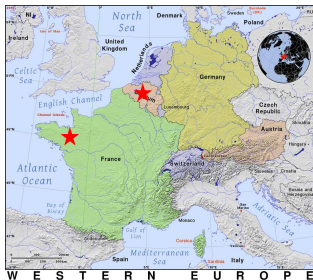


## About me: Ph.D., 2021-2024



**Jean-François Flot**  
*Université Libre de Bruxelles*  
Focus: assembling wild genomes



**Dominique Laveneir**  
*Université de Rennes*  
Focus: computational methods

## Metagenome assembly from long reads

## About me: postdoc, since 2025



Institut Pasteur, Paris



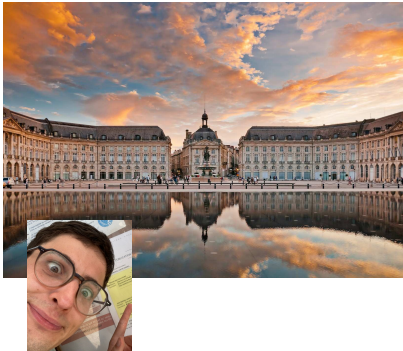
**Rayan Chikhi**

*Institut Pasteur, Paris*

Focus: massive genomics

## Index & Search the Logan database

## About me: future in Bordeaux??



*Inria*??  
cnrs??

# The Logan project: indexing and querying **all** the (meta)genomic data ever published

Roland Faure<sup>1</sup>

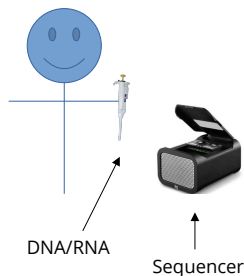
<sup>1</sup>Institut Pasteur

January 2026

Slides available (CC-BY) at: [rolandfaure.github.io](https://rolandfaure.github.io)

# The Logan database

# J. Doe sequenced something



Sequencing  
Data

ACAGCAGT  
CAGTCGTA  
CGTATAA

Article

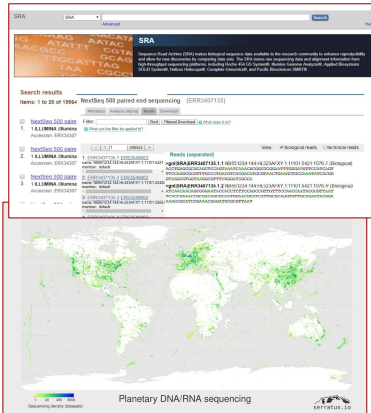


Supporting data  
e.g. SRR12345



# The SRA database

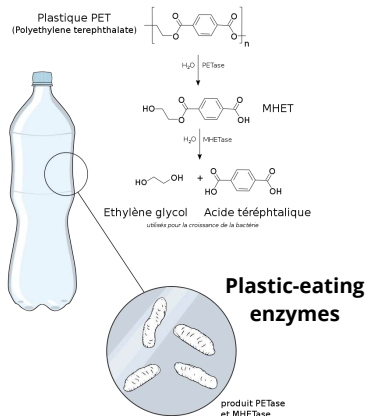
SRA: All public sequencing reads, 80 PB of data



