# Studying the transcriptome using RNA-seq

Cecilia Coimbra Klein

Computational Biology of RNA Processing, CRG
Departament de Genètica, IBUB, UB

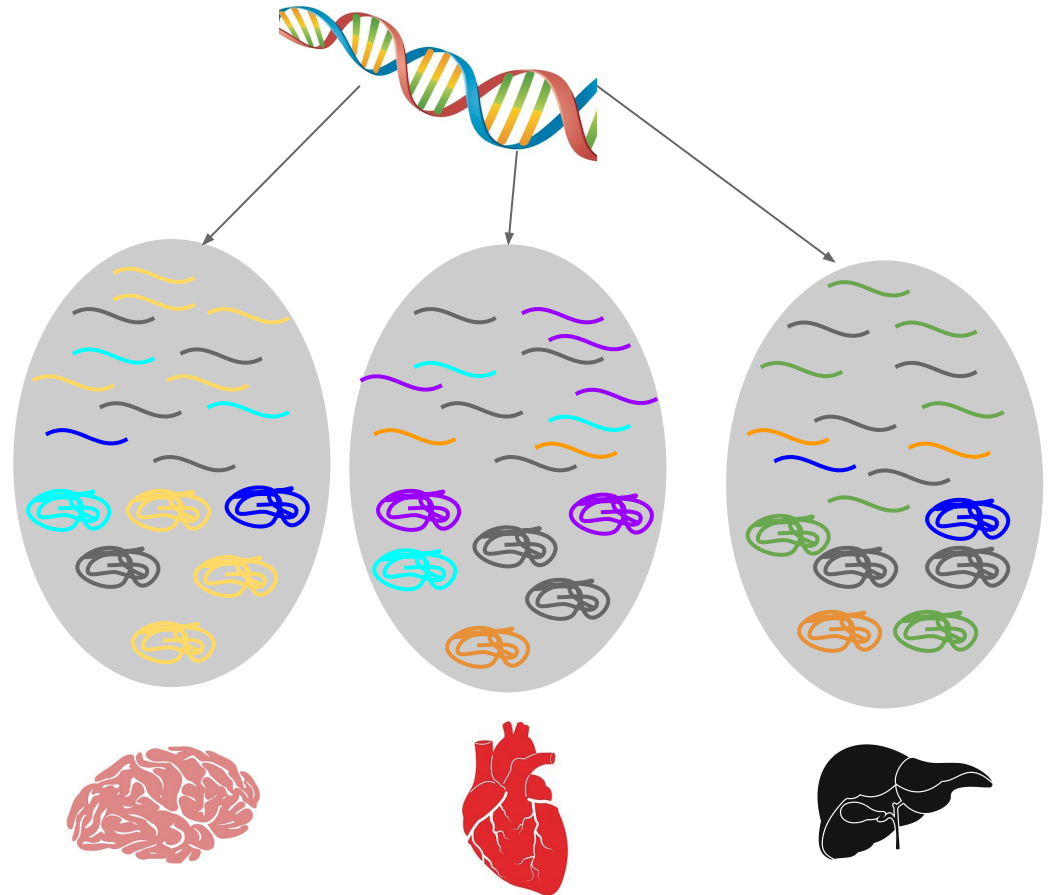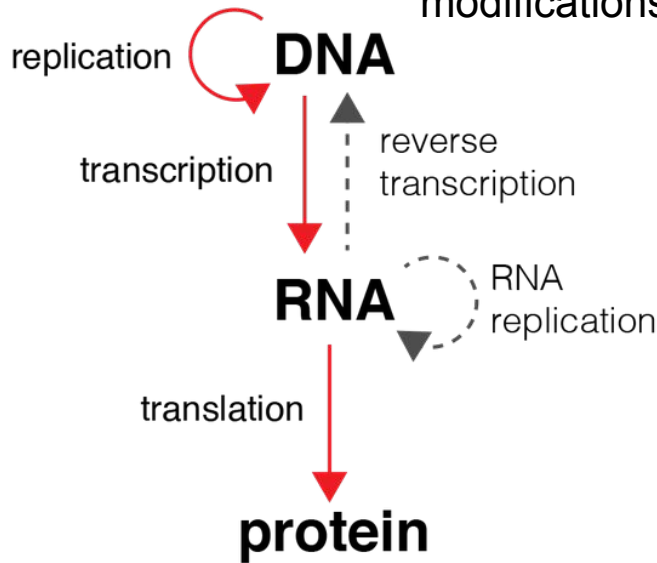Master in Omics Data Analysis
Jan. 2019

# Outline

# Outline

- Summary of the course
  - Day 1: RNA-seq introduction and processing
  - Day 2: RNA-seq analysis (clustering, differential gene expression, GO enrichment)
  - Day 3: RNA-seq analysis (splicing)
- MultiOmics
  - ChIP-seq (processing and data analysis)
  - ATAC-seq (visualization)
- Hands-on MultiOmics
  - ChIP-seq and ATAC-seq signal in the UCSC genome browser
  - promoter regions of differentially expressed genes
  - promoter regions of differentially spliced genes
  - omics portals
- Multiple-choice exercise

Cecilia Coimbra Klein

# Day 1: RNA-seq introduction and processing

# Molecular biology dogma

epigenetic modifications

replication DNA

transcription ← reverse transcription

RNA ← RNA replication

translation

protein

- The genome is identical in all cell types, however not all cell types have the same function. That's why the transcriptome (and the epigenome) becomes also relevant.
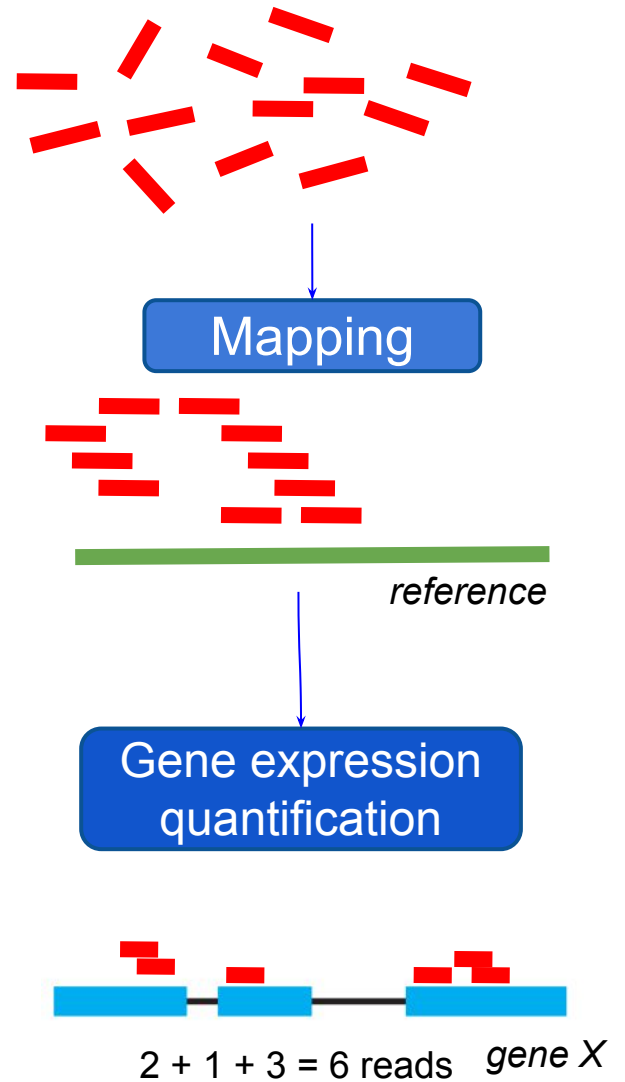
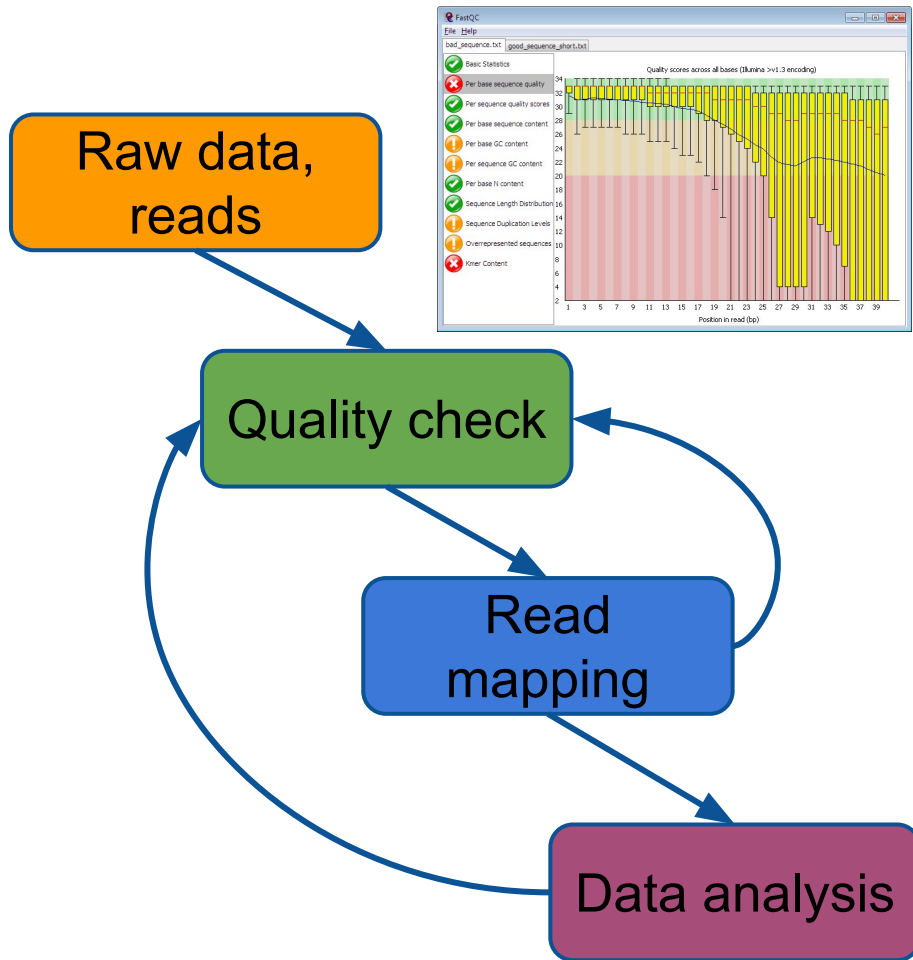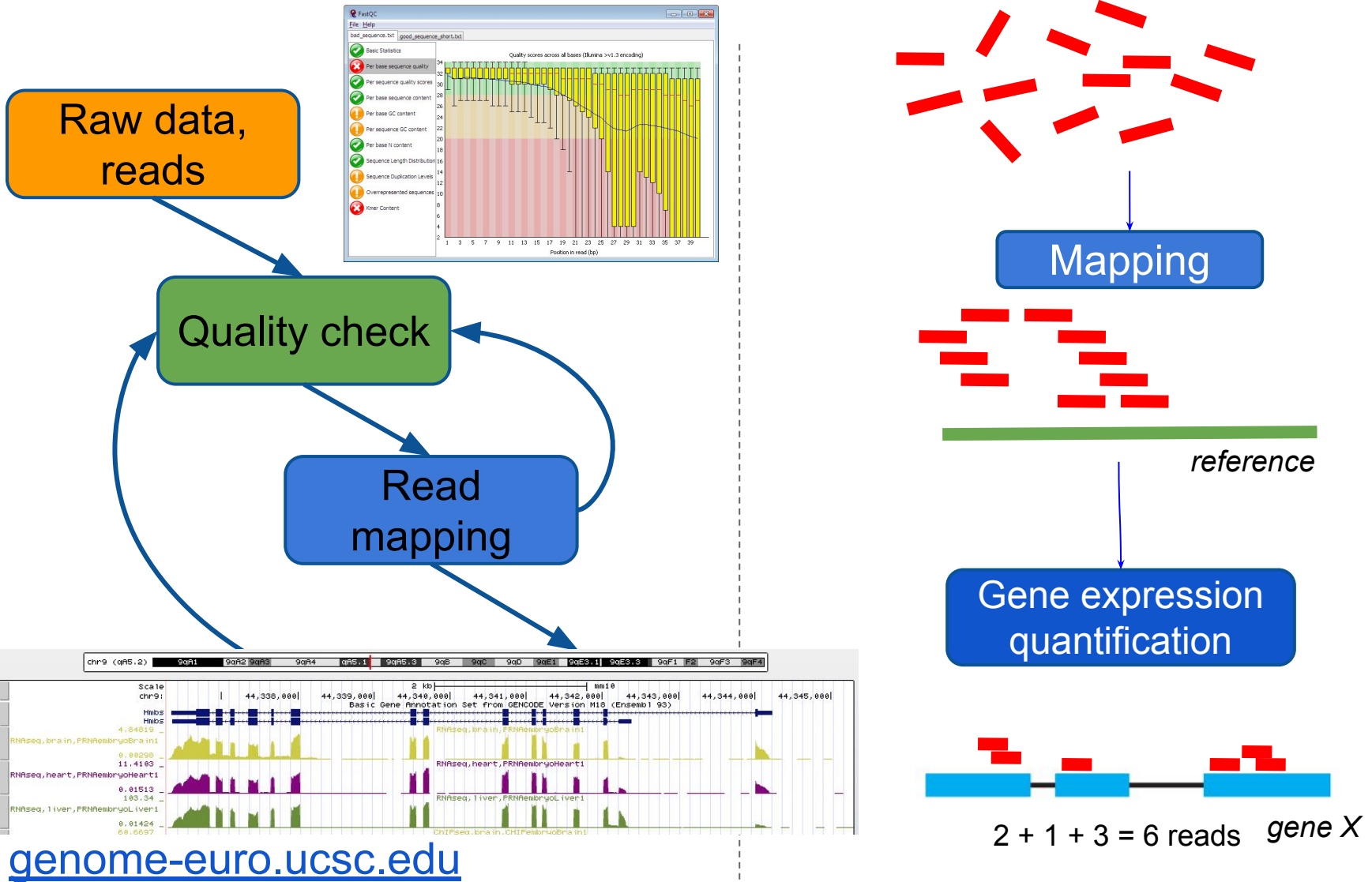Cecilia Coimbra Klein

