# Biological data analysis

## Single cell data analysis

Juan D. Montenegro

January 15th, 2024
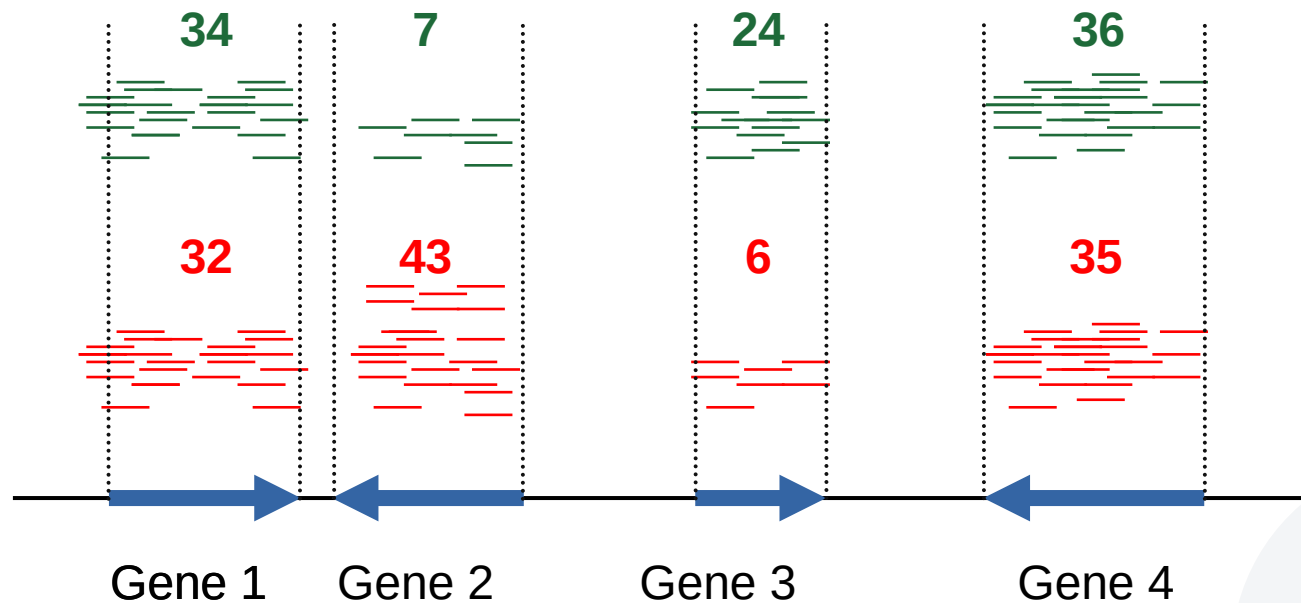
LibreOffice®

# Summarising read counts
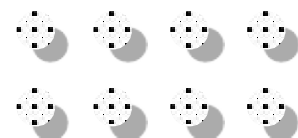
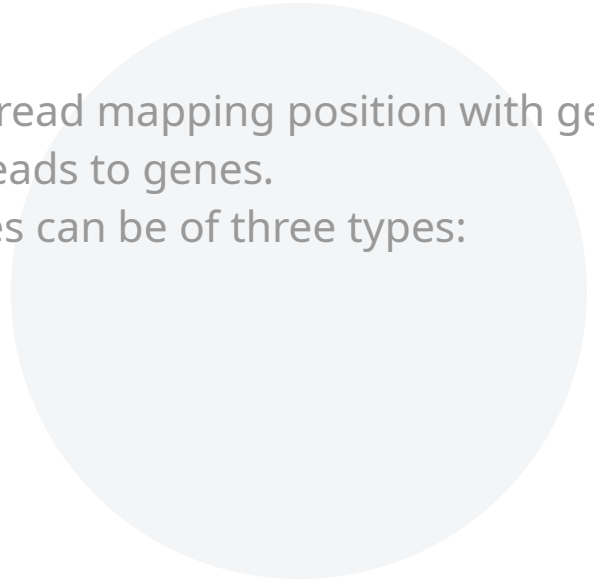|  | A | B | C |
|---|---|---|---|
| Gene 1 | 34 | 32 | 36 |
| Gene 2 | 2 | 0 | 16 |
| Gene 3 | 0 | 2 | 0 |
| Gene 4 | 4 | 6 | 10 |
| Gene 5 | 2 | 28 | 32 |
| Gene 6 | 7 | 32 | 33 |
| Gene 7 | 90 | 16 | 17 |
| Gene 8 | 13 | 0 | 13 |
| **TOTAL** | **152** | **116** | **157** |



