

FRACTIONAL HITTING SETS & SUPERSAMPLER

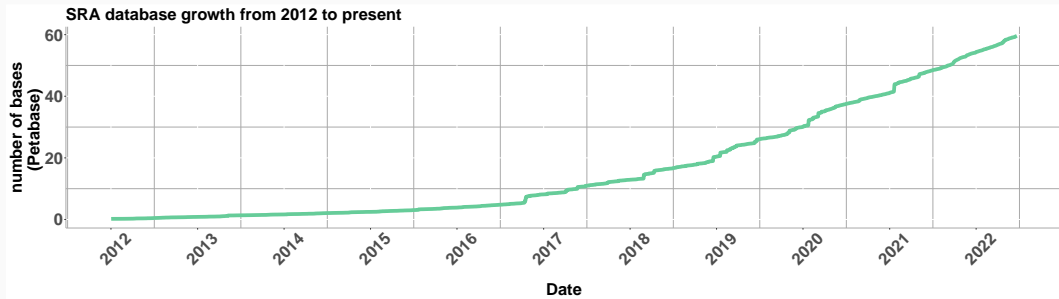
EFFICIENT AND LIGHTWEIGHT GENOMIC DATA SKETCHING

Timothé ROUZÉ, Igor MARTAYAN, Camille MARCHET & Antoine LIMASSET

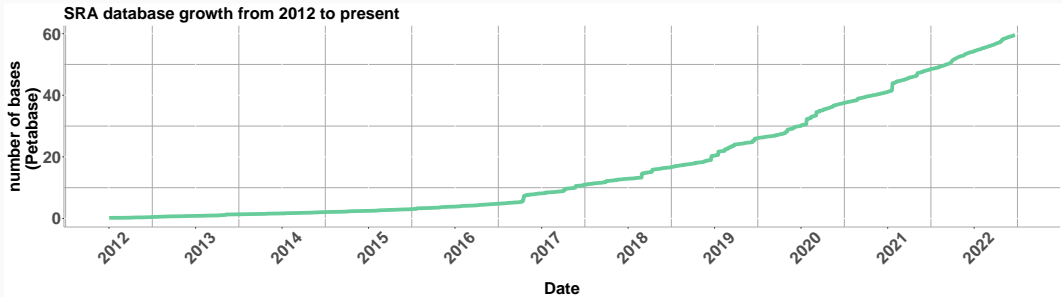
July 25, 2023



INTRODUCTION



INTRODUCTION

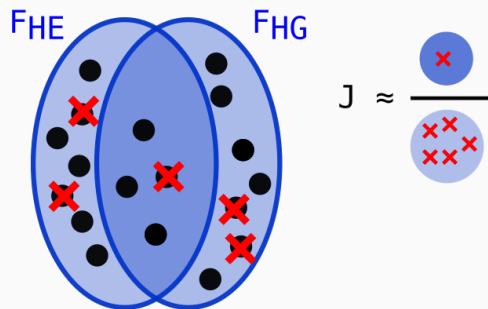
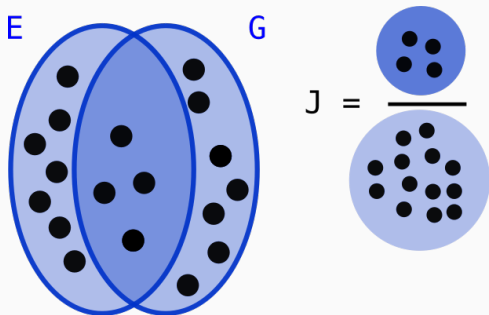


Problems

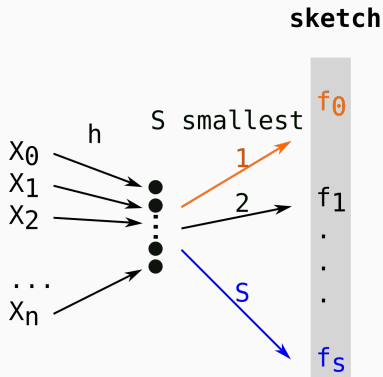
- keeps increasing
- large scale analysis not possible