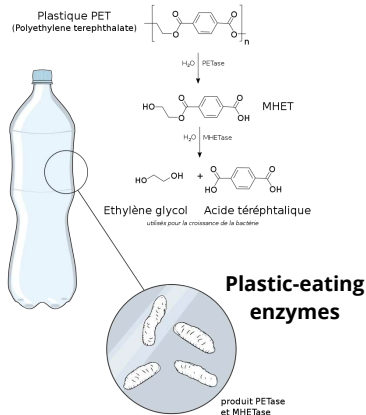
“Library of Alexandria” for genetics

Slide Credits: Rayan Chikhi

# Plastic-eating enzymes: PETases



# Plastic-eating enzymes: PETases

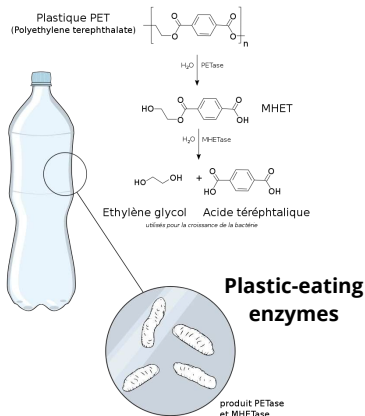


We know of 200 such enzymes but there is much we ignore, and they are hard to find



Artem Babaian

# Plastic-eating enzymes: PETases



We know of 200 such enzymes but there is much we ignore, and they are hard to find

There must be more of them in the SRA!



Artem Babaian

# The SRA is not queryable

SRA

CGACTCGTCGCTCGCATG

Search

[Create alert](#) [Advanced](#) [Help](#)

 The following term was not found in SRA: CGACTCGTCGCTCGCATG.

 No items found.

Search details

(CGACTCGTCGCTCGCATG[All Fields])

## The SRA now



Slide Credits: Teo Lemane

► 27 millions accessions

Quiz: At 1 Gbit/s, how much time to download the SRA?

A: 20 days

B: 20 weeks

C: 20 months

D: 20 years

Quiz: At 1 Gbit/s, how much time to download the SRA?

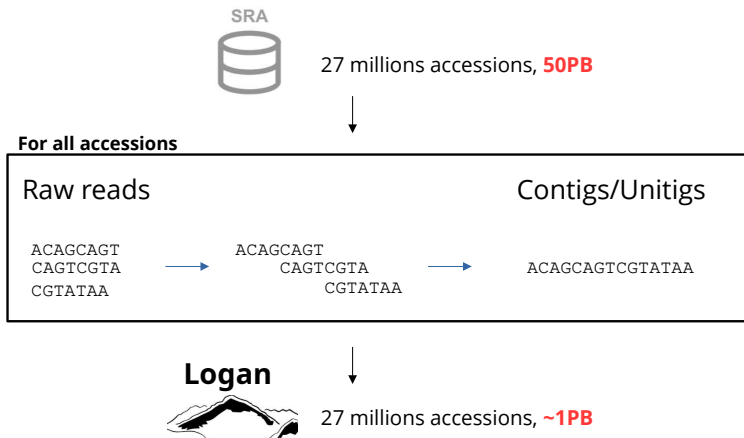
A: 20 days

B: 20 weeks

C: 20 months

D: 20 years

# The Logan database



How much did it cost to assemble Logan?

A: \$5,000

B: \$50,000

C: \$500,000

D: \$5,000,000

How much did it cost to assemble Logan?

A: \$5,000

B: \$50,000

C: \$500,000

D: \$5,000,000

# The Logan database

- ▶ 2.18 million parallel CPUs, 30h wall-clock time
- ▶ All the assemblies are available online

# The Logan database

- ▶ 2.18 million parallel CPUs, 30h wall-clock time
- ▶ All the assemblies are available online

## Downloading

To download one accession, type:

```
wget https://s3.amazonaws.com/logan-pub/c/[accession]/[accession].contigs.fa.zst
```



# The Logan database

- ▶ 2.18 million parallel CPUs, 30h wall-clock time
- ▶ All the assemblies are available online

## Downloading

To download one accessi

```
wget https://s3.amazo
```

Let's look for  
homologs of my  
enzyme in Logan!



Artem Babaian

a.zst

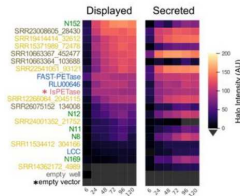


## Searching through Logan

- ▶ Downloaded the 1PB of contigs
- ▶ Aligned the known enzymes against the contigs (DIAMOND)

## Searching through Logan

- ▶ Downloaded the 1PB of contigs
- ▶ Aligned the known enzymes against the contigs (DIAMOND)
- ▶ 1.12 billion hits, 215 million clusters 90% identity
- ▶ Some discovered enzyme have better activity than known ones



How much did it cost to align on Logan?

