



UNIVERSITÄT
DES
SAARLANDES



CLEANIFIER: Removing human DNA contamination with a pangenomic gapped k -mer index

Jens Zentgraf, Johanna Elena Schmitz, Sven Rahmann

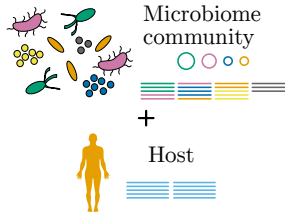
Saarland University

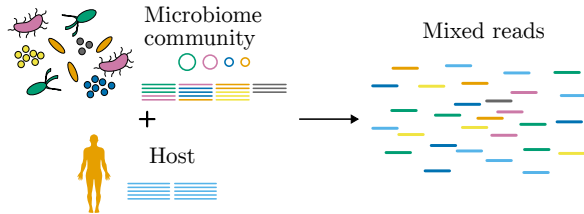
GCB 2025

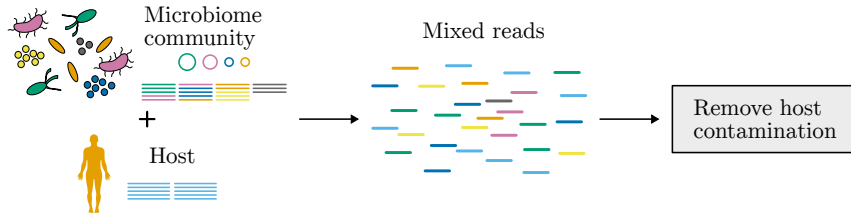


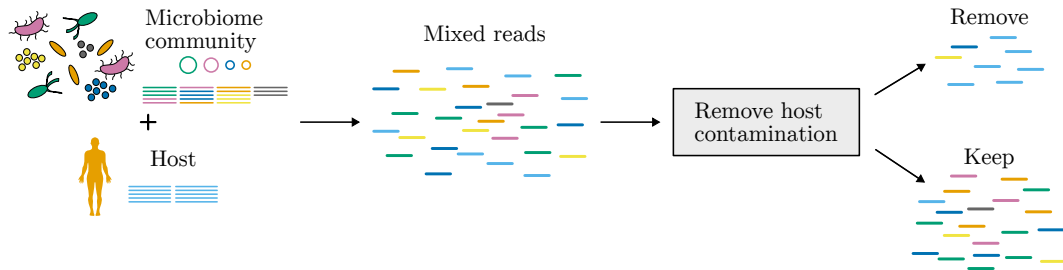
BIOCONDA
Reproducible research.
Install the software from bioconda:
> conda install -c bioconda cleaner
See bioconda.github.io

