

FRACTIONAL HITTING SETS

EFFICIENT AND LIGHTWEIGHT GENOMIC DATA SKETCHING

Timothé ROUZÉ, Igor MARTAYAN, Camille MARCHET & Antoine LIMASSET

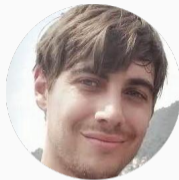
September 5, 2023



BEFORE WE START



Timothé Rouzé



Antoine LIMASSET



Camille MARCHET



slides



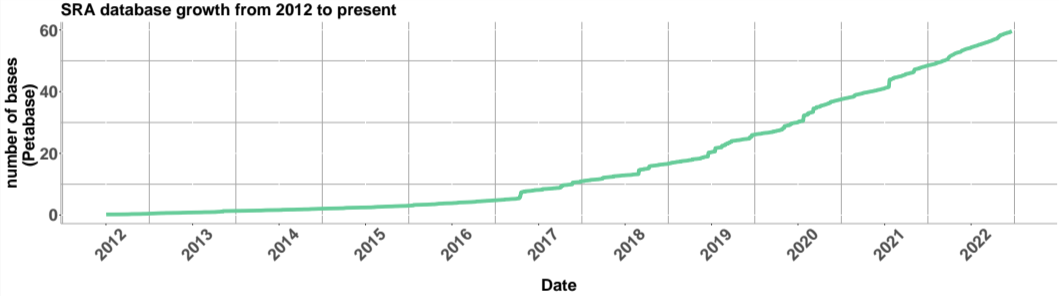
@IgorMartayan



@imartayan@genomic.social

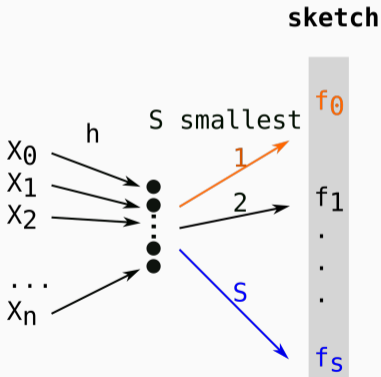
- Reminders on sketching & minimizers
- Fractional Hitting Sets
- SuperSampler, a sketching tool based on super- k -mers
- Experimental results
- Take home messages

BIONFORMATICIAN'S MOORE'S LAW

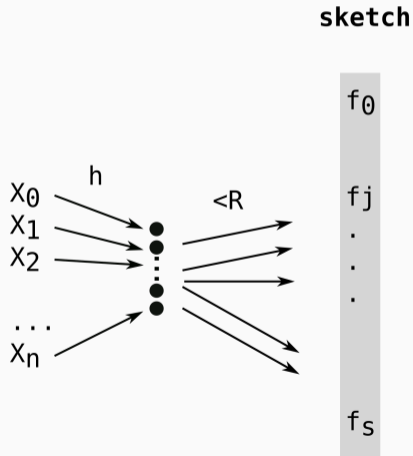


SKETCHING WITH MINHASH / FRACMINHASH

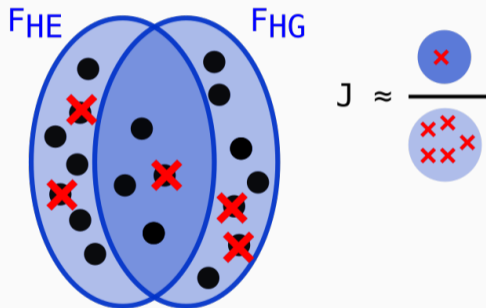
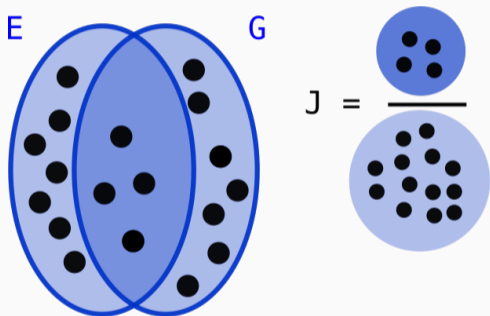
Bottom Minhash in MASH