# Typical pipeline



Raw data, reads

Quality check

Read mapping

Data analysis

Mapping

*reference*

Gene expression quantification

2 + 1 + 3 = 6 reads  *gene X*

Cecilia Coimbra Klein

# Typical pipeline



Raw data, reads

Quality check

Read mapping

genome-euro.ucsc.edu

Mapping

reference

Gene expression quantification
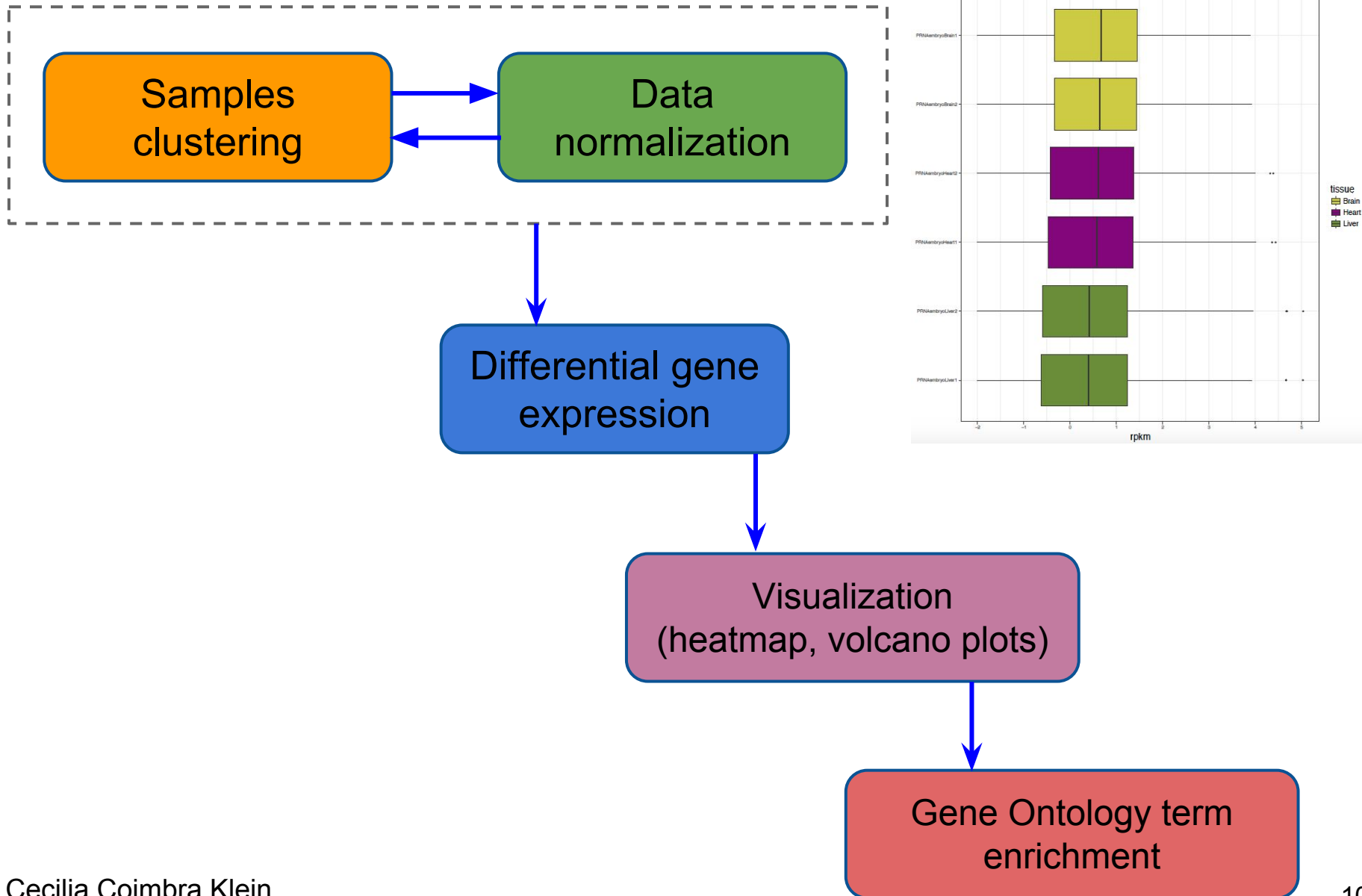
2 + 1 + 3 = 6 reads    *gene X*

Cecilia Coimbra Klein

7

# Day 2: RNA-seq analysis (clustering, differential gene expression, GO enrichment)
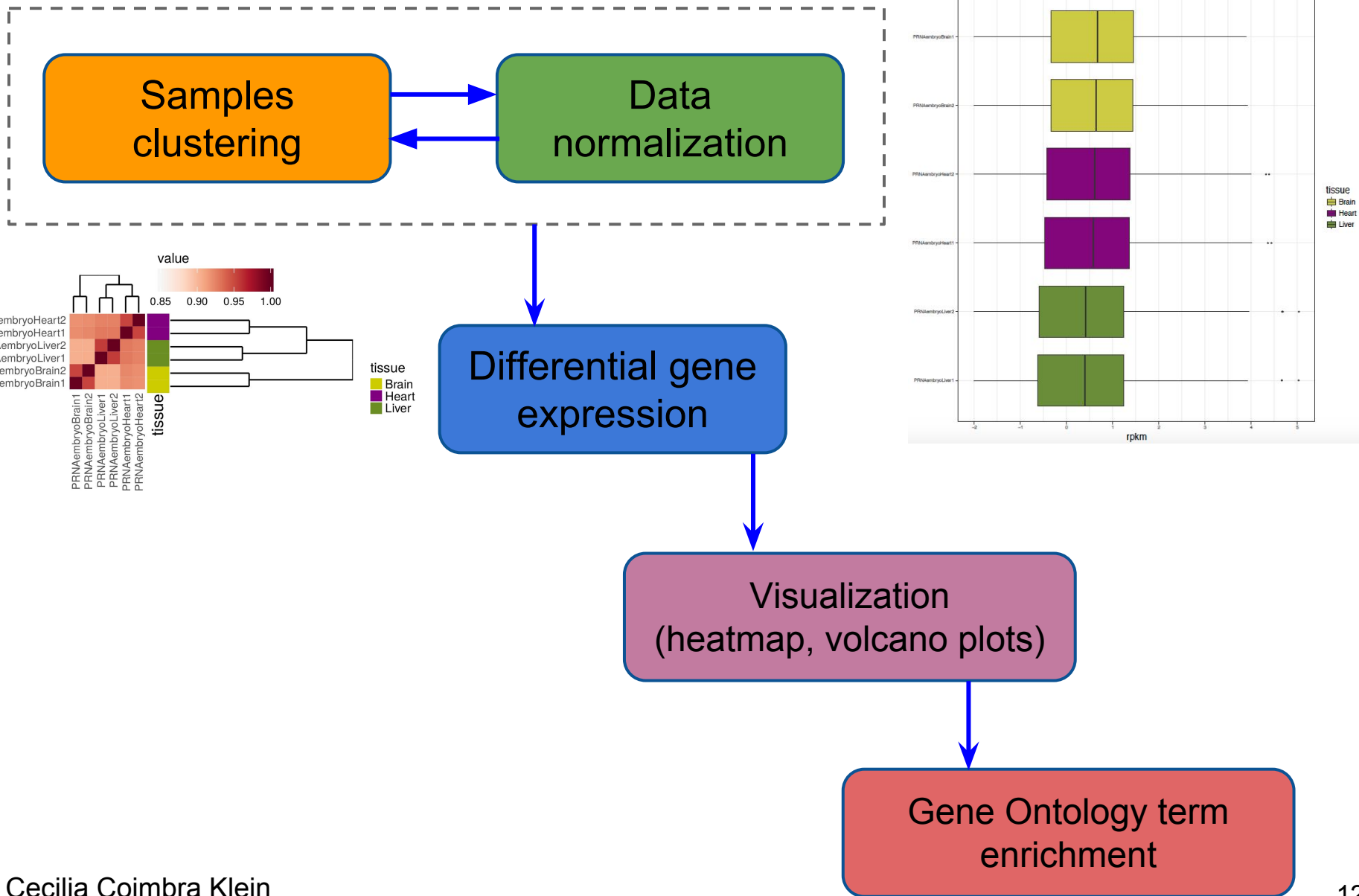
# Analysis pipeline