Illustrations by Pixeltrue on
icons8

# Read assignment

Process that compares read mapping position with gene
locations and assigns reads to genes.
Reads assigned to genes can be of three types:
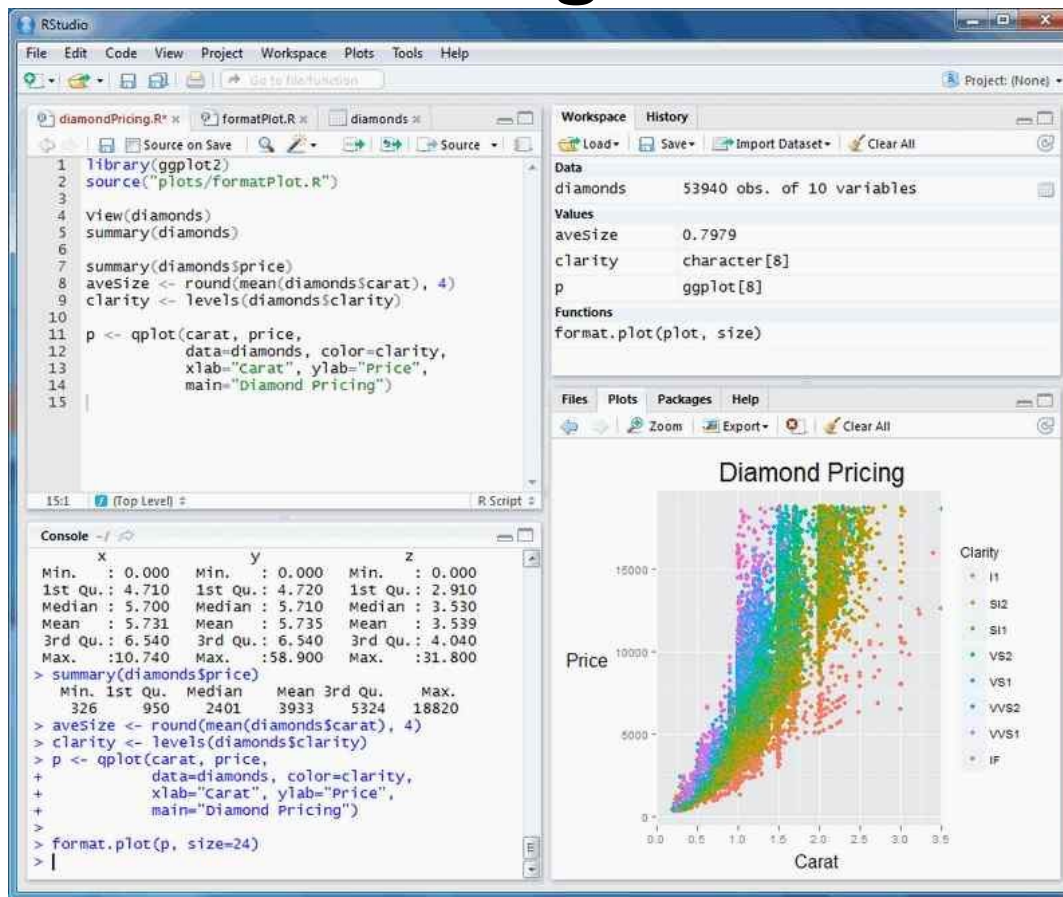    CDS
    UTR
    Intron

# Seventh Exercise

- Run featureCounts:
- Package Subread/ command featureCounts
- Write a slurm script and submit it.
- Review results:
  - Summary
  - Count Matrix

# Analising read counts
## Using R

# Eighth Exercise

- Open Rstudio on the LiSC
  - http://rstudio.lisc.univie.ac.at
- Basic commands:
  - getwd() / setwd()
  - c()
  - read.table() / read_tsv()
- Main libraries:
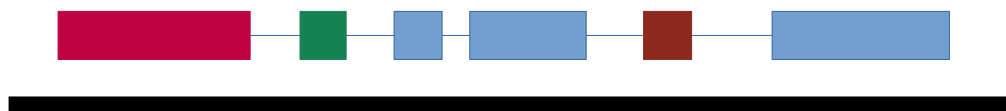  - tidyverse
  - ggplot2

# Eighth Exercise

- Plot basic statistics from the featureCounts results

- Write an R script in Rstudio, save it and share it on GitHub.

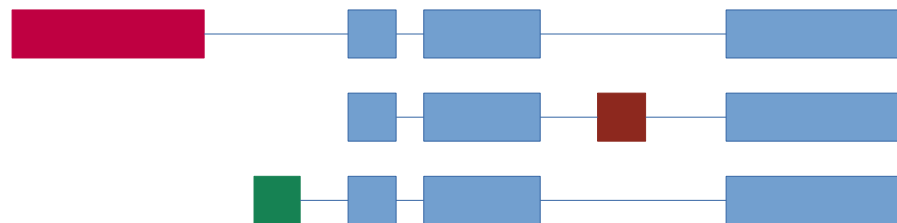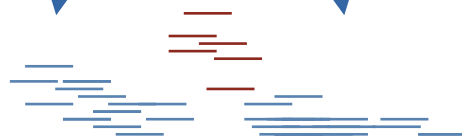- Perform PCA analysis to identify biological replicates

Reads map to multiple places in the transcriptome

# Eigth Exercise

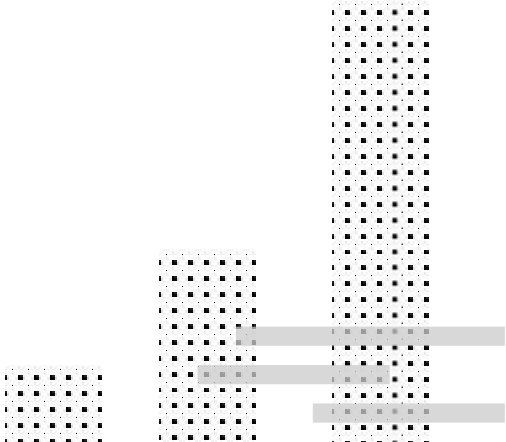- Download a reference transcriptome of Nematostella vectensis from the European Nucleotide Archive (ENA)