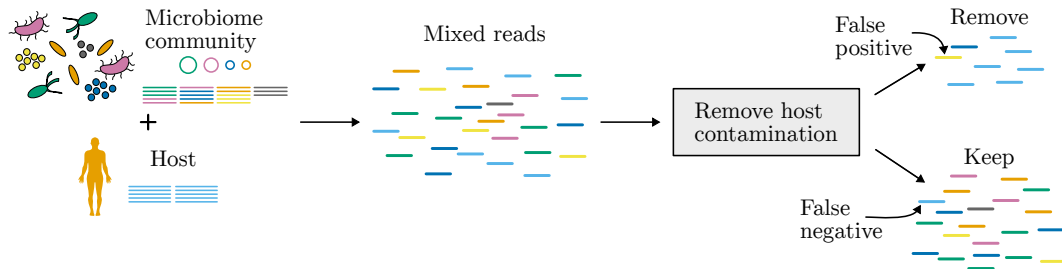


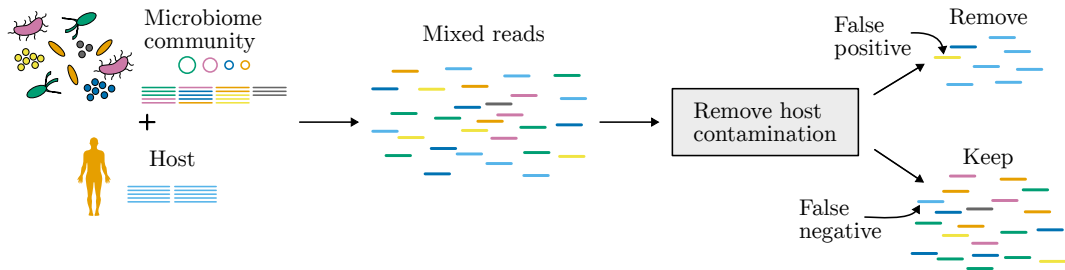






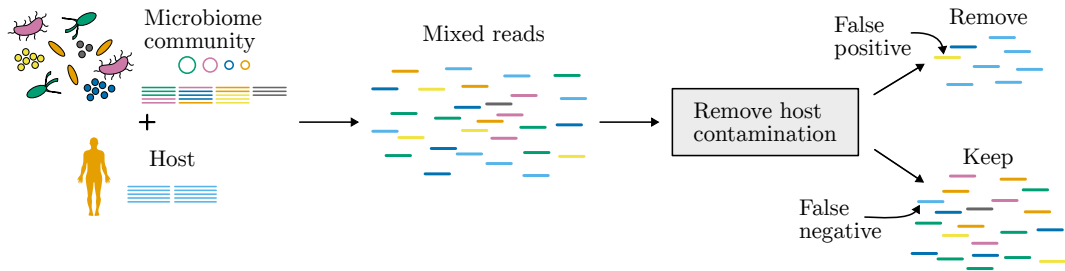






Privacy

- Human data cannot be made public



Privacy

- Human data cannot be made public

Downstream analysis

- Reduce problems in binning and assembly

Definitions

- k -mers
 - Substrings of length k

CGATCGACTAGCATCGAACGTACG . . .

k -mer

rc

canonical

Definitions

- k -mers
 - Substrings of length k

CGATCGACTAGCATCGAACGTACG . . .

k -mer rc canonical

CGATC

GATCG

ATCGA

TCGAC

CGACT

.

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$

CGATCGACTAGCATCGAACGTACG . . .

k -mer rc canonical

CGATC

GATCG

ATCGA

TCGAC

CGACT

.

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$

CGATCGACTAGCATCGAACGTACG . . .

k -mer rc canonical

CGATC GATCG

GATCG CGATC

ATCGA TCGAT

TCGAC GTCGA

CGACT AGTCG

.

.

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$
- Canonical k -mer
 - Maximum of k -mer and $rc(k\text{-mer})$

k -mer	rc	canonical
CGATC	GATCG	
GATCG	CGATC	
ATCGA	TCGAT	
TCGAC	GTCGA	
CGACT	AGTCG	
.....	

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$
- Canonical k -mer
 - Maximum of k -mer and $rc(k\text{-mer})$

k -mer	rc	canonical
CGATC	GATCG	GATCG
GATCG	CGATC	GATCG
ATCGA	TCGAT	TCGAT
TCGAC	GTCGA	TCGAC
CGACT	AGTCG	CGACT
.....

Definitions

- *k*-mers
 - Substrings of length *k*
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$
- Canonical *k*-mer
 - Maximum of *k*-mer and rc(*k*-mer)
- Gapped *k*-mers (spaced seeds)
 - *k* significant positions (#)
 - Window size *w*
 - *w* – *k* insignificant positions (_)
 - ##_#_##

CGATCGACTAGCATCGAACGTACG . . .		
<i>k</i> -mer	rc	canonical
CGATC	GATCG	GATCG
GATCG	CGATC	GATCG
ATCGA	TCGAT	TCGAT
TCGAC	GTCGA	TCGAC
CGACT	AGTCG	CGACT
.

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$
- Canonical k -mer
 - Maximum of k -mer and rc(k -mer)
- Gapped k -mers (spaced seeds)
 - k significant positions (#)
 - Window size w
 - $w - k$ insignificant positions (_)
 - ##_#_##

```
CGATCGACTAGCATCGAACGTACG . . .  
##_#_##  
CG T GA  
GA C AC  
AT G CT  
TC A TA  
CG C AG  
. . .
```

Definitions

- k -mers
 - Substrings of length k
- Reverse complements (rc)
 - Reverse order
 - $A \leftrightarrow T, C \leftrightarrow G$
- Canonical k -mer
 - Maximum of k -mer and $rc(k\text{-mer})$
- Gapped k -mers (spaced seeds)
 - k significant positions (#)
 - Window size w
 - $w - k$ insignificant positions (_)
 - ##_#_##

CGATCGACTAGCATCGAACGTACG . . .
##_#_##
CG T GA
GA C AC
AT G CT
TC A TA
CG C AG
. . .