Cecilia Coimbra Klein

# Analysis pipeline



Samples clustering → Data normalization

Differential gene expression

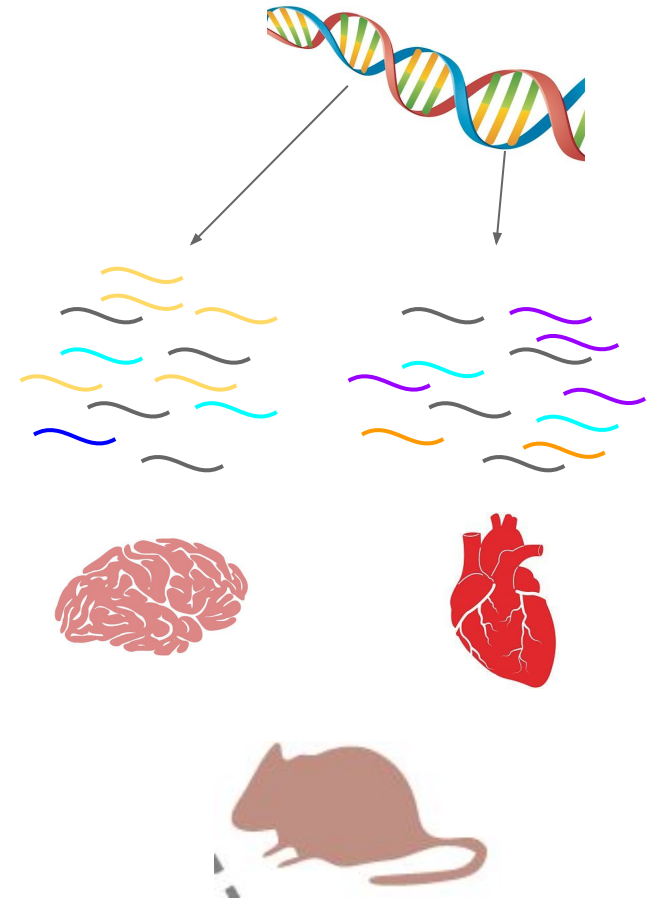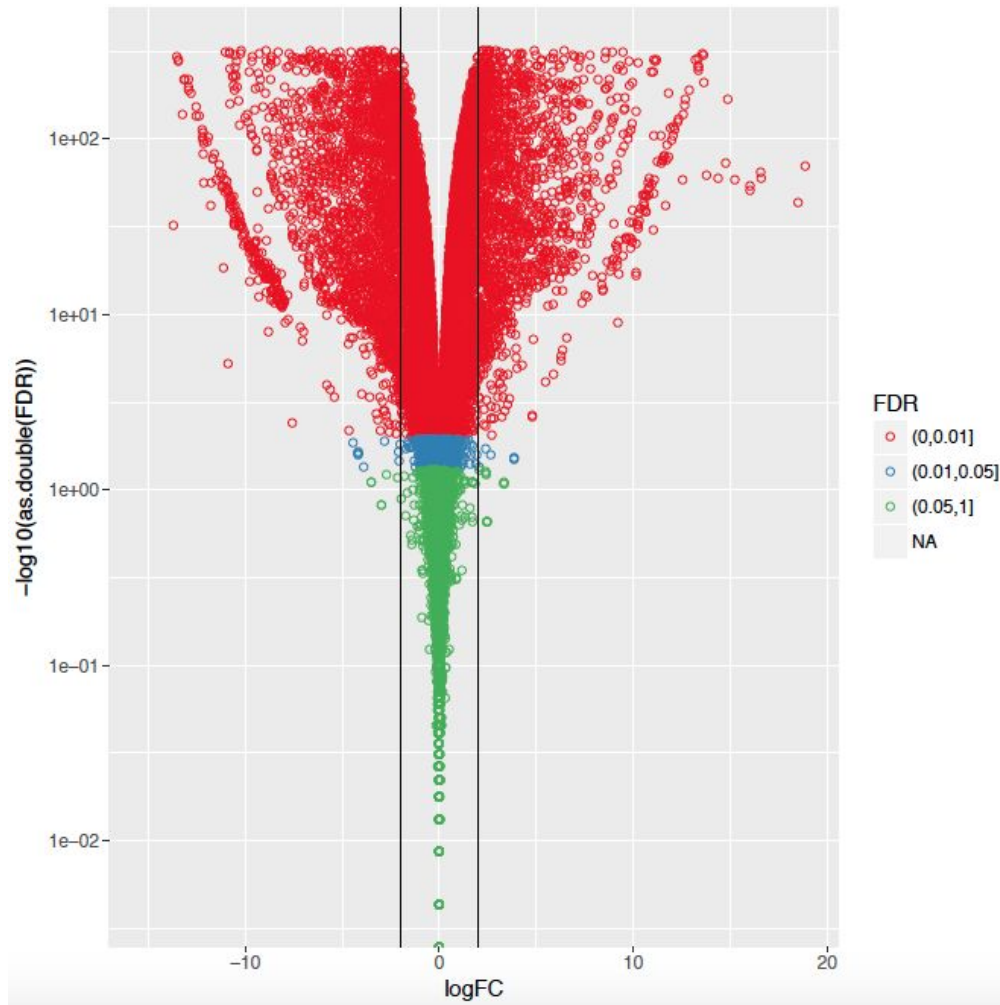Visualization (heatmap, volcano plots)

Gene Ontology term enrichment

Cecilia Coimbra Klein

10

# Samples clustering

# Analysis pipeline



Cecilia Coimbra Klein

# Differential Gene Expression

# Differential Gene Expression

# Analysis pipeline



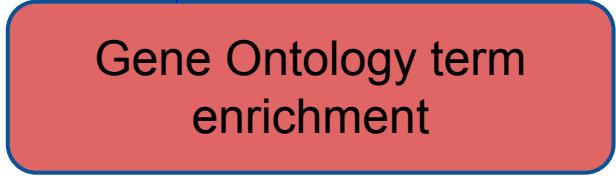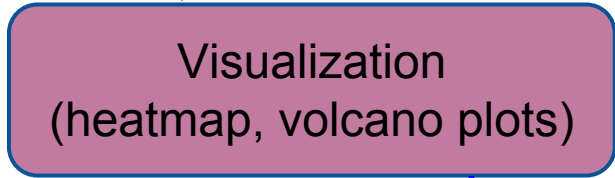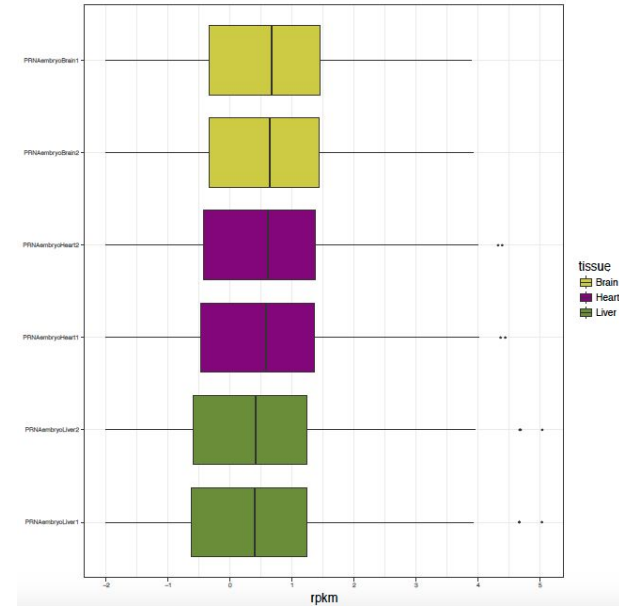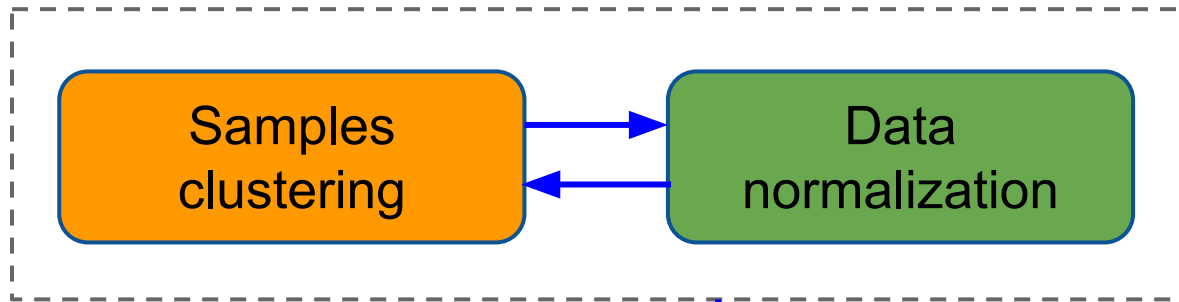Samples clustering ⇄ Data normalization → Differential gene expression → Visualization (heatmap, volcano plots) → Gene Ontology term enrichment
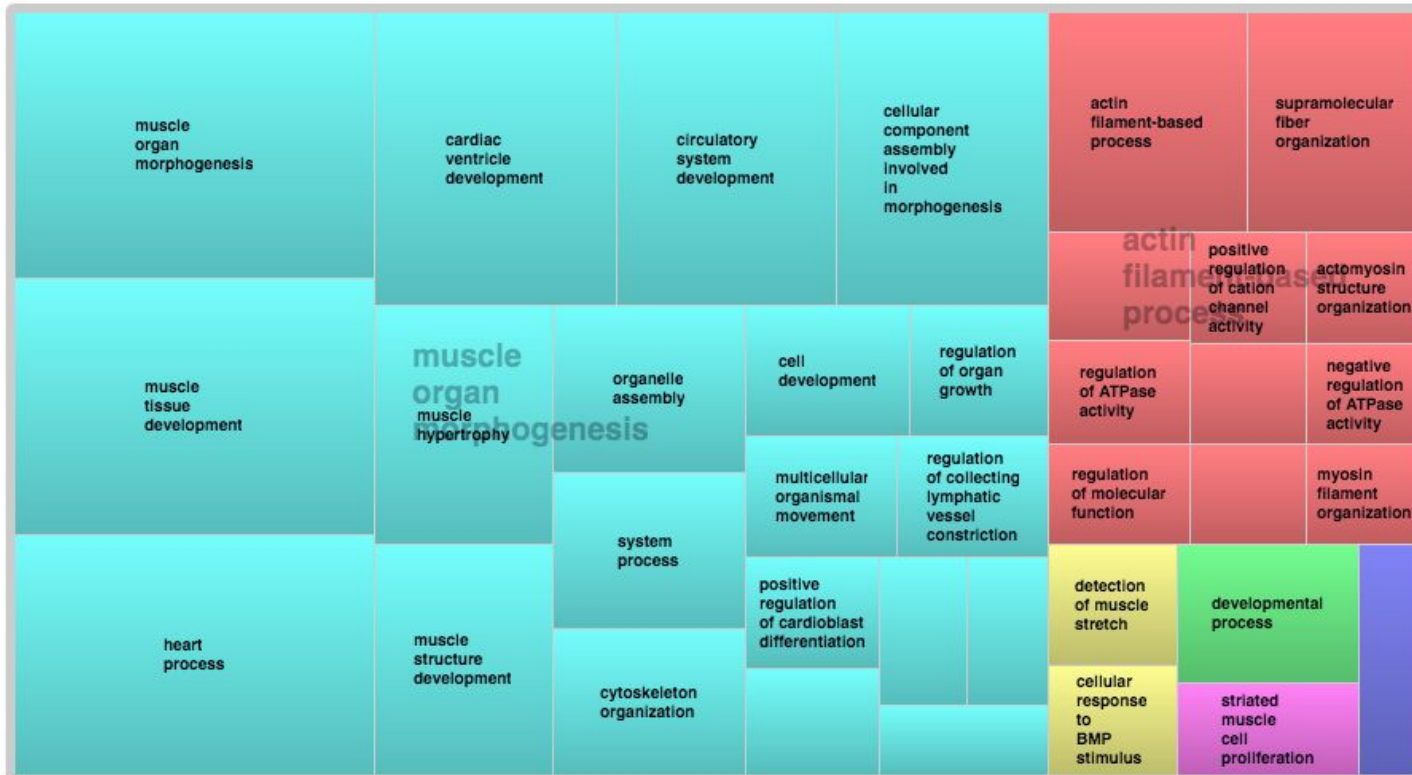
Cecilia Coimbra Klein

