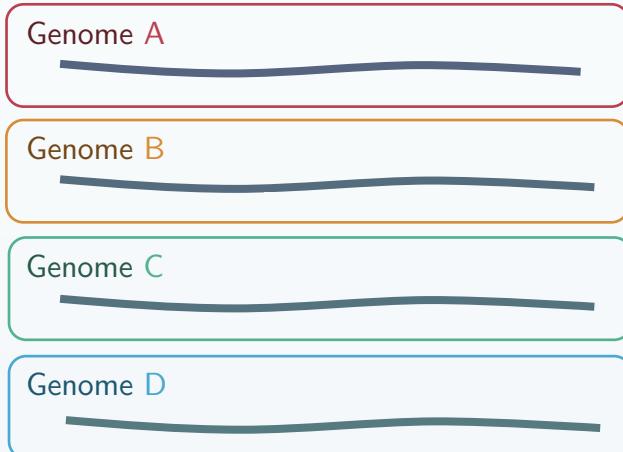


# Developing a standard interface for $k$ -mer-based index structures (update)

Andreas Rempel · June 13, 2022

# Software for Computational Pangenomics

## ❖ Data structures

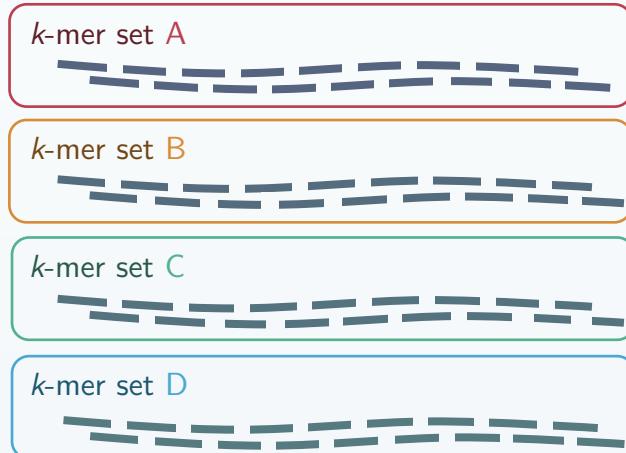


### Computational Pangenomics

- Thousands of input sequences
- Process large number of strings
  - in a reasonable time
  - within the limits of RAM
- Construct a compact index that allows fast queries

