# Conway-Bromage-Lyndon (CBL):

## an exact, dynamic representation of $k$-mer sets

Igor MARTAYAN, Bastien CAZAUX, Antoine LIMASSET & Camille MARCHET
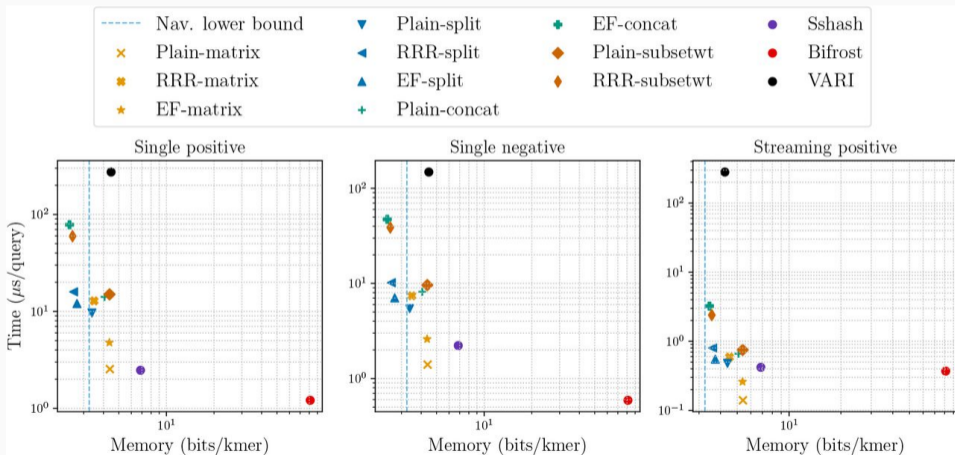University of Lille

July 14, 2024

ISMB 2024 — Montreal

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. 23]

Goal: designing a dynamic index of $k$-mers
with fast queries and relatively good compression

- membership
- enumeration
- insertion
- deletion
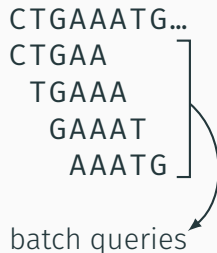- set operations $(\cup, \cap, \setminus)$

```
CTGAAATG…
CTGAA
 TGAAA
  GAAAT
   AAATG
```

batch queries

Use case: build an index incrementally, merge/intersect multiple indexes...

necklace:
smallest cyclic rotation of a word

CGAACT

CGAACT (0)
GAACTC (1)
AACTCG (2)
ACTCGA (3)
CTCGAA (4)
TCGAAC (5)

**Amortized necklace computation**
Consecutive necklaces can be
computed in $\mathcal{O}(\log k)$ amortized time

In practice: $\sim$ 10 ns / necklace

$x \longmapsto (\langle x \rangle, \text{rotation index})$ is reversible

*k*-mer view

GTCGTTCTTCCT**A**ACGTCATCTCTCATTCTG
TCGTTCTTCCT**A**ACGTCATCTCTCATTCTGT
CGTTCTTCCT**A**ACGTCATCTCTCATTCTGTG
GTTCTTCCT**A**ACGTCATCTCTCATTCTGTGA
TTCTTCCT**A**ACGTCATCTCTCATTCTGTGAC
TCTTCCT**A**ACGTCATCTCTCATTCTGTGACA
CTTCCT**A**ACGTCATCTCTCATTCTGTGACAC
TTCCT**A**ACGTCATCTCTCATTCTGTGACACG
TCCT**A**ACGTCATCTCTCATTCTGTGACACGC
CCT**A**ACGTCATCTCTCATTCTGTGACACGCA
CT**A**ACGTCATCTCTCATTCTGTGACACGCAG
T**A**ACGTCATCTCTCATTCTGTGACACGCAGG
**A**ACGTCATCTCTCATTCTGTGACACGCAGGG
ACGTCATCTCTCATTCTGTG**A**CACGCAGGGT