# GO term enrichment

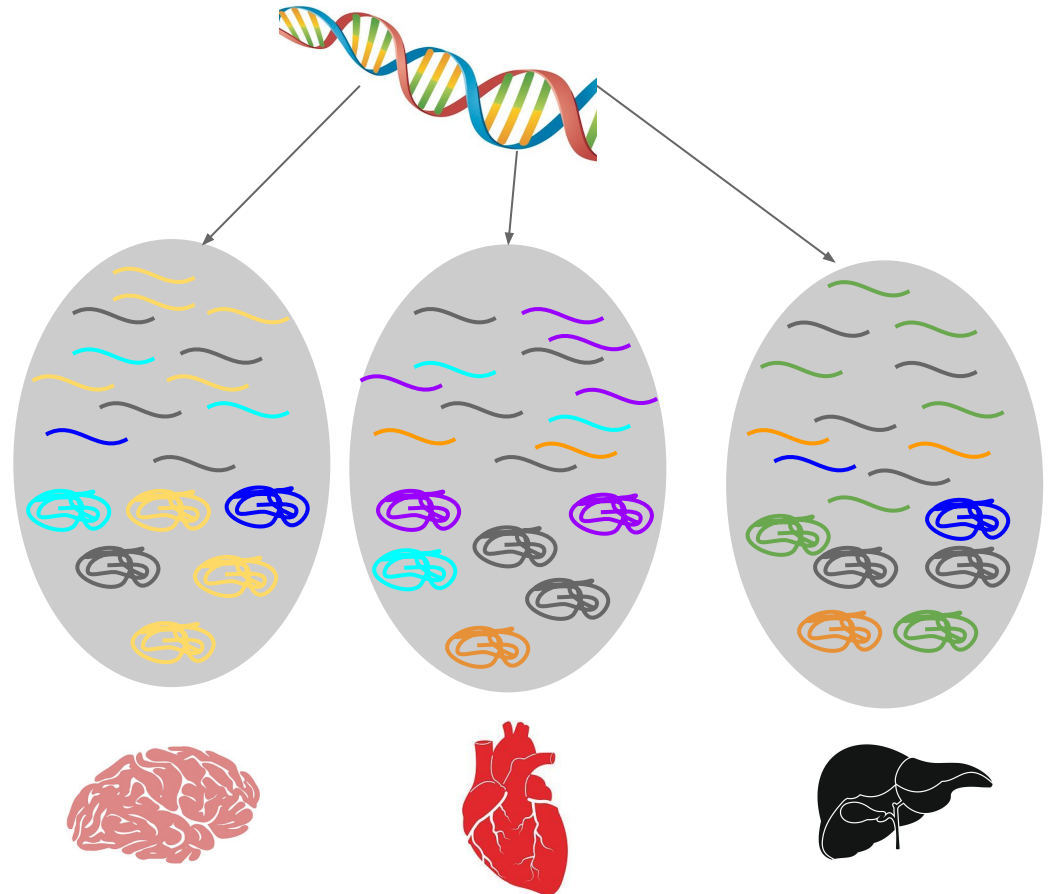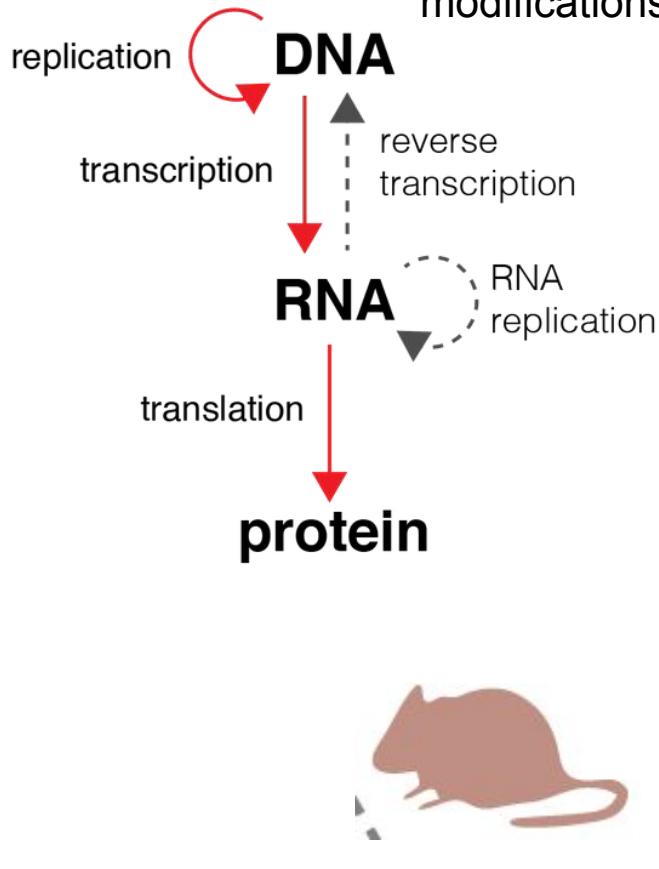# GO term enrichment

# Day 3: RNA-seq analysis (splicing)

# Molecular biology dogma



- The genome is identical in all cell types, however not all cell types have the same function. That's why the transcriptome (and the epigenome) becomes also relevant.

Cecilia Coimbra Klein

19

# Isoform usage

# SUPPA: generate events based on gene annotation

https://bitbucket.org/regulatorygenomicsupf/suppa

# SUPPA: generate events based on gene annotation



```
# number of alternative events with PSI values
ggbarplot.R -i input.tsv -o number_of_events.pdf --title "Events"
--y_title "Number of AS Events" --x_title "Event types" --palette_fill
/tutorial/palettes/cbbPalette.8.txt  -f 1
```

Cecilia Coimbra Klein

# SUPPA: Quantify event inclusion levels (PSIs)



```
# number of alternative single exon skipping
ls *.SE.psi|while read f;do echo -e $f"\t"$(grep -v nan $f |wc -l);done |
ggbarplot.R -i stdin -o number_of_SE.pdf -f 1 --palette_fill
/tutorial/palettes/palTissue.txt  --title "Skipping Exon (SE)" --y_title
"Number of Events" --x_title "Tissues"
```

Cecilia Coimbra Klein

# SUPPA: Quantify event inclusion levels (PSIs)



Skipping exon — SE

Alternative first exon — AF

## Skipping Exon (SE)

Number of Events / Tissues

V1
- Brain.SE.psi
- Heart.SE.psi
- Liver.SE.psi

## Alternative First Exon (AF)

Number of Events / Tissues

V1
- Brain.AF.psi
- Heart.AF.psi
- Liver.AF.psi

```
# number of alternative first exons
ls *.AF.psi|while read f;do echo -e $f"\t"$(grep -v nan $f |wc -l);done |
ggbarplot.R -i stdin -o number_of_AF.pdf -f 1 --palette_fill
/tutorial/palettes/palTissue.txt  --title "Alternative First Exon (AF)"
--y_title "Number of Events" --x_title "Tissues"
```

Cecilia Coimbra Klein

# SUPPA: compare conditions



- SUPPA calculates the magnitude of splicing change (ΔPSI) and their significance across multiple biological conditions, using two or more replicates per condition.

- Statistical significance is calculated by comparing the observed ΔPSI between conditions with the distribution of the ΔPSI between replicates as a function of the gene expression (measured as the expression of the transcripts defining the events).

https://bitbucket.org/regulatorygenomicsupf/suppa

Cecilia Coimbra Klein

# Skipping Exon (SE)

# **Skipping Exon (SE)**



- Select top events from pairwise comparison of Brain and heart
- p-value < 0.05
- $\Delta$PSI > 0.5
  - PSI = 1 inclusion of exon
  - PSI = 0 exclusion of exon

```
# prepare input for heatmap
event=SE; awk 'BEGIN{FS=OFS="\t"}NR>1 && $2!="nan" && ($2>0.5 || $2<-0.5)
&& $3<0.05{print}' DS.${event}.dpsi|cut -f1 > top-examples-SE.txt
selectMatrixRows.sh top-examples-SE.txt DS.SE.psivec >
matrix.top-examples-SE.tsv

# heatmap SE
ggheatmap.R -i matrix.top-examples-SE.tsv -o heatmap_top-examples-SE.pdf
--matrix_palette /tutorial/palettes/palSequential.txt --row_dendro
--matrix_fill_limits "0,1" -B 8
```

Cecilia Coimbra Klein

# Alternative First exon (AF)



```
# prepare input for heatmap
event=AF; awk 'BEGIN{FS=OFS="\t"}NR>1 && $2!="nan" && ($2>0.5 || $2<-0.5)
&& $3<0.05{print}' DS.${event}.dpsi|cut -f1 > top-examples-AF.txt
selectMatrixRows.sh top-examples-AF.txt DS.AF.psivec >
matrix.top-examples-AF.tsv

# heatmap alternative first exons top examples
ggheatmap.R -i matrix.top-examples-AF.tsv -o heatmap_top-examples-AF.pdf
--matrix_palette /tutorial/palettes/palSequential.txt --row_dendro
--matrix_fill_limits "0,1" -B 8
```
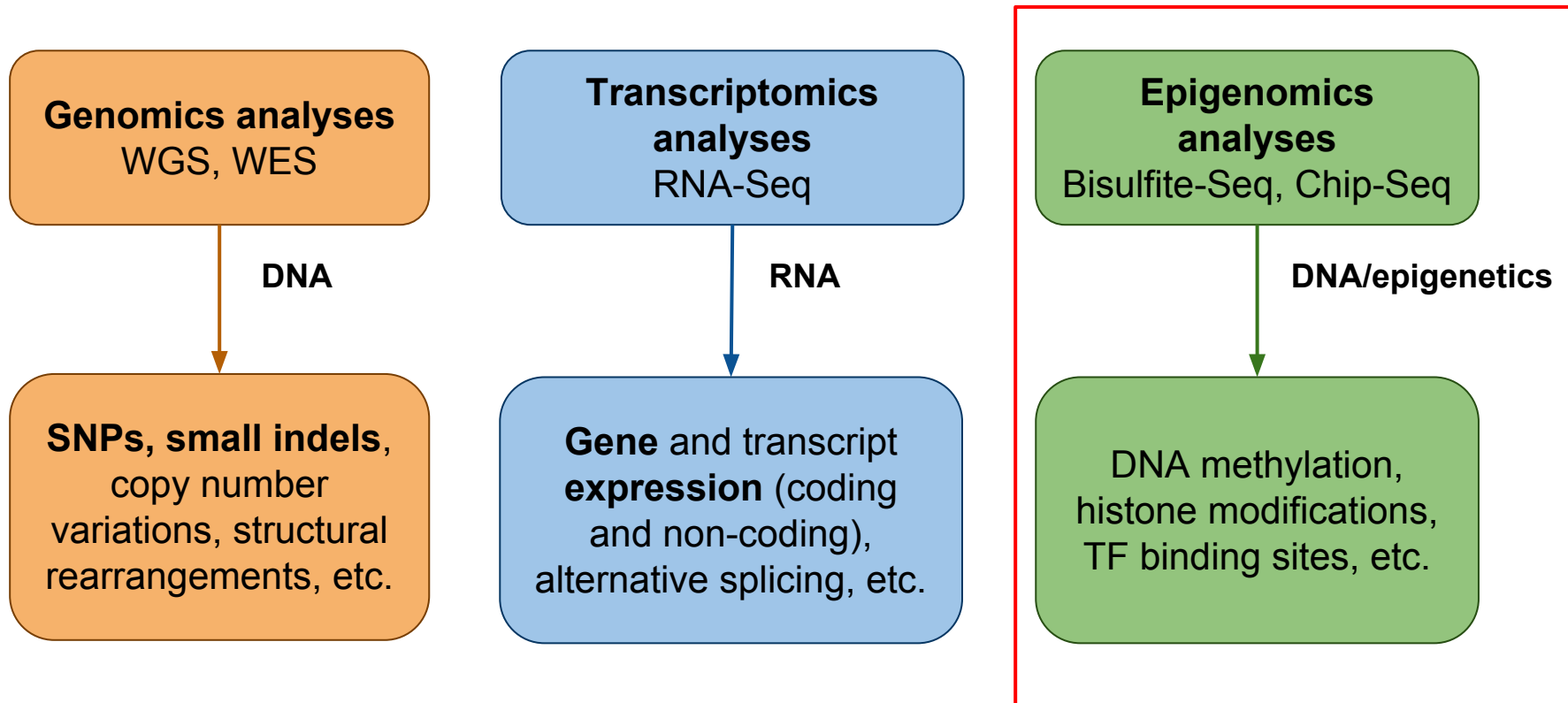
Cecilia Coimbra Klein