# Software for Computational Pangenomics

## ❖ Data structures



### Computational Pangenomics

- Requires advanced algorithms and data structures
- **colored *k*-mer sets**

Idea: *k*-mers & colors

Split each genomic sequence into overlapping substrings of length *k*

# Software for Computational Pangenomics

## ❖ Data structures

Genome A  
ATGTCA**G**CTA

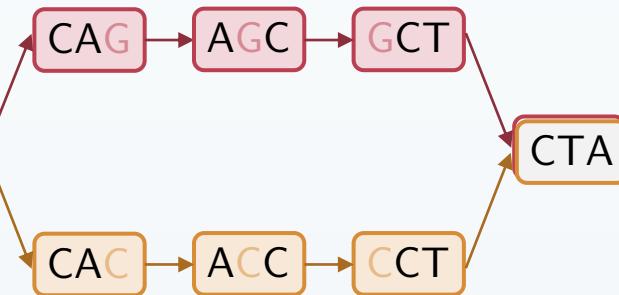
$k=3$

Genome B  
ATGTCA**C**CTA



symmetric “bubbles”  
→ SNPs, MNPs

colored De Bruijn Graph



# Software for Computational Pangenomics

## ❖ Data structures

Genome A  
ATGTCA**G**CTA

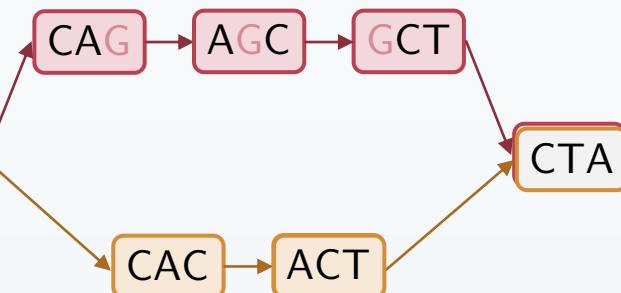
$k=3$

Genome B  
ATGTCA**-**CTA



asymmetric “bubbles”  
→ insertions, deletions

colored De Bruijn Graph

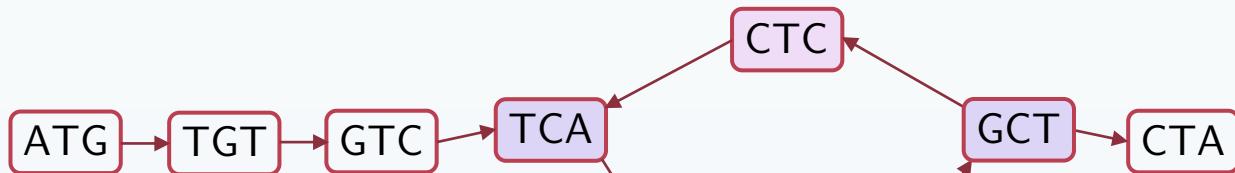


# Software for Computational Pangenomics

## ❖ Data structures

Genome A  
ATGT**CAGCTCAGCTA**

colored De Bruijn Graph



“loops”  
→ repeats

and many more complex structures...

# Software for Computational Pangenomics

## ❖ Data structures

### Computational Pangenomics

- Requires advanced algorithms and data structures  
→ **colored  $k$ -mer sets**

Idea:  $k$ -mers & colors

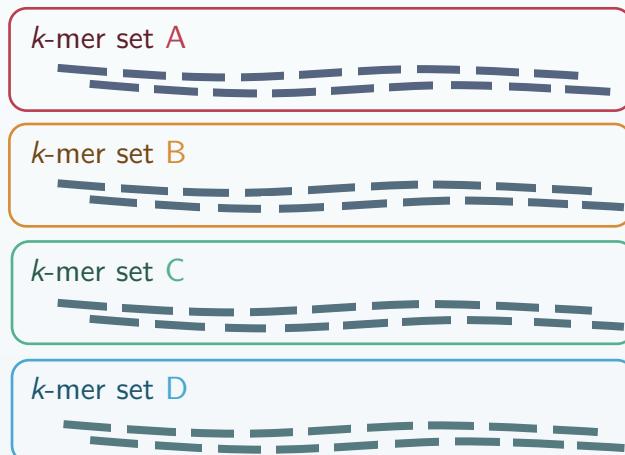
Split each genomic sequence into overlapping substrings of length  $k$

AACG	CAAG
Colors: {A}	Colors: {A, B}
AACT	CACG
Colors: {B}	Colors: {C, D}
ACCG	CAGG
Colors: {C}	Colors: {B, C, D}
ACCT	CCGG
Colors: {D}	Colors: {A, B, C, D}

# Software for Computational Pangenomics

## ❖ Data structures

For each input, store the set of  $k$ -mers



<b>AACG</b> Colors: {A}	<b>CAAG</b> Colors: {A, B}
<b>AACT</b> Colors: {B}	<b>CACG</b> Colors: {C, D}
<b>ACCG</b> Colors: {C}	<b>CAGG</b> Colors: {B, C, D}
<b>ACCT</b> Colors: {D}	<b>CCGG</b> Colors: {A, B, C, D}