- Align the reads to the reference transcriptome

- Assign reads to genes

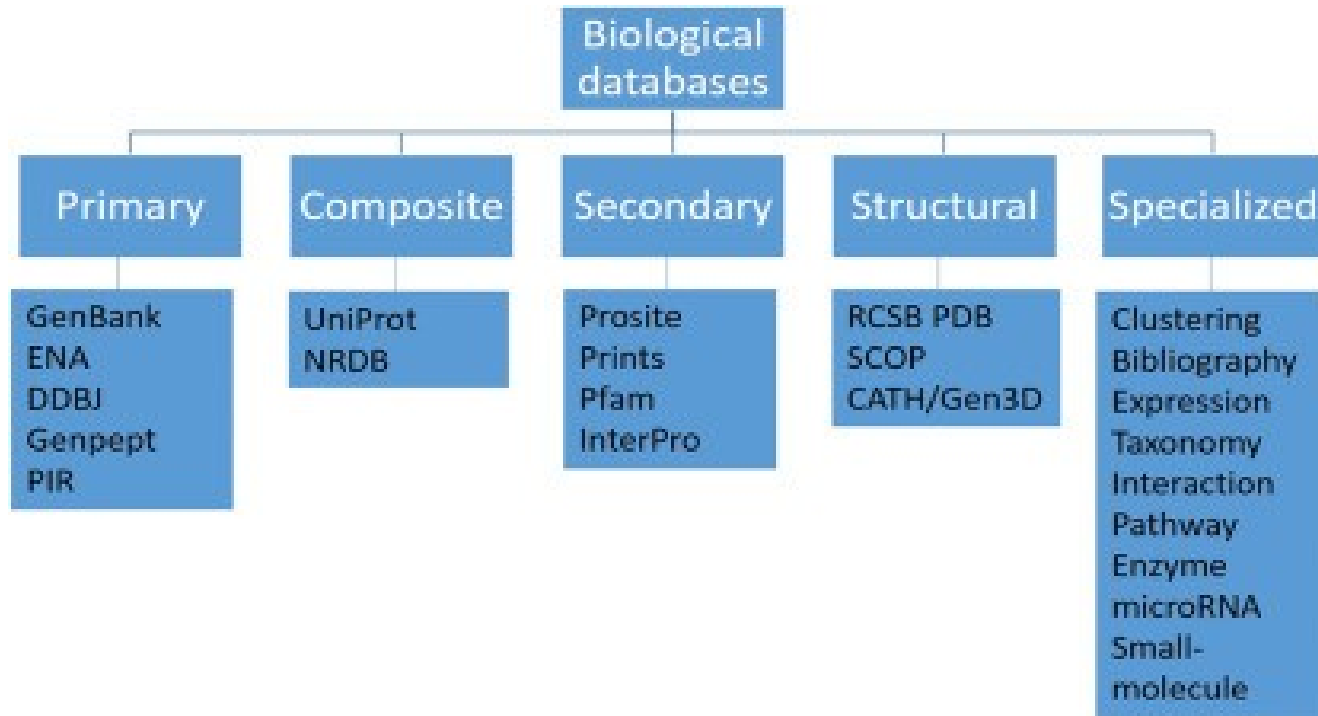- Compare mapping efficiency and assignment efficiency between transcriptome and genome using R

# Functional annotation of genes

Process of assigning functions to genes
Relies on high throughput comparison to large annotated databases

Sharma & Yadav (2022) Biological databases and their application;
Bioinformatics: Methods and applications:17-31

# InterProScan



## InterProScan 5: Large scale protein function classification

ft Nuka, Simon Potter, Siew-Yit Yong, Maxim Scheremetjew, Alex Mitchell, Matthew
aser and Rob Finn

ropean Bioinformatics Institute (EMBL-EBI), United Kingdom

**EMBL-EBI**

### Introduction
InterPro (http://www.ebi.ac.uk/interpro/) is a freely available resource that is used to classify sequences into protein families and to predict the presence of important domains and sites.

**InterProScan** (http://www.ebi.ac.uk/interpro/interproscan.html) is the underlying software that allows both protein and nucleic acid sequences to be scanned against InterPro's predictive models (signatures), which are provided by the resource's member databases.
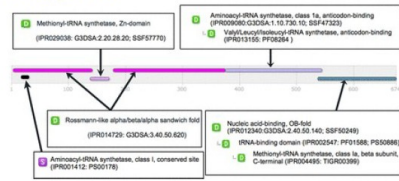
**Figure 1. InterProScan matches for UniProtKB protein Q3JCG5 showing predicted protein family membership, domains and sites.**

### Structure-Function Linkage Database (SFLD)
SFLD's hidden Markov models that offer structure-function mapping have also been incorporated in InterProScan. SFLD models allow evolutionary classification of related enzymes according to shared chemical functions to determine conserved active sites.
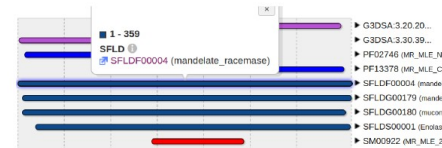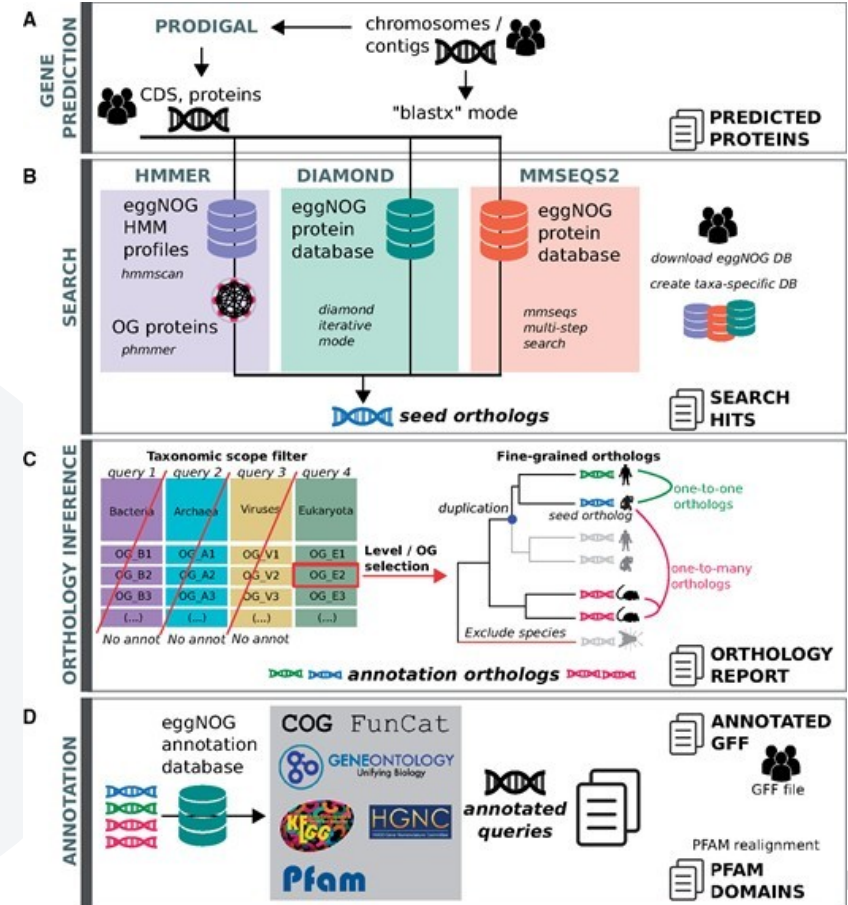
**Figure 3. InterProScan matches for UniProtKB protein T2HDW6. The matches include hits to SFLD signatures (SFLDF00004, SFLDG00179, SFLDG00180, SFLDS00001).**

### Performance improvements
Optimisation in the pipeline filters and database query refinements have improved throughput for large-scale protein sequence analysis and accelerated InterProScan domain searches by several orders of magnitude.