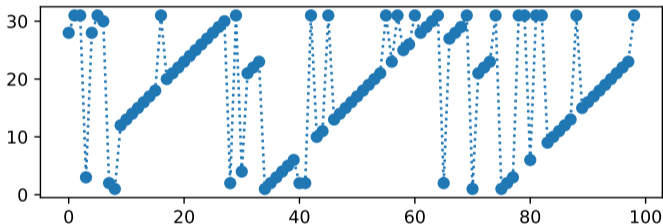| necklace view | *k*-mer view |
|---|---|
| **A**ACGTCATCTCTCATTCTG GTCGTTCTTCCT | GTCGTTCTTCCT**A**ACGTCATCTCTCATTCTG |
| **A**ACGTCATCTCTCATTCTGT TCGTTCTTCCT | TCGTTCTTCCT**A**ACGTCATCTCTCATTCTGT |
| **A**ACGTCATCTCTCATTCTGTG CGTTCTTCCT | CGTTCTTCCT**A**ACGTCATCTCTCATTCTGTG |
| **A**ACGTCATCTCTCATTCTGTGA GTTCTTCCT | GTTCTTCCT**A**ACGTCATCTCTCATTCTGTGA |
| **A**ACGTCATCTCTCATTCTGTGAC TTCTTCCT | TTCTTCCT**A**ACGTCATCTCTCATTCTGTGAC |
| **A**ACGTCATCTCTCATTCTGTGACA TCTTCCT | TCTTCCT**A**ACGTCATCTCTCATTCTGTGACA |
| **A**ACGTCATCTCTCATTCTGTGACAC CTTCCT | CTTCCT**A**ACGTCATCTCTCATTCTGTGACAC |
| **A**ACGTCATCTCTCATTCTGTGACACG TTCCT | TTCCT**A**ACGTCATCTCTCATTCTGTGACACG |
| **A**ACGTCATCTCTCATTCTGTGACACGC TCCT | TCCT**A**ACGTCATCTCTCATTCTGTGACACGC |
| **A**ACGTCATCTCTCATTCTGTGACACGCA CCT | CCT**A**ACGTCATCTCTCATTCTGTGACACGCA |
| **A**ACGTCATCTCTCATTCTGTGACACGCAG CT | CT**A**ACGTCATCTCTCATTCTGTGACACGCAG |
| **A**ACGTCATCTCTCATTCTGTGACACGCAGG T | T**A**ACGTCATCTCTCATTCTGTGACACGCAGG |
| **A**ACGTCATCTCTCATTCTGTGACACGCAGGG | **A**ACGTCATCTCTCATTCTGTGACACGCAGGG |
| **A**CACGCAGGGT ACGTCATCTCTCATTCTGTG | ACGTCATCTCTCATTCTGTG**A**CACGCAGGGT |

necklace view

**A**ACGTCATCTCTCATTCTG GTCGTTCTTCCT
**A**ACGTCATCTCTCATTCTGT TCGTTCTTCCT
**A**ACGTCATCTCTCATTCTGTG CGTTCTTCCT
**A**ACGTCATCTCTCATTCTGTGA GTTCTTCCT
**A**ACGTCATCTCTCATTCTGTGAC TTCTTCCT
**A**ACGTCATCTCTCATTCTGTGACA TCTTCCT
**A**ACGTCATCTCTCATTCTGTGACAC CTTCCT
**A**ACGTCATCTCTCATTCTGTGACACG TTCCT
**A**ACGTCATCTCTCATTCTGTGACACGC TCCT
**A**ACGTCATCTCTCATTCTGTGACACGCA CCT
**A**ACGTCATCTCTCATTCTGTGACACGCAG CT
**A**ACGTCATCTCTCATTCTGTGACACGCAGG T
**A**ACGTCATCTCTCATTCTGTGACACGCAGGG
**A**CACGCAGGGT ACGTCATCTCTCATTCTGTG

$k$-mer view

GTCGTTCTTCCT**A**ACGTCATCTCTCATTCTG
TCGTTCTTCCT**A**ACGTCATCTCTCATTCTGT
CGTTCTTCCT**A**ACGTCATCTCTCATTCTGTG
GTTCTTCCT**A**ACGTCATCTCTCATTCTGTGA
TTCTTCCT**A**ACGTCATCTCTCATTCTGTGAC
TCTTCCT**A**ACGTCATCTCTCATTCTGTGACA
CTTCCT**A**ACGTCATCTCTCATTCTGTGACAC
TTCCT**A**ACGTCATCTCTCATTCTGTGACACG
TCCT**A**ACGTCATCTCTCATTCTGTGACACGC
CCT**A**ACGTCATCTCTCATTCTGTGACACGCA
CT**A**ACGTCATCTCTCATTCTGTGACACGCAG
T**A**ACGTCATCTCTCATTCTGTGACACGCAGG
**A**ACGTCATCTCTCATTCTGTGACACGCAGGG
ACGTCATCTCTCATTCTGTG**A**CACGCAGGGT