For each  $k$ -mer, indicate the input data in which the  $k$ -mer is present

# Software for Computational Pangenomics

## ❖ Data structures

For each input, store the set of  $k$ -mers

Genome A  
ATGTCAGCTA

$k=3$        $h=3$

ATG

TGT

GTC

TCA

CAG

AGC

GCT

$h_1$        $h_2$        $h_3$

set of hash functions

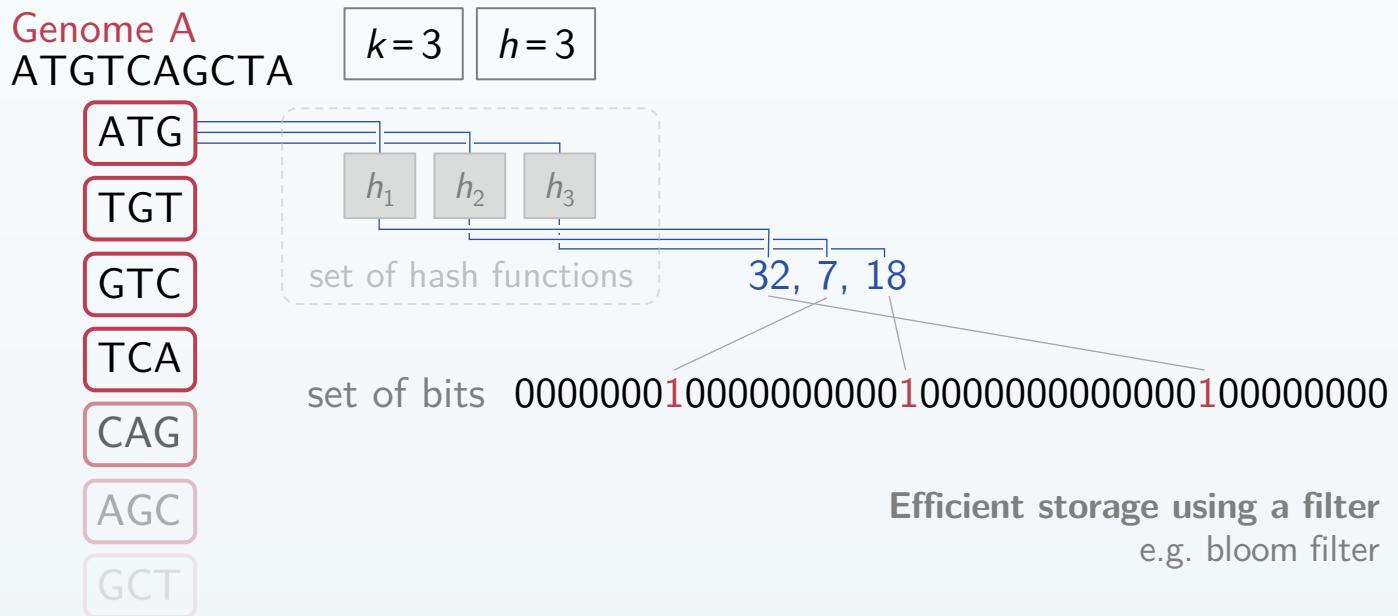
set of bits 000

Efficient storage using a filter  
e.g. bloom filter

# Software for Computational Pangenomics

## ❖ Data structures

For each input, store the set of  $k$ -mers



# Software for Computational Pangenomics

## ❖ Data structures

For each input, store the set of  $k$ -mers

Genome A

ATGTCAGCTA

00111101101010001111110001001001100010000

Genome B

ACGTTAGCTG

01010111011010110100000011110111000111100

Genome C

CGGTAAACGAG

00010111010011011001001011000011110010011

Genome D

GTATTACCAT

00101100110011110011011101110001000010100

Efficient storage using a filter  
e.g. bloom filter

# Software for Computational Pangenomics

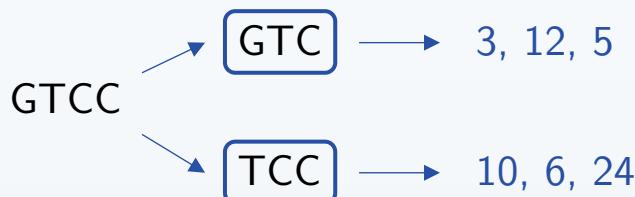
## ❖ Data structures

For each input, store the set of  $k$ -mers

Index

Genome A / SRA Exp.	00111101101010001111110001001001100010000
Genome B / SRA Exp.	01010111011010110100000011110111000111100
Genome C / SRA Exp.	00010111010011011001001011000011110010011
Genome D / SRA Exp.	00101100110011110011011101110001000010100

