

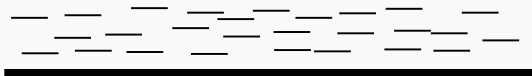


Generalization of the minimizers schemes

Guillaume Marçais, Dan DeBlasio, Carl Kingsford

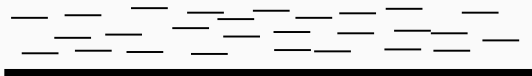
Carnegie Mellon University

Computing read overlaps

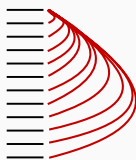


Roberts, *et al.*
(2004). Reducing
storage
requirements for
biological
sequence
comparison.

Computing read overlaps



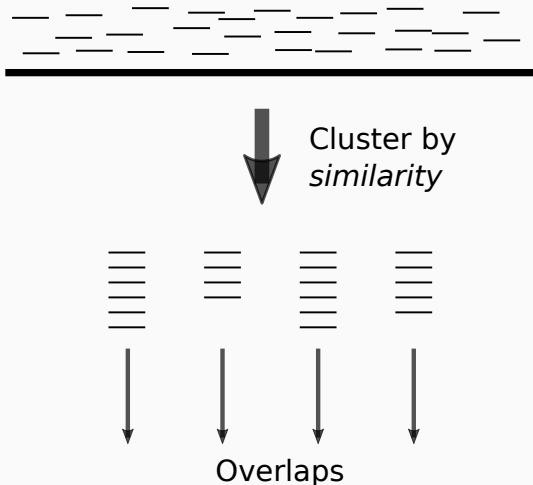
Roberts, *et al.*
(2004). Reducing
storage
requirements for
biological
sequence
comparison.



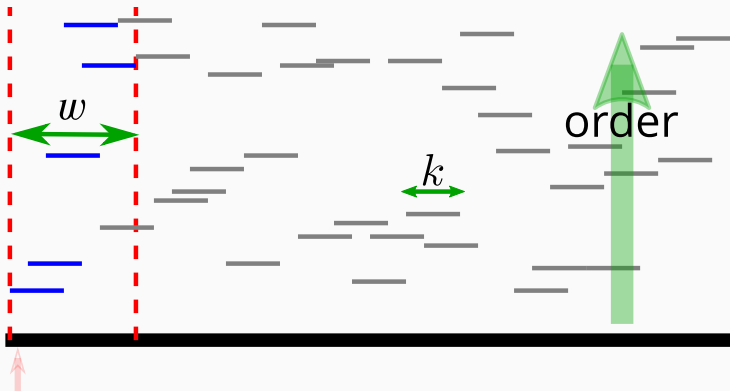
$O(n^2)$ alignments

Computing read overlaps

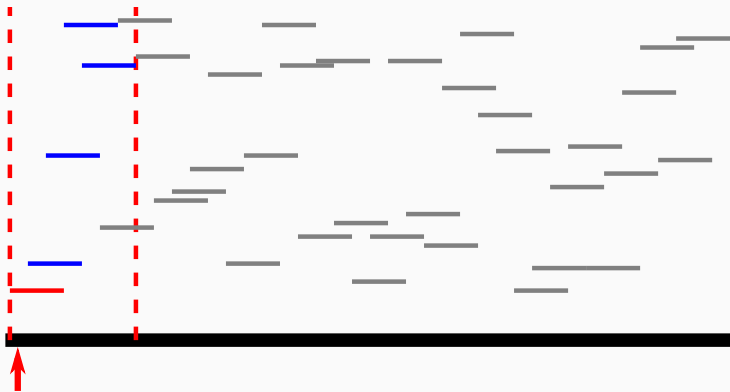
Roberts, *et al.*
(2004). Reducing
storage
requirements for
biological
sequence
comparison.



