

# USING NECKLACES TO BUILD A LOCALITY-PRESERVING AND DYNAMIC INDEX FOR $K$ -MERS

---

Igor MARTAYAN, Bastien CAZAUX, Camille MARCHET, Antoine LIMASSET

December 14, 2023

Seminar on Lyndon words — Rouen



# DNA SEQUENCING & TOKENIZATION WITH $k$ -MERS

DNA samples  →



→ CTGAAATG...

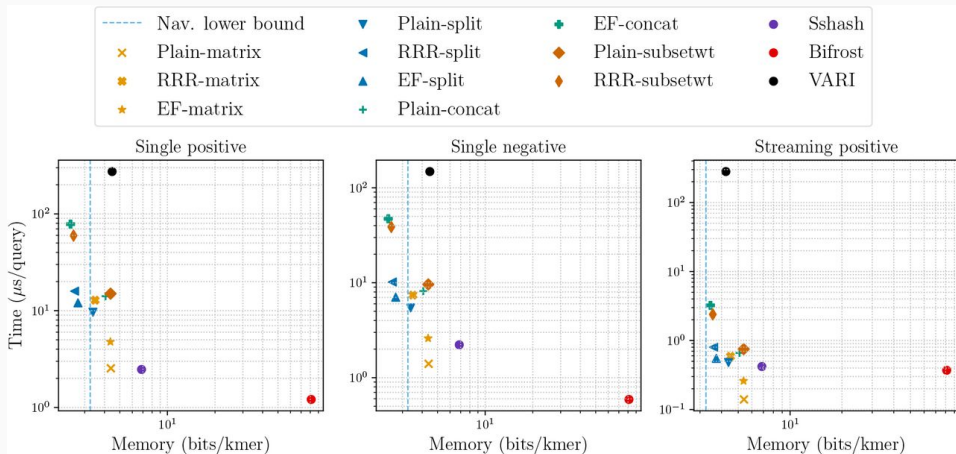
We typically **index the words of size  $k$**  ( $k$ -mers) instead of the sequence itself.

In practice, we usually consider  $k \leq 63$  so that each  $k$ -mer fits inside a machine word.

```
CTGAA
TGAAA
GAAAT
AAATG
```

# MOTIVATION OF THIS WORK

Plenty of compact data structures for storing  $k$ -mers ...but most of them are **static**



Query time and memory usage of some efficient data structures, taken from [Alanko et al. 22]

[Conway & Bromage 11]

- we can see  $k$ -mers as integers in  $\llbracket 4^k \rrbracket$   
 $A \rightarrow 00$     $C \rightarrow 01$     $G \rightarrow 10$     $T \rightarrow 11$
- since they're usually very sparse, we can use a sparse bitvector to store them

Limitations

- it's not really dynamic
  - it's not cache-efficient
    - $\text{index}(\text{ATAACGCCA}) = 49,556$
    - $\text{index}(\text{TAACGCCAT}) = 198,227$
- average distance of  $4^k/2$

[Conway & Bromage 11]

- we can see  $k$ -mers as integers in  $\llbracket 4^k \rrbracket$   
 $A \rightarrow 00$     $C \rightarrow 01$     $G \rightarrow 10$     $T \rightarrow 11$
- since they're usually very sparse, we can use a sparse bitvector to store them

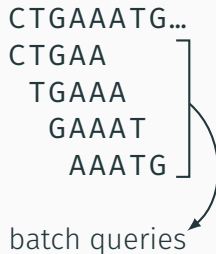
Limitations

- it's not really dynamic
  - it's not cache-efficient
    - $\text{index}(\text{ATAACGCCA}) = 49,556$
    - $\text{index}(\text{TAACGCCAT}) = 198,227$
- average distance of  $4^k/2$

*How can we improve this approach?*

# WISH LIST FOR AN IDEAL DATA STRUCTURE

- **space-efficient**: few bits /  $k$ -mer
- **dynamic**: support insertion and deletion after construction
- **efficient queries**:
  - membership
  - enumeration
  - insertion
  - deletion
- **locality-preserving**: reduce cache misses when querying consecutive  $k$ -mers



## PRESERVING LOCALITY WITH NECKLACES

---

## A LOCALITY-PRESERVING ENCODING OF K-MERS





## A LOCALITY-PRESERVING ENCODING OF K-MERS



Alternative encoding based on necklaces

The necklace of  $x$  is its smallest cyclic rotation  $\langle x \rangle = \min_{0 \leq i < k} x^{(i)}$

