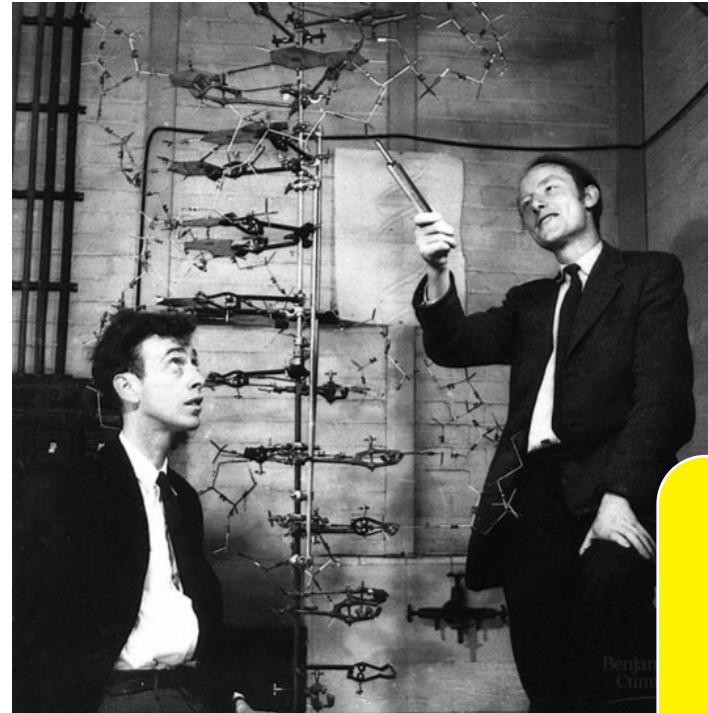# howest
## university of applied sciences

# Bioinformatics for dummies
# MB&C2024 Workshop
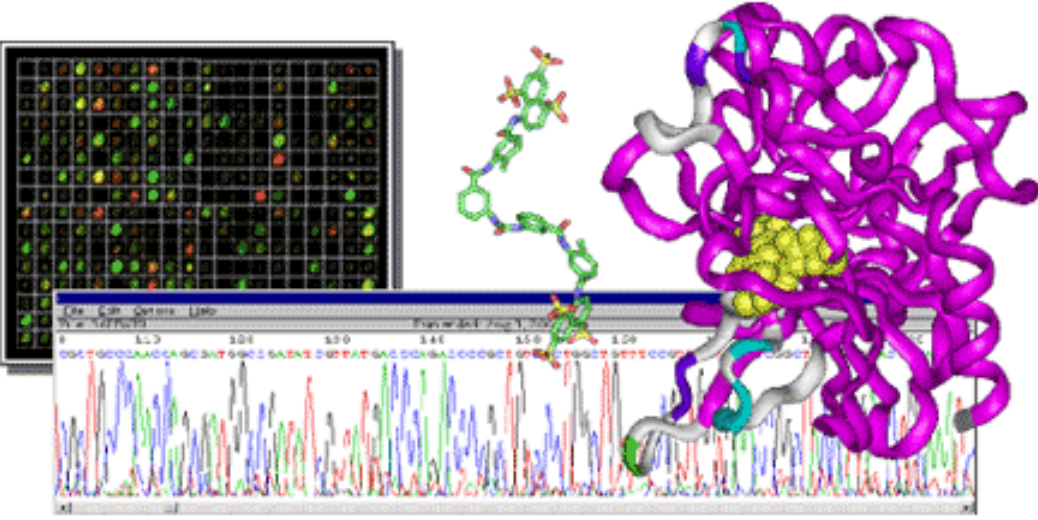
Cedric Hermans

Paco Hulpiau

# Introduction



**Molecular biology**

**Information technologies**

**Bioinformatics**

Combine:

- New insights and technologies in molecular biology

- Advances in information technologies

howest
university of applied sciences

# Introduction

**Informatics**
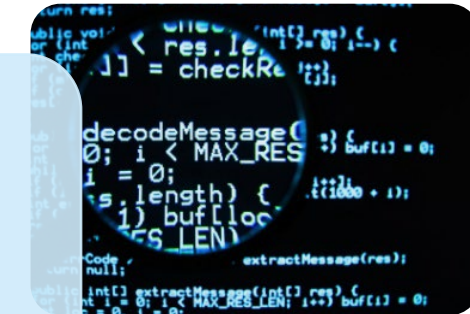To store, organize and share molecular biological data in database systems



**Bioinformatics**
To process and analyse biological data by using bioinformatics tools in a "dry lab"
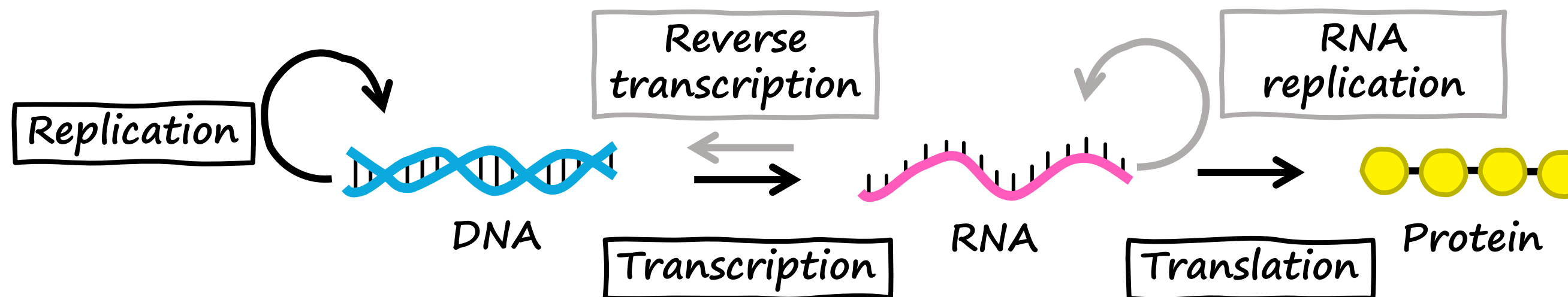


**Programming**
To integrate the different tools by means of scripting into a bioinformatics pipeline

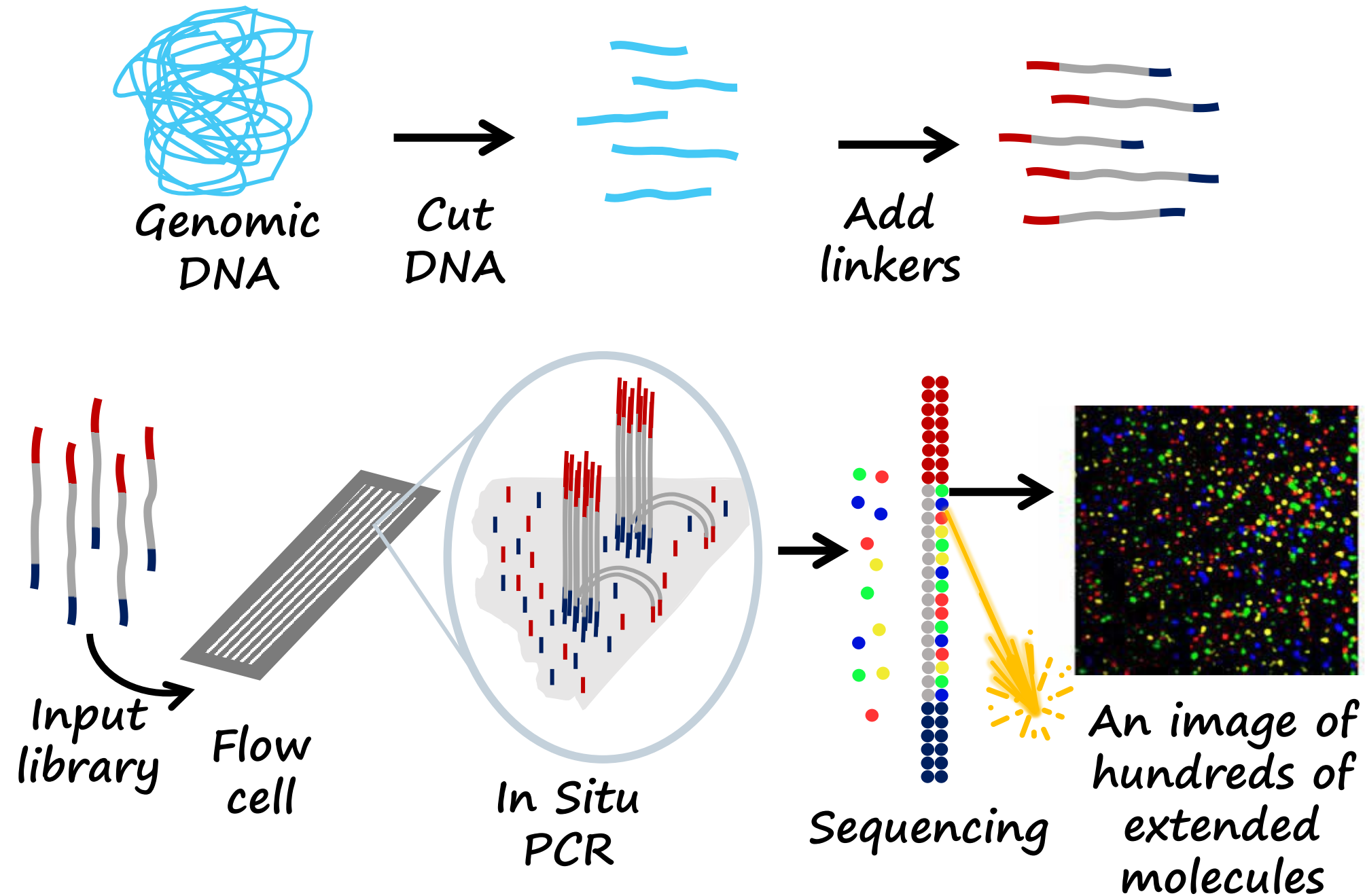**howest**
university of applied sciences