PRELIMINARIES

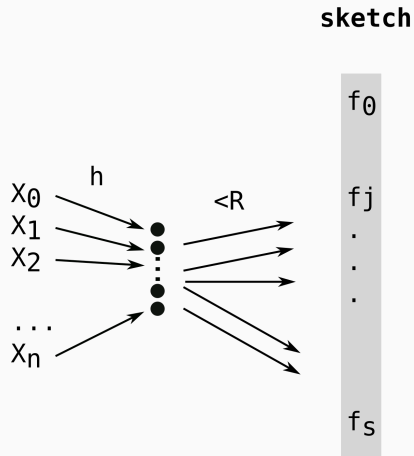
JACCARD INDEX



Bottom Minhash in MASH

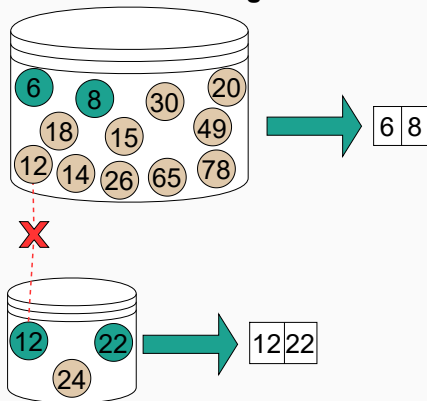


Scaled MinHash in Sourmash

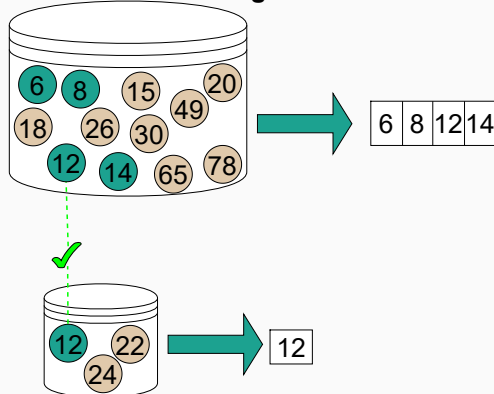


FIXED VS SCALED SIZE SKETCHING

Fixed size sketching



Scaled size sketching



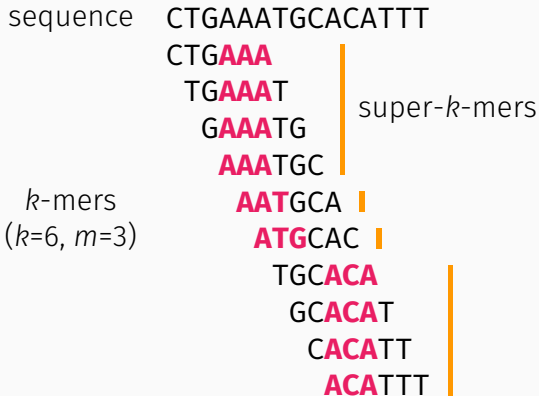
THEORETICAL BACKGROUND

MINIMIZERS & SUPER-K-MERS

Minimizer

smallest m -mer of a k -mer according to some order (e.g. lexicographical)

width parameter: $w = k - m + 1$



MINIMIZERS & SUPER-K-MERS

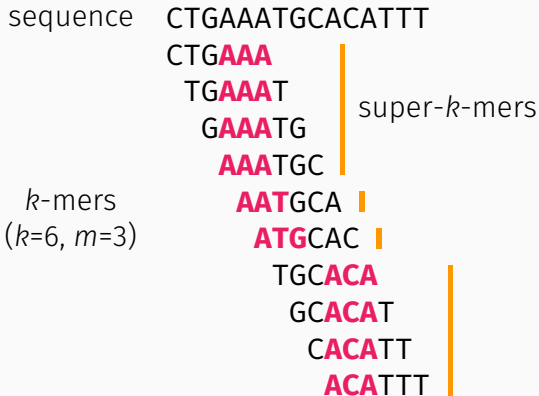
Minimizer

smallest m -mer of a k -mer according to some order (e.g. lexicographical)

width parameter: $w = k - m + 1$

Super- k -mer

run of consecutive k -mers sharing the same minimizer



We use minimizers as a footprint for selecting super- k -mers

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

$$d = \frac{\text{\#selected minimizers}}{\text{\#}m\text{-mers}}$$

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

$$d = \frac{\text{\#selected minimizers}}{\text{\#m-mers}}$$

sequence TGTGTGCTATT

T**GTT**GT

GTT**GTG**

TT**GTGC**

TGT**GCT**

GTG**CTA**

TG**CTAT**

GCT**ATT**

k-mers
(*k*=6, *m*=3)

selected
minimizers

. * . . * . * * . *

high density
(lexicographical order)