# Which *-Seq do I need?

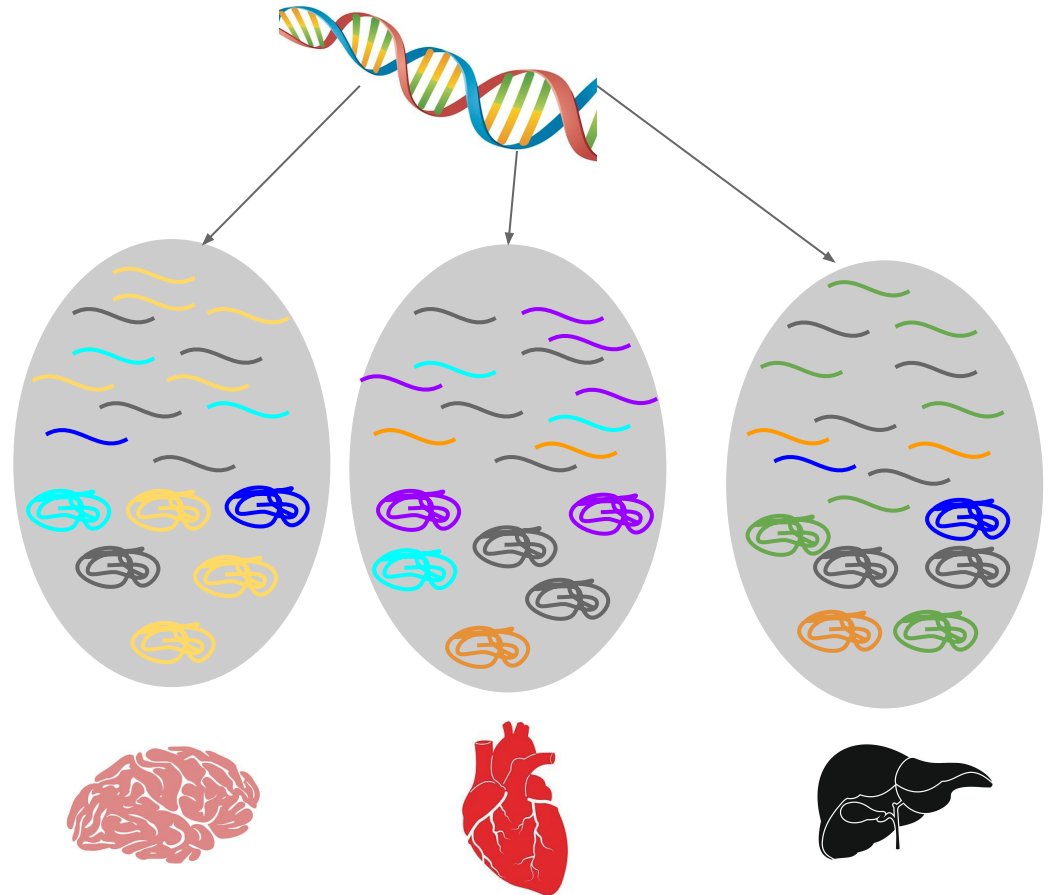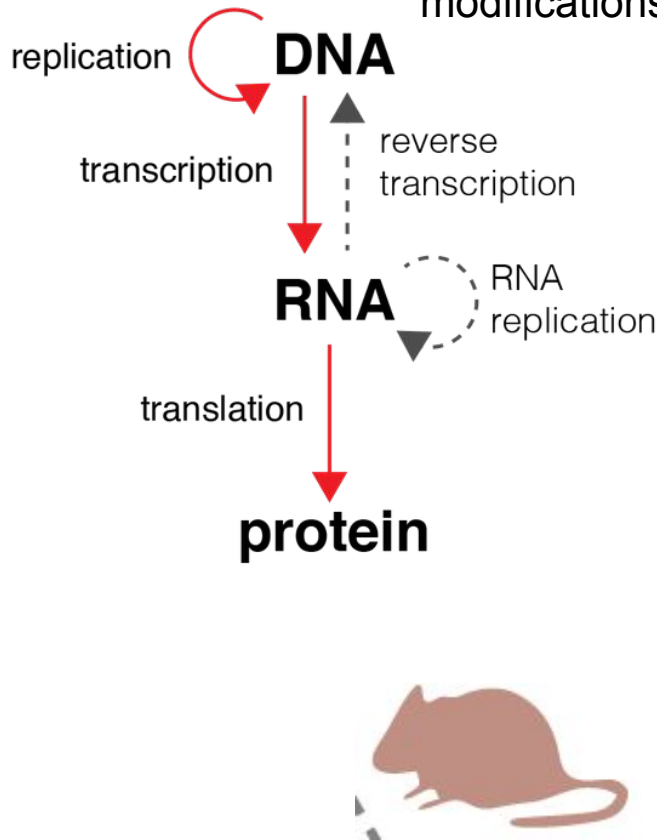| Genomics analyses<br>WGS, WES | Transcriptomics analyses<br>RNA-Seq | Epigenomics analyses<br>Bisulfite-Seq, Chip-Seq |
|---|---|---|
| ↓ **DNA** | ↓ **RNA** | ↓ **DNA/epigenetics** |
| **SNPs, small indels**, copy number variations, structural rearrangements, etc. | **Gene** and transcript **expression** (coding and non-coding), alternative splicing, etc. | DNA methylation, histone modifications, TF binding sites, etc. |

- Learn more about your favourite *-Seq here!

- Note that we are always talking about *re-sequencing*, which is something different from *de novo sequencing* (what is done for a new genome assembly)

Cecilia Coimbra Klein

# Basic concepts

# Molecular biology dogma

epigenetic modifications



- The genome is identical in all cell types, however not all cell types have the same function. That's why the transcriptome (and the epigenome) becomes also relevant.

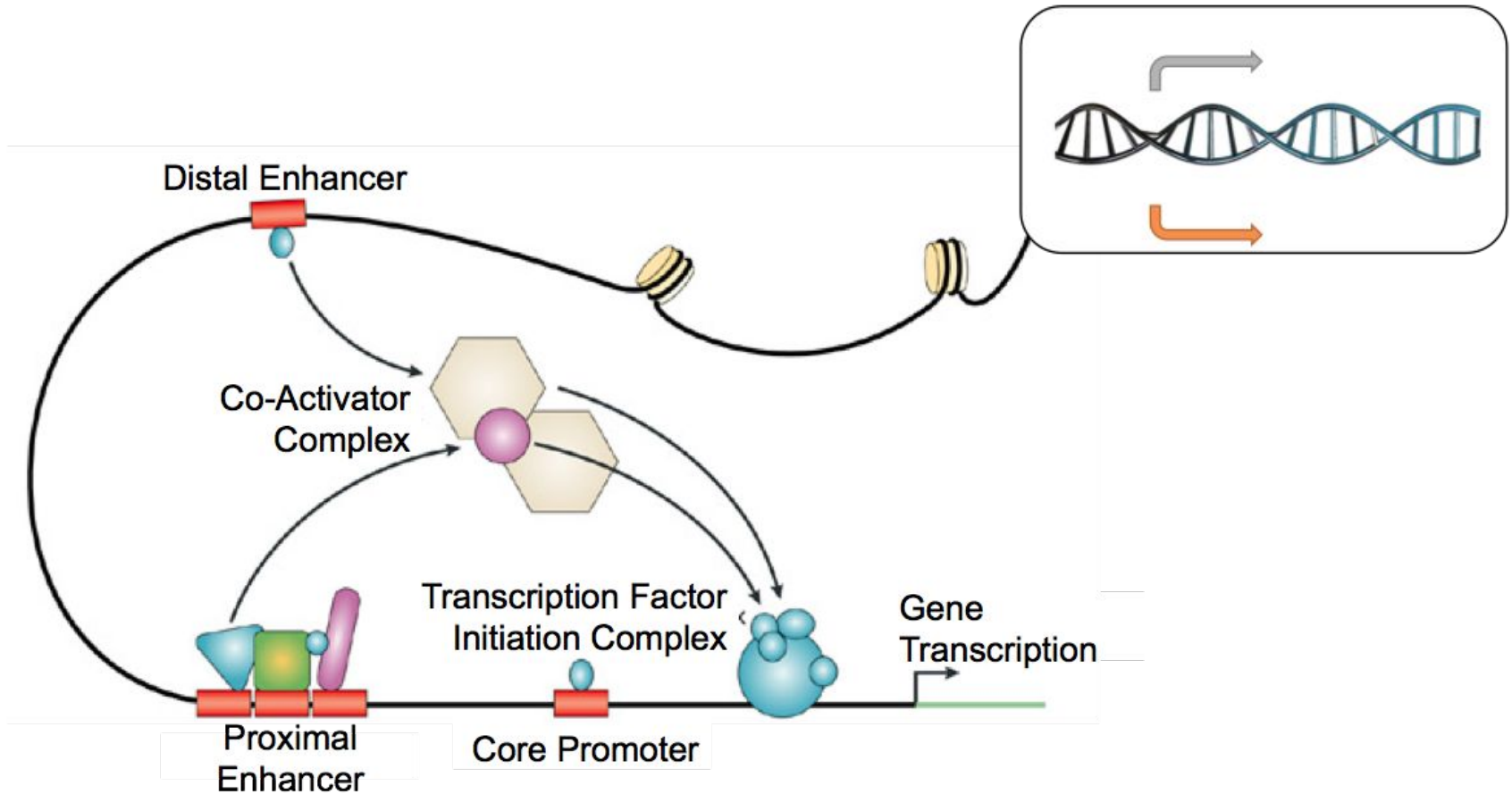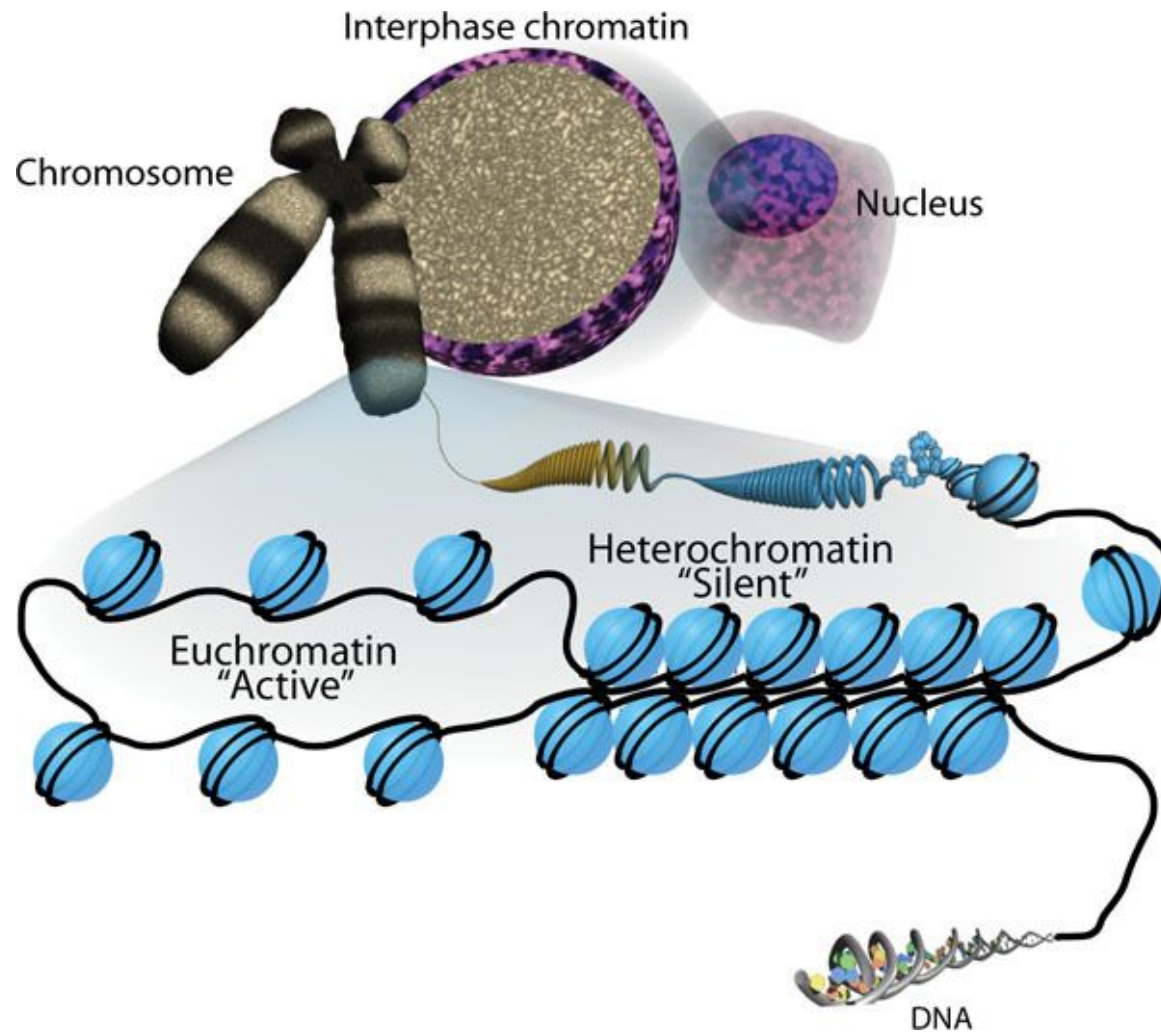Cecilia Coimbra Klein