Query



is probably in Genome B  
definitely not in A, C, D

false positive rate due to collisions,  
is controlled by  $h$  & number of bits

# Software for Computational Pangenomics

## ❖ Available tools

<b>SeqOthello</b>	<b>Cortex</b>	<b>kmtricks</b>
<b>Mantis</b>	<b>McCortex</b>	<b>kcollections</b>
<b>SBT/SSBT</b>	<b>BIGSI</b>	<b>Cosmo VARI</b>
<b>AllSomeSBT</b>	<b>COBS</b>	<b>Rainbowfish</b>
<b>HowDeSBT</b>	<b>RAMBO</b>	<b>Cuttlefish</b>
<b>BFT</b>	<b>Raptor</b>	<b>REINDEER</b>
<b>Bifrost</b>	<b>Themisto</b>	<b>GATB-core</b>

... and many more

# Software for Computational Pangenomics

## ❖ Available tools

**SeqOthello**

PreProcess

Group

Build

Query

**Mantis**

count

build

mst

query

**SBT/SSBT**

hashes

count

build

query

**HowDeSBT**

makebf

cluster

build

query

**BIGSI**

build

bloom

build

search

# Software for Computational Pangenomics

## ❖ Available tools

**SeqOthello**    **Mantis**    **SBT/SSBT**    **HowDeSBT**    **BIGSI**



# Developing a common standard interface

## ❖ Software for colored $k$ -mer sets (SOCKS)



Knut Reinert



Nadia Pisanti



Carl Kingsford



Brad Solomon



Gaurav Gupta



Pierre Peterlongo



Rayan Chikhi



Yoann Dufresne



Paul Medvedev



Sven Rahmann

# Software for colored $k$ -mer sets (SOCKS)

## ❖ SOCKS interface

- construct index
  - input: set of sequences, each with a distinct color
  - output: index

input: plain text file containing the sequence file names

- list of assembled genomes/contigs, multiple fasta
- list of read sets, single-end or paired-end fastq
- one line per color, e.g. whitespace-separated

output: one or multiple files for the constructed index

- proprietary/binary format, compact, fast read/write
- interoperable format, e.g. gfa | fasta |  $k$ -mer matrix

# Software for colored $k$ -mer sets (SOCKS)

## ❖ SOCKS interface

- **build** : construct an index from a set of sequences
  - **input**: plain text file containing the sequence file names, e.g.:

```
COLOR_NAME_1: /PATH/T0/GENOME.FASTA
```

```
COLOR_NAME_2: /PATH/T0/READ_1.FASTQ /PATH/T0/READ_2.FASTQ
```

- **output**: index in binary or interoperable format, e.g. [kmer file format](#)  
*(at least one of the options, binary or interoperable format, should be provided)*  
*(if both options are provided, it should be possible to switch using a parameter)*

# Software for colored $k$ -mer sets (SOCKS)

## ❖ SOCKS interface

- query sequence
  - input:  $k$ -mer(s)
  - output: set of colors

input: plain text file containing the  $k$ -mers

- one line per  $k$ -mer query

output: plain text file listing the color sets

- list of positive hits, e.g. whitespace-separated
- hit/miss vector, e.g. present = '1', absent = '0'
- one line per  $k$ -mer query result

# Software for colored $k$ -mer sets (SOCKS)

## ❖ SOCKS interface

- **lookup-kmer** : find the color sets for a list of  $k$ -mers
  - **input:** plain text file containing the  $k$ -mers, one per line, e.g.:

```
ACGTACGT  
ACCTAGGT
```

- **output:** plain text file listing the color set for each  $k$ -mer, e.g.:

(as list of positive hits)

```
ACGTACGT: COLOR_NAME_1 COLOR_NAME_4 COLOR_NAME_7 ...  
ACCTAGGT: COLOR_NAME_1 COLOR_NAME_5 COLOR_NAME_8 ...
```

# Software for colored $k$ -mer sets (SOCKS)

## ❖ SOCKS interface

- **lookup-kmer** : find the color sets for a list of  $k$ -mers
  - **input:** plain text file containing the  $k$ -mers, one per line, e.g.:

```
ACGTACGT  
ACCTAGGT
```

- **output:** plain text file listing the color set for each  $k$ -mer, e.g.:

(or as a binary vector)

```
ACGTACGT: 10010010...  
ACCTAGGT: 10001001...
```

# Software for colored $k$ -mer sets (SOCKS)

## ❖ Available tools

<b>SeqOthello</b>	<b>Cortex</b>	<b>kmtricks</b>
<b>Mantis</b>	<b>McCortex</b>	<b>kcollections</b>
<b>SBT/SSBT</b>	<b>BIGSI</b>	<b>Cosmo VARI</b>
<b>AllSomeSBT</b>	<b>COBS</b>	<b>Rainbowfish</b>
<b>HowDeSBT</b>	<b>RAMBO</b>	<b>Cuttlefish</b>
<b>BFT</b>	<b>Raptor</b>	<b>REINDEER</b>
<b>Bifrost</b>	<b>Themisto</b>	<b>GATB-core</b>

... and many more

UNIVERSITÄT  
BIELEFELD  
Technische Fakultät

localhost:8000/tools/

▶ Genome Informatics   ▶ Research   ▶ Studies   ▶ Internal

You are here: home » research » [panbench](#)

## Pangenomics Benchmark & Workbench

≡

**bifrost**  
Bifrost: Highly parallel construction and indexing of colored and compacted de Bruijn graphs

**BIGSI**  
Bitsliced Genomic Signature Index - Efficient indexing and search in very large collections of WGS data

**BloomFilterTrie**  
An alignment-free, reference-free and incremental data structure for colored de Bruijn graph with application to pan-genome indexing.

**bloomtree**  
No description provided.

**bloomtree-allsome**  
Sequence Bloom Trees with All/Some split

**cohs**

seqan/raptor × +

localhost:8000/tools/raptor/

UNIVERSITÄT  
BIELEFELD  
Technische Fakultät



▶ Genome Informatics   ▶ Research   ▶ Studies   ▶ Internal

You are here: home » research » [panbench](#)

## Raptor

 Raptor CI passing  codecov 100% [Install with bioconda](#) [Install with brew](#)

from seqan/raptor Last update: May 18, 2022

**Required**

names ?  files ?

[Browse...](#) No file selected. [Browse...](#) No file selected.

**Optional**

kmer ?  size ?

hash ?  compressed ?

**Privacy**

I would like to be notified when my job status changes. Please contact me via the following email address: \_\_\_\_\_

I want my browser to keep track of my current jobs and parameter settings. This requires the use of cookies.

[Reset](#) [Submit](#)

---

A fast and space-efficient pre-filter for querying very large collections of nucleotide sequences

Download and Installation

seqan/raptor

localhost:8000/tools/

Open + names.txt /home/example Save ...

UNIVERSITÄT  
BIELEFELD  
Technische Fakultät

► Genome Informatics ► Research

Raptor Raptor CI passing codecov 10 from seqan/raptor

**Required**  
names ?   
[Browse...](#) No file selected.

**Optional**  
kmer ?  20  
hash ?  2  
compressed ?  false

Plain Text ▾ Tab Width: 4 ▾ Ln 1, Col 1 ▾ INS

1 536: 536.fasta  
2 APEC01: APEC01.fasta  
3 ATCC8739: ATCC8739\_1.fasta ATCC8739\_2.fasta  
4 B18BS512: B18BS512.fasta  
5 B45b227: B45b227.fasta  
6 BW2952: BW2952\_1.fasta BW2952\_2.fasta BW2952\_3.fasta  
7 CB9615: CB9615.fasta  
8 CFT073: CFT073.fasta  
9 D15d197: D15d197.fasta  
10 DH10B: DH10B\_1.fasta DH10B\_2.fasta DH10B\_3.fasta DH10B\_4.fasta  
11 E234869: E234869.fasta  
12 E24377A: E24377A.fasta  
13 ED1a: ED1a\_1.fasta ED1a\_2.fasta ED1a\_3.fasta  
14 EDL933: EDL933.fasta  
15 F2a2457T: F2a2457T\_1.fasta F2a2457T\_2.fasta  
16 F2a301: F2a301.fasta  
17 F5b8401: F5b8401.fasta  
18 HS: HS\_1.fasta HS\_2.fasta HS\_3.fasta  
19 IAI1: IAI1.fasta  
20 IAI39: IAI39.fasta  
21 MG1655: MG1655.fasta  
22 S88: S88\_1.fasta S88\_2.fasta S88\_3.fasta S88\_4.fasta  
23 Sakai: Sakai.fasta  
24 SE11: SE11.fasta  
25 SMS35: SMS35\_1.fasta SMS35\_2.fasta SMS35\_3.fasta  
26 SSSs046: SSSs046.fasta  
27 UMN026: UMN026\_1.fasta UMN026\_2.fasta  
28 UTI89: UTI89.fasta  
29 W3110: W3110.fasta