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

$$d = \frac{\text{\#selected minimizers}}{\text{\#m-mers}}$$

sequence TGT**T**TGTGCTATT
T**GTT**GT
GTT**G**TG
TT**G**TGC
TGT**G**CT
GTG**C**TA
TG**C**TAT
GCT**A**TT

k-mers
(*k*=6, *m*=3)

selected
minimizers . * . . * . * * . *

high density
(lexicographical order)

sequence TGT**T**TGTGCTATT
TGT**TGT**
GT**TGT**G
TT**TGT**GC
TGTGCT
GTG**C**TA
TG**C**TAT
GCT**A**TT

k-mers
(*k*=6, *m*=3)

selected
minimizers . . . * . . . * . .

low density
(TGT < CTA < ...)

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

$$d = \frac{\# \text{selected minimizers}}{\# m\text{-mers}}$$

sequence TGT TGT GCT ATT
T **GTT** GT
GTT **GTG**
TT **GTG** C
TGT **GCT**
GTG **CTA**
TG **CTA** T
GCT **ATT**

k-mers
(*k*=6, *m*=3)

selected
minimizers . * . . * . * . *

high density
(lexicographical order)

sequence TGT TGT GCT ATT
TGT **TGT**
GTT **GTG**
TT **GTG** C
TGT GCT
GTG **CTA**
TG **CTA** T
GCT **ATT**

k-mers
(*k*=6, *m*=3)

selected
minimizers . . . * . . . * . .

low density
(TGT < CTA < ...)

low density \iff long super-*k*-mers

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

$$d = \frac{\text{\#selected minimizers}}{\text{\#m-mers}}$$

Optimal density: $d = 1/w$

When using a random order,
the **expected** density is $\frac{2}{w+1}$

sequence TGTTGTGCTATT
 T**GTT**GT
 GTT**G**TG
 TT**G**TGC
 TGT**G**CT
 GTG**C**TA
 TG**C**TAT
 GCT**A**TT

k-mers
($k=6, m=3$)

selected
minimizers . * . . * . * . *

high density
(lexicographical order)

sequence TGTTGTGCTATT
 TGT**TGT**
 GTT**G**TG
 TT**G**TGC
 TGTGCT
 GTG**C**TA
 TG**C**TAT
 GCT**A**TT

k-mers
($k=6, m=3$)

selected
minimizers . . . * . . . * . .

low density
($\text{TGT} < \text{CTA} < \dots$)

low density \iff long super- k -mers

UNIVERSAL HITTING SETS & DENSITY LOWER BOUND

Universal Hitting Set (UHS)

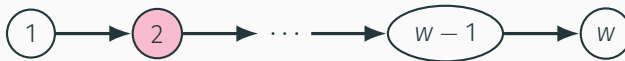
set S of m -mers s.t. every run of w consecutive m -mers has ≥ 1 element in S



UNIVERSAL HITTING SETS & DENSITY LOWER BOUND

Universal Hitting Set (UHS)

set S of m -mers s.t. every run of w consecutive m -mers has ≥ 1 element in S



Density lower bound

In any UHS, the density is $\geq \frac{1.5}{w+1}$ (i.e. the density factor is ≥ 1.5)

UNIVERSAL HITTING SETS & DENSITY LOWER BOUND

Universal Hitting Set (UHS)

set S of m -mers s.t. every run of w consecutive m -mers has ≥ 1 element in S



Density lower bound

In any UHS, the density is $\geq \frac{1.5}{w+1}$ (i.e. the density factor is ≥ 1.5)

Can we cross this lower bound by relaxing some constraints?