Scaled MinHash in Sourmash

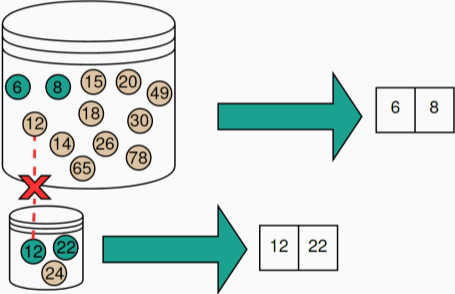


APPROXIMATING JACCARD INDEX

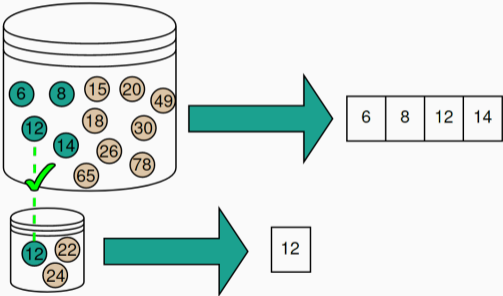


FIXED-SIZE VS SCALED-SIZE SKETCHING

Fixed size sketch



Scaled size sketch

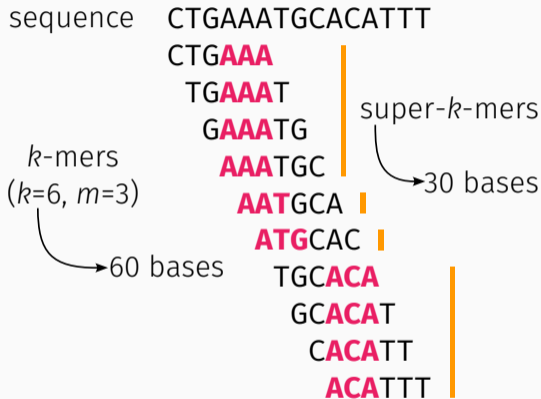


MINIMIZERS & SUPER-K-MERS

Minimizer

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

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



MINIMIZERS & SUPER-K-MERS

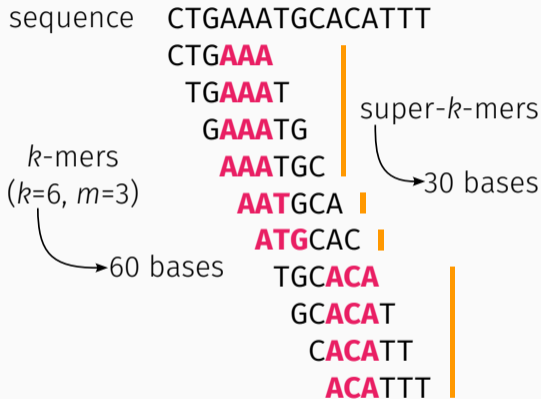
Minimizer

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

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

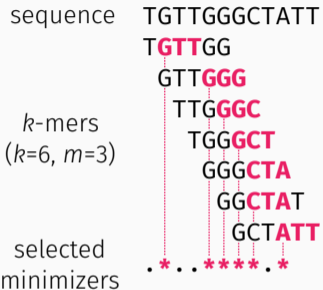
$$d = \frac{\text{\#selected minimizers}}{\text{\#m-mers}}$$

DENSITY OF MINIMIZERS

We want a **sparse** minimizer set

Density

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



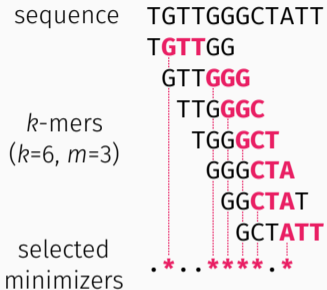
high density
(lexicographic order)

DENSITY OF MINIMIZERS

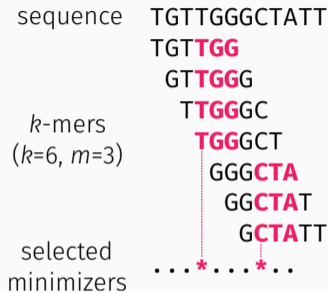
We want a **sparse** minimizer set

Density

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



high density
(lexicographic order)



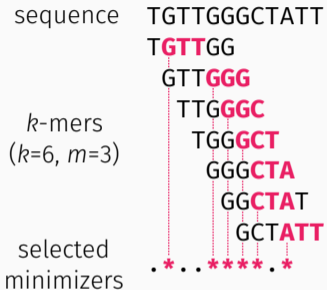
low density
(TGG < CTA < ...)

DENSITY OF MINIMIZERS

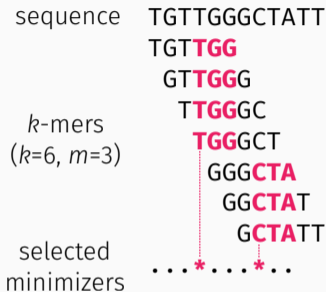
We want a **sparse** minimizer set

Density

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



high density
(lexicographic order)



low density
(TGG < CTA < ...)