In Figure 4, we look at the performance of InteProScan since version 5.1-44.0, the first official release of InteProScan 5. We run InteProScan on over 120 million proteins from the UniProt

# EggNogMapper

# Ninth Exercise

- Use gffread to extract protein sequences from genome

- Use blastp to align proteins to uniref

- Use interProScan and EggNogMapper to add GO term annotation and identify conserved motifs

-

# Assessing the mapping tool

Using external and internal clues to determine how useful a mapping tool is.
1) Contiguity: Fragmentation, N50, NG50, AuN curve
2) Coherence: Size estimation, proper mapping of reads
3) Completeness: Mapping efficiency of different data sources (cloned genes, ESTs, RNAseq, BAC ends, proteins)
4) Correctness: DNA sequence variation compared to actual known sequences

# Contiguity:
## Assembly fragmentation depends on technology



DNA sequencing history

Maxam-Gilbert and Sanger sequencing

PacBio and Nanopore sequencing

1st generation   2nd generation   3rd generation
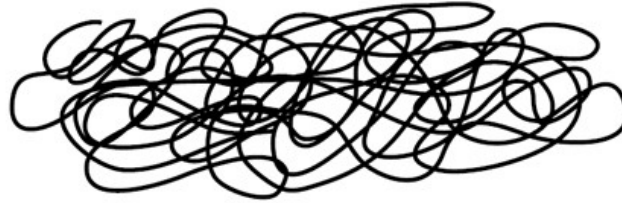
Pyro- and Illumina sequencing

Creator: Werner, Anina
https://www.integra-biosciences.com/canada/en/blog/article/dna-sequencing-methods-sanger-ngs

# Hierarchical shotgun sequencing
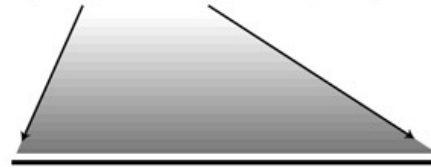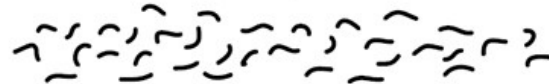
Genomic DNA