FRACTIONAL HITTING SETS

Instead of covering every k -mer, we cover a fraction f of them

Fractional Hitting Set (FHS)

set S of m -mers s.t. every run of w consecutive m -mers
has ≥ 1 element in S with probability $\geq f$

Instead of covering every k -mer, we cover a fraction f of them

Fractional Hitting Set (FHS)

set S of m -mers s.t. every run of w consecutive m -mers
has ≥ 1 element in S with probability $\geq f$

In practice, we select minimizers smaller than a certain threshold t

$$t = \left[1 - (1 - f)^{1/w} \right] \cdot 4^m$$

minimizers $\leq t$ are called **small minimizers**

Density upper bound

Given a covering fraction f , assuming $m > (3 + \varepsilon) \log_4 w$,

$$d \leq \frac{2f}{w+1} + o(1/w)$$

Density upper bound

Given a covering fraction f , assuming $m > (3 + \varepsilon) \log_4 w$,

$$d \leq \frac{2f}{w+1} + o(1/w)$$

⊕ simple, consistent with known results for $f = 1$

Density upper bound

Given a covering fraction f , assuming $m > (3 + \varepsilon) \log_4 w$,

$$d \leq \frac{2f}{w+1} + o(1/w)$$

- ⊕ simple, consistent with known results for $f = 1$
- ⊖ not very meaningful as $f \rightarrow 0$
(since most k -mers are not covered)

Density upper bound

Given a covering fraction f , assuming $m > (3 + \varepsilon) \log_4 w$,

$$d \leq \frac{2f}{w+1} + o(1/w)$$

- ⊕ simple, consistent with known results for $f = 1$
- ⊖ not very meaningful as $f \rightarrow 0$
(since most k -mers are not covered)

Is there a more meaningful metric?

Restricted density upper bound

Given a covering fraction f , assuming $m > (3 + \varepsilon) \log_4 w$,
when restricting to k -mers containing small minimizers,

$$d \leq 2 \cdot \frac{f + (1 - f) \ln(1 - f)}{f^2(w + 1)} + o(1/w)$$

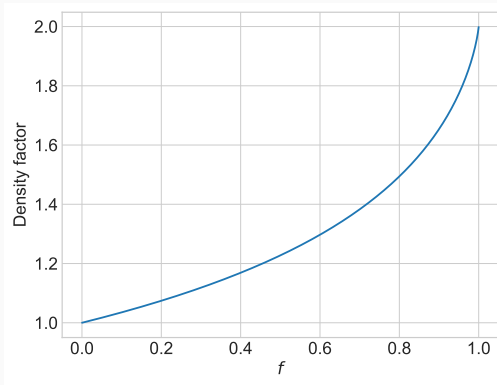
RESTRICTED DENSITY UPPER BOUND FOR SMALL MINIMIZERS

Restricted density upper bound

Given a covering fraction f , assuming $m > (3 + \epsilon) \log_4 w$,
when restricting to k -mers containing small minimizers,

$$d \leq 2 \cdot \frac{f + (1-f) \ln(1-f)}{f^2(w+1)} + o(1/w)$$

- below the $\frac{1.5}{w+1}$ barrier for $f \leq 0.8$
- approaches optimal density as $f \rightarrow 0$



Proportion of maximal super- k -mers

The average proportion of maximal super- k -mers is

$$\left[\left(1 - \frac{1}{w} \right) \frac{f}{1+f} \right]^2 + \frac{1 - f(1 - 2/w)}{1+f}$$

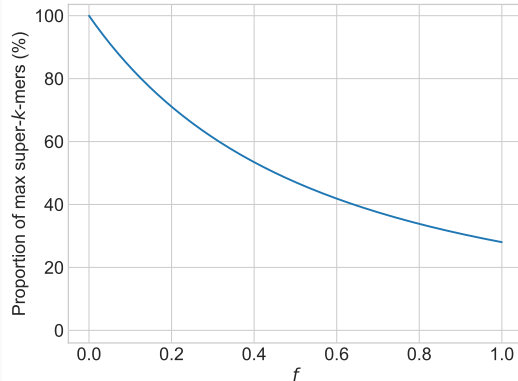
PROPORTION OF MAXIMAL SUPER-K-MERS

Proportion of maximal super- k -mers

The average proportion of maximal super- k -mers is

$$\left[\left(1 - \frac{1}{w} \right) \frac{f}{1+f} \right]^2 + \frac{1 - f(1 - 2/w)}{1+f}$$