## A LOCALITY-PRESERVING ENCODING OF $k$ -MERS



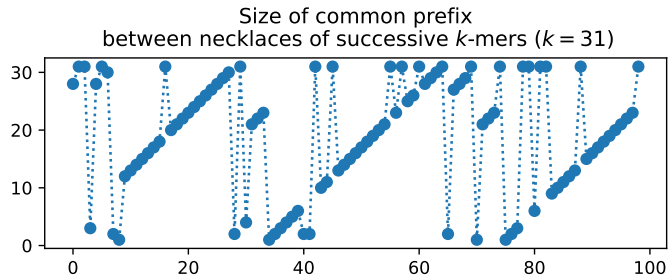
Alternative encoding based on necklaces

The necklace of  $x$  is its **smallest cyclic rotation**  $\langle x \rangle = \min_{0 \leq i < k} x^{(i)}$

- $x \mapsto (\langle x \rangle, \text{rotation index})$  is a **reversible** transformation
- necklaces of consecutive  $k$ -mers **share long prefixes**

## A CLOSER LOOK AT THE LOCALITY OF NECKLACES

AACGTCATCTCTCATTCTGGTCGTTCTTCCT  
AACGTCATCTCTCATTCTGTTTCGTTCTTCCT  
AACGTCATCTCTCATTCTGTGCGTTCTTCCT  
AACGTCATCTCTCATTCTGTGAGTTCTTCCT  
AACGTCATCTCTCATTCTGTGACTTCTTCCT  
AACGTCATCTCTCATTCTGTGACATCTTCCT  
AACGTCATCTCTCATTCTGTGACACCCTTCCT  
AACGTCATCTCTCATTCTGTGACACGTTCTTCCT  
AACGTCATCTCTCATTCTGTGACACGCTCCT  
AACGTCATCTCTCATTCTGTGACACGCACCT  
AACGTCATCTCTCATTCTGTGACACGCAGCT  
AACGTCATCTCTCATTCTGTGACACGCAGGT  
AACGTCATCTCTCATTCTGTGACACGCAGGG  
ACACGCAGGGTACGTCATCTCTCATTCTGTG



# PRACTICAL USE OF NECKLACES

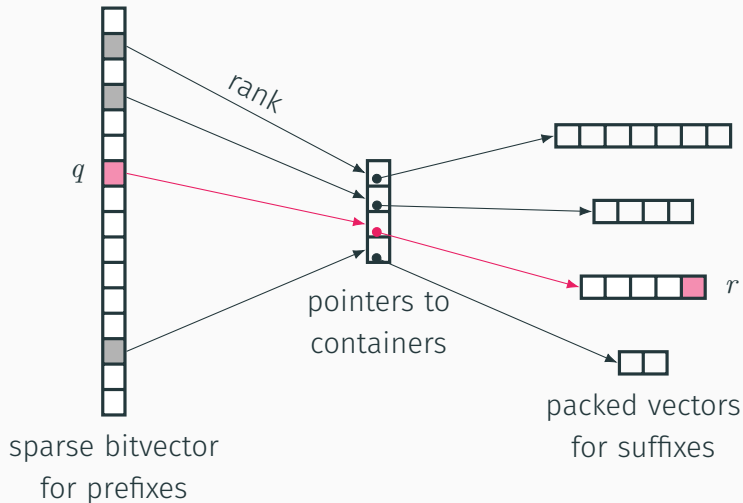
---

# OVERVIEW OF OUR DATA STRUCTURE (CBL)

Quotiented  
data structure

Query  $x$ :

1. compute  $\langle x \rangle$
2. split  $\langle x \rangle$  as  $q \parallel r$
3. look for  $(q, r)$



## ACCELERATING THE COMPUTATION OF CONSECUTIVE NECKLACES

Basic approach: compute every cyclic rotation and select the smallest in  $\mathcal{O}(k)$ .  
→  $\mathcal{O}(nk)$  for  $n$  queries

Better approach: amortize the computation cost for consecutive queries.

### Key observation

Given a fixed  $m$ , if  $\langle x \rangle$  does not start at one of the  $m - 1$  last positions of  $x$ , its prefix of size  $m$  is the smallest factor of size  $m$  in  $x$ .

Good news: we can keep track of the smallest factors of size  $m$  in  $\mathcal{O}(1)$  amortized time using a monotone queue.

# ACCELERATING THE COMPUTATION OF CONSECUTIVE NECKLACES

## Faster necklace computation

Only consider the cyclic rotations that start:

- at one of the smallest factors of size  $m$
- at one of the  $m - 1$  last positions

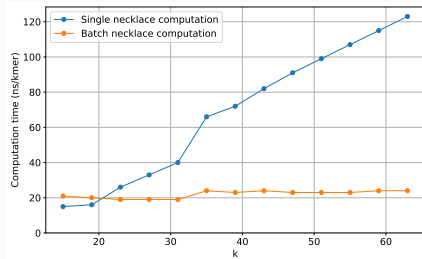
## Useful property [Zheng et al. 20]

Assuming  $m = \Omega(\log k)$ , the probability that a  $k$ -mer contains **duplicate  $m$ -mers** is  $o(1/k)$ .