# Molecular biology and bioinformatics

Important (high-throughput) technologies:

- Next Generation Sequencing
  - ➢ Sequencing and expression analysis

- Microarray
  - ➢ Expression and genetic variation analysis

- Mass spectrometry
  - ➢ Protein (sequence) identification

# Next generation sequencing



Genomic DNA → Cut DNA → Add linkers

Input library → Flow cell → In Situ PCR → Sequencing → An image of hundreds of extended molecules
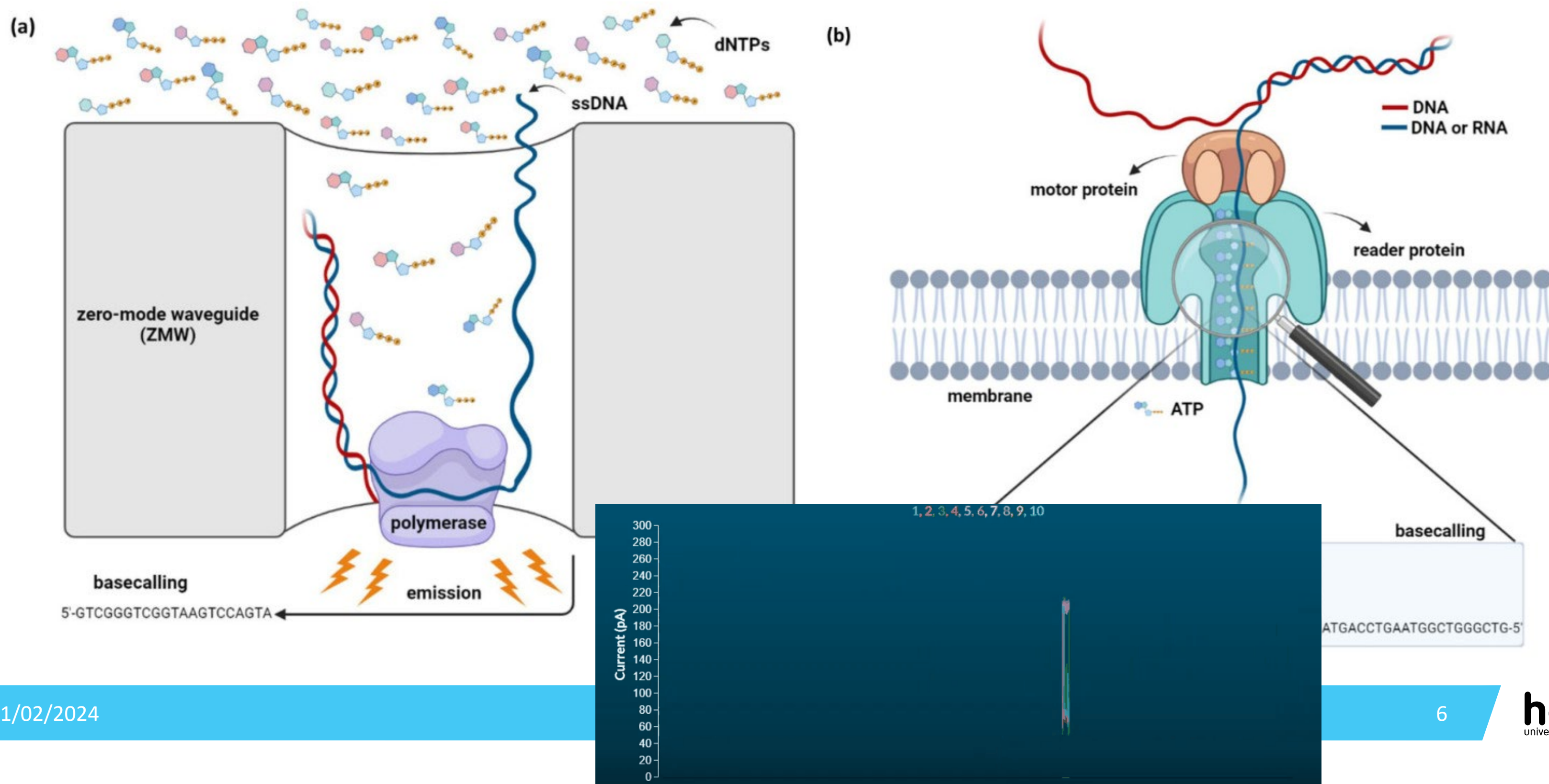
Short-read NGS

- 2 approaches:
  - Sequencing by synthesis
  - Sequencing by ligation

- 35-700 bp read length

- High accuracy (~ 99,99%)

- Complex assembly

howest
university of applied sciences
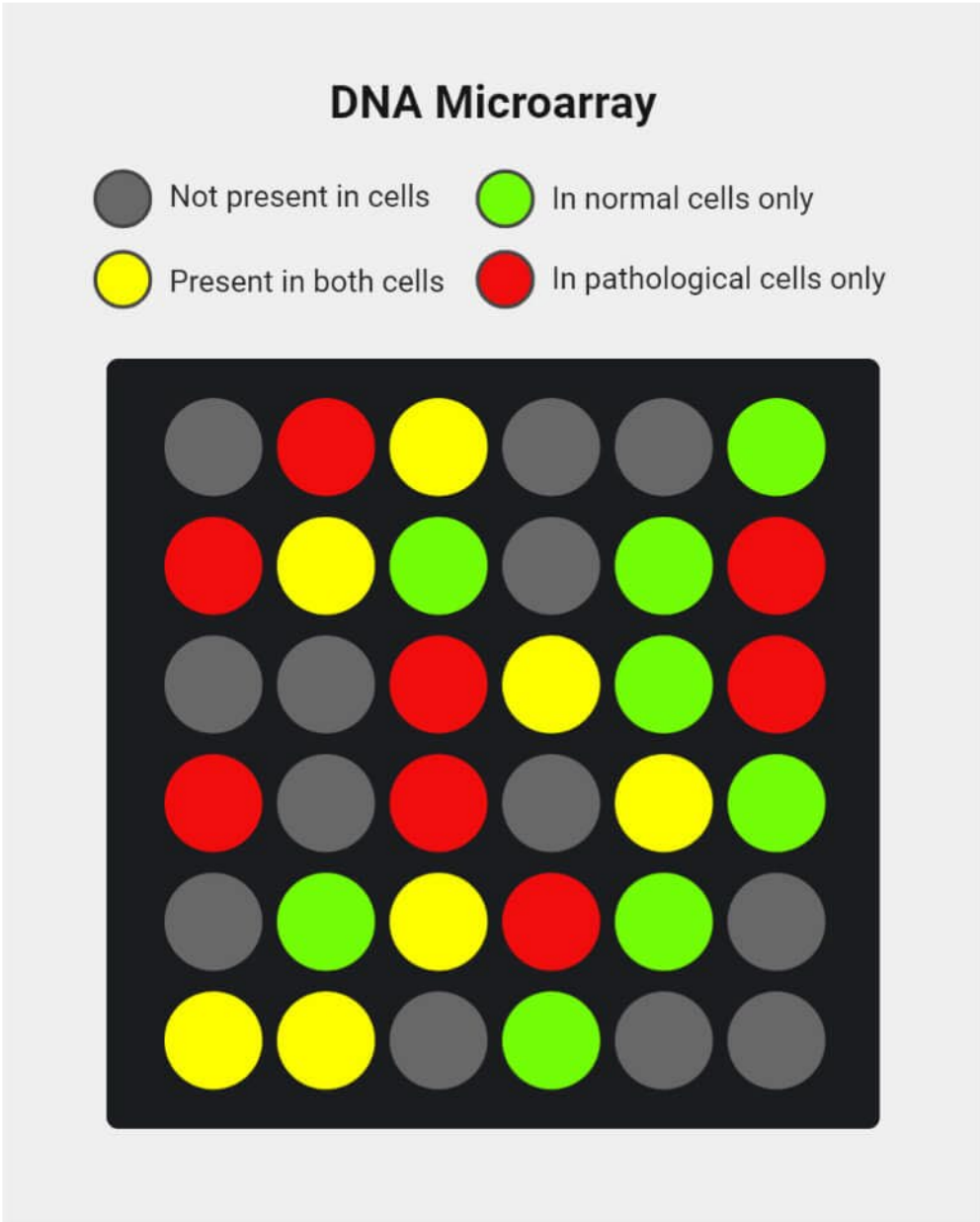
# Next next generation sequencing