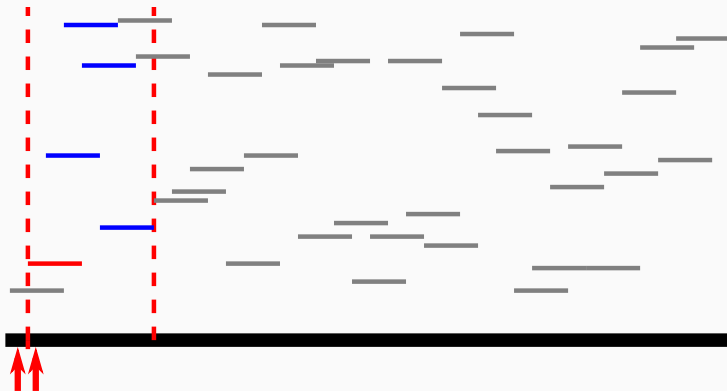
Computing minimizers



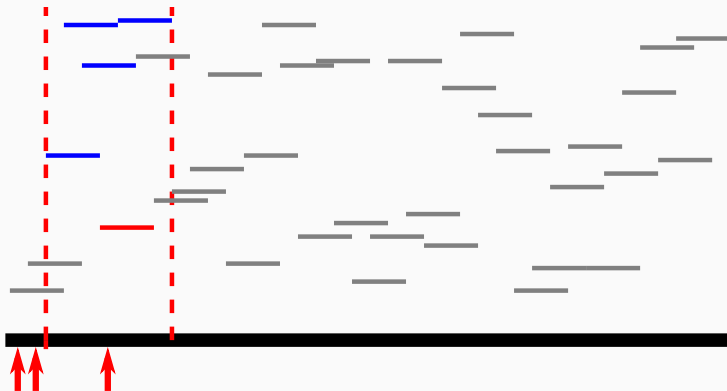
Computing minimizers



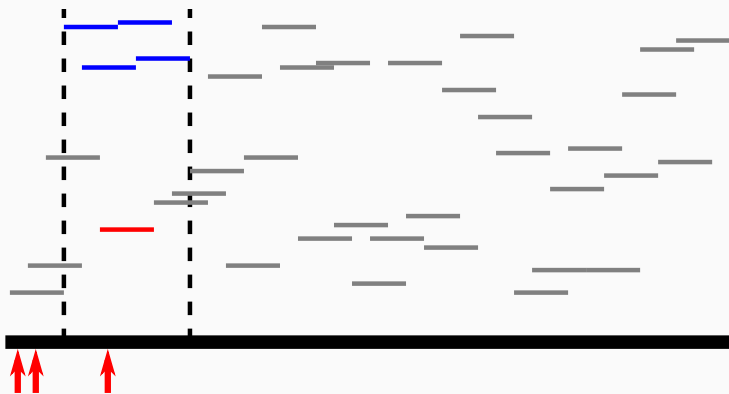
Computing minimizers



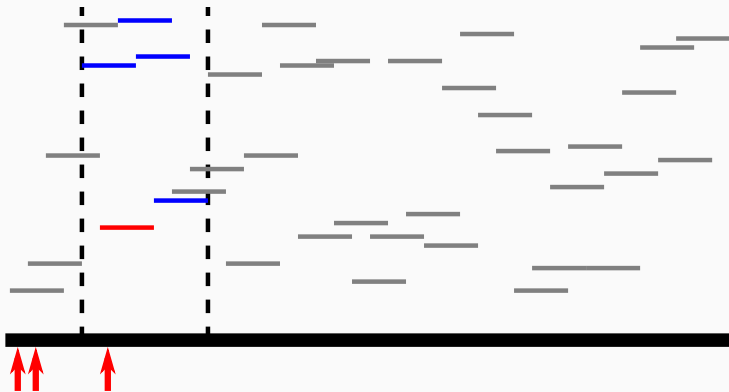
Computing minimizers



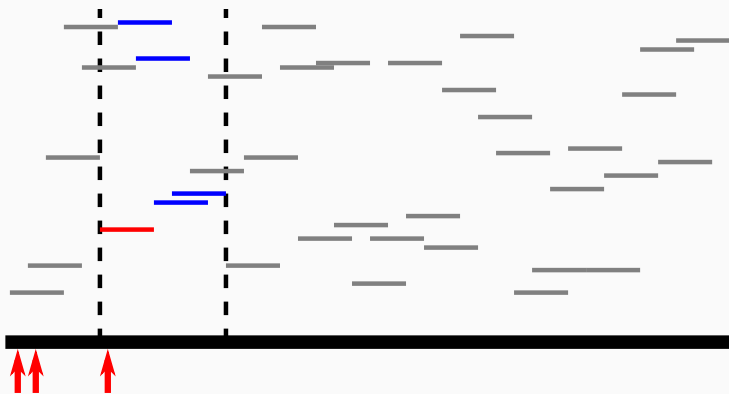
Computing minimizers



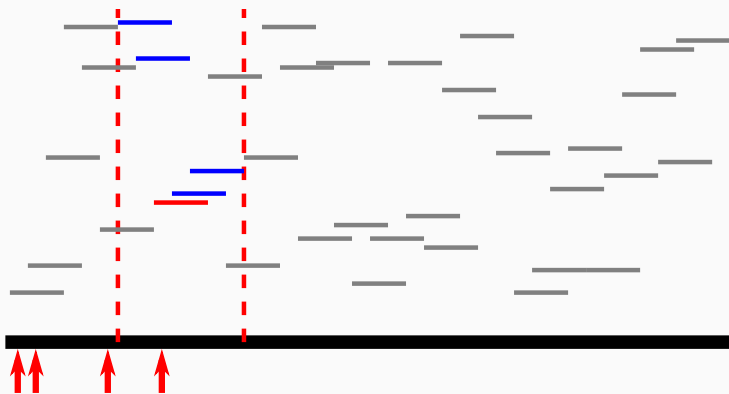
Computing minimizers



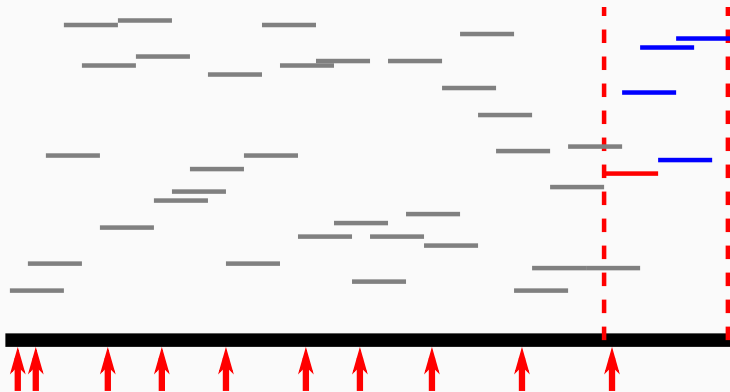
Computing minimizers



Computing minimizers