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

DENSITY OF MINIMIZERS

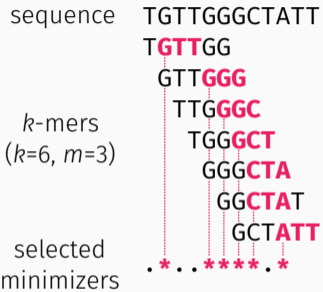
We want a **sparse** minimizer set

Density

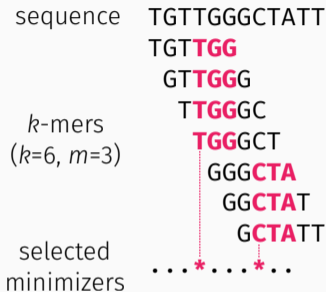
$$d = \frac{\# \text{selected minimizers}}{\# m\text{-mers}}$$

Optimal density: $d = 1/w$

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



high density
(lexicographic order)



low density
($TGG < CTA < \dots$)

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



e.g. Decycling sets (Pellow & al., 2022), Miniception (Zheng & al., 2020)

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



e.g. Decycling sets (Pellow & al., 2022), Miniception (Zheng & al., 2020)

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



e.g. Decycling sets (Pellow & al., 2022), Miniception (Zheng & al., 2020)

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,

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

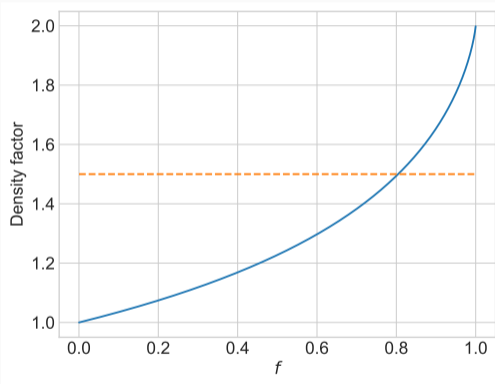
RESTRICTED DENSITY UPPER BOUND FOR SMALL MINIMIZERS

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

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

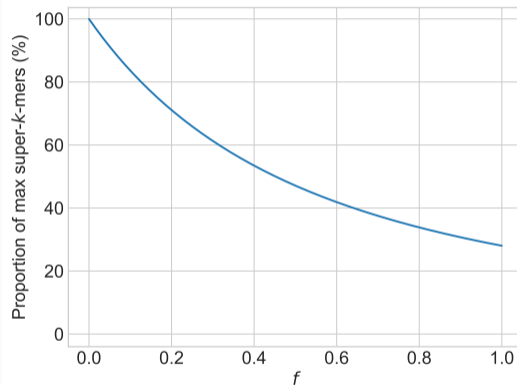
PROPORTION OF MAXIMAL SUPER-K-MERS

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 $w = 17$)