31

# Dynamics of gene regulation



Cecilia Coimbra Klein

# Chromatin organization



The chromatin signature of pluripotent cells, StemBook, NCBI
https://www.ncbi.nlm.nih.gov/books/NBK27041/figure/thechromatinsignature.F1/

# Histone modifications

N-terminal histone tail

double-stranded DNA

H3   H4

H2A   H2B

histone proteins

nucleosomes

chromatin

chromosome

Acetyl Group → Ac

Histone Tail →

DNA →

H2A   H2B

H3   H4

Methyl Group → Me
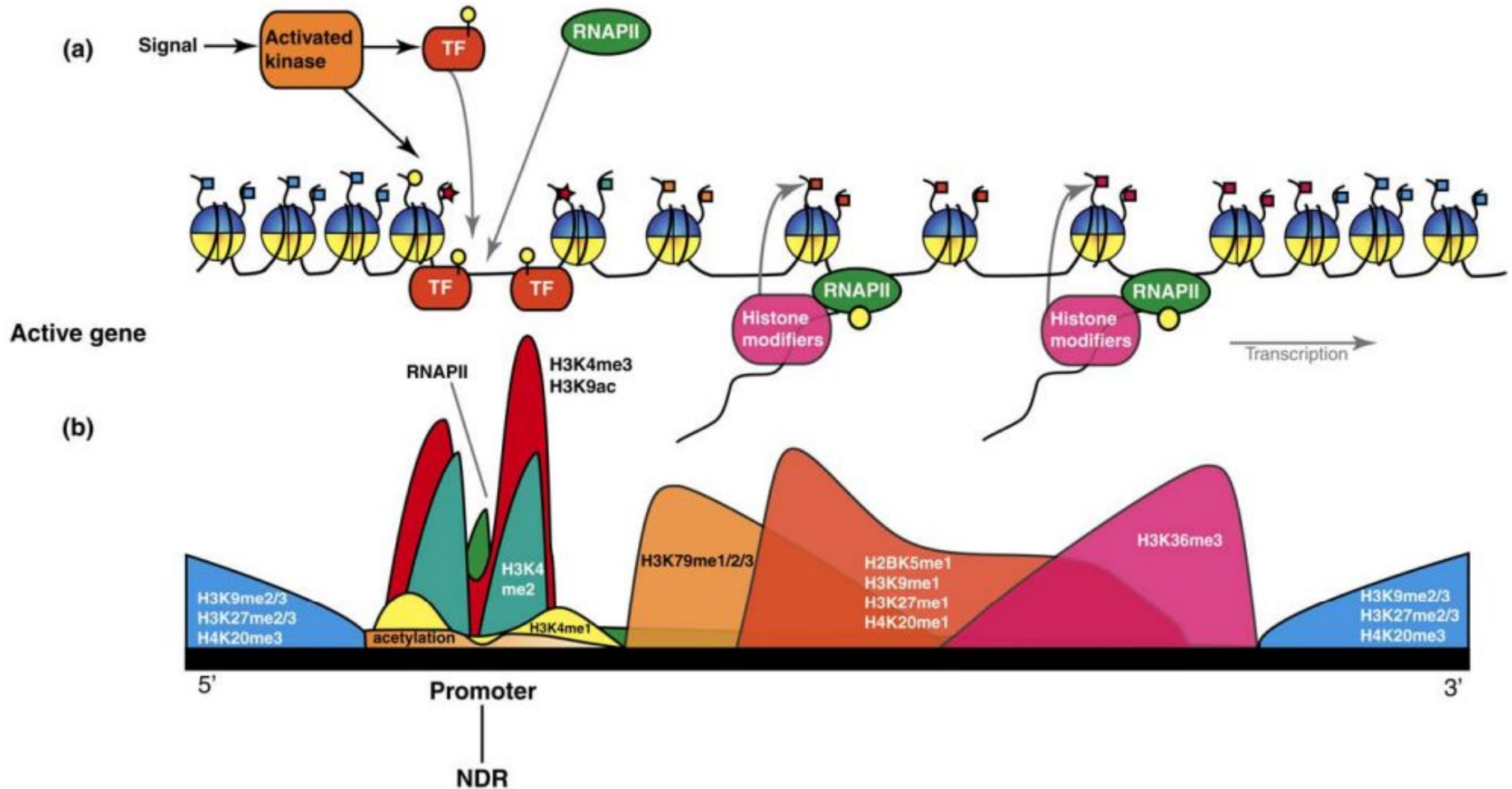
Histone Octamer

- methylation, acetylation and phosphorylation
- involve covalent post-translational modifications mostly to the residues at the positively charged N-terminal tails of histones.
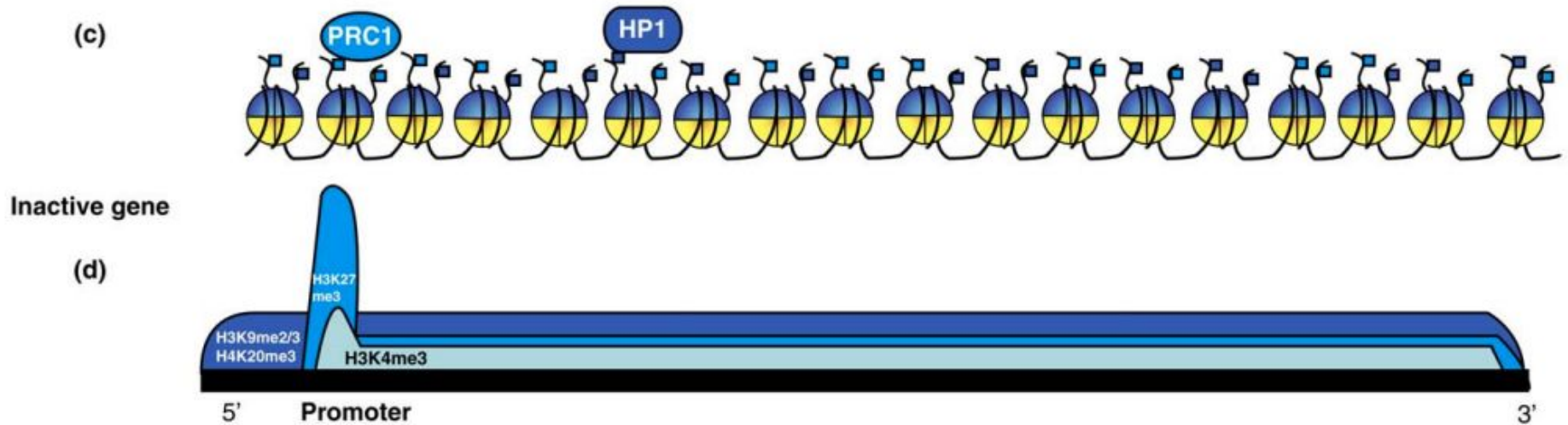
https://www.whatisepigenetics.com/histone-modifications/
https://www.epigentek.com/catalog/advanced-epigenetic-overview-of-histone-modifications-n-5.html?currency=es
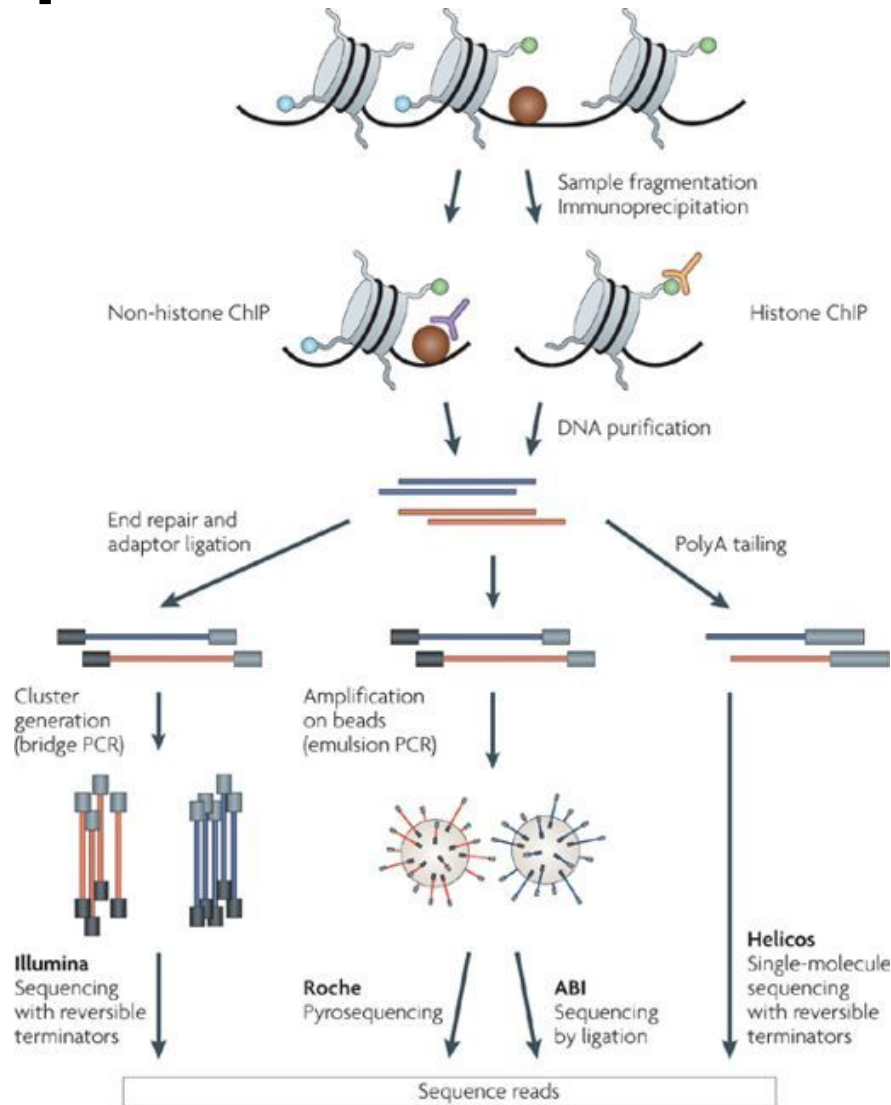
Cecilia Coimbra Klein

