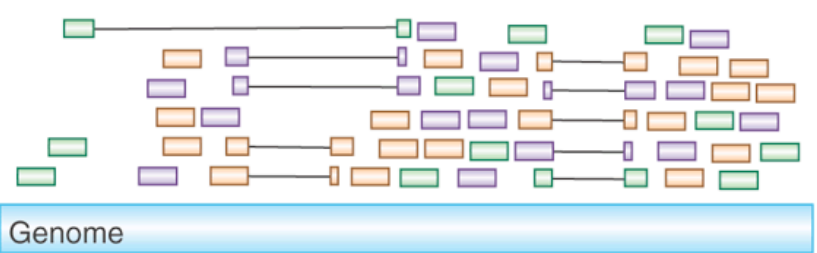# The RNA seq course

- 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

# Sequencing platforms



|  | ABI 3730xl Sanger Sequencing | 454 Life Sciences pyrosequencing | SOLiD + Illumina | Pacific Biosciences, Oxford Nanopore etc Single-molecule sequencing |
|---|---|---|---|---|
| **Length/read** | 800 bp | 400 bp | 100 bp | 20 000+ bp |
| **Reads/run** | 96 | 1 million | 2 billion | 5 million |
| **Bases/run** | 60 kbp | 400 Mbp | 500 Gbp | 100 Gbp |
| **Speed** | 10 years/HG | 1 month/HG | 1 day/HG | 10 min/HG |
|  | "Old school" | "2ⁿᵈ gen" | | "3ʳᵈ gen" |

"Old school"          "2nd gen"          "3rd gen"

# Promises and pitfalls

## Sanger

- 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
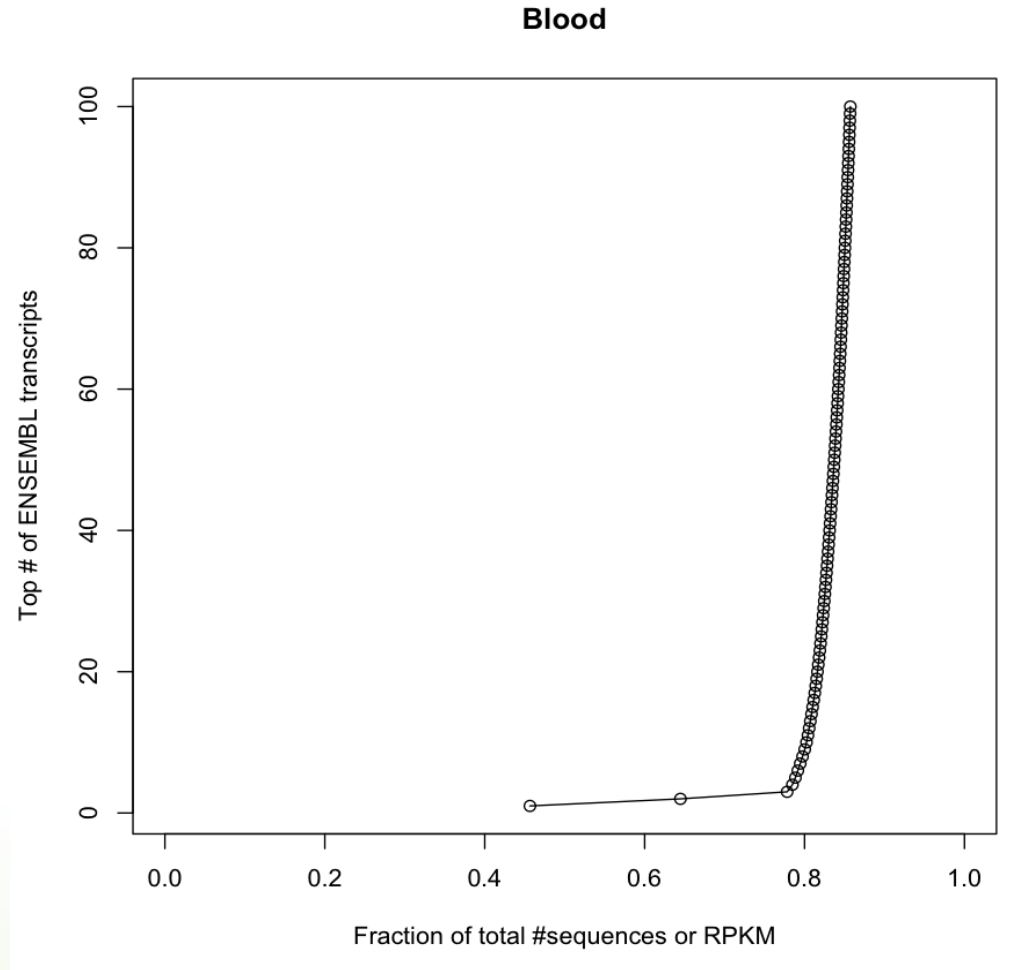
RNA->