BAC library

Organized
mapped large
clone contigs

BAC to be
sequenced

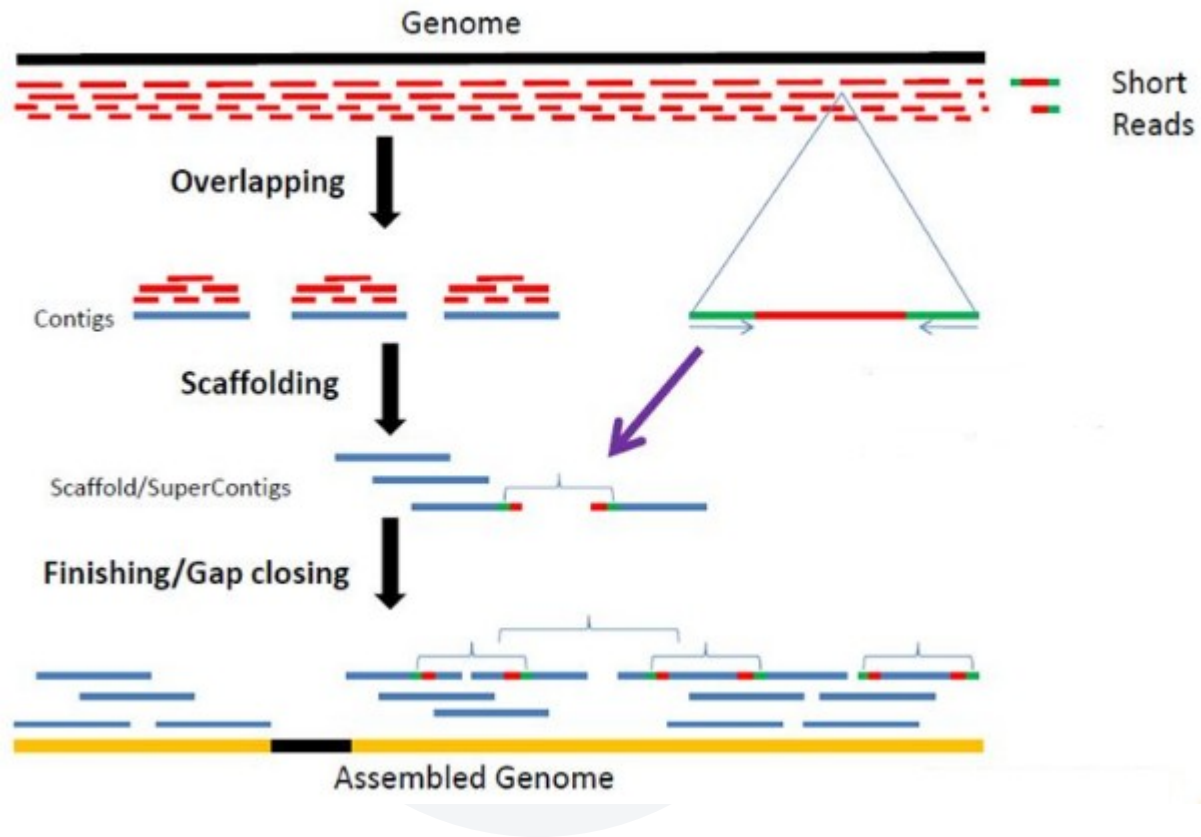Shotgun
clones

Shotgun
sequence

...ACCGTAAATGGGCTGATCATGCTTAAA
            TGATCATGCTTAAACCCTGTGCATCCTACTG...
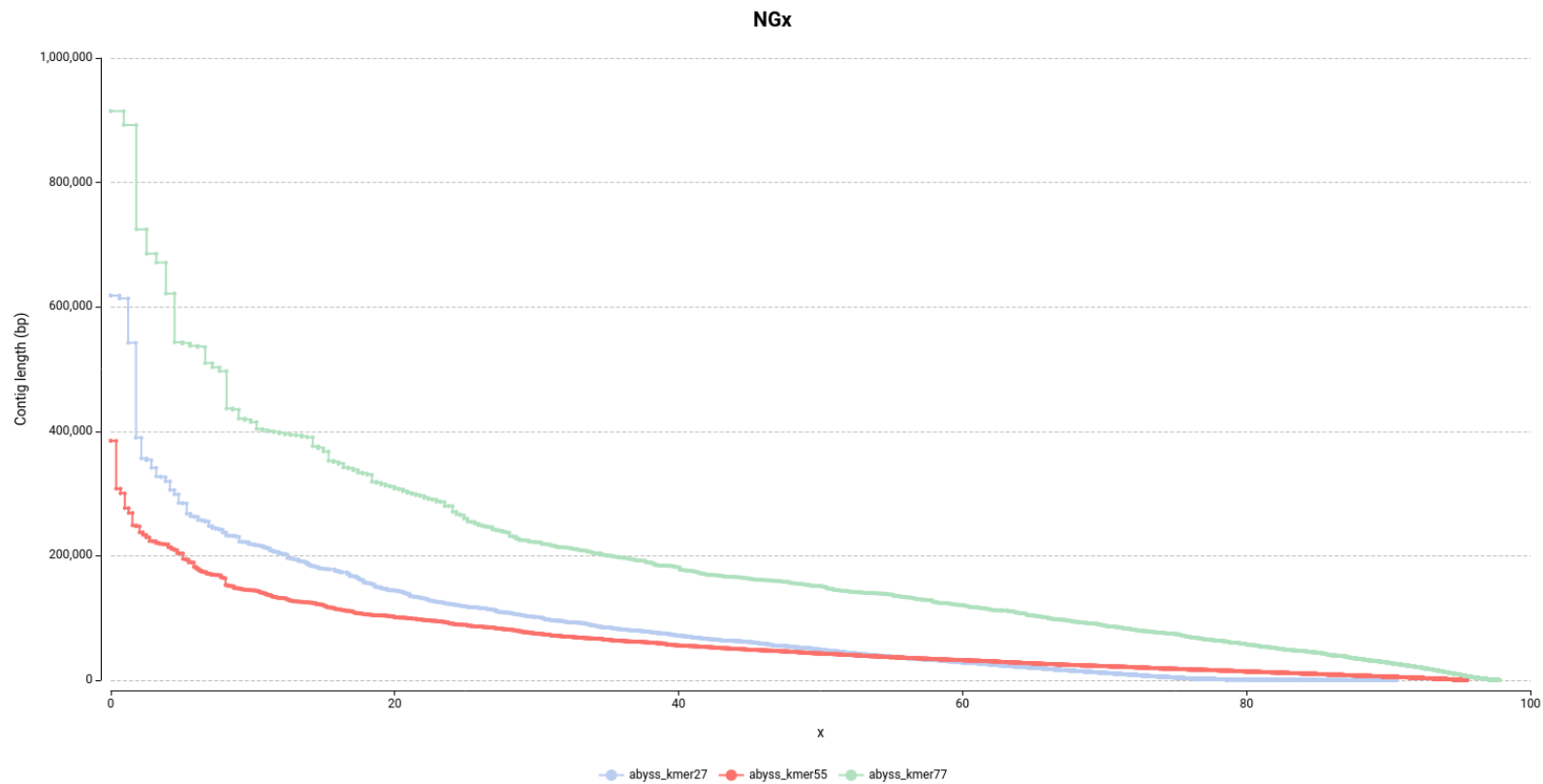
Assembly  ...ACCGTAAATGGGCTGATCATGCTTAAACCCTGTGCATCCTACTG...