Computing minimizers



Minimizers definition and properties

Minimizers (k, w, o)

In each window of w consecutive k -mers, select the smallest k -mer according to order o .

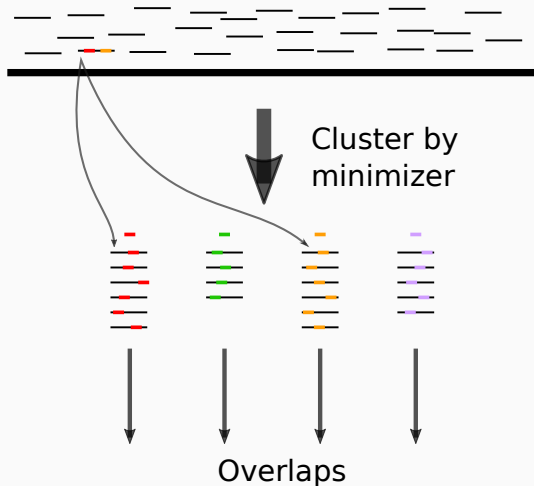
1. **Uniform:** distance between selected k -mers is $\leq w$
2. **Deterministic:** two strings matching on w consecutive k -mers select the same minimizer

Computing read overlaps

1. **Uniform:** no sequence ignored

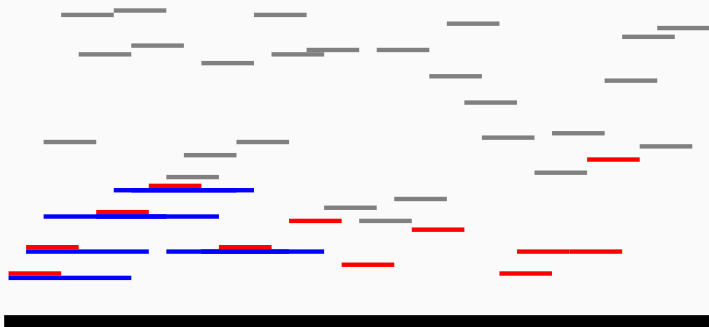
2.