# Microarrays
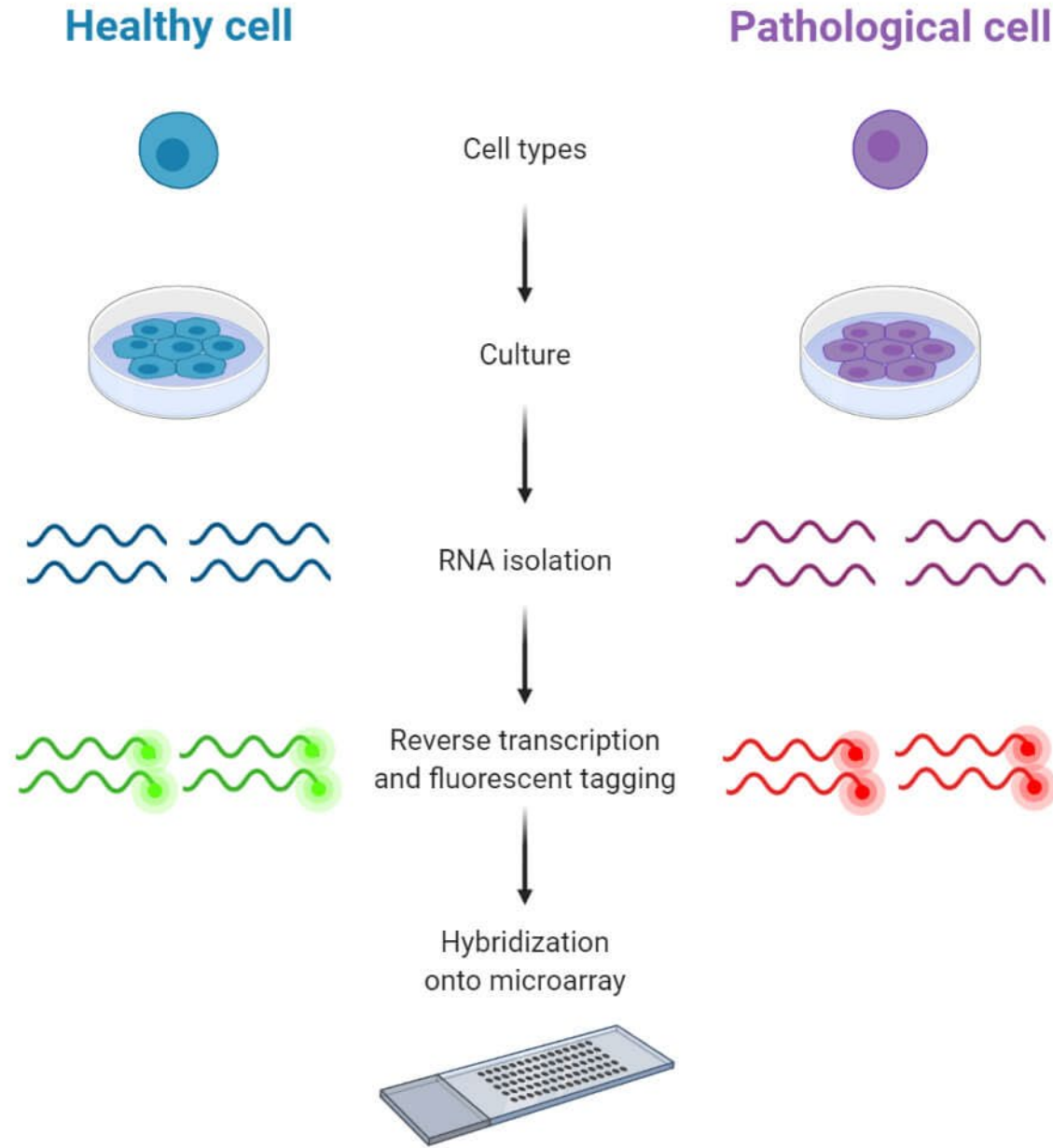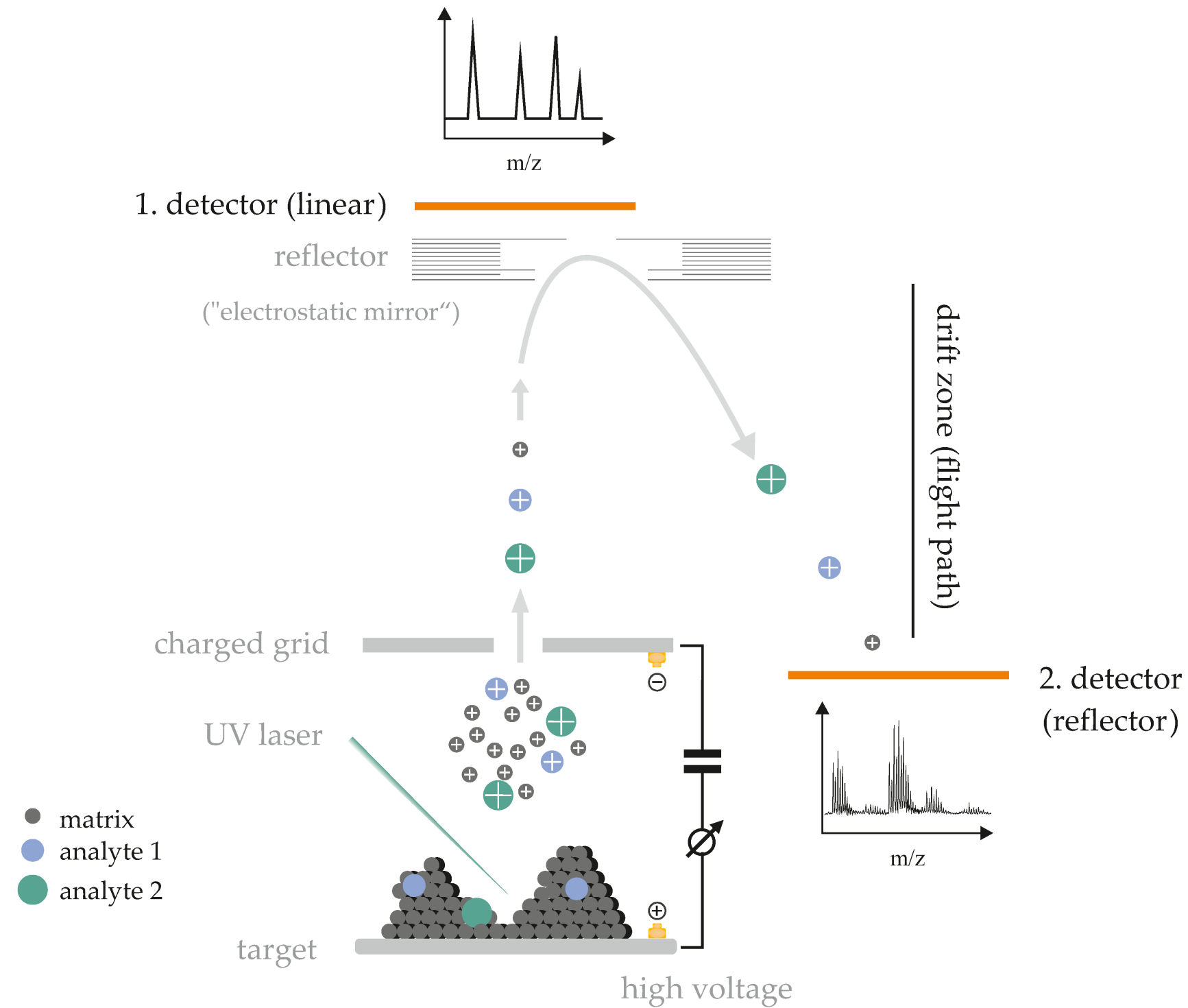


Image By Sagar Aryal, created using biorender.com

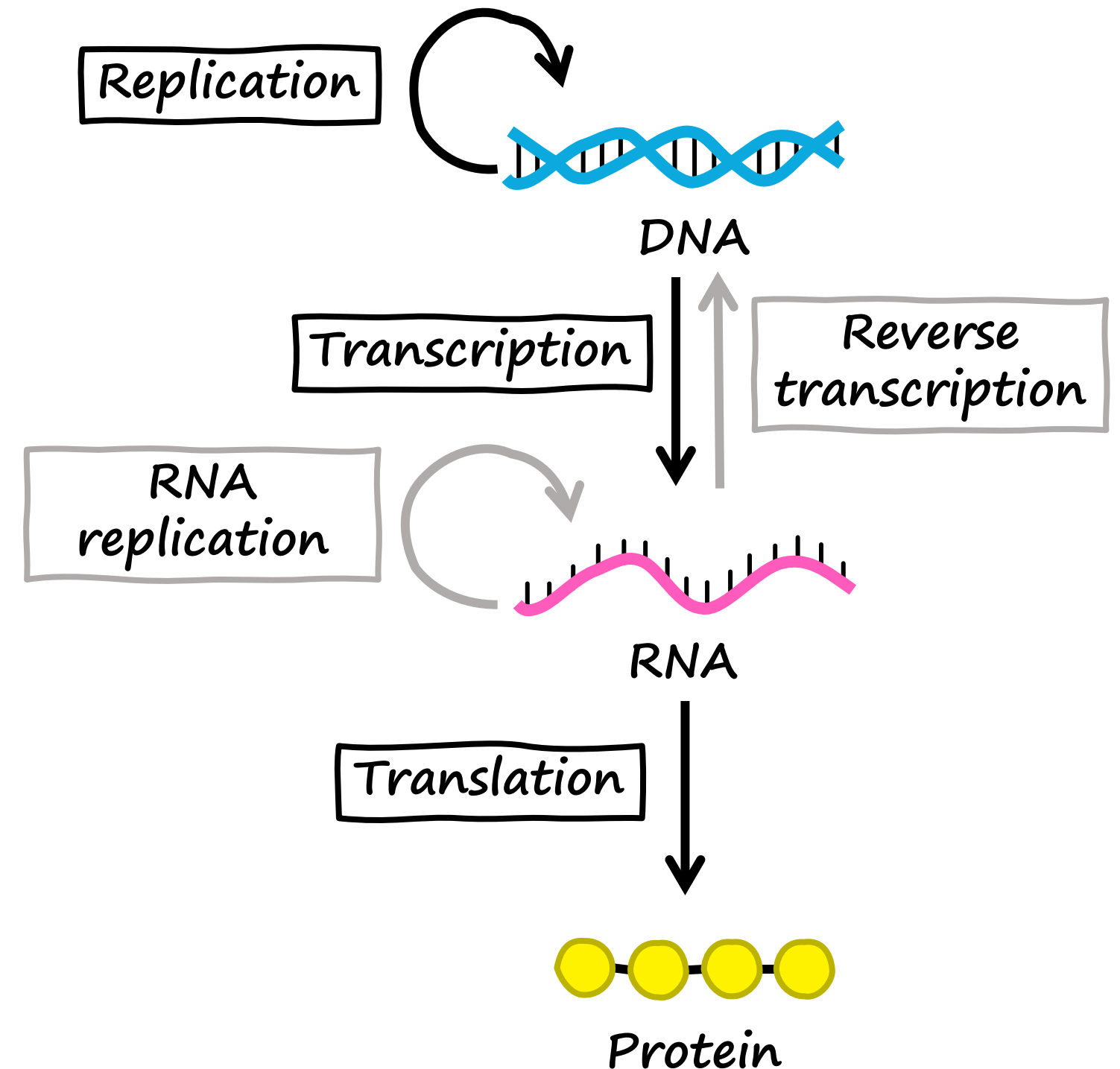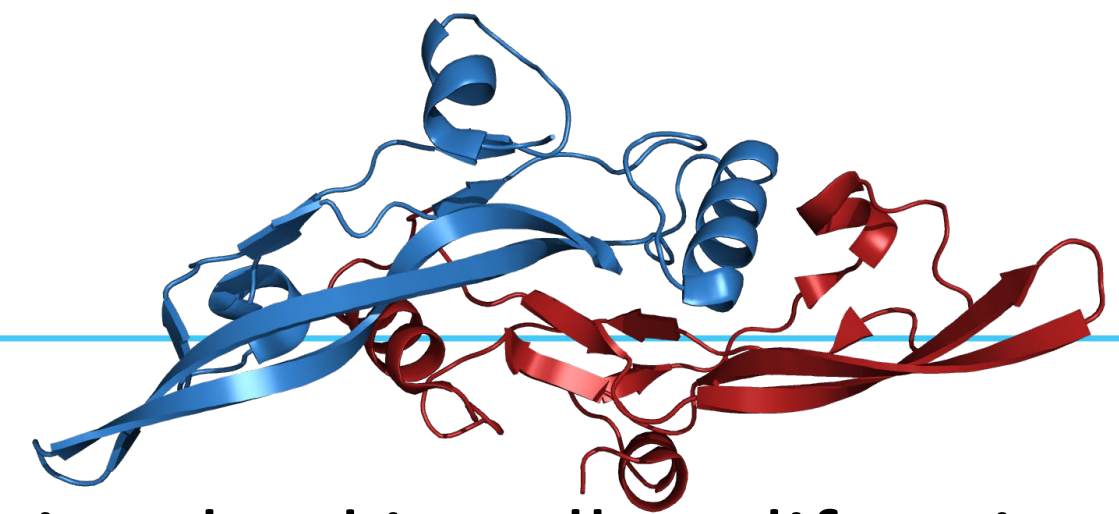# Mass spectrometry

# Molecular biology and bioinformatics

Biological databases:

- DNA
  - ➢ Sequence and loci
  - ➢ (Natural) genetic variation
- RNA
  - ➢ Transcripts (and variants)
  - ➢ Gene expression
- Protein
  - ➢ Sequence and function
  - ➢ Phenotype (and diseases)



Replication

DNA

Transcription

Reverse transcription

RNA replication

RNA

Translation

Protein

# (Sequence) repositories

Exploratory example: **TGF beta 1** – an important protein involved in cell proliferation, differentiation and growth

| NCBI Gene | https://www.ncbi.nlm.nih.gov/gene/7040<br>General and integrated sequence and locus information |
|---|---|
| NCBI Nucleotide | https://www.ncbi.nlm.nih.gov/nuccore/?term=TGFB1+AND+"Homo+sapiens"[Organism]<br>All available (partial) TGF beta 1 nucleotide sequences → ± 138 records (!) |
| Ensembl | https://www.ensembl.org/Homo_sapiens/Gene/Summary?db=core;g=ENSG00000105329<br>General information + detailed transcripts and gene expression |
| UniProt or NCBI Protein | http://www.uniprot.org/uniprot/P01137<br>High-quality recourse of protein sequence and functional information |

howest
university of applied sciences