(for $w = 17$)



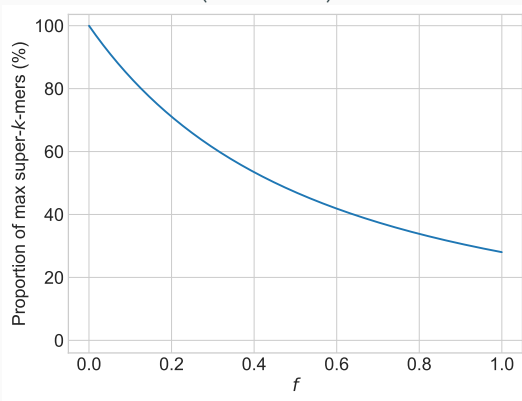
PROPORTION OF MAXIMAL SUPER-K-MERS

Proportion of maximal super- k -mers

The average proportion of maximal super- k -mers is

$$\left[\left(1 - \frac{1}{w} \right) \frac{f}{1+f} \right]^2 + \frac{1 - f(1 - 2/w)}{1+f}$$

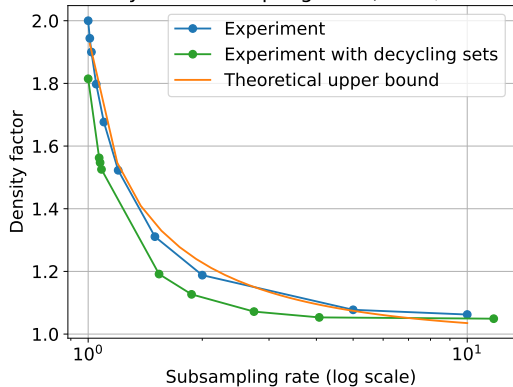
(for $w = 17$)



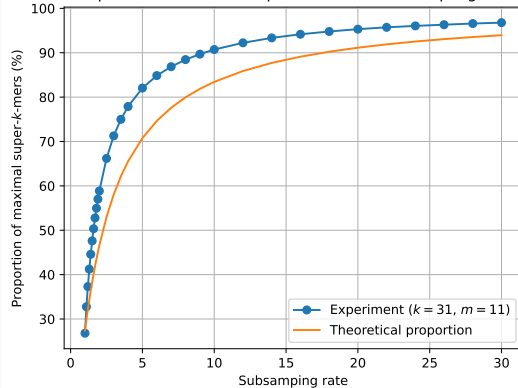
How accurate is it in practice?

COMPARISON WITH EXPERIMENTAL RESULTS

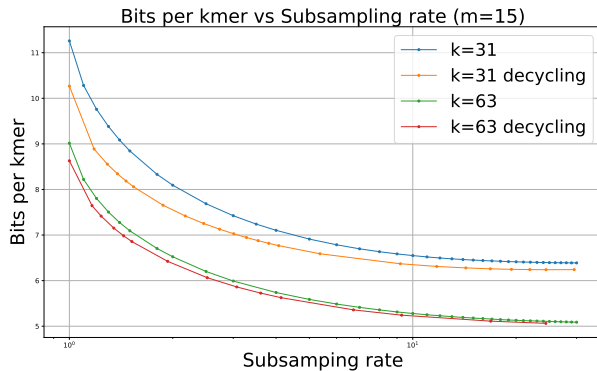
Density vs Subsampling rate ($k=31, m=15$)