- More robust against substitutions
- *Design of Worst-Case-Optimal Spaced Seeds*
at WABI 2025

■ HOSTILE

- Alignment based
- Human reference
- BOWTIE2 or MINIMAP2

- HOSTILE
 - Alignment based
 - Human reference
 - BOWTIE2 or MINIMAP2
- KRAKEN2
 - *k*-mer based
 - Metagenomic classification
 - Default database (very large)

- HOSTILE
 - Alignment based
 - Human reference
 - BOWTIE2 or MINIMAP2
- KRAKEN2
 - *k*-mer based
 - Metagenomic classification
 - Default database (very large)
- NOHUMAN
 - KRAKEN2 wrapper
 - Custom Human database (smaller)

■ HOSTILE

- Alignment based
- Human reference
- BOWTIE2 or MINIMAP2

■ KRAKEN2

- *k*-mer based
- Metagenomic classification
- Default database (very large)

■ NoHUMAN

- KRAKEN2 wrapper
- Custom Human database (smaller)

■ HRRT

- *k*-mer based
- Min-hash based (Jaccard Similarity)
- Include all *k*-mers from human derived eukaryotic species
- Exclude all *k*-mers in non-eukaryotic species

■ HOSTILE

- Alignment based
- Human reference
- BOWTIE2 or MINIMAP2

■ KRAKEN2

- *k*-mer based
- Metagenomic classification
- Default database (very large)

■ NOHUMAN

- KRAKEN2 wrapper
- Custom Human database (smaller)

■ HRRT

- *k*-mer based
- Min-hash based (Jaccard Similarity)
- Include all *k*-mers from human derived eukaryotic species
- Exclude all *k*-mers in non-eukaryotic species

■ DEACON

- Minimizer based
- Pangenome approach

■ HOSTILE

- Alignment based
- Human reference
- BOWTIE2 or MINIMAP2

■ KRAKEN2

- *k*-mer based
- Metagenomic classification
- Default database (very large)

■ NOHUMAN

- KRAKEN2 wrapper
- Custom Human database (smaller)

■ HRRT

- *k*-mer based
- Min-hash based (Jaccard Similarity)
- Include all *k*-mers from human derived eukaryotic species
- Exclude all *k*-mers in non-eukaryotic species

■ DEACON

- Minimizer based
- Pangenome approach

■ CLEANIFIER

- Gapped *k*-mer based
- Pangenome approach

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants
(Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants
(Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%
- Human Pangenome Reference Consortium
 - 47 assemblies

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants
(Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%
- Human Pangenome Reference Consortium
 - 47 assemblies
- IPD-IMGT/HLA database
 - HLA is highly variable

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants
(Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%
- Human Pangenome Reference Consortium
 - 47 assemblies
- IPD-IMGT/HLA database
 - HLA is highly variable
- cDNA transcripts
 - Deal with RNA

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants (Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%
- Human Pangenome Reference Consortium
 - 47 assemblies
- IPD-IMGT/HLA database
 - HLA is highly variable
- cDNA transcripts
 - Deal with RNA

Data structures

- 1 Bucketed Cuckoo hash table
 - Exact data structure
 - Stores the k -mers
 - Size 13.85 GB

#####_#####_###_#####_##### $k = 29, w = 33$

Human references

- T2T reference
- 1000 Genome Project
 - Extract all variants (Substitution, Insertion and Deletion)
 - Allele frequency of at least 1%
- Human Pangenome Reference Consortium
 - 47 assemblies
- IPD-IMGT/HLA database
 - HLA is highly variable
- cDNA transcripts
 - Deal with RNA

Data structures

- 1 Bucketed Cuckoo hash table
 - Exact data structure
 - Stores the k -mers
 - Size 13.85 GB
- 2 Windowed Cuckoo filter
 - Probabilistic set membership data structure
 - Store a fingerprint (of p bits) instead of the k -mer
 - False positive rate of $2^{-p} = 2^{-14}$
 - Size 6.9 GB

Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2, \ell = 4$ fill rate of ≈ 0.9989

Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2, \ell = 4$ fill rate of ≈ 0.9989

Slots



Windows



Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2, \ell = 4$ fill rate of ≈ 0.9989

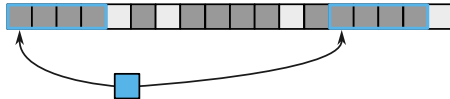
Slots



Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2$, $\ell = 4$ fill rate of ≈ 0.9989