# (Sequence) repositories

**Example 1:** Look for the nucleotide sequence of PSA

- https://www.ncbi.nlm.nih.gov/nucleotide/

- NCBI nucleotide query: "(prostate specific antigen)" restricted to humans

# (Sequence) repositories

**Example 2:** the Genome Projects

- 1000 Genomes Project (2008-2015)
  - Goal: to find most genetic variants with frequencies of at least 1% in the populations studied

- 100 000 Genomes Project (2013-2018)
  - Goal: focus on rare diseases, some common types of cancer and infectious diseases

- 1+ Million Genomes (2018-2027)

```
ACGTACGTACGTACGTACGTACGT
ACGTACCTACGTACGTACGTACGT
ACGTACCTACGTATGTTCGTACGT
ACGTACGTACGTATGTTCGTACGT
```

howest
university of applied sciences

# (Sequence) repositories

**Example 2:** the Genome Projects

- 100 000 Genomes Project (2013-2018)



The NEW ENGLAND JOURNAL of MEDICINE

**The 100,000 Genomes Pilot on Rare-Disease Diagnosis**

U.K. PATIENTS WITH RARE DISEASES AND NO DIAGNOSIS — PRELIMINARY REPORT

**2183** Probands with 161 undiagnosed disorders

**Diagnostic yield** ⟩ 25% of probands received a genetic diagnosis

**Diagnostic pipeline**

- **86%** of diagnoses were identified through automated pipeline
- **14%** of diagnoses required additional research

**Novel discoveries**

- **3** new disease genes discovered
- **19** new disease–gene associations identified

25% of genetic diagnoses had immediate ramifications for clinical decision making.