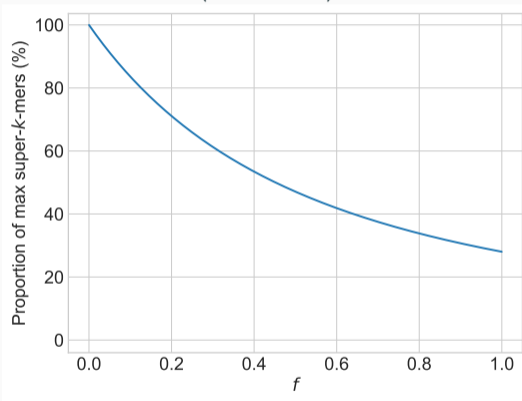
PROPORTION OF MAXIMAL SUPER-K-MERS

Proportion of maximal super- k -mers

The average proportion of maximal super- k -mers is

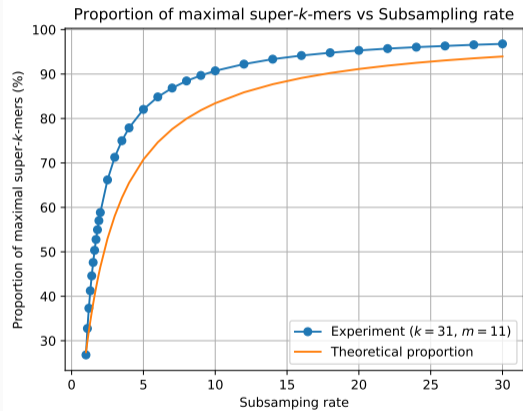
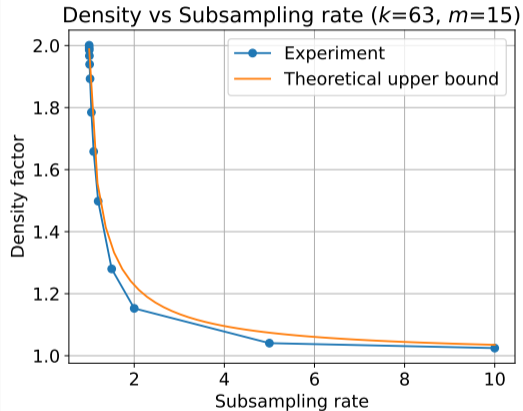
$$\left[\left(1 - \frac{1}{w} \right) \frac{f}{1+f} \right]^2 + \frac{1 - f(1 - 2/w)}{1+f}$$

(for $w = 17$)



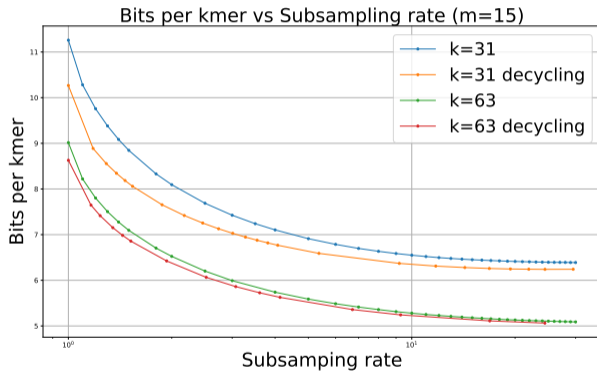
How accurate is it in practice?

COMPARISON WITH EXPERIMENTAL RESULTS



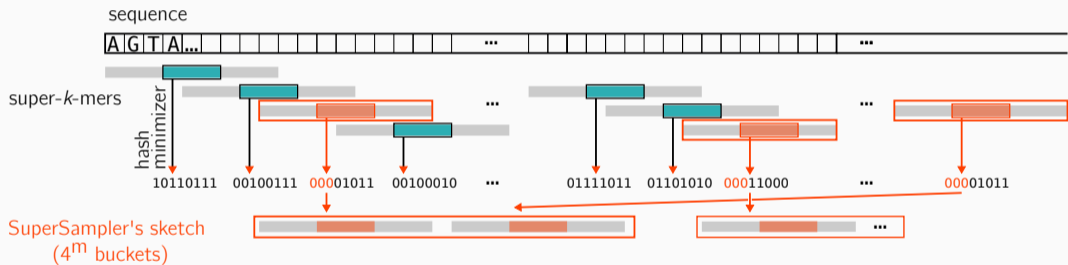
Instead of applying a threshold on minimizers, we can:

1. build a universal hitting set S (e.g. a decycling set)
2. sample elements from S (by hashing elements and applying a threshold)



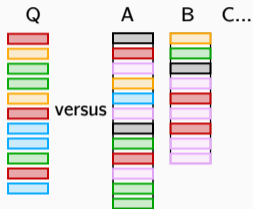
SUPERSAMPLER

SUPERAMPLER'S SKETCHES



SKETCH COMPARISON

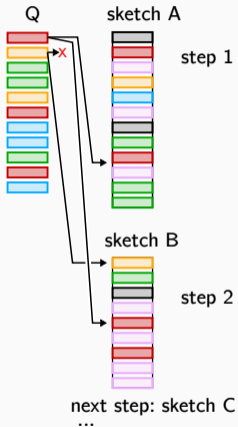
Sketches to compare



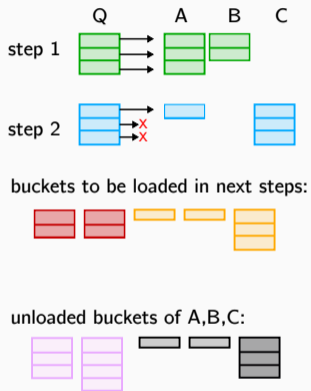
legend

- fingerprint
- 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



EXPERIMENTAL RESULTS VS SOURMASH

PERFORMANCE COMPARISON ON DISSIMILAR DATA (REFSEQ)