A: \$10

B: \$100

C: \$1,000

D: \$10,000

How much did it cost to align on Logan?

A: \$10

B: \$100

C: \$1,000

D: \$10,000

## Indexing nucleotides

## Let's index nucleotides

- ▶ Reminder: **1PB** of data, 27 million datasets
- ▶ Query: sequence
- ▶ Answer:
  - ▶ Difficulty level 1: datasets containing similar sequences
  - ▶ Difficulty level 2: the actual similar sequences



Téo Lemane

## Indexing k-mers efficiently: bloom filters

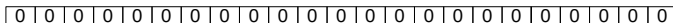
## Dataset

CACTCTGACTGA

cut in k-mers (6-mers here)

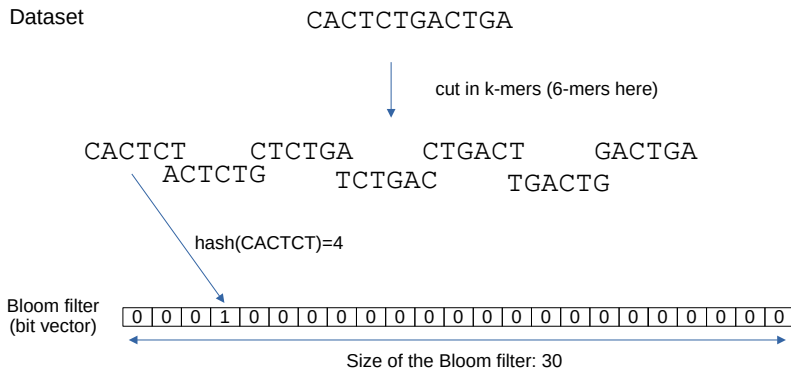
CACTCT      CTCTGA      CTGACT      GACTGA  
          ACTCTG      TCTGAC      TGACTG

Bloom filter  
(bit vector)



Size of the Bloom filter: 30

## Indexing k-mers efficiently: bloom filters



# Indexing k-mers efficiently: bloom filters

Dataset

CACTCTGACTGA

cut in k-mers (6-mers here)

CACTCT   CTCTGA   CTGACT   GACTGA  
ACTCTG   TCTGAC   TGACTG

Bloom filter  
(bit vector)

0 1 0 1 1 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 1 1 0 0 0 0 0

Size of the Bloom filter: 30

# Indexing k-mers efficiently: bloom filters

Is CTCTGA in my dataset ?

Bloom filter  
(bit vector)

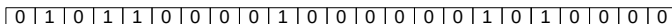
|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Indexing k-mers efficiently: bloom filters

Is CTCTGA in my dataset ?

hash(CTCTGA)=5

Bloom filter  
(bit vector)

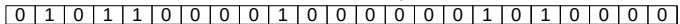


# Indexing k-mers efficiently: bloom filters

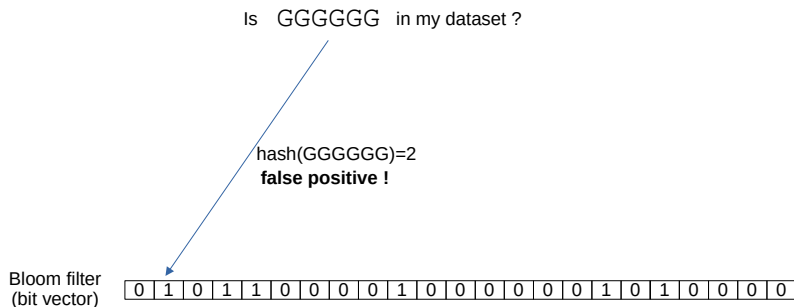
Is **AAAAAA** in my dataset ?

hash(AAAAA)=14

Bloom filter  
(bit vector)



# Indexing k-mers efficiently: bloom filters



# Indexing k-mers efficiently: bloom filters

## Trick: query (k+s)-mers

Is GGGGGGAT in my dataset ?