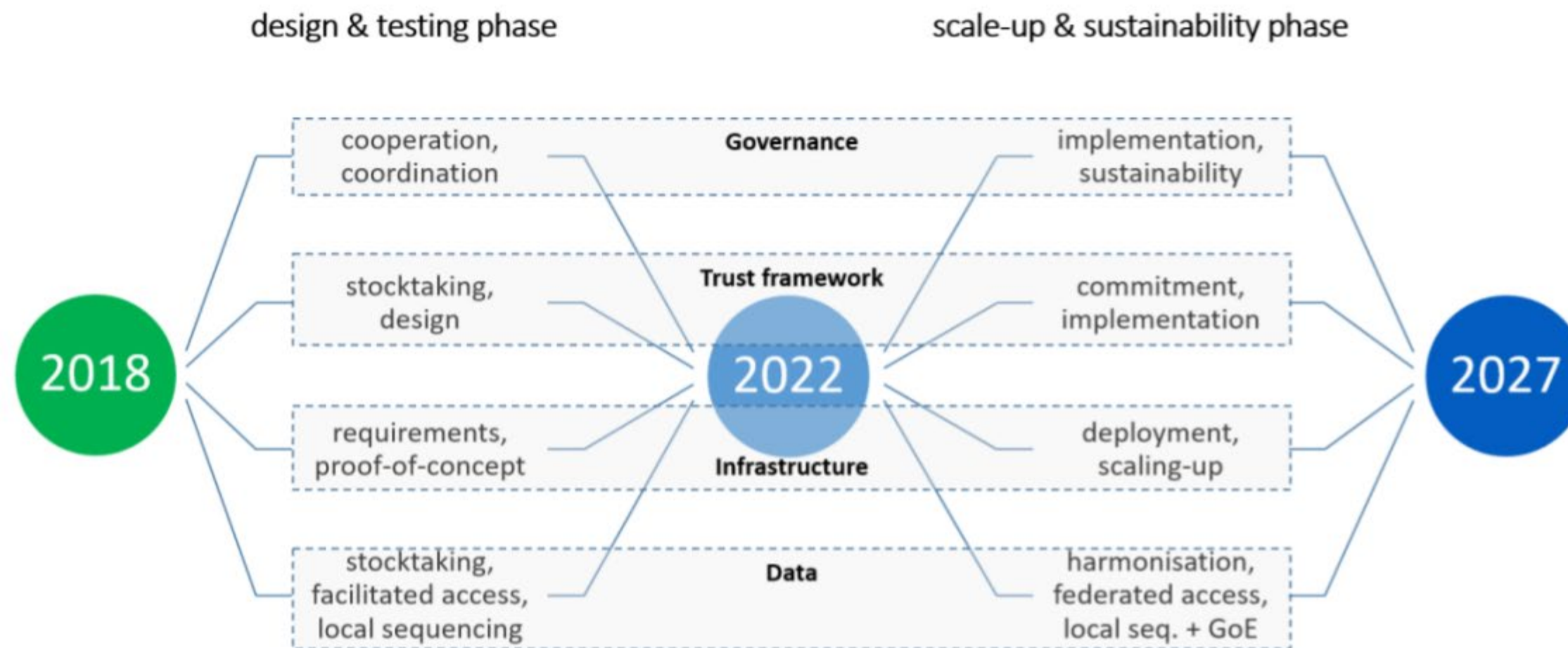The 100,000 Genomes Project Pilot Investigators    10.1056/NEJMoa2035790    Copyright © 2021 Massachusetts Medical Society

# (Sequence) repositories

**Example 2:** the Genome Projects

- 1+ Million Genomes

# (Sequence) repositories

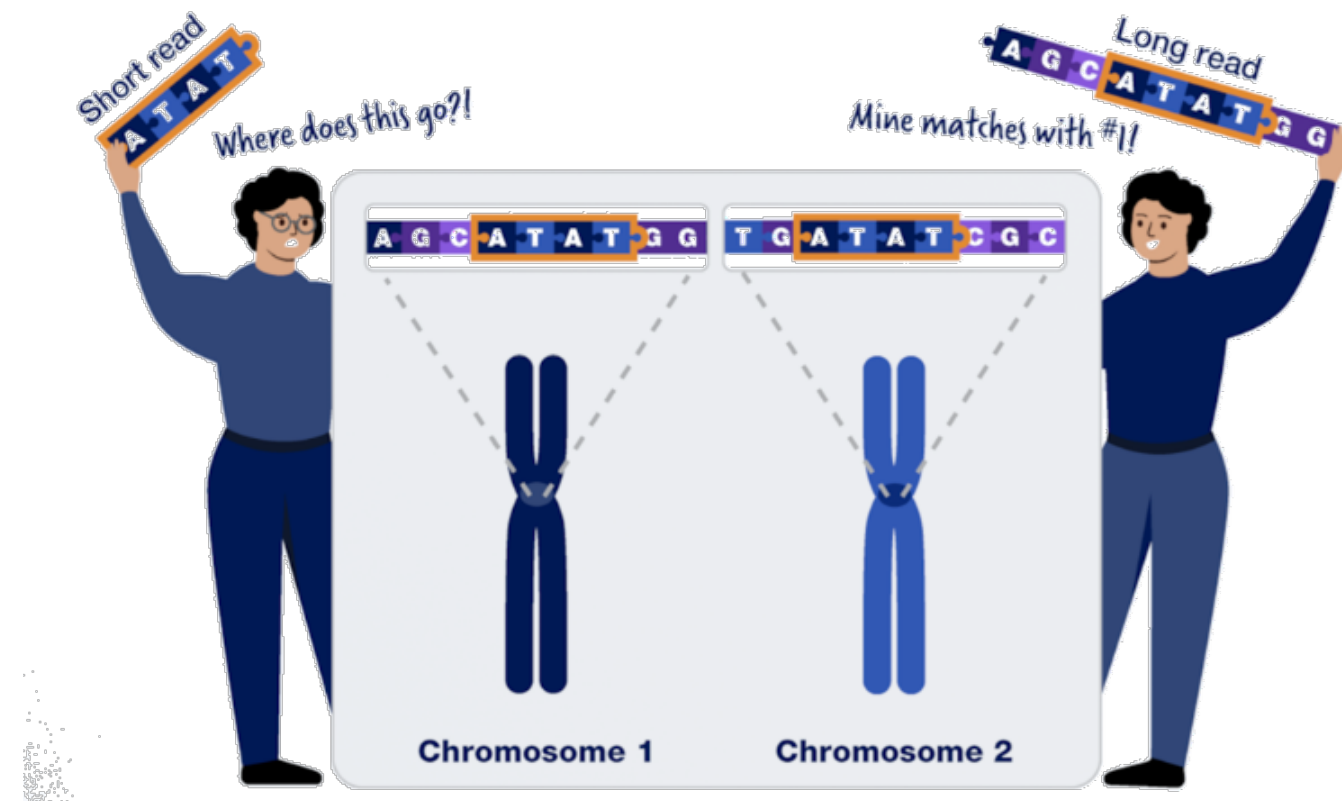Lots of diversity between genomes => solution is needed

- **Genome Reference Consortium (GRC)**
  - Goal: create the best possible reference assembly for humans
  - -> latest major release: GRCh38 (also know as hg38) version 14
  - https://www.ncbi.nlm.nih.gov/grc/human
- **NCBI Reference Sequence Database (RefSeq)**
  - Non-redundant, well annotated set of reference sequences including genomes, transcripts and proteins
  - https://www.ncbi.nlm.nih.gov/refseq/
  - One gene/transcript/protein = one sequence

# (Sequence) repositories

Note: GRCh38 is not complete…

- Telomere-to-Telomere (T2T) consortium: https://www.genome.gov/t2t

- Data: https://github.com/marbl/CHM13

# (Sequence) repositories



A good place to start searching for a (reference) sequence!