Computational time

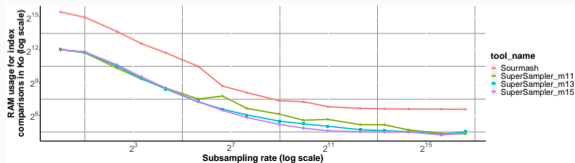
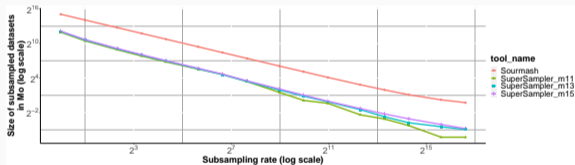
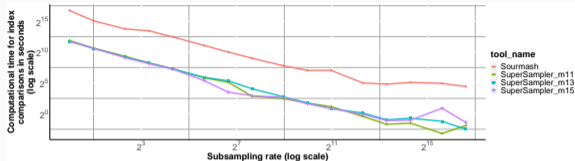
≈ 40× faster

Disk usage

≈ 15× lighter

RAM usage

≈ 5× less RAM



PERFORMANCE COMPARISON ON SIMILAR DATA (SALMONELLAS)

Computational time

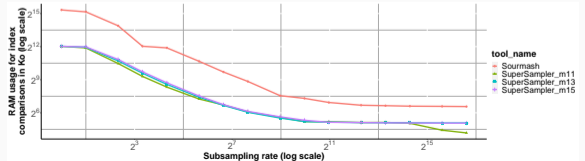
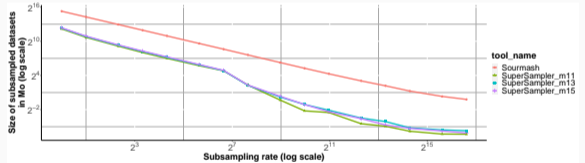
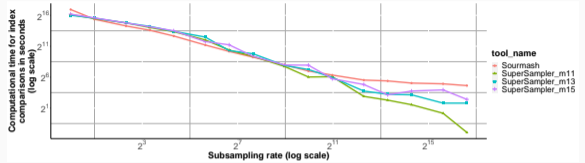
similar w/ low sampling,
 $\approx 2\times$ faster w/ high sampling

Disk usage

$\approx 40\times$ lighter

RAM usage

$\approx 5\times$ less RAM



CONCLUSION

Fractional Hitting Sets:

- unify UHS and sketching problems
- lead to lower density / longer super- k -mers
- can benefit from existing UHS building techniques

Super- k -mers:

- provide a space-efficient representation
- speed-up genome sketching & comparison

paper



SuperSampler



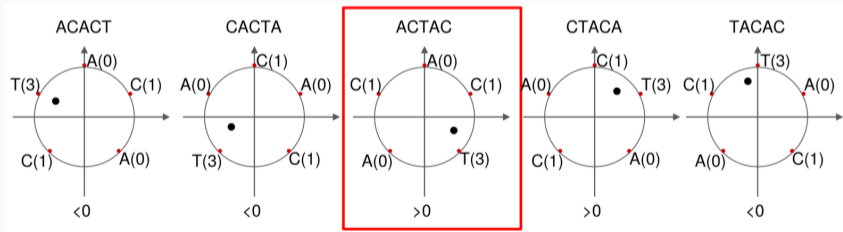
APPENDIX

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



Pellow & al., 2022

DENSITY UPPER BOUND: SKETCH OF THE PROOF

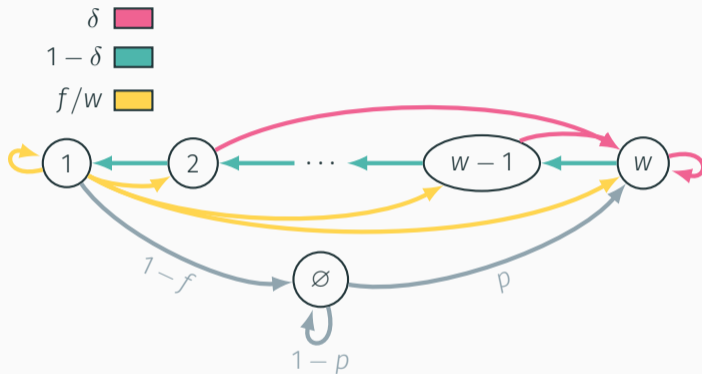
Key property (from Zheng & al., 2020)

Assuming $m > (3 + \varepsilon) \log_{\sigma} w$, the probability of having duplicate m -mers in a k -mer is negligible

We consider two consecutive k -mers,
the density is equal to the probability that they have different minimizers,
which is the expectation of $\frac{\#\text{small boundary } m\text{-mers}}{\#\text{small } m\text{-mers}}$.

The \ln factor in the restricted density bound comes from a Taylor expansion.

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