GGGGGG   GGGGGA   GGGGAT

Bloom filter  
(bit vector)

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

# Indexing strategy

- ▶ Index 26-mers of all datasets in Bloom filters
- ▶ At query time, query 31-mers

|          |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SRR00001 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |   |
| SRR00002 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 |
| SRR00003 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| SRR00004 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| SRR00005 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| SRR00006 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| SRR00007 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| .        |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| .        |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
| .        |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |

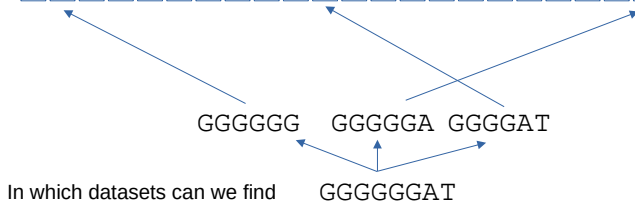
- ▶ In total, 1PB

# kmindex

|          |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SRR00001 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00002 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 |
| SRR00003 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00004 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00005 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| SRR00006 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| SRR00007 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |

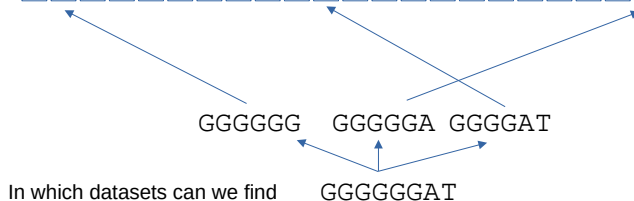
# kmindex

|          |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SRR00001 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00002 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| SRR00003 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00004 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| SRR00005 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| SRR00006 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 |
| SRR00007 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |



# kmindex

|          |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|----------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| SRR00001 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00002 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| SRR00003 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 |
| SRR00004 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| SRR00005 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| SRR00006 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| SRR00007 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 |



## ► Index of 1PB

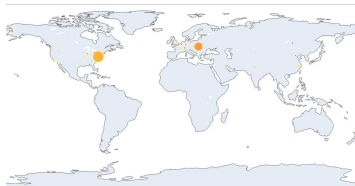
# Logan-search.org

- “Find all datasets that share x% of 31-mers with my query”

Table Map Plot BLAST-like alignment Help

kmer\_coverage = 5.7 AND assay\_type IN (WGS, WGA)