# Homology searching



## Next (next) Generation Sequencing

- Result: unknown nucleotide sequences

Determination of sequence ≠ simple keyword search strategy

⇒ Usage of <u>evolutionary models</u> to determine <u>homology</u> between nucleotide (or protein) sequences



- Based on sequence alignment

- **BLAST**: <u>B</u>asic <u>L</u>ocal <u>A</u>lignment <u>S</u>earch <u>T</u>ool

# Homology searching

**Homology**

- Derived from a common ancestor

- 2 types:
    - Orthologs = speciation event (different species)
    - Paralogs = duplication event (same species)

- Typically based on morphological characteristics

- Use "molecular phylogeny" to determine homology

$\Rightarrow$ Phylogenetic tree



https://lifemap-ncbi.univ-lyon1.fr/#

# Homology searching

**Homology**

- Derived from a common ancestor

- 2 types:
  - Orthologs
  - Paralogs
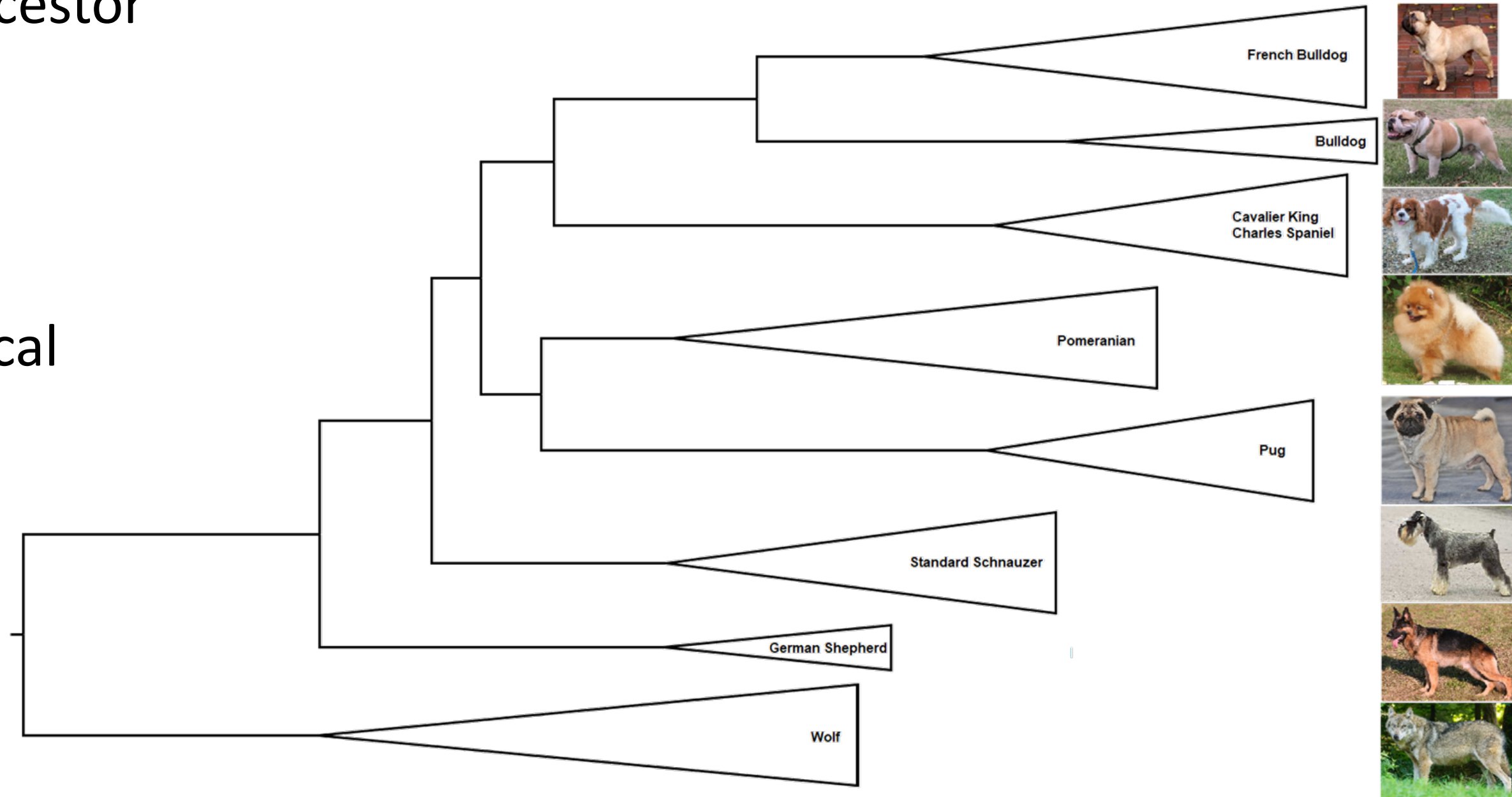
Typically based on morphological characteristics

- Use "molecular phylogeny" to determine homology

⇒ Phylogenetic tree

howest
university of applied sciences

# Homology searching

**File format**

Fasta-file:

- Header line starting with ">"

- One or more lines containing the sequence

Multifasta-file:

- Mutiple sequences in Fasta format below one another.

- A new sequence is recognized by the ">" in front of each header

```
>1
GGCCGGTAAAACTCGTGCCAGCCACCGCGGTTAAACGAGAGGCCCTAGTTGATAA
>2
GGCCGGTAAAACTCGTGCCAGCCACCGCGGTTAAACGAGAGGCCCTAGTTGATAT
>3
GGTCGGTTAAACTCGTGCCAGCCACCGCGGTTATACGAGAGACCCTAGTTGACTCA
>4
GGCCGGTAAATTCGCGTGCCAGCAACCGCGGTTAGACGTACATAGGCCTAAGTTG
>5
GGCCGGTAAAACTCGTGCCAGCCACCGCGGTTAAACGAGAGGCCCTAGTTGATAG
>6
GGCCGGTAAAACTCGTGCCAGCCACCGCGGTTAGACGAGAGGCCCTAGTTGATAT
```

**howest**
university of applied sciences

# Homology searching

```
>unknown human nucleotide sequence
CAAGGCTGTCCCCCCAAGACGTGCTCCCAGGACGAGTTTCGCTGCCACGATGGGAAGTGCATCTCTCGGCAGTTCGTCTGTGACTC
AGACCGGGACTGCTTGGACGGCTCAGACGAGGCCTCCTGCCCGGTGCTCACCTGTGGTCCCGCCAGCTTCCAGTGCAACAGCTCCA
CCTGCATCCCCCAGCTGTGGGCCTGCGACAAC
```

- Given: an unknown human nucleotide Fasta sequence
-> https://www.bio-informatica.be/workshops/
  > "unknown human nucleotide sequence.fasta"



- To determine the identity -> use BLAST (https://blast.ncbi.nlm.nih.gov/Blast.cgi)
  - Settings:
    - Organism: *Homo sapiens*
    - Database: refseq_rna
    - Exclude: models

howest
university of applied sciences

# Homology searching

**Results:**

- Identity
- Bits score
- Expect value
- Gaps

# Homology searching

**BLAST**

- Not a simple keyword search strategy

- 3 steps
  - LIST
  - SCAN
  - EXTEND

- Based on a model of evolution and scoring system

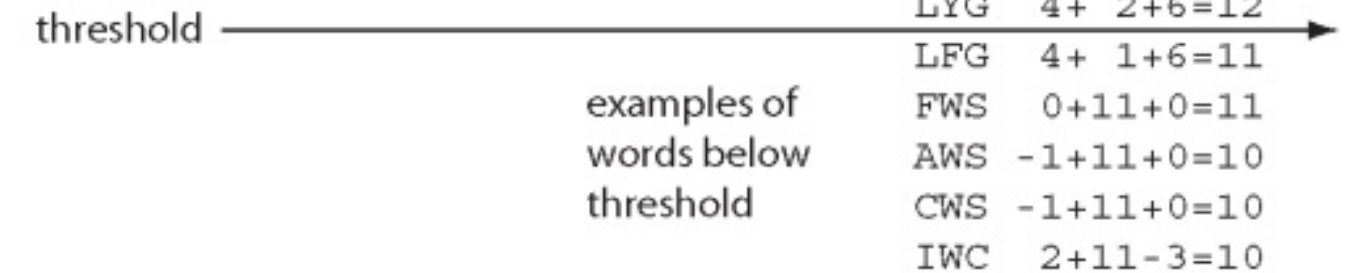Phase 1: Setup: compile a list of words (w=3) above threshold T

• Query sequence: human beta globin NP_000509.1 (includes ...VTALWGKVNVD...).
This sequence is read; low complexity or other filtering is applied; a "lookup" table is built.