necklace view

AACGTCATCTCTCATTCTG GTCGTTCTTCCT
AACGTCATCTCTCATTCTGT TCGTTCTTCCT
AACGTCATCTCTCATTCTGTG CGTTCTTCCT
AACGTCATCTCTCATTCTGTGA GTTCTTCCT
AACGTCATCTCTCATTCTGTGAC TTCTTCCT
AACGTCATCTCTCATTCTGTGACA TCTTCCT
AACGTCATCTCTCATTCTGTGACAC CTTCCT
AACGTCATCTCTCATTCTGTGACACG TTCCT
AACGTCATCTCTCATTCTGTGACACGC TCCT
AACGTCATCTCTCATTCTGTGACACGCA CCT
AACGTCATCTCTCATTCTGTGACACGCAG CT
AACGTCATCTCTCATTCTGTGACACGCAGG T
AACGTCATCTCTCATTCTGTGACACGCAGGG
ACACGCAGGGT ACGTCATCTCTCATTCTGTG



Size of common prefix
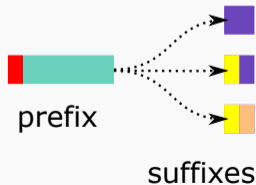between necklaces of successive $k$-mers ($k = 31$)

necklace view

```
AACGTCATCTCTCATTCTG GTCGTTCTTCCT
AACGTCATCTCTCATTCTGT TCGTTCTTCCT
AACGTCATCTCTCATTCTGTG CGTTCTTCCT
AACGTCATCTCTCATTCTGTGA GTTCTTCCT
AACGTCATCTCTCATTCTGTGAC TTCTTCCT
AACGTCATCTCTCATTCTGTGACA TCTTCCT
AACGTCATCTCTCATTCTGTGACAC CTTCCT
AACGTCATCTCTCATTCTGTGACACG TTCCT
AACGTCATCTCTCATTCTGTGACACGC TCCT
AACGTCATCTCTCATTCTGTGACACGCA CCT
AACGTCATCTCTCATTCTGTGACACGCAG CT
AACGTCATCTCTCATTCTGTGACACGCAGG T
AACGTCATCTCTCATTCTGTGACACGCAGGG
ACACGCAGGGT ACGTCATCTCTCATTCTGTG
```



Size of common prefix
between necklaces of successive $k$-mers ($k = 31$)

*How to exploit these common prefixes?*

necklaces

quotienting

prefix

suffixes

necklaces

quotienting

prefix

suffixes

Conway-Bromage-Lyndon

aaaa
ac..
cc..
cg..
gg..
gt..
tt..

tiered vector
storing pointers

rank

ac  cg  cc
         ct

bitvector
storing prefixes

vectors/tries
storing suffixes

Conway-Bromage-Lyndon

Main query steps:

1. compute $\langle x \rangle$
2. split $\langle x \rangle$ as $q \parallel r$
3. query $r$ in the bucket of $q$

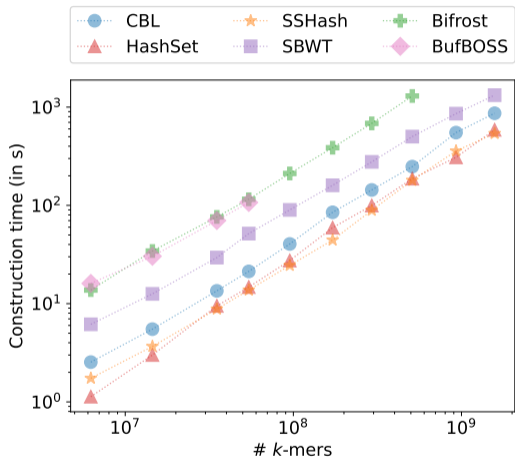$\rightarrow$ faster for consecutive $k$-mers (likely in the same bucket)

tiered vector storing pointers

rank

aaaa
ac..
cc..
cg..
gg..
gt..
tt..