enrichments ->

library ->

reads ->

extraction of poly-A RNAs

conversion into ds-cDNA and shearing

amplification and adapter ligation

sequencing

single end (SET)

paired-end (PET)

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)

RNA-Seq reads

Align reads to genome

Assemble transcripts *de novo*

Mapping reads to reference (Johan)

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

Transcriptome assembly using reference (Estelle)

More abundant

Less abundant

RNA-Seq reads

Align reads to genome

Transcriptome assembly without reference (Estelle)

*de novo*

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant
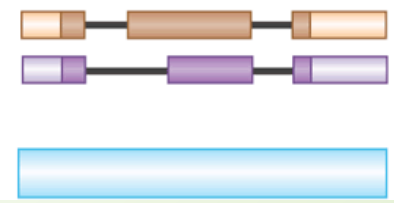
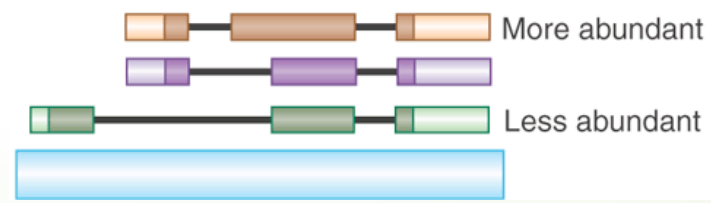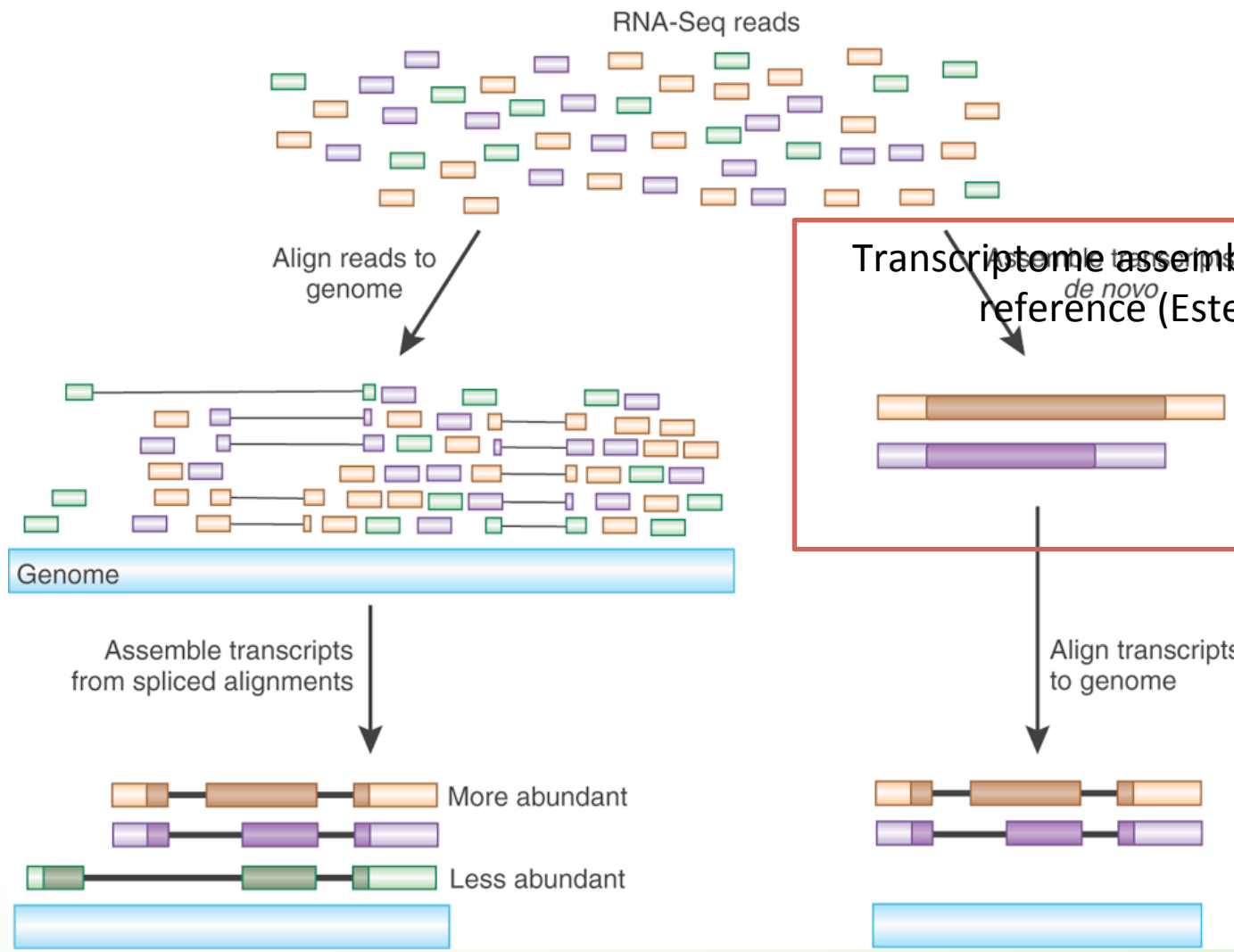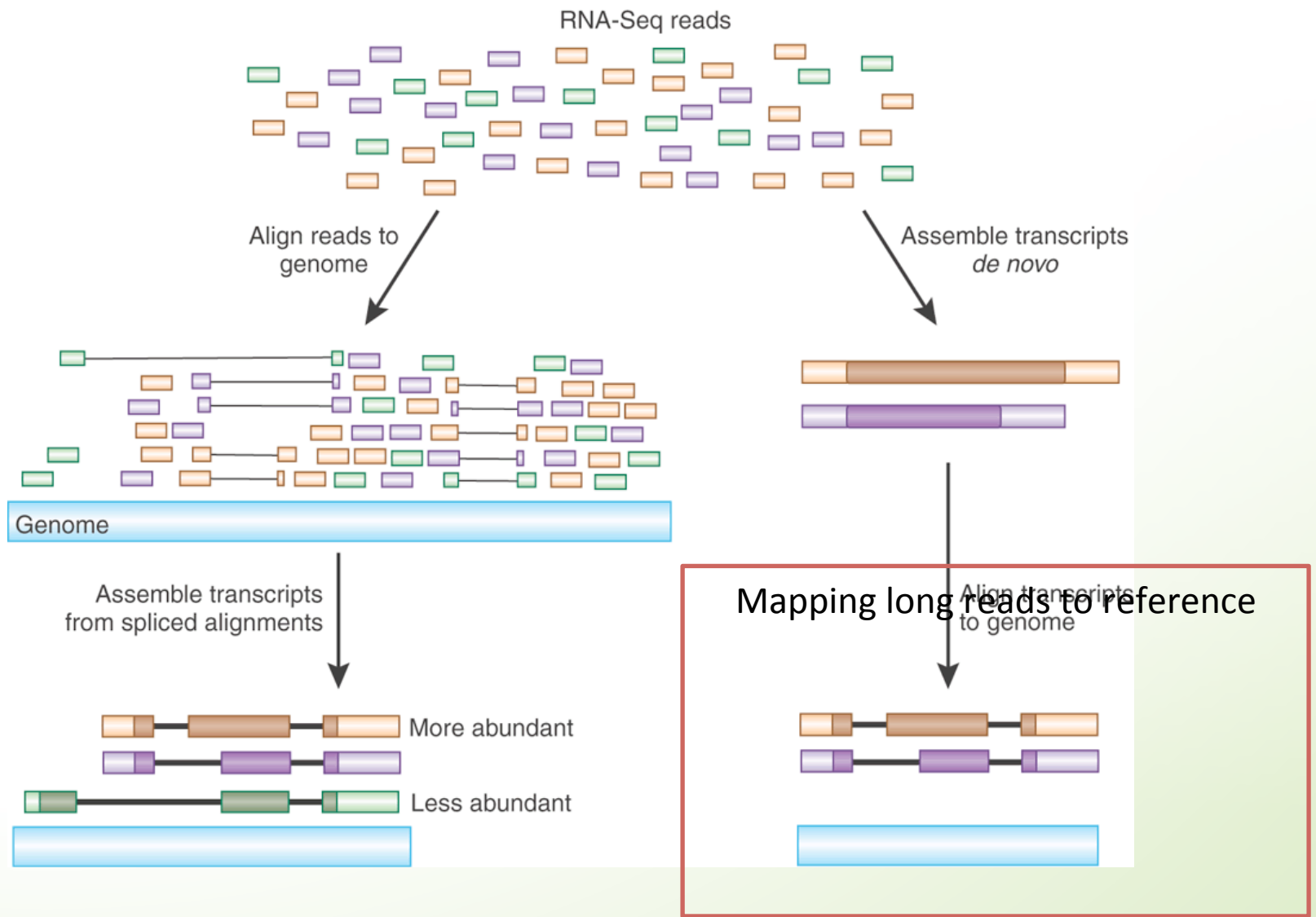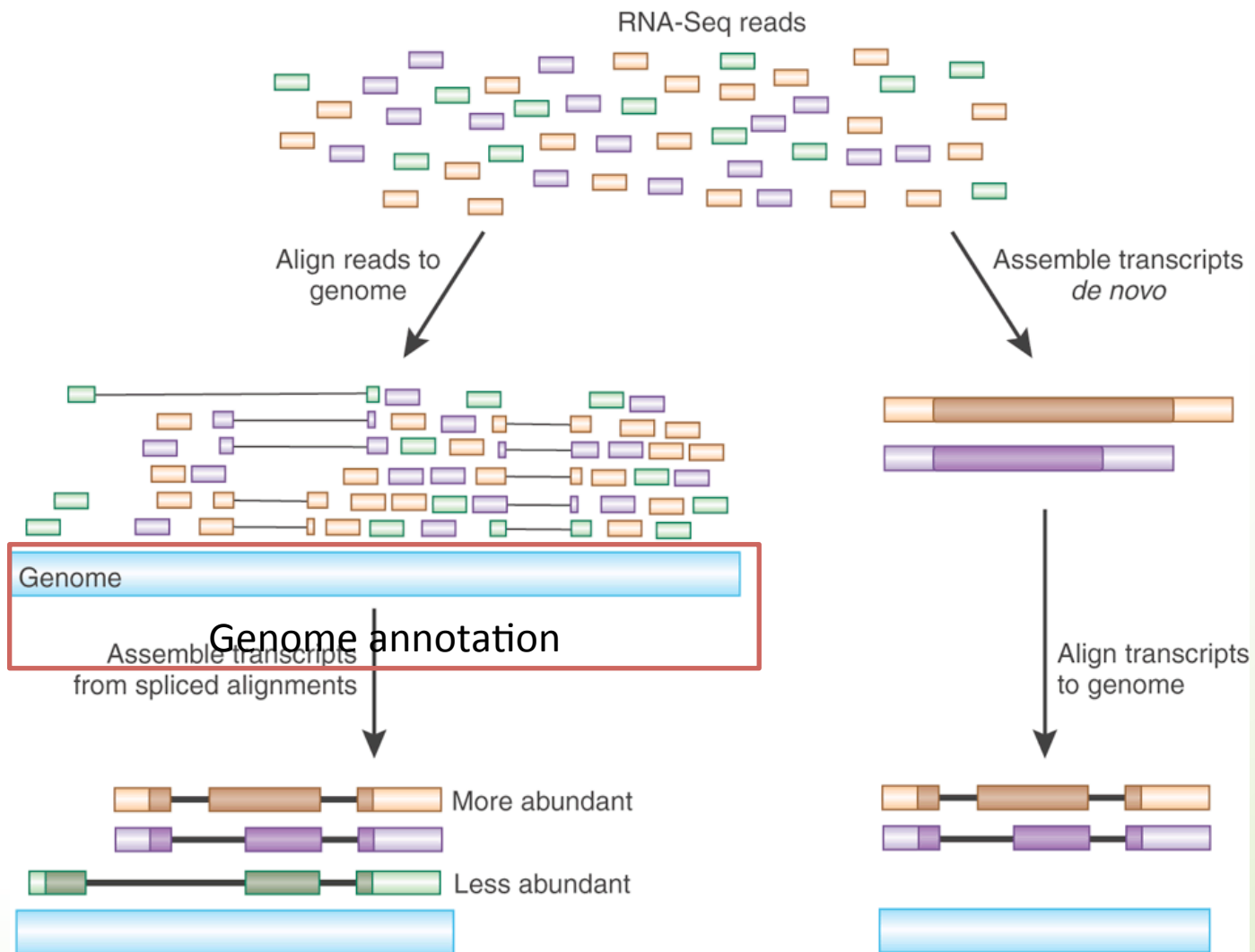Less abundant