Deterministic: reads with overlap in same bin



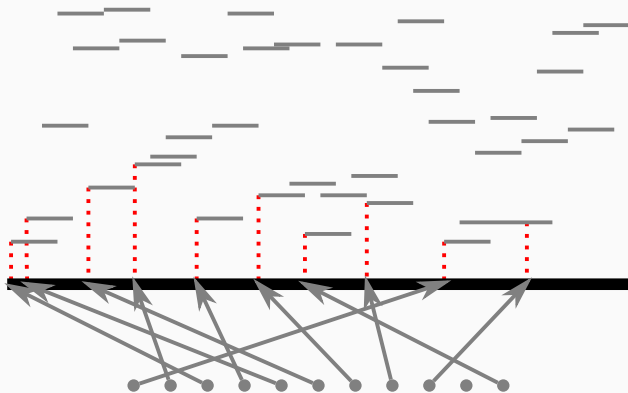
Counting k -mers: super-windows

- Minimizers: m -mer ($m = 7$)
- k -mers: window with $k = w + m - 1$



Count k -mers in parallel in buckets of blue super-windows

Sparse suffix-array

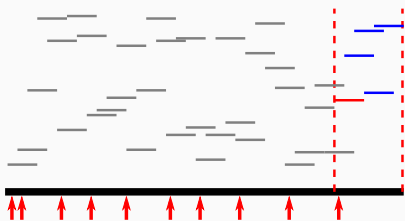


Check only for minimizers in query sequence in SA

Many applications of minimizers

- **UMDOverlapper (Roberts, 2004)**: bin sequencing reads by shared minimizers to compute overlaps
- **MSPKmerCounter (Li, 2015), KMC2 (Deorowicz, 2015), Gerbil (Erber, 2017)**: bin input sequences based on minimizer to count k -mers in parallel
- **SparseAssembler (Ye, 2012), MSP (Li, 2013), DBGFM (Chikhi, 2014)**: reduce memory footprint of de Bruijn assembly graph with minimizers
- **SamSAMi (Grabowski, 2015)**: sparse suffix array with minimizers
- **MiniMap (Li, 2016), MashMap (Jain, 2017)**: sparse data structure for sequence alignment
- **Kraken (Wood, 2014)**: taxonomic sequence classifier
- **Schleimer *et al.* (2003)**: winnowing, text fingerprinting

Improving minimizers by lowering density

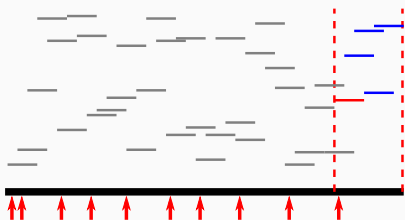


Density

Density of a scheme is the expected proportion of selected k -mer in a random sequence:

$$d = \frac{\# \text{ of selected } k\text{-mers}}{\text{length of sequence}}$$

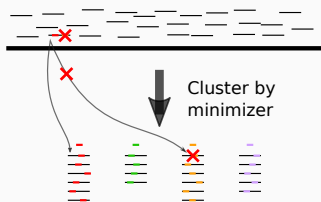
Improving minimizers by lowering density



Density

Density of a scheme is the expected proportion of selected k -mer in a random sequence:

$$d = \frac{\# \text{ of selected } k\text{-mers}}{\text{length of sequence}}$$



Lower density

⇒ smaller bins

⇒ less computation

Minimizers density minimizing problem

For fixed k and w :

- Properties “Uniform” & “Deterministic” unaffected by order
- Density changes with ordering o
- Lower density \implies sparser data structures and/or less computation
- Benefit existing and new applications

Density minimization problem

For fixed w, k , find k -mer **order** o giving the lowest expected **density**

Minimizers density minimizing problem

For fixed k and w :

- Properties “Uniform” & “Deterministic” unaffected by order
- Density changes with ordering o
- Lower density \implies sparser data structures and/or less computation
- Benefit existing and new applications

Density minimization problem

For fixed w, k , find k -mer **order** o giving the lowest expected **density**

Expected and bound on density

For an *idealized random*
order σ :

$$d = \frac{2}{w+1} \quad df = 2$$

Expect ≈ 2 minimizers per
window

For any order σ :

$$df \geq 1.5 + \frac{1}{2w}$$

Requires ≥ 1.5 minimizers
per window

Expected and bound on density

For an *idealized random*
order o :

$$d = \frac{2}{w+1} \quad df = 2$$

Expect ≈ 2 minimizers per
window

For any order o :

$$df \geq 1.5 + \frac{1}{2w}$$

Requires ≥ 1.5 minimizers
per window

Is density factor constant?