| ID                                  | kmer_coverage | bioproject                           | biosample                             |
|-------------------------------------|---------------|--------------------------------------|---------------------------------------|
| SRR16173100 <a href="#">GSA:Q10</a> | 1             | PRJNA768258 <a href="#">GSA:Q10</a>  | SAMRQ22020249 <a href="#">GSA:Q10</a> |
| SRR18416387 <a href="#">GSA:Q10</a> | 1             | PRJNA727098 <a href="#">GSA:Q10</a>  | SAMR19006852 <a href="#">GSA:Q10</a>  |
| SRR15166888 <a href="#">GSA:Q10</a> | 1             | PRJNA746342 <a href="#">GSA:Q10</a>  | SAMRQ20203734 <a href="#">GSA:Q10</a> |
| ERR33935101 <a href="#">GSA:Q10</a> | 1             | PRJEB27179 <a href="#">GSA:Q10</a>   | SAMR5390111 <a href="#">GSA:Q10</a>   |
| SRR3740047 <a href="#">GSA:Q10</a>  | 1             | PRJNA327431 <a href="#">GSA:Q10</a>  | SAMRQ5335350 <a href="#">GSA:Q10</a>  |
| ERR2398873 <a href="#">GSA:Q10</a>  | 1             | PRJEB22684 <a href="#">GSA:Q10</a>   | SAMR104883220 <a href="#">GSA:Q10</a> |
| SRR1465031 <a href="#">GSA:Q10</a>  | 1             | PRJNA732827 <a href="#">GSA:Q10</a>  | SAMR13221634 <a href="#">GSA:Q10</a>  |
| SRR13300852 <a href="#">GSA:Q10</a> | 1             | PRJNA687219 <a href="#">GSA:Q10</a>  | SAMR177141291 <a href="#">GSA:Q10</a> |
| SRR18416320 <a href="#">GSA:Q10</a> | 1             | PRJNA727098 <a href="#">GSA:Q10</a>  | SAMR19006764 <a href="#">GSA:Q10</a>  |
| ERR2399820 <a href="#">GSA:Q10</a>  | 1             | PRJEB22684 <a href="#">GSA:Q10</a>   | SAMR104885216 <a href="#">GSA:Q10</a> |
| SRR13011679 <a href="#">GSA:Q10</a> | 1             | PRJNA668406 <a href="#">GSA:Q10</a>  | SAMR16710006 <a href="#">GSA:Q10</a>  |
| ERR2395390 <a href="#">GSA:Q10</a>  | 1             | PRJEB22684 <a href="#">GSA:Q10</a>   | SAMR104881707 <a href="#">GSA:Q10</a> |
| ERR2394034 <a href="#">GSA:Q10</a>  | 1             | PRJEB22684 <a href="#">GSA:Q10</a>   | SAMR104888353 <a href="#">GSA:Q10</a> |
| SRR25617432 <a href="#">GSA:Q10</a> | 1             | PRJNA1004049 <a href="#">GSA:Q10</a> | SAMRQ6020974 <a href="#">GSA:Q10</a>  |
| SRR25609779 <a href="#">GSA:Q10</a> | 1             | PRJNA1006165 <a href="#">GSA:Q10</a> | SAMRQ37013649 <a href="#">GSA:Q10</a> |
| SRR6069137 <a href="#">GSA:Q10</a>  | 1             | PRJNA511728 <a href="#">GSA:Q10</a>  | SAMRQ0318067 <a href="#">GSA:Q10</a>  |
| SRR11318504 <a href="#">GSA:Q10</a> | 1             | PRJNA6612988 <a href="#">GSA:Q10</a> | SAMR14389256 <a href="#">GSA:Q10</a>  |
| SRR25507401 <a href="#">GSA:Q10</a> | 1             | PRJNA1001958 <a href="#">GSA:Q10</a> | SAMR13623017 <a href="#">GSA:Q10</a>  |



How much does it cost to query Logan-search?

A: \$1

B: \$10

C: \$100

D: \$1,000

How much does it cost to query Logan-search?

A: \$1

B: \$10

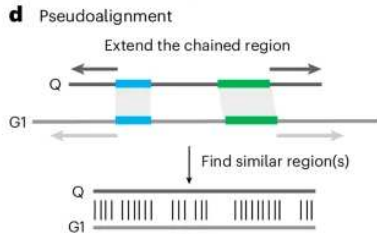
C: \$100

D: \$1,000

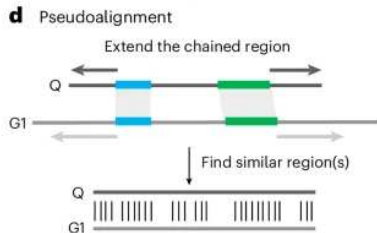
## Indexing sequences: limits & future

- ▶ Limit: slow to get the actual sequences
- ▶ Limit: query and target need to share 31-mers