RNA-Seq reads

Align reads to genome

Assemble transcripts *de novo*

Genome

Assemble transcripts from spliced alignments

Align transcripts to genome

Mapping long reads to reference

More abundant

Less abundant

RNA-Seq reads

Align reads to genome

Assemble transcripts *de novo*

RNA seq Long reads

Genome

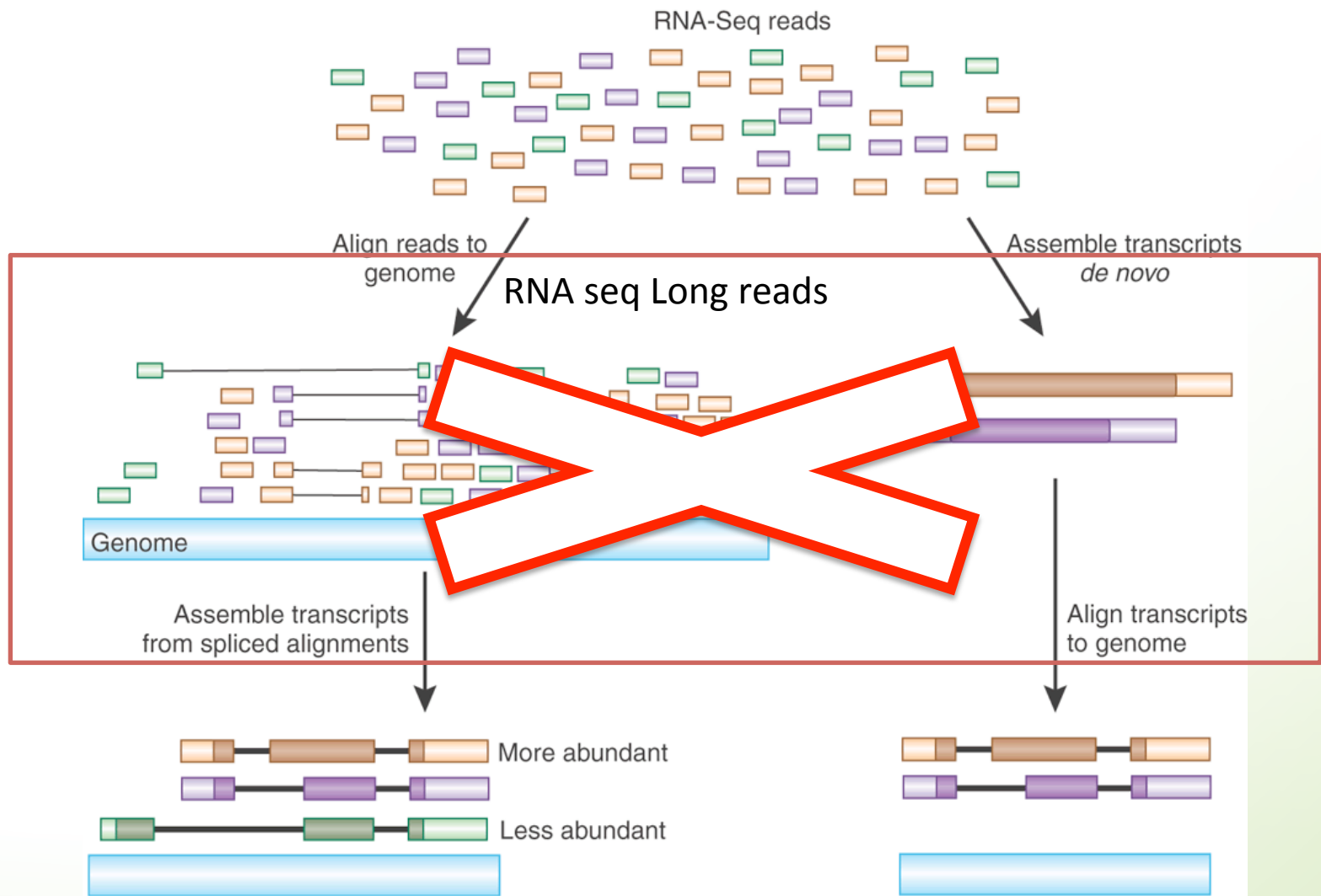Assemble transcripts from spliced alignments

Align transcripts to genome

More abundant

Less abundant

# microRNA analysis (Jakub)



(Berezikov et al. Genome Research, 2011.)

**ANOTHER WAY OF LOOKING AT IT**

# 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 (Åsa)



Read quality

Mapping statistics

Transcript quality

# Compare expression between different samples (Åsa)

# Differential expression analysis using univariate analysis (Åsa)

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

# Multi variate differential expression analysis (Sanela)

Multivariate methods such as PCA (unsupervised) or PLS (supervised) can be used to obtain loadings for features (genes/transcripts/...) that contribute to separation of groups



The loading scores can be used as a different kind of measure of which genes are interesting

# Welcome to WABI RNA-seq tutorial packages

This page contains links to different tutorials that are used in the RNA-seq course. Some of the tutorials are well documented and should be easy to follow. We also supply more beta versions of labs that requires more from the user and may contain errors.

## Covered labs in the course

- Introduction to the RNA seq data provided
- Short introduction to R
- Short introduction to IGV
- Mapping reads to a reference and convering them to the BAM format
- isoform-visualisation
- Tutorial for reference guided assembly
- Tutorial for de novo assembly
- Tutorial concerning RNA seq Quality Controll
- Tutorial for small RNA analysis
- Tutorial for differential expression analysis
- Tutorial for multi variate analysis

## Beta labs

- Differential expression analysis using kallisto

We will try to keep these tutorials up to date. If you find any errors or things that you think should be updated please contact Johan (johan.reimegard@scilifelab.se)

Karolinska Institutet
KTH VETENSKAP OCH KONST — ROYAL INSTITUTE OF TECHNOLOGY
Stockholm University
UPPSALA UNIVERSITET

SciLifeLab