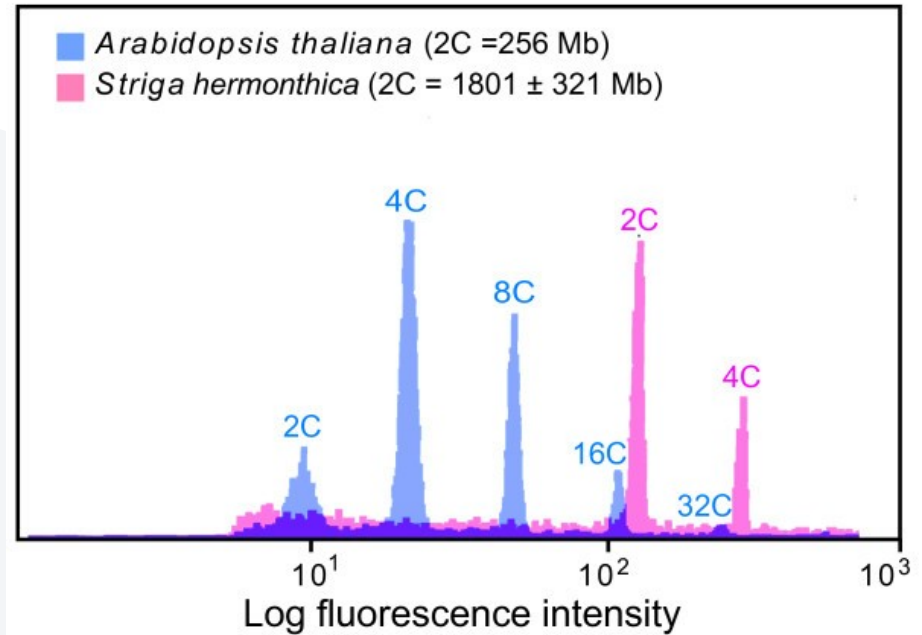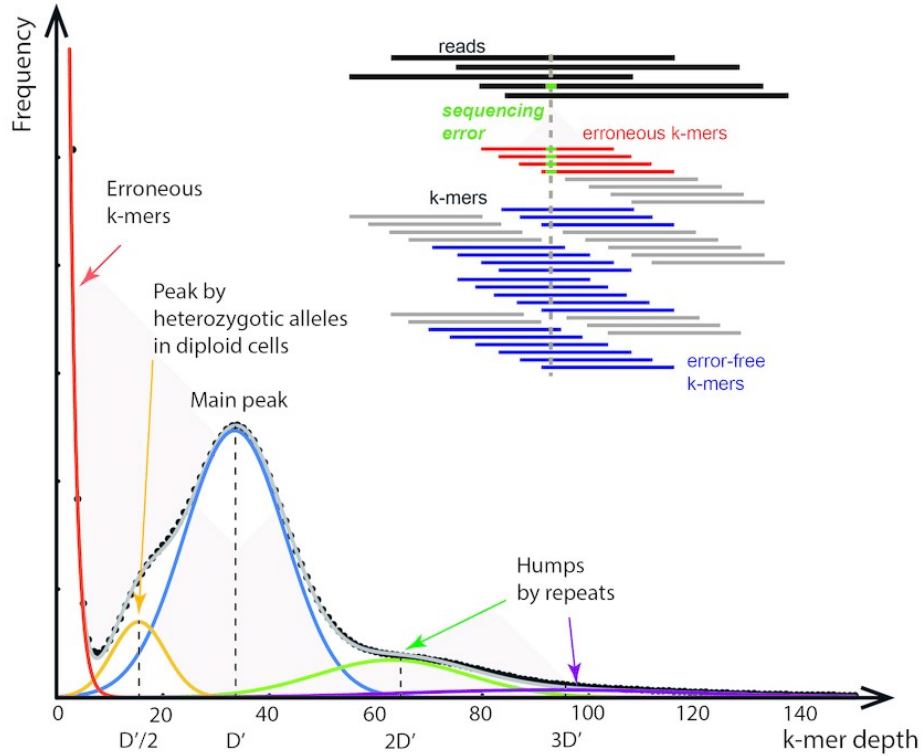# Whole genome shotgun sequencing

# Level of fragmentation (AuN)

**NGx**

# Coherence:
## Estimating genome size

### In silico: kmer analysis

### Flow cytometry

# Coherence:
## Read distance and orientation

Mapped reads

A) Properly paired reads (PPR)

5'　　　　　　　　　　　　　　3'

Insert size distance
and
Mates orientation ✔ ✔

B) Improperly paired reads (IPPR)

5'　　　　　　　　　　　　　　3'

5'　　　　　　　　　　　　　　3'

Insert size distance ✖
and/or
Mates orientation ✖

C) Both mates map, but on different scaffolds

Broken reads

scaffold 10　　　　　scaffold 203

5'　　　　3'　　　5'　　　　3'

D) One of the mate does not map

scaffold 10

5'　　　　3'

scaffold 10

5'　　　　3'

Category B – potential misassembled regions
Category C – scaffolding purposes (if read coverage is high)

# Completeness
## Mapping efficiency



Alignment of:
- Assembly reads
- RNAseq
- Cloned genes
- Genomic surveillance sequences (EST, BES, Fosmids, flcDNA)

# Completeness
## Gene Space



**BUSCO Assessment Results**

Complete (C) and single-copy (S)  Complete (C) and duplicated (D)
Fragmented (F)  Missing (M)

genome1: C:1097 [S:1090, D:7], F:137, M:285, n:1519

genome2: C:1459 [S:1453, D:6], F:17, M:43, n:1519

genome3: C:1045 [S:1012, D:33], F:72, M:402, n:1519

%BUSCOs

# Single cell data analysis

Library preparation
Indexing and alignment
Exploring the results
Data analysis

# Library preparation



Adapted from 10X genomics

# sc-RNAseq are enriched in the 3' end



Ma et al. (2019). BMC Genomics 20(1).