## Another strategy to index sequences: LexicMap



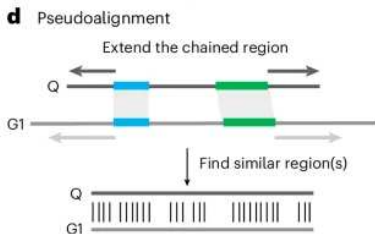
## Another strategy to index sequences: LexicMap



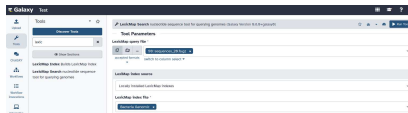
- Available on Galaxy soon



## Another strategy to index sequences: LexicMap



- Available on Galaxy soon



- More sensitive but slow (several hours)

## Indexing proteins

# Obtaining all the proteins of SRA

- ▶ Ran prodigal on all assemblies: 100 billion proteins
- ▶ Clustered with MMseqs2 at 50% identity: 3 billion representative proteins

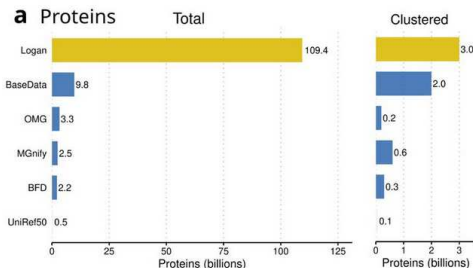
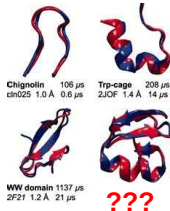


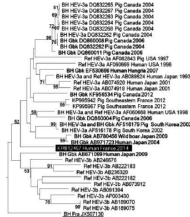
Figure from the Logan preprint

# What can we do with these proteins?

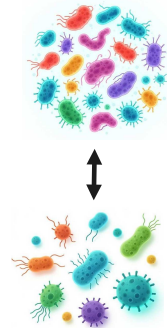
Discover new proteins



Improve protein  
phylogenies



Compare datasets



## My work: looking for proteins by similarity

- ▶ Query: a protein
- ▶ Answer: all similar proteins in Logan

Let's look for  
homologs of my  
enzyme in Logan!



Artem Babaian

# How to compare (3 billion) proteins?

## Strategy 1: sequence comparisons

```
MRIFGFFITLVAIIQ  
  |||||  
MRIKGFFITLAIIFQ
```

# How to compare (3 billion) proteins?

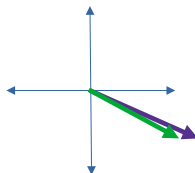
Strategy 1: sequence  
comparisons

MRIF**GFFITLVA**II**GQ**  
 |||||  
 MR**IKGFFITL****IA**II**FQ**

Strategy 2: embedding  
comparisons

**MR**IK**GFFITL**IA**II**F**Q**  
 MR**I**F**GFFITL**V**AI**I**GQ**

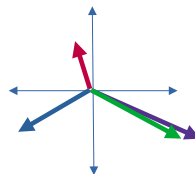
↓ Protein Language  
Model embedding



# Indexing 3 billion proteins

MRIGGFFITLIAIIFQ  
MRIFGFFITLVAIIGQ  
MSIYHMKVRTITGKDMTLQP  
MTFFLYISPMISILIGFK  
.....

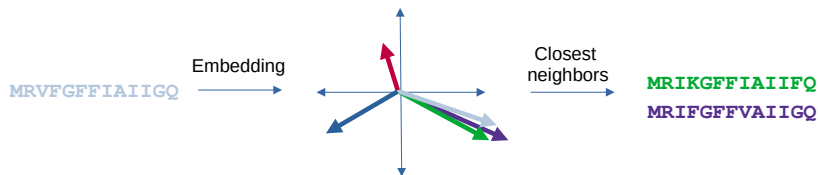
Embedding



Vector database

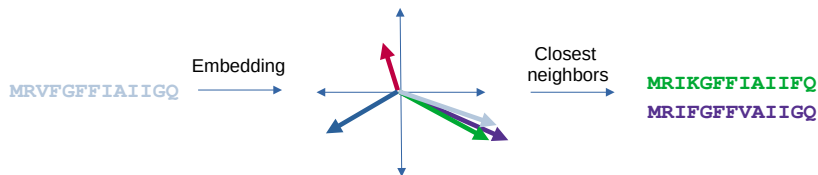
- ▶ Protein Language Model: gLM2
- ▶ 512-dimension vectors
- ▶ 3k GPU.hours
- ▶ Space taken by final database: 1.5TB