# Histone modifications



Barth & Imhof (2010) Trends in Biochemical Sciences
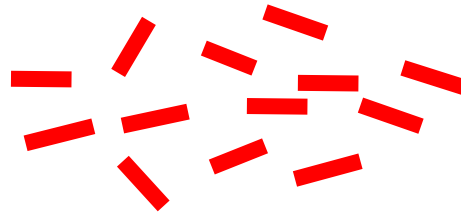
Cecilia Coimbra Klein

# Histone modifications
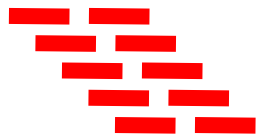
# ChIP-seq: Chromatin ImmunoPrecipitation
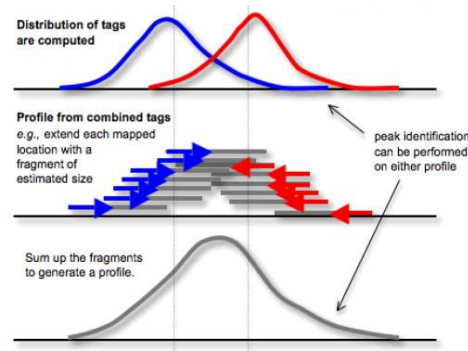


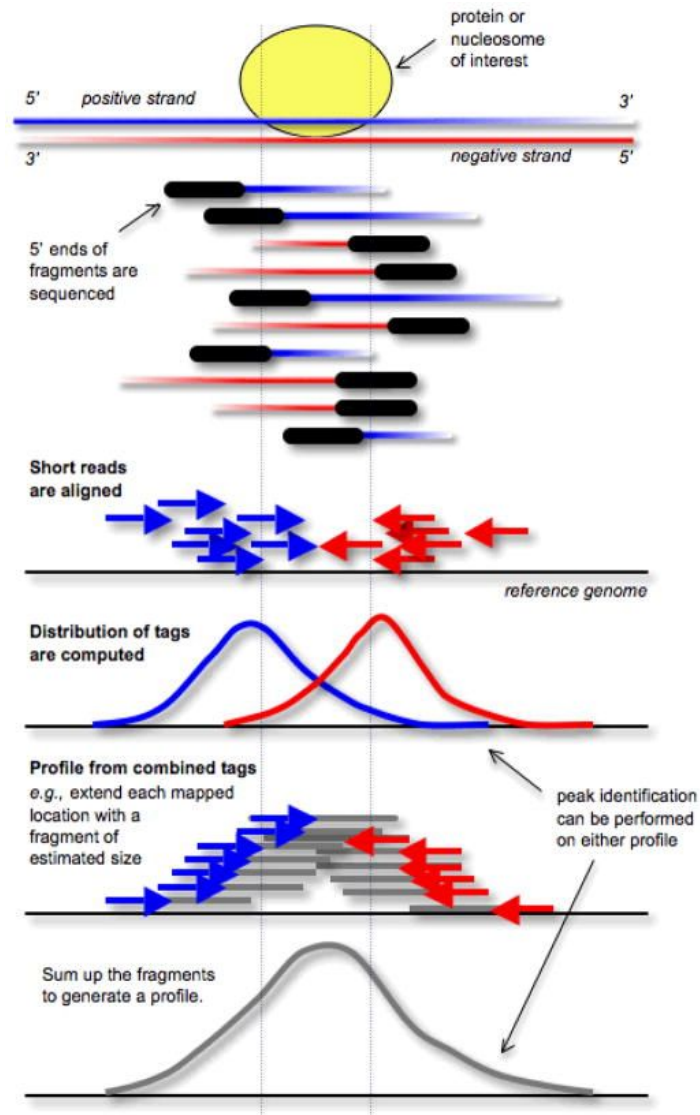Park (2009) doi: 10.1038/nrg2641

# Processing

# Mapping

**Mapping**

**Find a correspondence between the query sequences (ChIP-seq reads) and our prior knowledge (reference genome sequence).**
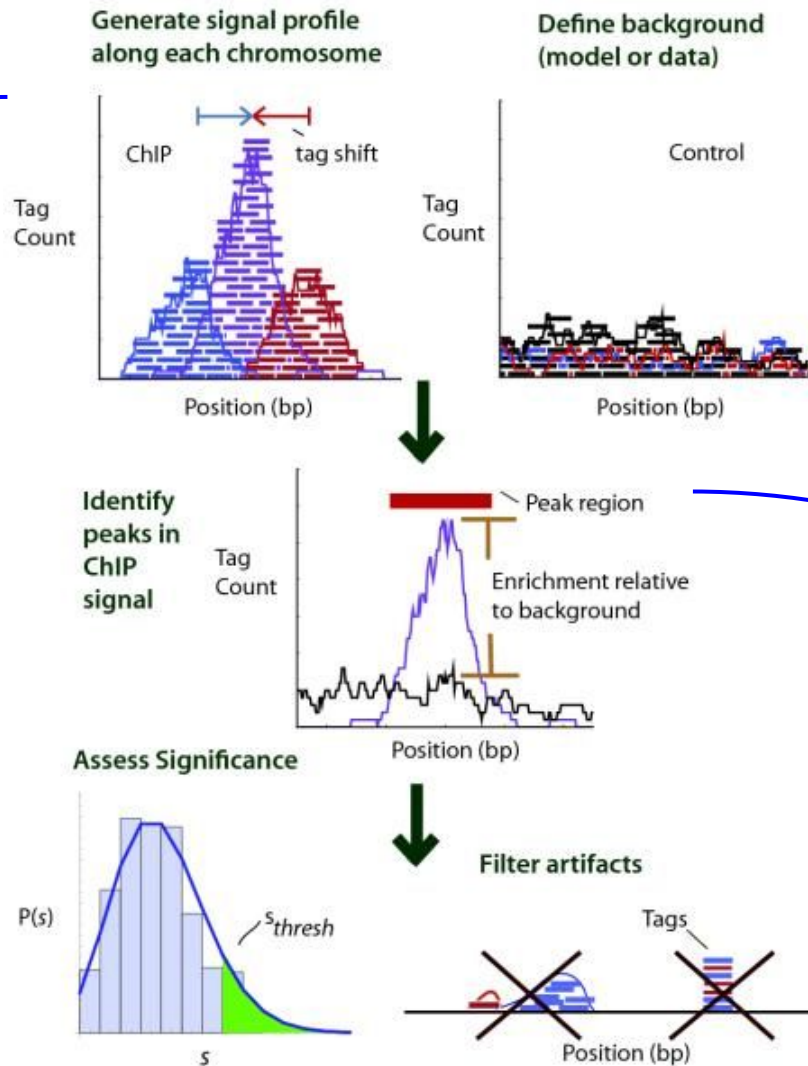
*reference*

**Peak calling**

Distribution of tags are computed

Profile from combined tags
e.g., extend each mapped location with a fragment of estimated size

peak identification can be performed on either profile

Sum up the fragments to generate a profile.

Cecilia Coimbra Klein

# Peak calling

Cecilia Coimbra Klein

# Peak calling: MACS2
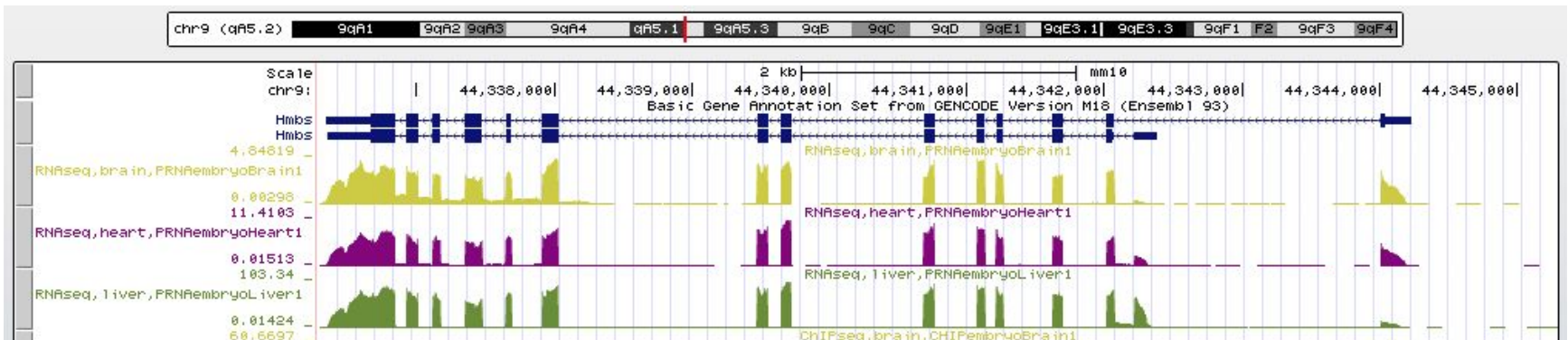


Signal profile (bigWig)

Peak regions (BED)

Pepke et al. (2009) doi: 10.1038/nmeth.1371
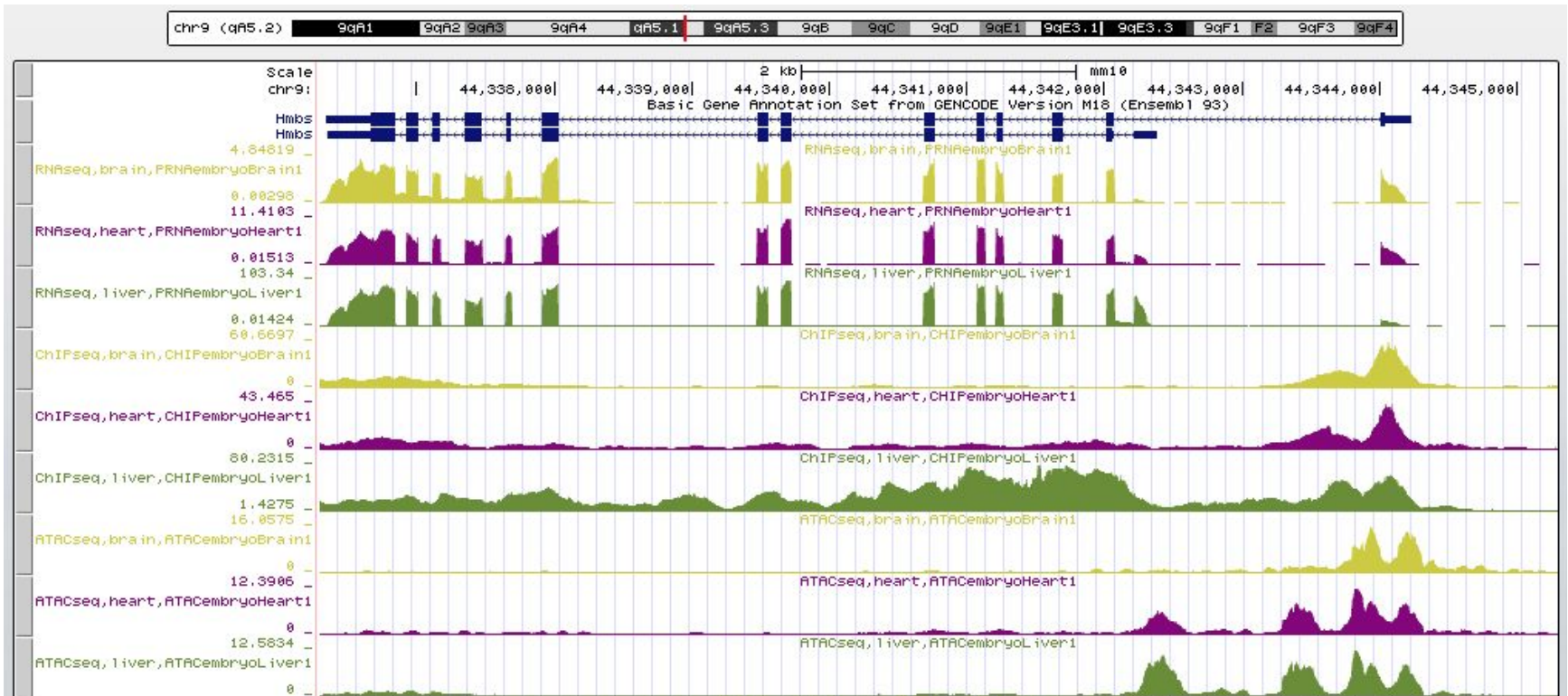
# Visualization