By choosing  $m = \Theta(\log k)$ ,

the smallest factor of size  $m$  is unique w.h.p.

→  $\mathcal{O}(nm) = \mathcal{O}(n \log k)$  for  $n$  queries (on average)



## DENSIFYING THE SPACE OF NECKLACES

---



## DENSIFYING THE SPACE OF NECKLACES BY RANKING

The **number of necklaces** of size  $k$  on an alphabet with  $\sigma$  letters is

$$N(k) = \frac{1}{k} \sum_{d|k} \varphi\left(\frac{k}{d}\right) \sigma^d \sim \frac{\sigma^k}{k}$$

so only a fraction  $\frac{1}{k}$  **of the universe** is actually used

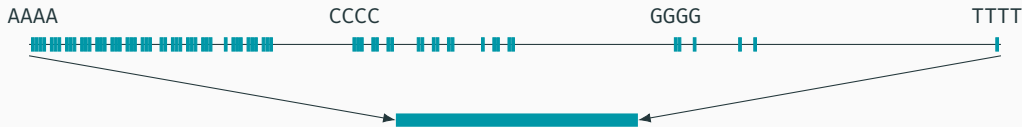


## DENSIFYING THE SPACE OF NECKLACES BY RANKING

The **number of necklaces** of size  $k$  on an alphabet with  $\sigma$  letters is

$$N(k) = \frac{1}{k} \sum_{d|k} \varphi\left(\frac{k}{d}\right) \sigma^d \sim \frac{\sigma^k}{k}$$

so only a fraction  $\frac{1}{k}$  of the universe is actually used



**Ranking:** given a necklace  $\langle x \rangle$ , find  $i$  s.t.  $\langle x \rangle$  is the  $i$ -th smallest necklace of size  $k$   
We can compute the rank in  $\mathcal{O}(k^2)$  time [Sawada & Williams 17]

**Tradeoff:** better locality + compression vs  $\mathcal{O}(k^2)$  queries

## CAN WE DO BETTER FOR CONSECUTIVE NECKLACES? (I DON'T KNOW YET)

Ranking in  $\mathcal{O}(k^2)$  is generally too expensive for our use case, but it might be faster to rank necklaces of consecutive  $k$ -mers.

Since most necklaces of consecutive words share the same starting position, they only **differ by a single letter**.

```
AACGTCATCTCTCATTCTGGTCGTTCTTCCT  
AACGTCATCTCTCATTCTGTTCGTTCTTCCT
```

**Formulation in the binary case ( $\sigma = 2$ )**

How does the rank of  $\langle x \rangle$  change if we flip its  $i$ -th bit?

## CONCLUSION

---

## TAKE-HOME MESSAGES & OPEN QUESTIONS

Indexing  $k$ -mers with their necklaces:

- preserves locality
- improves compression
- fits in well with a quotiented data structure
- combines easily with dynamic operations

Future questions:

- What is the average distance between necklaces of consecutive  $k$ -mers?
- Can we rank necklaces in subquadratic time?
- Can we accelerate ranking for necklaces of consecutive  $k$ -mers?

## TAKE-HOME MESSAGES & OPEN QUESTIONS

Indexing  $k$ -mers with their necklaces:





- preserves locality
- improves compression
- fits in well with a quotiented data structure
- combines easily with dynamic operations

Future questions:

- What is the average distance between necklaces of consecutive  $k$ -mers?
- Can we rank necklaces in subquadratic time?
- Can we accelerate ranking for necklaces of consecutive  $k$ -mers?

*Thank you!*

## REFERENCES I

-  Alanko, Jarno N, Simon J Puglisi & Jaakko Vuohtoniemi (2022). “**Succinct k-mer sets using subset rank queries on the spectral burrows-wheeler transform**”. In: *bioRxiv*, pp. 2022–05.
-  Conway, Thomas C & Andrew J Bromage (2011). “**Succinct data structures for assembling large genomes**”. In: *Bioinformatics* 27.4, pp. 479–486.
-  Sawada, Joe & Aaron Williams (2017). “**Practical algorithms to rank necklaces, Lyndon words, and de Bruijn sequences**”. In: *Journal of Discrete Algorithms* 43, pp. 95–110.
-  Zheng, Hongyu, Carl Kingsford & Guillaume Marçais (2020). “**Improved design and analysis of practical minimizers**”. In: *Bioinformatics* 36.Supplement\_1, pp. i119–i127.