bitvector storing prefixes

ac cg cc ct

vectors/tries storing suffixes

# Comparison with some *k*-mer set data structures

| category | data structure | membership | insert | delete | ∪ ∩ \ |
|----------|----------------|------------|--------|--------|-------|
| BWT | FM-index | ✓ | ✗ | ✗ | ✗ |
| — | SBWT | ✓ | ✗ | ✗ | ✗ |
| — | dynamic BOSS | ✓ | ✓ | ✓ | ✗ |
| hashing | SSHash | ✓ | ✗ | ✗ | ✗ |
| — | Bifrost | ✓ | ✓ | ✗ | ✗ |
| — | Bloom filter | approx | ✓ | ✗ | union |
| — | Quotient filter | approx* | ✓ | ✗ | union |
| other | Conway-Bromage | ✓ | ✓ | ✓ | ✓ |
| — | CBL | ✓ | ✓ | ✓ | ✓ |

*exact if a perfect hash function is used

TLDR: almost as fast as a hash table, $\sim$40–50 bits / $k$-mer ($k = 31$)

TLDR: $4\times$ faster and $3\times$ smaller than a hash table when merging a billion $k$-mers

Improving CBL's memory usage:

- suffixes among the same bucket are similar and can be compressed
- better layout of the tries (e.g. adaptive radix tries)

Extending the data structure:

- associate data (e.g. abundance) to each $k$-mer $\rightarrow$ CBL Map
- concurrent version (distribute suffix buckets between threads)

Using CBL to enumerate $k$-mers satisfying a given constraint
e.g. find $k$-mers present in ref $A$ and $B$ but not in $C$ $\longrightarrow$ preprint:

- new dynamic structure based on necklaces
- very fast queries, cache efficient
- limited memory usage ($\sim$ 40 bpk for $k$=31)
- supports fast insertion, deletion & set ops
- available as a CLI and a Rust library

*Thank you!*

github.com/imartayan/CBL

Paper

## References

📄 Alanko, Jarno N., Simon J. Puglisi & Jaakko Vuohtoniemi (2023). "Small Searchable $\kappa$-Spectra via Subset Rank Queries on the Spectral Burrows-Wheeler Transform". In: *ACDA23*, pp. 225–236.

📄 Bille, Philip et al. (2017). "Fast Dynamic Arrays". In: LIPIcs 87, 16:1–16:13. ISSN: 1868-8969.

📄 Conway, Thomas C & Andrew J Bromage (2011). "Succinct data structures for assembling large genomes". In: *Bioinformatics* 27.4, pp. 479–486.

📄 Marchini, Stefano & Sebastiano Vigna (2020). "Compact Fenwick trees for dynamic ranking and selection". In: *Software: Practice and Experience* 50.7, pp. 1184–1202.

📄 Pibiri, Giulio Ermanno & Shunsuke Kanda (2021). "Rank/select queries over mutable bitmaps". In: *Information Systems* 99, p. 101756.

📄 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.

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

Better: amortize the computation for consecutive $k$-mers.

### Key observation

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.

$m$

```
A T A A C G T C
T A A C G T C A
A A C G T C A T
A C G T C A T A
C G T C A T A A
G T C A T A A C
T C A T A A C G
C A T A A C G T
```
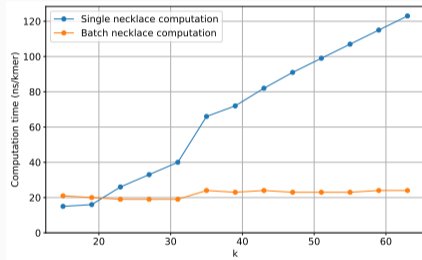
## Faster necklace computation

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.
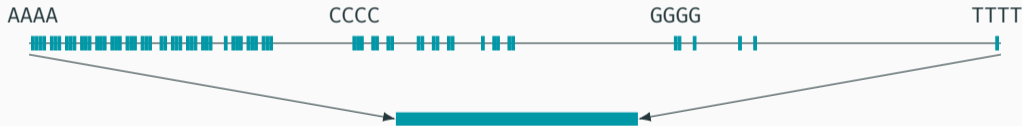$\rightarrow \mathcal{O}(nm) = \mathcal{O}(n \log k)$ for $n$ necklaces (on avg)

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