## Querying 3 billion proteins



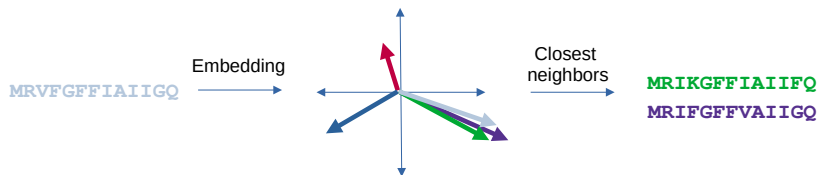
- Embed query & compare to existing vectors

## Querying 3 billion proteins



- ▶ Embed query & compare to existing vectors
- ▶ Compare to 3 billion vectors using industrial vector databases

## Querying 3 billion proteins



- ▶ Embed query & compare to existing vectors
- ▶ Compare to 3 billion vectors using industrial vector databases
- ▶ Actually just brute force comparison

## Protein search: performance

- ▶ Query: 301 known papillomavirus proteins

## Protein search: performance

- ▶ Query: 301 known papillomavirus proteins
- ▶ 4h computation, 120 GB RAM (~ \$10)

## Protein search: performance

- ▶ Query: 301 known papillomavirus proteins
- ▶ 4h computation, 120 GB RAM (~ \$10)
- ▶ 50k homologs in Logan

## Protein search: performance

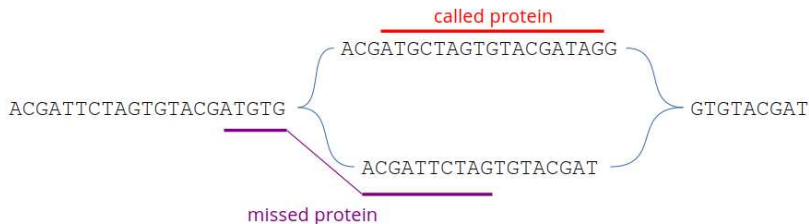
- ▶ Query: 301 known papillomavirus proteins
- ▶ 4h computation, 120 GB RAM (~ \$10)
- ▶ 50k homologs in Logan
- ▶ Soon available online on Galaxy and downloadable (~ 5TB)
- ▶ Contact me if you are interested now

# Limits

- ▶ Only full proteins match

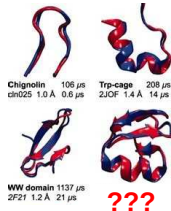
## Limits

- ▶ Only full proteins match
- ▶ 90% proteins missing in the database because of the protein calling

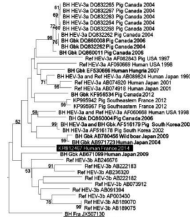


# Future developments of Logan

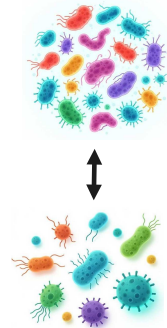
Discover new proteins



Improve protein  
phylogenies



Compare datasets



## Take-home message

- ▶ Logan centralizes all public datasets
- ▶ Nucleotide and proteins

## Take-home message

- ▶ Logan centralizes all public datasets
- ▶ Nucleotide and proteins



- ▶ Contact us if you need help using Logan tools
- ▶ Contact us if you want some features
- ▶ You can propose features too