**Step**

**Alignment and molecular counting**

CTGCAGG
CTGTATG
CTACATG
CTGCATG

**Cell filtering and quality control**

No. UMIs in cell — 100, 1,000, 10,000

Real cells

Cell rank — 1, 10, 100, 1,000

**Doublet scoring**

Simulated and real doublets
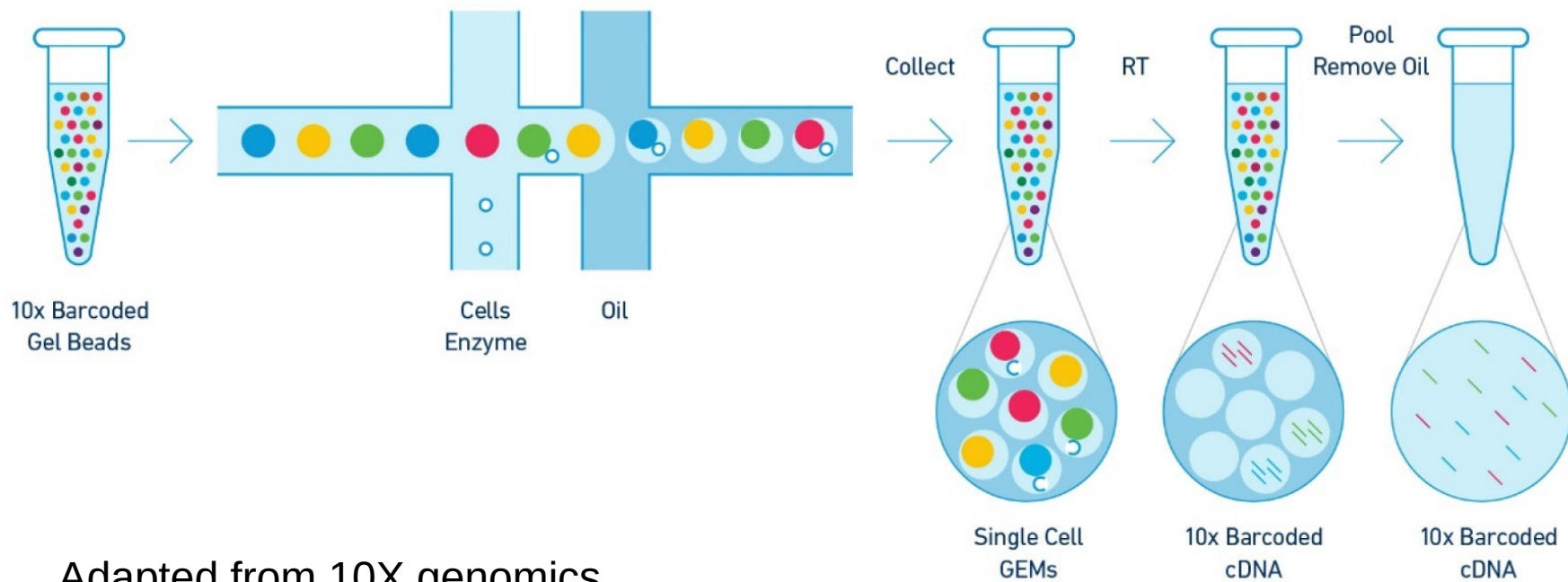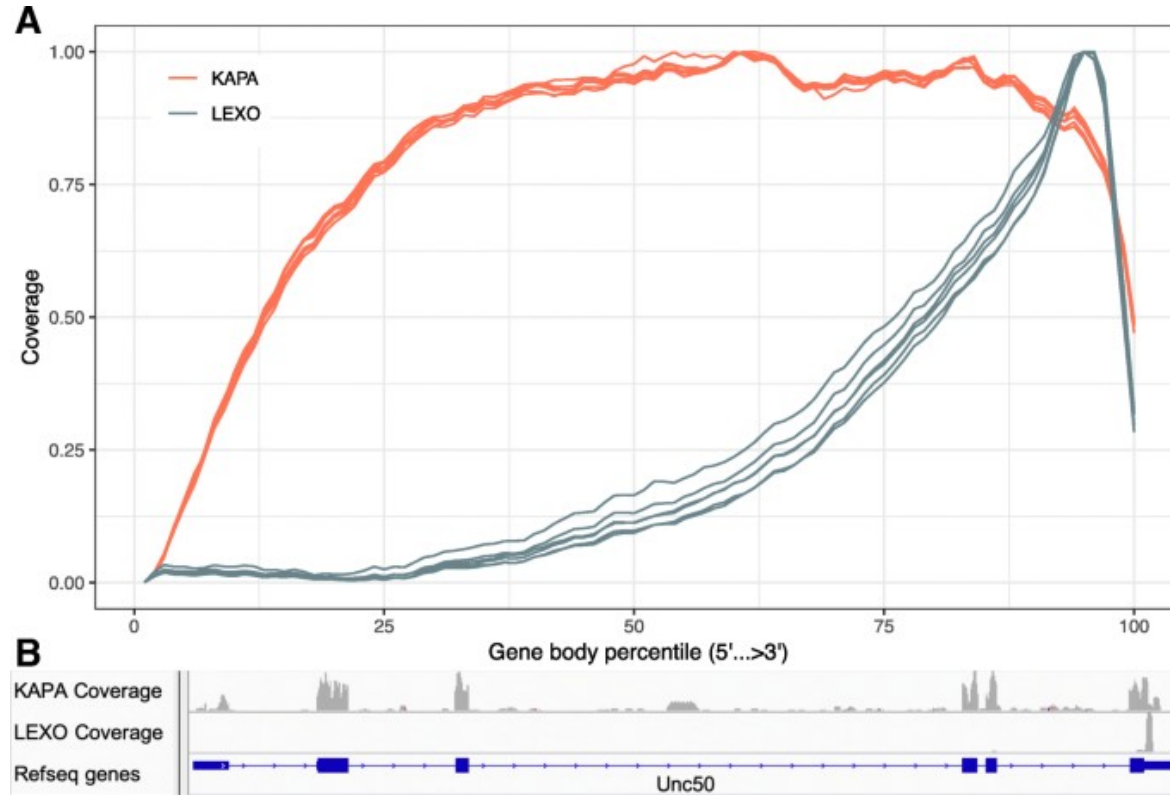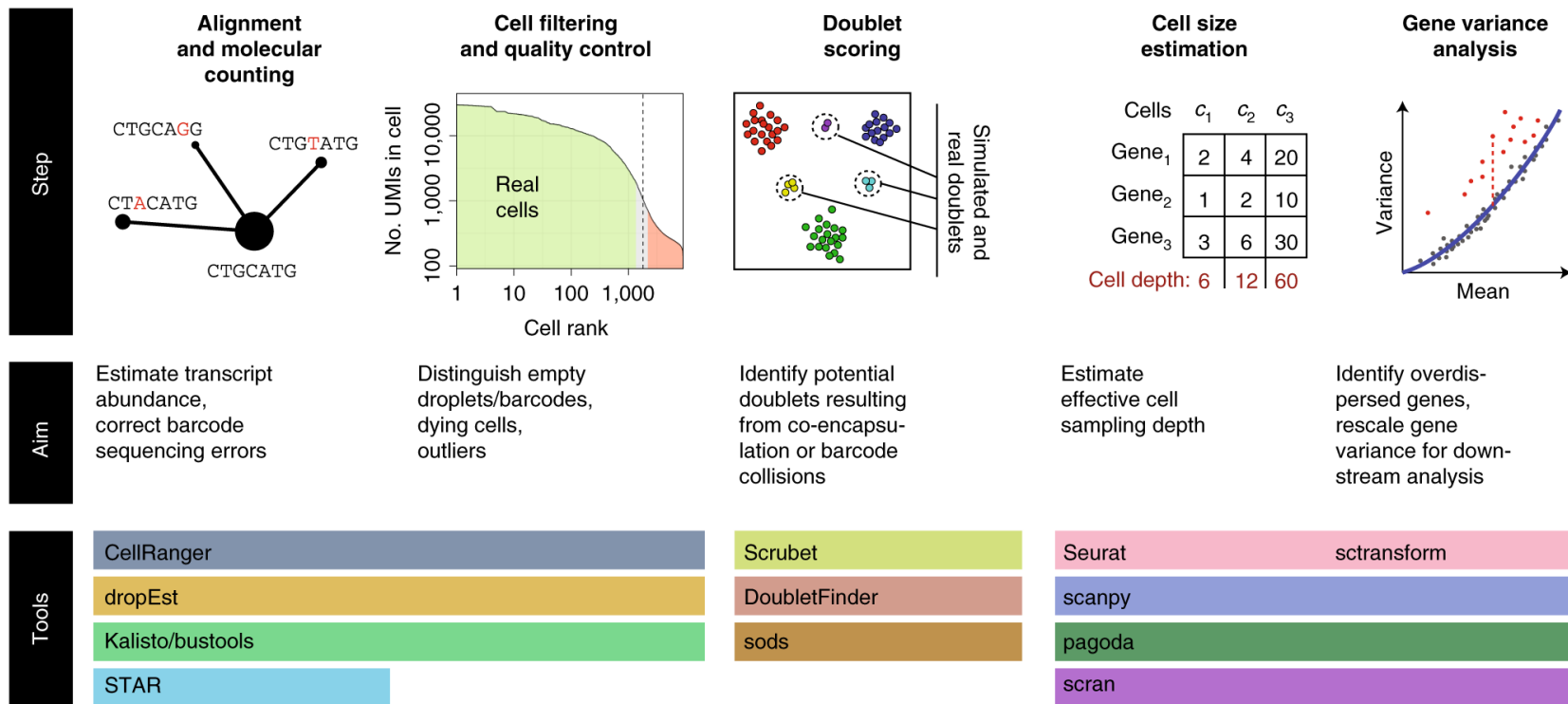
**Cell size estimation**