For an *idealized random*
order o :

$$d = \frac{2}{w+1} \quad df = 2$$

Expect ≈ 2 minimizers per
window

For any order o :

$$df \geq 1.5 + \frac{1}{2w}$$

Requires ≥ 1.5 minimizers
per window

Is density factor constant?

For an *idealized random*
order o :

$$d = \frac{2}{w+1} \quad df = 2$$

Expect ≈ 2 minimizers per
window

Not valid for $w \gg k$

For any order o :

$$df \geq 1.5 + \frac{1}{2w}$$

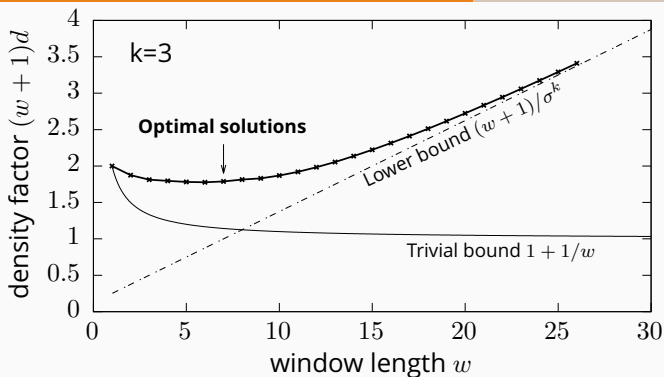
Requires ≥ 1.5 minimizers
per window

Valid only for $w \gg k$

What is the best ordering possible when:

- k is fixed and $w \rightarrow \infty$
- w is fixed and $k \rightarrow \infty$

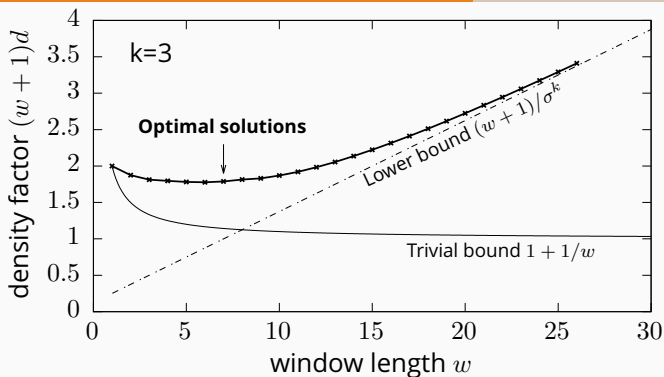
Asymptotic behavior in w



$$df \geq \frac{w+1}{\sigma^k}$$

Density factor is $\Omega(w)$, not constant

Asymptotic behavior in w



$$df \geq \frac{w+1}{\sigma^k}$$

Density factor is $\Omega(w)$, not constant

Asymptotic behavior in k

Asymptotically optimal minimizers schemes

There exists a sequence of orders $(o_k)_{k \in \mathbb{N}}$ which are asymptotically optimal:

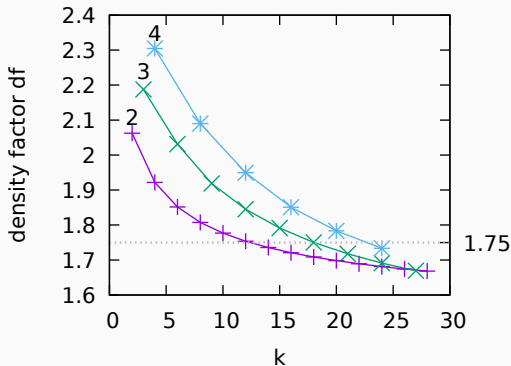
$$df_{o_k} \xrightarrow{k \rightarrow \infty} 1 + \frac{1}{W}$$