Privacy

I would like to be notified when my job status changes. Please contact me via the following email address: \_\_\_\_\_  
 I want my browser to keep track of my current jobs and parameter settings. This requires the use of cookies.

Reset Submit

A fast and space-efficient pre-filter for querying very large collections of nucleotide sequences

Download and Installation

seqan/raptor

localhost:8000/tools/raptor/6807e7e7-566c-47a8-921b-2465b1e8fcfd7

UNIVERSITÄT  
BIELEFELD  
Technische Fakultät

▶ Genome Informatics   ▶ Research   ▶ Studies   ▶ Internal

You are here: home » research » [panbench](#)

## Raptor

Raptor CI passing codecov 100% [Install with bioconda](#) [Install with brew](#)

from seqan/raptor Last update: May 18, 2022

```
6807e7e7-566c-47a8-921b-2465b1e8fcfd7$ ./build/bin/raptor socks build ./names.txt --output ./index --kmer 20 --hash 2  
--size 1k

Process finished with exit code 0
```

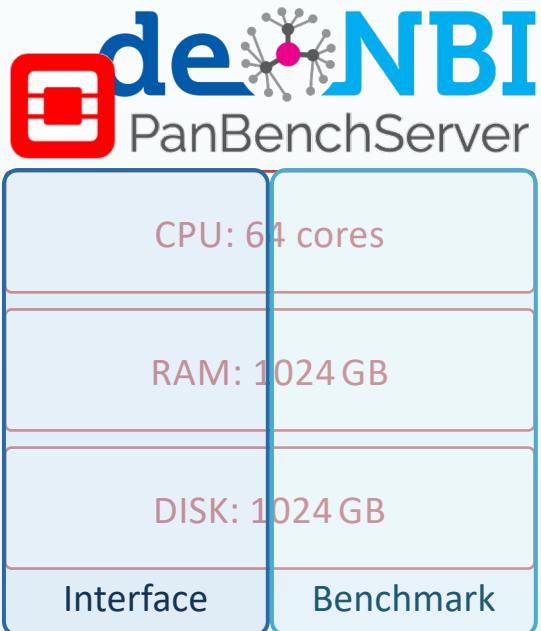
• COMPLETE

◀ BACK CANCEL DOWNLOAD

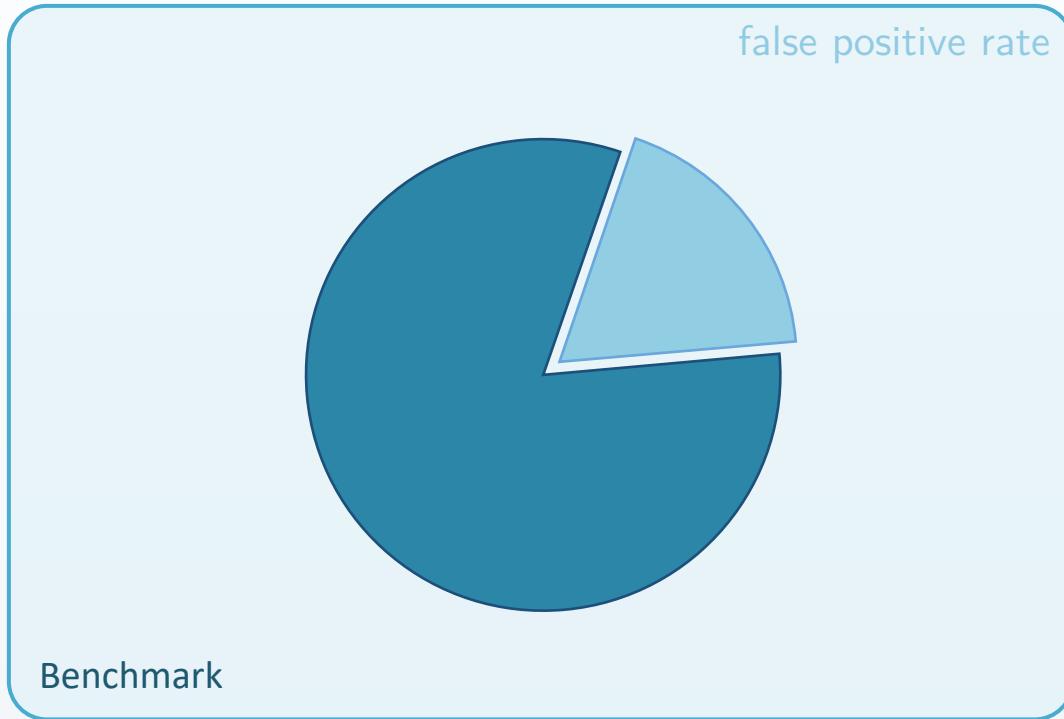
A fast and space-efficient pre-filter for querying very large collections of nucleotide sequences