| Cells | $c_1$ | $c_2$ | $c_3$ |
|-------|-------|-------|-------|
| $Gene_1$ | 2 | 4 | 20 |
| $Gene_2$ | 1 | 2 | 10 |
| $Gene_3$ | 3 | 6 | 30 |
| Cell depth: | 6 | 12 | 60 |

**Gene variance analysis**

Variance vs Mean

**Aim**

Estimate transcript abundance, correct barcode sequencing errors

Distinguish empty droplets/barcodes, dying cells, outliers

Identify potential doublets resulting from co-encapsulation or barcode collisions

Estimate effective cell sampling depth

Identify overdispersed genes, rescale gene variance for downstream analysis

**Tools**

CellRanger
dropEst
Kalisto/bustools
STAR

Scrubet
DoubletFinder
sods

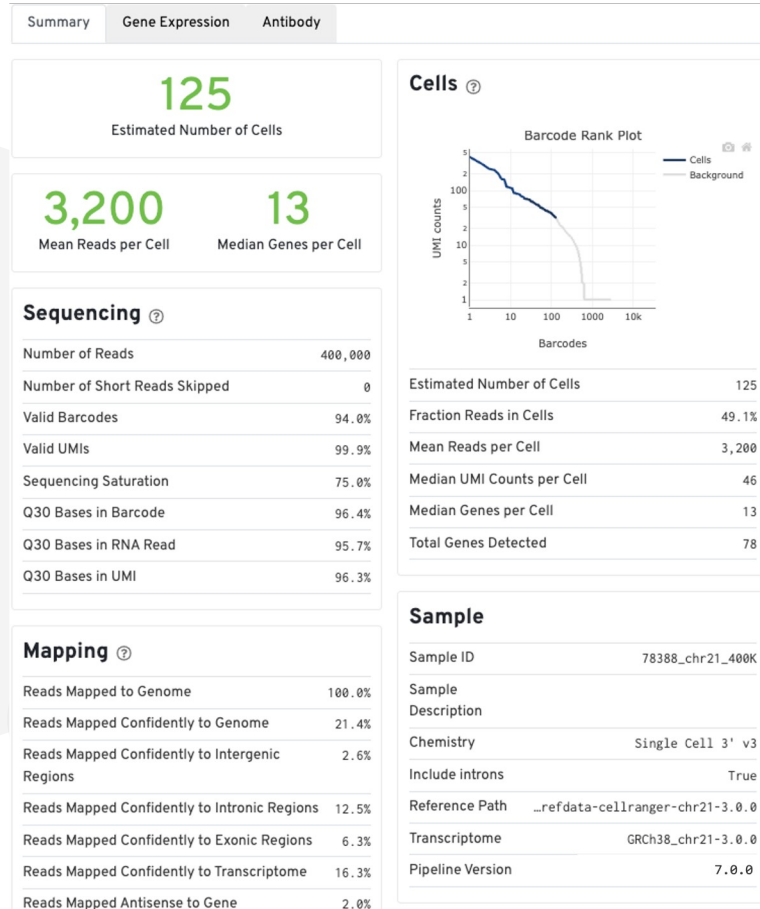Seurat | sctransform
scanpy
pagoda
scran

Karchenko (2021) Nature Methods 18:723–732

# Tenth Exercise

- Find and load cellranger software
- Index the reference genome
- Align reads to the reference genome
- Understand the results

# Cellranger html reports

# End of the second week

Thank you for your attention