• Words derived from query sequence (HBB):  VTA  TAL  ALW  LWG  WGK  GKV  KVN  VNV  NVD

• Generate a list of words matching query (both above and below T). Consider LWG in the query and the scores (derived from a BLOSUM62 matrix) for various words.
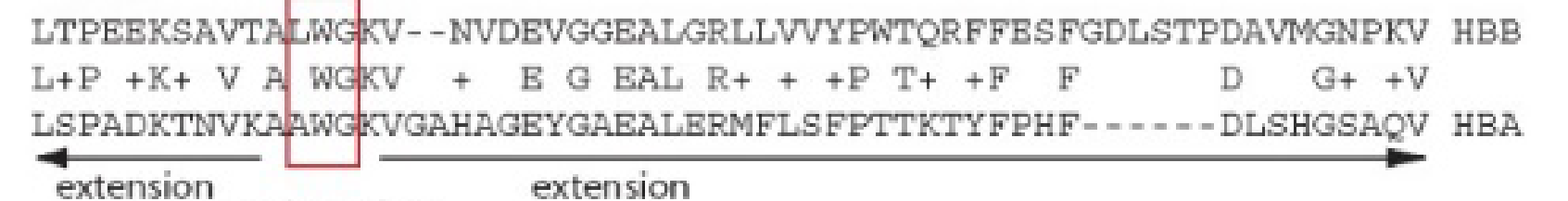
• Generate similar lists of words spanning the query (e.g. words for WGW, GWG, WGK...).

|       | LWG | 4+11+6=21 |
|-------|-----|-----------|
|       | IWG | 2+11+6=19 |
|       | MWG | 2+11+6=19 |
|       | VWG | 1+11+6=18 |
| examples of words >= threshold 12 | FWG | 0+11+6=17 |
|       | AWG | 0+11+6=17 |
|       | LWS | 4+11+0=15 |
|       | LWN | 4+11+0=15 |
|       | LWA | 4+11+0=15 |
|       | LYG | 4+ 2+6=12 |

threshold ———————————————————►

|       | LFG | 4+ 1+6=11 |
|-------|-----|-----------|
| examples of words below threshold | FWS | 0+11+0=11 |
|       | AWS | -1+11+0=10 |
|       | CWS | -1+11+0=10 |
|       | IWC | 2+11-3=10 |

Phase 2:  Scanning and extensions
• Select all the words above threshold T (LWG, IWG, MWG, VWG, FWG, AWG, LWS, LWN, LWA, LYG)
• Scan the database for entries ("hits") that match the compiled list
• Create a hash table index with the locations of all the hits for each word
• Perform gap free extensions
• Perform gapped extensions

```
LTPEEKSAVTALWGKV--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV HBB
L+P +K+ V A WGKV   +   E G EAL R+ + +P T+ +F   F        D    G+ +V
LSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHF------DLSHGSAQV HBA
```

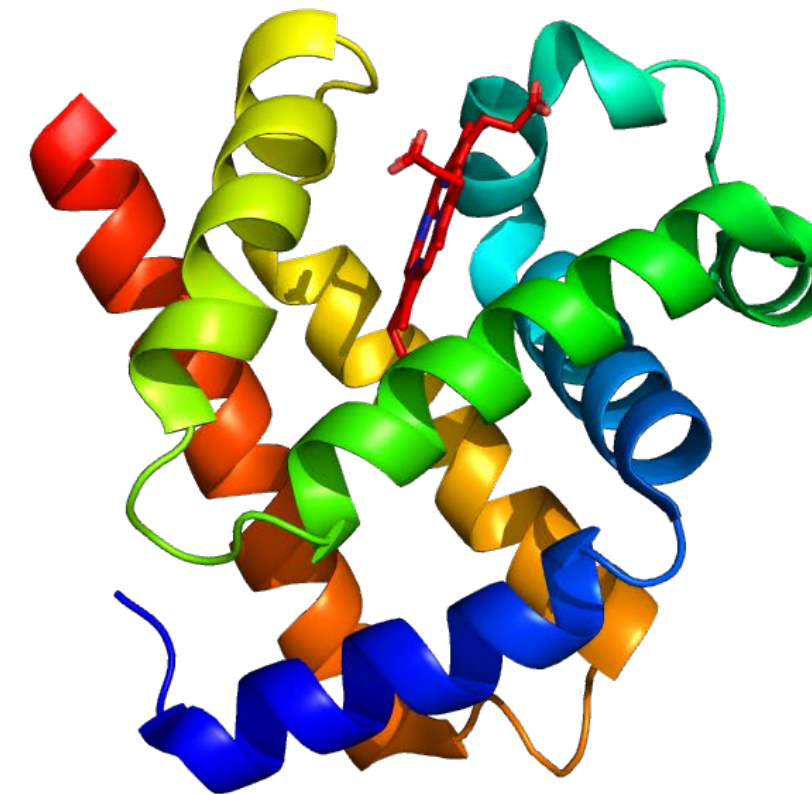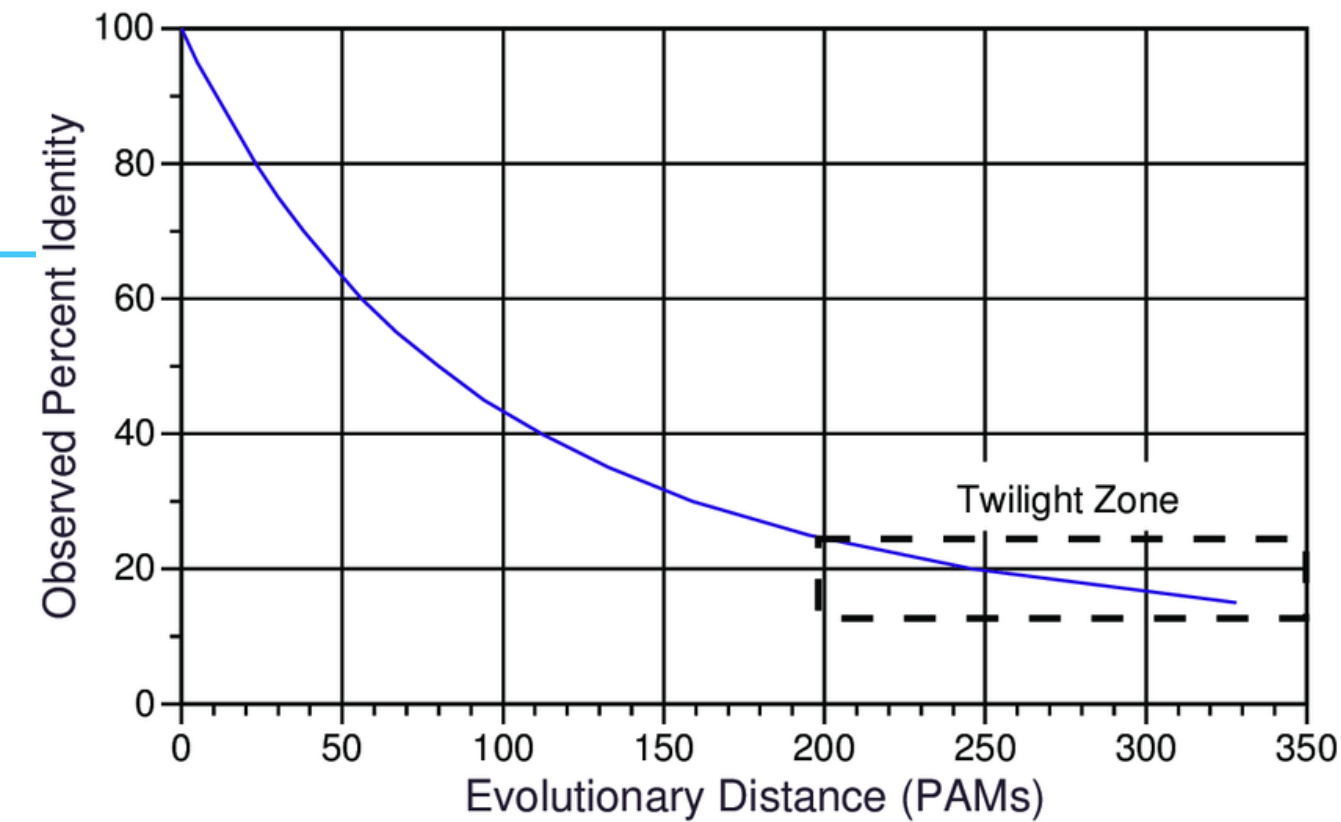◄——— extension          extension ———►