Download and Installation

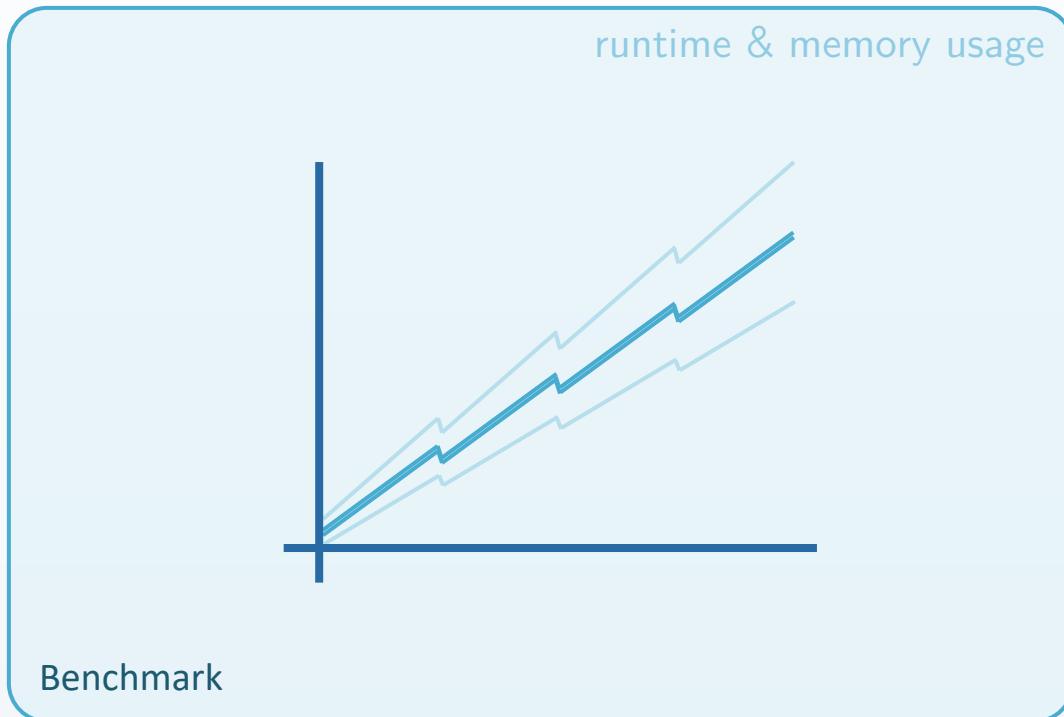
# Pangenomics Benchmark and Workbench



# Pangenomics Benchmark and Workbench



# Pangenomics Benchmark and Workbench



# Pangenomics Benchmark and Workbench

## ❖ Special Thanks

Prof. Dr. Jens Stoye  
Genome Informatics group  
Graduate School “DILS”

**Thank you for  
your attention!**



 [andreas\\_rempel](https://twitter.com/andreas_rempel)

 [andreas.rempel@uni-bielefeld.de](mailto:andreas.rempel@uni-bielefeld.de)