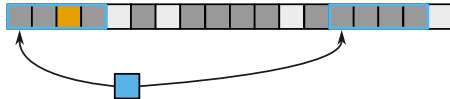
Slots



Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2$, $\ell = 4$ fill rate of ≈ 0.9989

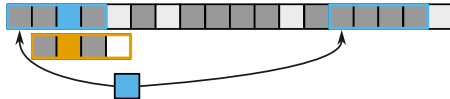
Slots



Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2$, $\ell = 4$ fill rate of ≈ 0.9989

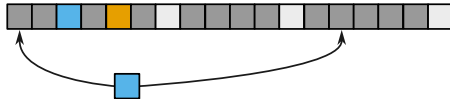
Slots



Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2$, $\ell = 4$ fill rate of ≈ 0.9989

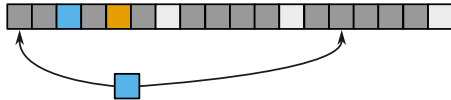
Slots



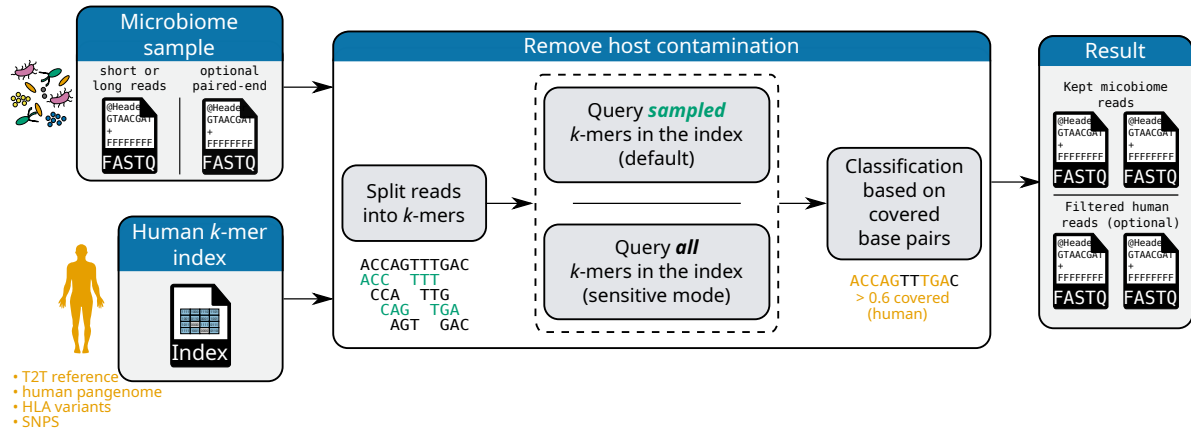
Data structure

- Probabilistic set membership data structure
- Store a fingerprint (of p bits) instead of the k -mer
- False positive rate of 2^{-p}
- d hash functions
- Window size of ℓ
- High fill rates possible
 - $d = 2, \ell = 4$ fill rate of ≈ 0.9989

Slots



Smaller and More Flexible Cuckoo Filters
at ALENEX 2026



- Query all gapped k -mers
- All k bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

- Query all gapped k -mers
- All k bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

- Query all gapped k -mers
- All k bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##

CGATCGACTAGCATCGAACGTACG . . .

CG T GA AG A CG

GA C AC GC T GA

AT G CT CA C AA

TC A TA AT G AC

CG C AG TC A CG

GA T GC CG A GT

AC A CA GA C TA

CT G AT AA G AC

TA C TC AC T CG

- Query all gapped k -mers
- All k bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

CG T GA AG A CG

GA C AC GC T GA

AT G CT CA C AA

TC A TA AT G AC

CG C AG TC A CG

GA T GC CG A GT

AC A CA GA C TA

CT G AT AA G AC

TA C TC AC T CG

- Query all gapped k -mers
- All k bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

CG T GA AG A CG

GA C AC GC T GA

AT G CT CA C AA

TC A TA AT G AC

CG C AG TC A CG

GA T GC CG A GT

AC A CA GA C TA

CT G AT AA G AC

TA C TC AC T CG

- Query every $\lfloor w/2 \rfloor$ gapped k -mer
- All w bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

- Query every $\lfloor w/2 \rfloor$ gapped k -mer
- All w bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

- Query every $\lfloor w/2 \rfloor$ gapped k -mer
- All w bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

CG T GA AT G AC

TC A TA GA C TA

AC A CA

AG A CG

- Query every $\lfloor w/2 \rfloor$ gapped k -mer
- All w bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

CGATCGACTAGCATCGAACGTACG . . .