word pair from
first phases of search
"hits" alpha globin,
triggers extension
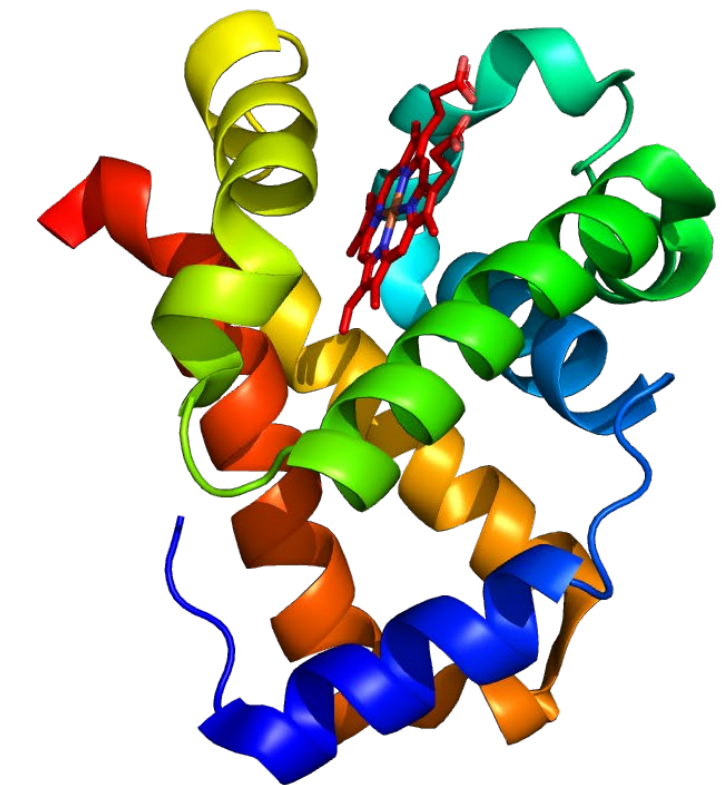
# Homology searching

- Are two sequences homologues?
    - Look at percent identity (quantitative) + expect value

- Problem: homology = YES/NO question

Example case:

Is it possible to predict that human **myoglobin** (NP_05359) and **beta hemoglobin** (NP_000509) are paralogs?



The Limits of Sequence Similarity

**Myoglobin**

**Hemoglobin**

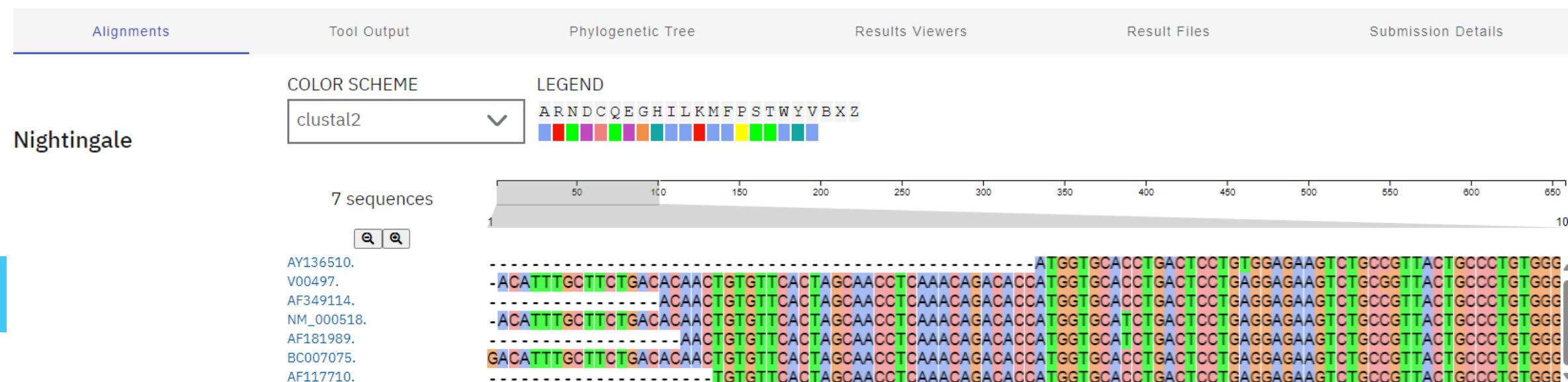howest
university of applied sciences

# DNA variant analysis

Compare nucleotide sequence with a reference sequence

- Nucleotide diversity -> DNA variant identification

- Example: nucleotide diversity in multiple hemoglobin beta variants
  - https://www.bio-informatica.be/workshops/
    > "HBB multiple sequence alignment.fasta"
  - Align sequences using MUSCLE software (Muscle < EMBL-EBI)
    - -> output: HTML

Multiple sequence alignment (MSA)

-> phylogenetic analysis

# DNA variant analysis

Browsing genetic variations

- Natural genetic variation -> Variation Viewer (https://www.ncbi.nlm.nih.gov/variation/view)

- Database of short genetic variations -> NCBI dbSNP (https://www.ncbi.nlm.nih.gov/snp/)

# DNA variant analysis

Database of variants with clinical significance: **ClinVar**
(https://www.ncbi.nlm.nih.gov/clinvar/)

# DNA variant analysis

Genetic variation -> effect on protein structure/function?

- Depends on the location of the mutation/variation



- Use PROVEAN or SIFT (**S**orts **I**ntolerant **F**rom **T**olerant) score for amino acid substitutions

| Variant ID | Chr: bp | Alleles | Global MAF | Class | Source | Evidence | Clin. Sig. | Conseq. Type | AA | AA co-ord | SIFT | Poly-Phen | CADD | REVEL | MetaLR | Mutation Assessor | Transcript |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| rs33954264 | 11:5225602 | T/A/C/G | - | SNP | dbSNP | | ⚠ ? | missense variant | H/L | 147 | 0 | 0.76 | 24 | 0.865 | 0.874 | 0.955 | ENST00000335295.4 |
| rs33954264 | 11:5225602 | T/A/C/G | - | SNP | dbSNP | | ⚠ ? | missense variant | H/R | 147 | 0.01 | 0.76 | 24 | 0.722 | 0.69 | 0.403 | ENST00000335295.4 |
| rs33954264 | 11:5225602 | T/A/C/G | - | SNP | dbSNP | | ⚠ ? | missense variant | H/P | 147 | 0 | 0.974 | 24 | 0.883 | 0.908 | 0.955 | ENST00000335295.4 |
| rs33961444 | 11:5225603 | G/A/C | - | SNP | dbSNP | | ⚠ ? | missense variant | H/Y | 147 | 0.02 | 0.146 | 23 | 0.809 | 0.801 | 0.938 | ENST00000335295.4 |

howest
university of applied sciences

# DNA variant analysis

Genetic variation → **effect** on protein structure/function?

- Variant Effect Predictor (https://www.ensembl.org/Homo_sapiens/Tools/VEP)



- Example: investigate rs13306510
  - Look up the SNP in the dbSNP database
  - Examine the SNP with the Variant Effect Predictor

howest
university of applied sciences

# Concluding remarks

Bioinformatics is more than sequence alignment, BLAST, variant calling…

➢ Interested in more? Be sure to check our offers of further training!

**Click us!**

**(advanced bachelor) Bioinformatics**

- Also in @Home version!

R for Data Analysis and Visualisation

- Also in @Home version!

**howest**
university of applied sciences

# howest
## university of applied sciences

# Bioinformatics for dummies
# MB&C2024 Workshop

Cedric Hermans

Paco Hulpiau