Asymptotic behavior in k

Asymptotically optimal minimizers schemes

There exists a sequence of orders $(o_k)_{k \in \mathbb{N}}$ which are asymptotically optimal:

$$df_{o_k} \xrightarrow{k \rightarrow \infty} 1 + \frac{1}{W}$$



Density factor of minimizers

Asymptotic behavior of minimizers is fully characterized:

- Minimizers scheme is optimal for large k : $df \xrightarrow[k \rightarrow \infty]{} 1 + \frac{1}{w}$
- Minimizers scheme is not optimal for large w : $df = \Omega(w)$
- Better lower bound on d

Density factor of minimizers

Asymptotic behavior of minimizers is fully characterized:

- Minimizers scheme is optimal for large k : $df \xrightarrow[k \rightarrow \infty]{} 1 + \frac{1}{w}$
- Minimizers scheme is not optimal for large w : $df = \Omega(w)$
- Better lower bound on d

Good:

- First example of optimal minimizers scheme
- Constructive proof

Not good:

- Large k less practical than large w
- Minimizers **don't** have **constant** density factor

Generalizing minimizers: local and forward schemes

Local scheme

Given $f : \Sigma^{w+k-1} \rightarrow [0, w-1]$, for each window ω , select k -mer at position $f(\omega)$.

Generalizing minimizers: local and forward schemes

Local scheme

Given $f : \Sigma^{w+k-1} \rightarrow [0, w-1]$, for each window ω , select k -mer at position $f(\omega)$.

Minimizers scheme with order o is a local scheme where

$$f = \arg \min_{i \in [0, w-1]} o(\omega[i : k])$$

Generalizing minimizers: local and forward schemes

Local scheme

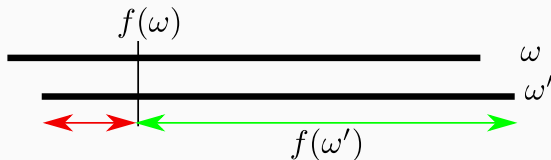
Given $f : \Sigma^{w+k-1} \rightarrow [0, w-1]$, for each window ω , select k -mer at position $f(\omega)$.

Minimizers scheme with order o is a local scheme where

$$f = \arg \min_{i \in [0, w-1]} o(\omega[i : k])$$

Forward scheme

Local scheme such that $f(\omega') \geq f(\omega) - 1$ if suffix of ω' equals prefix of ω



Minimizers \subsetneq Forward \subsetneq Local

- Properties “Uniform” & “Deterministic” also satisfied
- Drop-in replacement for minimizers
- Potential for lower density

Density factor overview

Density factor df		
	$k \rightarrow \infty$	$w \rightarrow \infty$
Scheme	Best	Bound
Minimizers		
Forward		
Local		

Density factor overview

Density factor df			
	$k \rightarrow \infty$	$w \rightarrow \infty$	
Scheme		Best	Bound
Minimizers	$1 + \frac{1}{w}$	$O(w)$	$\Omega(w)$
Forward			
Local			

Density factor overview

Density factor df			
	$k \rightarrow \infty$	$w \rightarrow \infty$	
Scheme		Best	Bound
Minimizers	$1 + \frac{1}{w}$	$O(w)$	$\Omega(w)$
Forward	$1 + \frac{1}{w}$	$O(\sqrt{w})$	$\sim 1.5 + \frac{1}{2w}$
Local			

Density factor overview

Density factor df			
	$k \rightarrow \infty$	$w \rightarrow \infty$	
Scheme		Best	Bound
Minimizers	$1 + \frac{1}{w}$	$O(w)$	$\Omega(w)$
Forward	$1 + \frac{1}{w}$	$O(\sqrt{w})$	$\sim 1.5 + \frac{1}{2w}$
Local	$1 + \frac{1}{w}$	$O(\sqrt{w})$	$1 + \frac{1}{w}$

Conclusion: the quest for optimal constant density factor

- Minimizers schemes:
 - **can't** achieve **constant** density factor
 - **can't** be optimal for large w
- Forward schemes
 - **may** achieve **constant** density factor
 - **can't** be optimal for large w
- Local schemes
 - **may** achieve **constant** density factor
 - **may** be optimal for large w

- Design of **optimal** orders or functions f still open



Carnegie
Mellon
University

GORDON AND BETTY
MOORE
FOUNDATION
GBMF4554



CCF-1256087
CCF-1319998



R01HG007104
R01GM122935