CG T GA AT G AC

TC A TA GA C TA

AC A CA

AG A CG

- Query every $\lfloor w/2 \rfloor$ gapped k -mer
- All w bases count as covered
- Check how many bases are covered by a human gapped k -mer
- Threshold to decide if human or not

##_#_##_

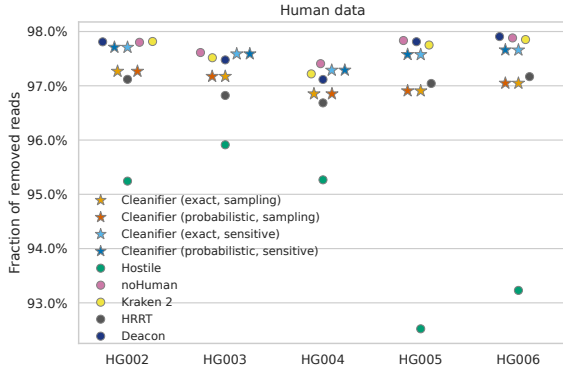
CGATCGACTAGCATCGAACGTACG . . .

CG T GA AT G AC

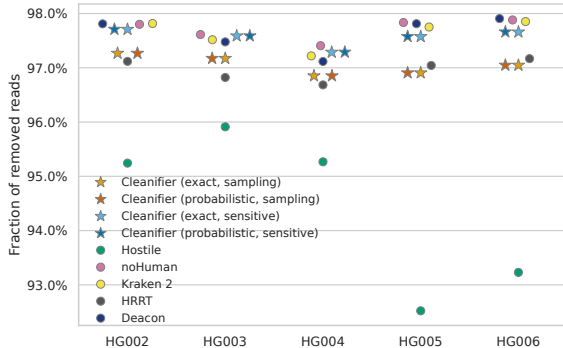
TC A TA GA C TA

AC A CA

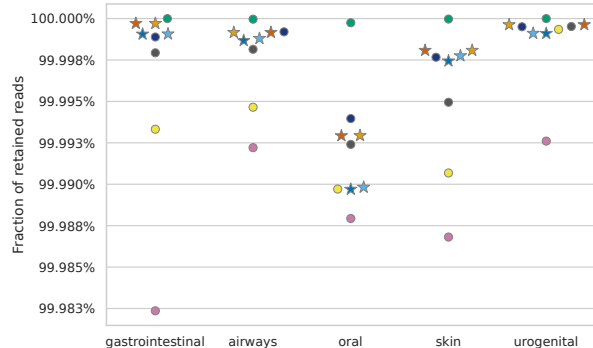
AG A CG

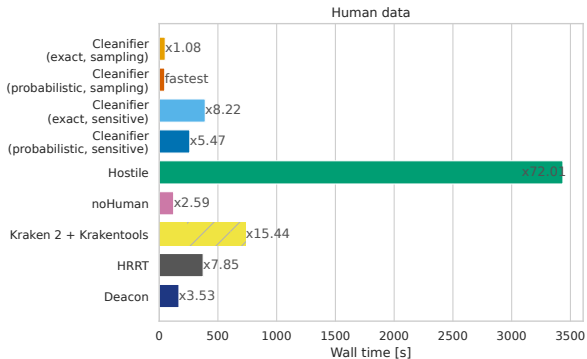


Human data

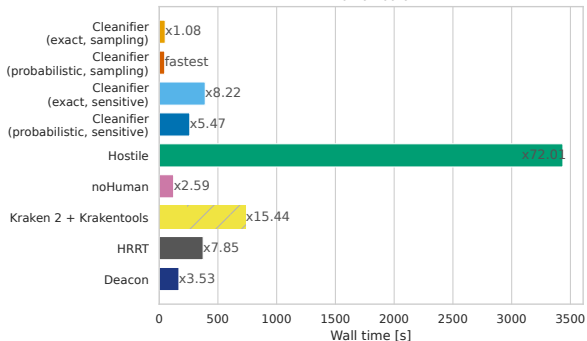


Microbiome data

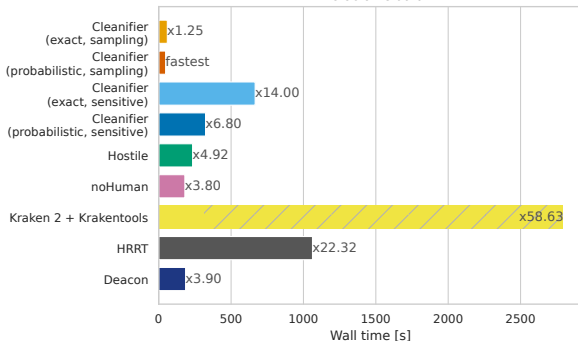




Human data

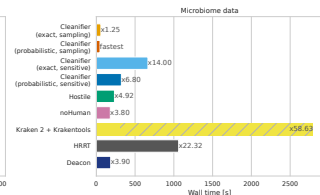
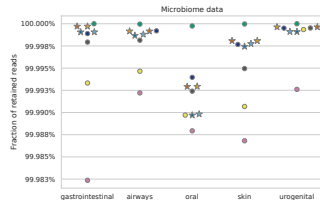
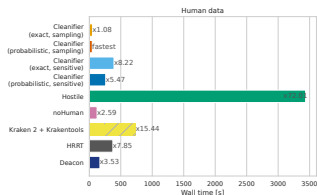
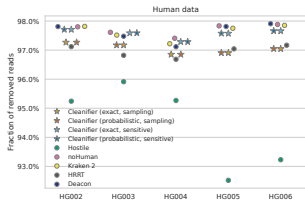


Microbiome data



CLEANIFIER

- High accuracy
- Low memory footprint (supports shared memory)
- Fast filtering
- Supports short and long reads

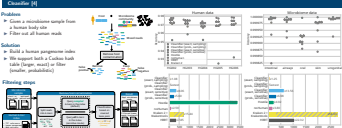
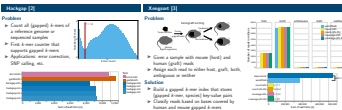
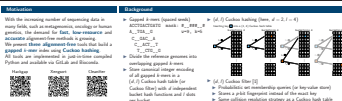


BIOCONDA
Reproducible research.
Install the software from bioconda:
> conda install -c bioconda cleanifier
See bioconda.github.io

P037



Fast and Low-Resource Alignment-free Methods for Sequence Analysis