Proportion of maximal super-k-mers vs Subsampling rate

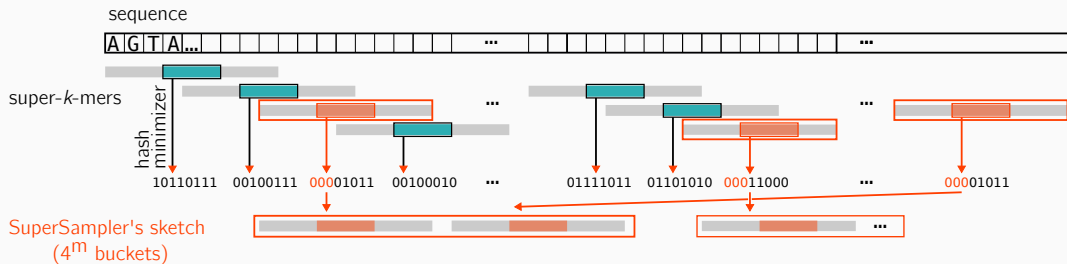


COMPARISON WITH EXPERIMENTAL RESULTS



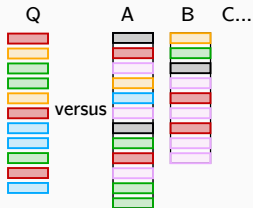
METHODS

SUPERAMPLER'S SKETCHES



SKETCH COMPARISON

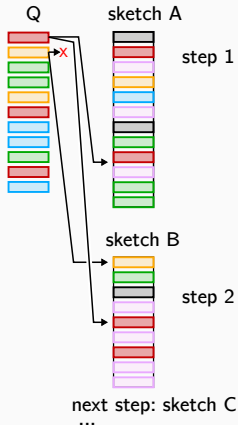
Sketches to compare



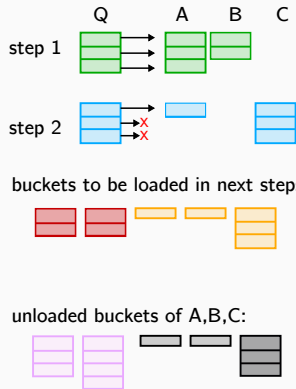
legend

- fingerprints
- find a match
- x no match found

Sourmash's Q vs all (A,B,C,...) comparison



SPSP's Q vs all (A,B,C,...) comparison



RESULTS

PERFORMANCE COMPARISON ON DISSIMILAR DATA (REFSEQ)

Computational time

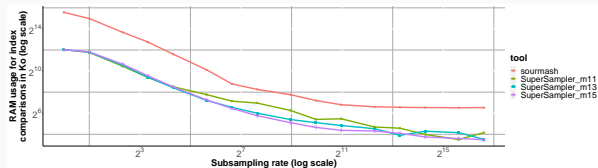
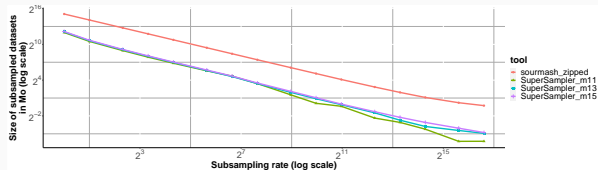
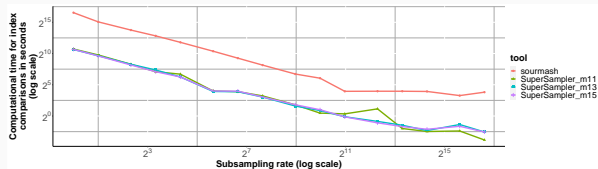
50× faster

Disk usage

20× lighter

RAM usage

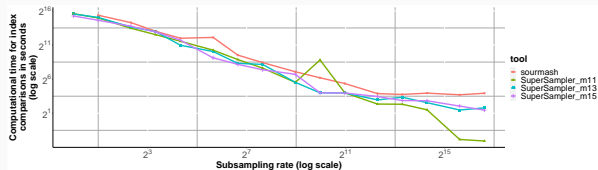
6× less RAM



PERFORMANCE COMPARISON ON SIMILAR DATA (SALMONELLAS)

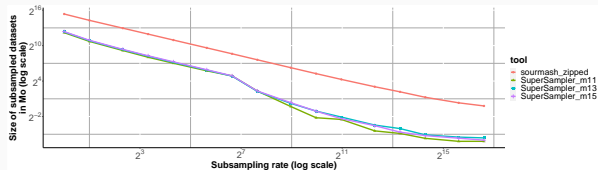
Computational time

7× faster



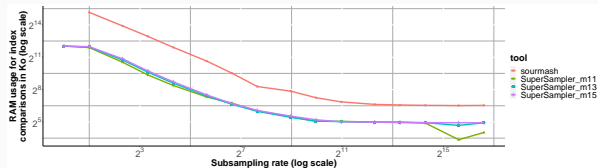
Disk usage

50× lighter



RAM usage

6× less RAM



CONCLUSION





TAKE HOME MESSAGES

- Super- k -mers
 - robust fingerprints
 - low memory cost
- Fractional Hitting Sets
 - generalization of UHS
 - lower density
 - longer super- k -mers
 - can be combined w/ existing UHS
- Go check out our posters!
 - **Igor Martayan's poster (#123)** on the theory behind FHS
 - **My poster about SuperSampler (#147)**, if you want to have a chat with me !
- If you want to dig deeper into FHS and SuperSampler, go checkout our preprint:



APPENDIX

REFERENCES I

-  Brian D Ondov, Todd J Treangen, Páll Melsted, Adam B Mallonee, Nicholas H Bergman, Sergey Koren, and Adam M Phillippy.
Mash: fast genome and metagenome distance estimation using minhash.
Genome biology, 17(1):1–14, 2016.
-  N Tessa Pierce, Luiz Irber, Taylor Reiter, Phillip Brooks, and C Titus Brown.
Large-scale sequence comparisons with sourmash.
F1000Research, 8, 2019.
-  Mahmudur Rahman Hera, N Tessa Pierce-Ward, and David Koslicki.
Debiasing fracminhash and deriving confidence intervals for mutation rates across a wide range of evolutionary distances.
bioRxiv, 2022.
-  Kenneth Katz, Oleg Shutov, Richard Lapoint, Michael Kimelman, J Rodney Brister, and Christopher O’Sullivan.
The sequence read archive: a decade more of explosive growth.
Nucleic acids research, 50(D1):D387–D390, 2022.

REFERENCES II



Andrei Z Broder.

On the resemblance and containment of documents.

In *Proceedings. Compression and Complexity of SEQUENCES 1997* (Cat. No. 97TB100171), pages 21–29. IEEE, 1997.



Hongyu Zheng, Carl Kingsford, and Guillaume Marçais.

Improved design and analysis of practical minimizers.

Bioinformatics, 36(Supplement_1):i119–i127, 2020.



David Pellow, Lianrong Pu, Baris Ekim, Lior Kotlar, Bonnie Berger, Ron Shamir, and Yaron Orenstein.

Efficient minimizer orders for large values of k using minimum decycling sets.

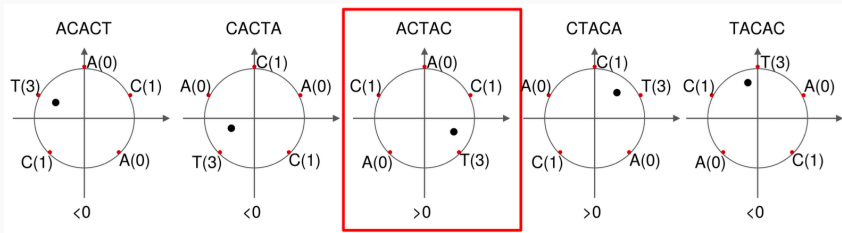
bioRxiv, pages 2022–10, 2022.

DECYCLING SETS

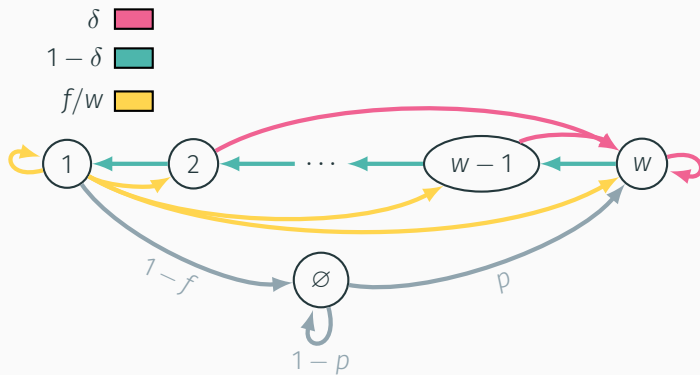
Decycling set

set S of m -mers whose removal make the De Bruijn graph **acyclic**

- if at least one m -mer is in S , take it in your UHS
- otherwise, use a random order to select a minimizer



SUPER-K-MERS' MARKOV CHAIN



- state i : small minimizer starts at i in the k -mer
- state \emptyset : no small minimizer in the k -mer