# RNA-seq signal

genome-euro.ucsc.edu



- expected read depth at each position in the genome
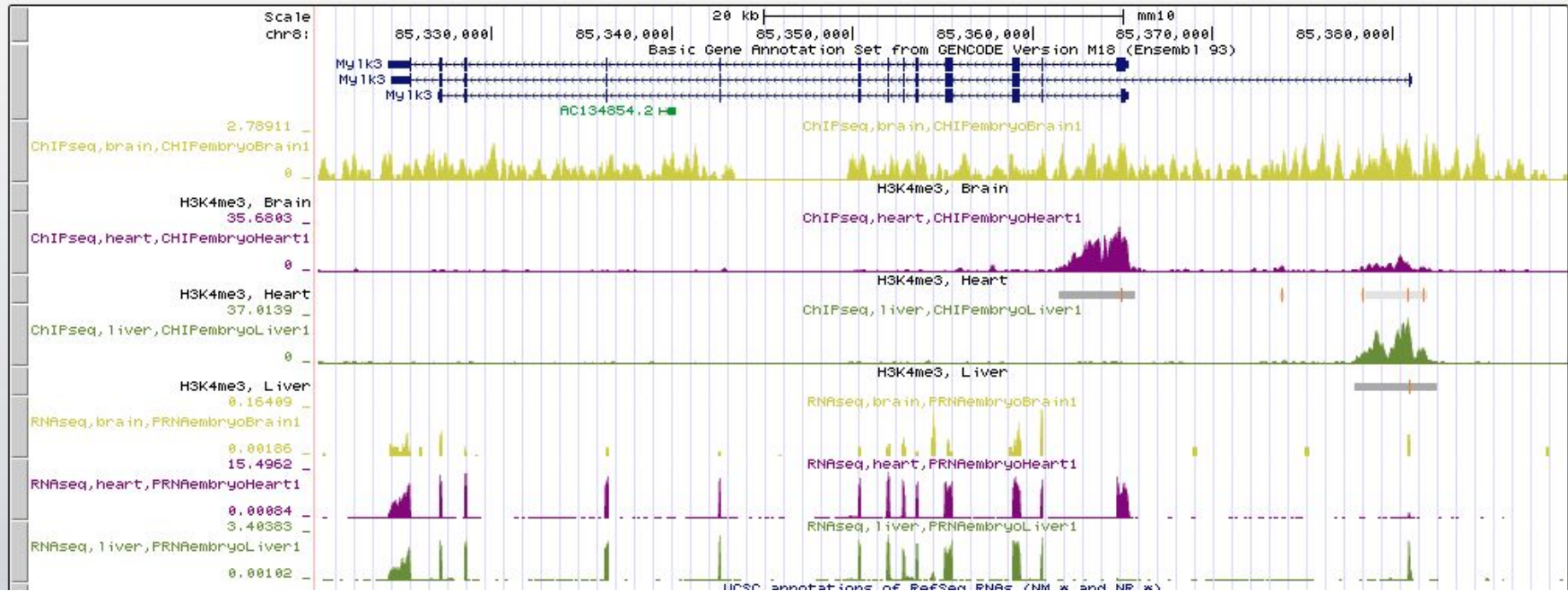- can be normalized (e.g. RPM, reads per million reads)
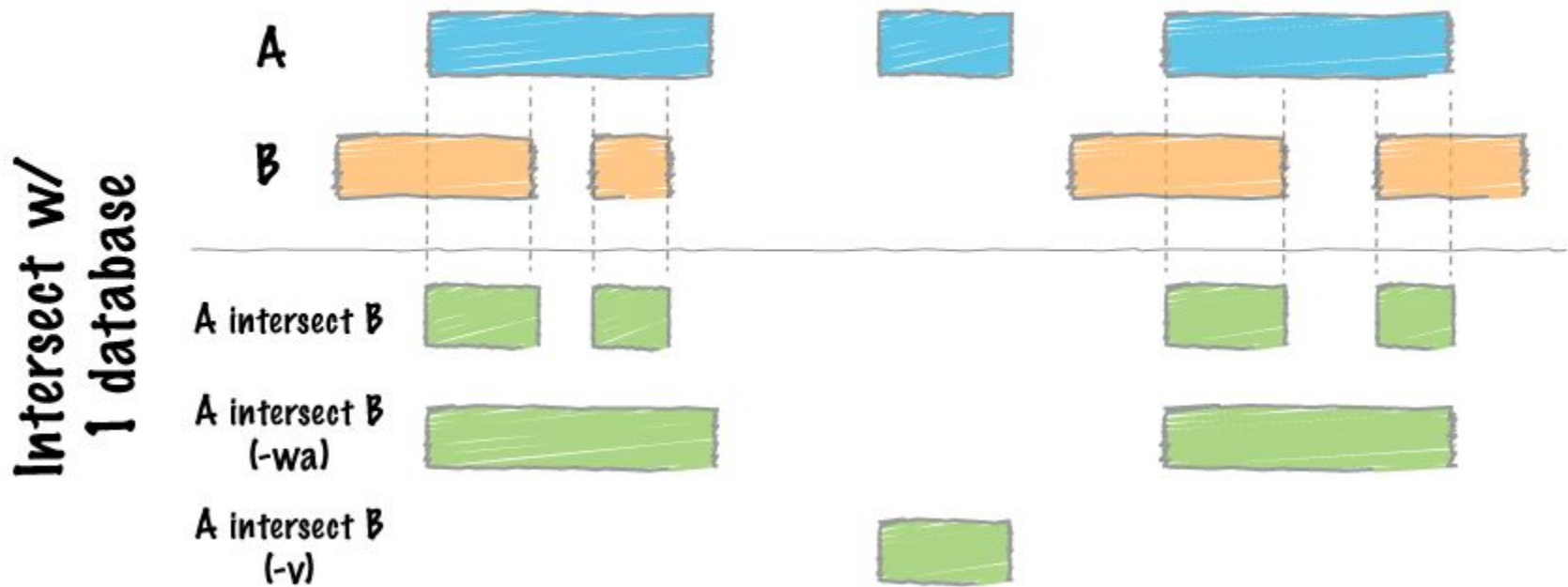
# UCSC: signal files

genome-euro.ucsc.edu

# UCSC

# Analysis

# BEDTools intersect



https://bedtools.readthedocs.io/en/latest/content/tools/intersect.html

Cecilia Coimbra Klein

# ATAC-seq

# Open chromatin

Cecilia Coimbra Klein
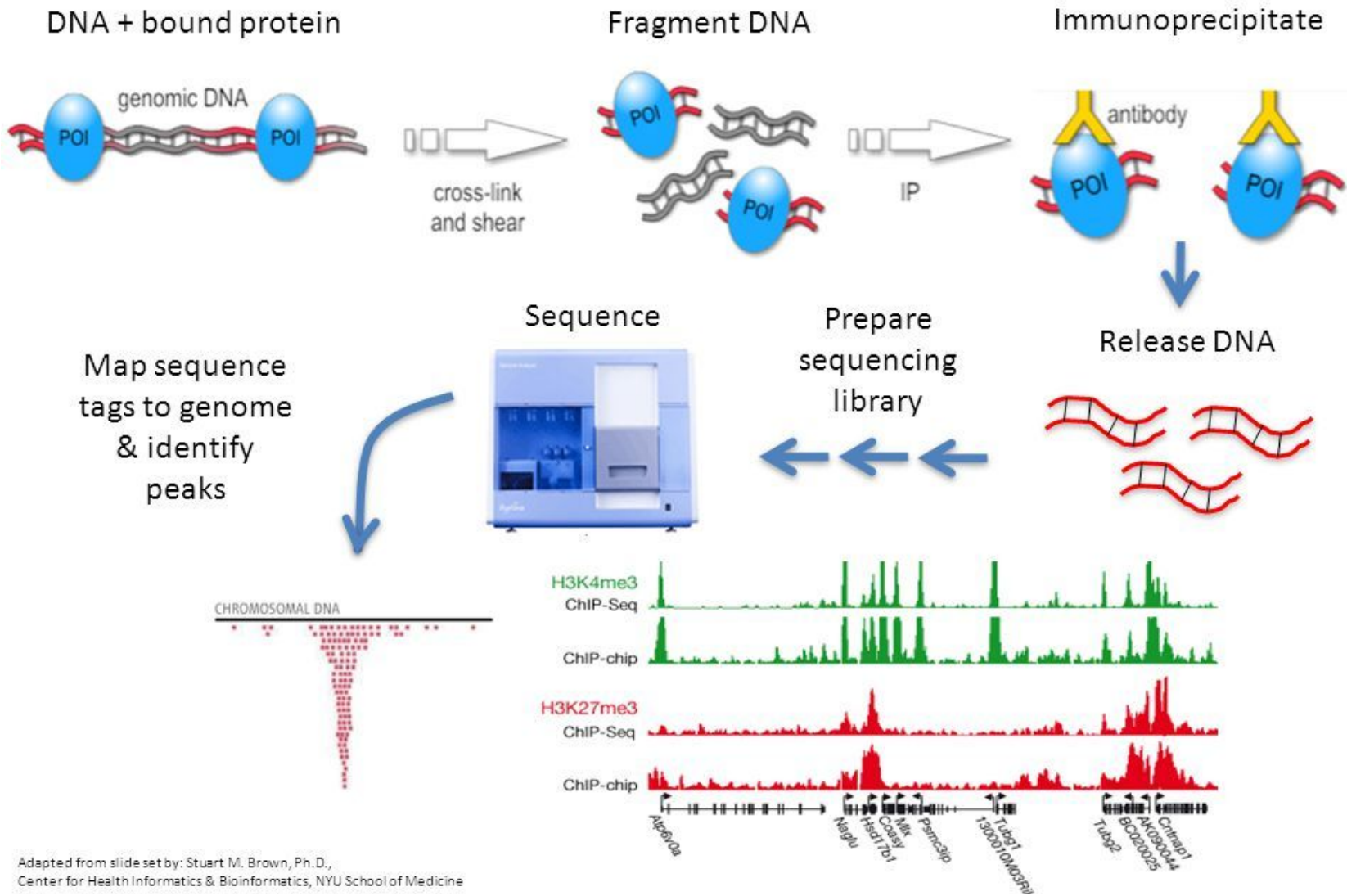
# Hands-on

**Setup environment 1**

**Integrative data analysis 7**

https://public_docs.crg.es/rguigo/Data/cklein/courses/UVIC/handsOn/

Cecilia Coimbra Klein

# **Additional slides**

# ChIP-seq: Chromatin ImmunoPrecipitation



Adapted from slide set by: Stuart M. Brown, Ph.D.,
Center for Health Informatics & Bioinformatics, NYU School of Medicine

Cecilia Coimbra Klein

# BEDTools intersect



https://bedtools.readthedocs.io/en/latest/content/tools/intersect.html

Cecilia Coimbra Klein

# UCSC



Cecilia Coimbra Klein