Jens Zenggraf^{1,2} Johanna Elena Schmitz^{1,2} Sven Rahmann¹¹ Algorithmic Bioinformatics, Saarland University and Center for Bioinformatics, Saarland Informatics Campus, Saarbrücken, Germany² Saarbrücken Graduate School of Computer Science, Saarland Informatics Campus, Saarbrücken, Germany

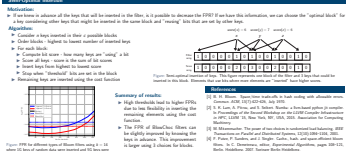
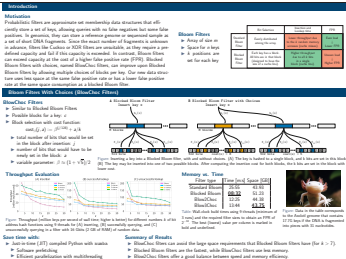
using this we can build a fast, memory-efficient and accurate

Johanna Elena Schmitz

P075



Blocked Bloom Filters with Choices

Johanna Elena Schmitz^{1,2} Jens Zenggraf^{1,2} Inês Alves Ferreira^{1,2} Sven Rahmann¹¹ Algorithmic Bioinformatics, Saarland University and Center for Bioinformatics, Saarland Informatics Campus, Saarbrücken, Germany² Saarbrücken Graduate School of Computer Science, Saarland Informatics Campus, Saarbrücken, Germany

using this we can build a fast, memory-efficient and accurate

Inês Alves Ferreira

P079



Cap k-mers: Simple but efficient and flexible seeds

Moein Karami^{1,2} Jens Zenggraf^{1,2} Sven Rahmann¹¹ Algorithmic Bioinformatics, Saarland University and Center for Bioinformatics, Saarland Informatics Campus, Saarbrücken, Germany² Saarbrücken Graduate School of Computer Science, Saarland Informatics Campus, Saarbrücken, Germany

using this we can build a fast, memory-efficient and accurate

Moein Karami