

A COMPRESSED, DYNAMIC AND LOCALITY-PRESERVING REPRESENTATION OF K-MER SETS FOR GENOMIC ANALYSIS

Igor MARTAYAN

November 10, 2023

RT MIA — Journée Réduction de Dimension — ENS Lyon



DNA SEQUENCING

DNA samples



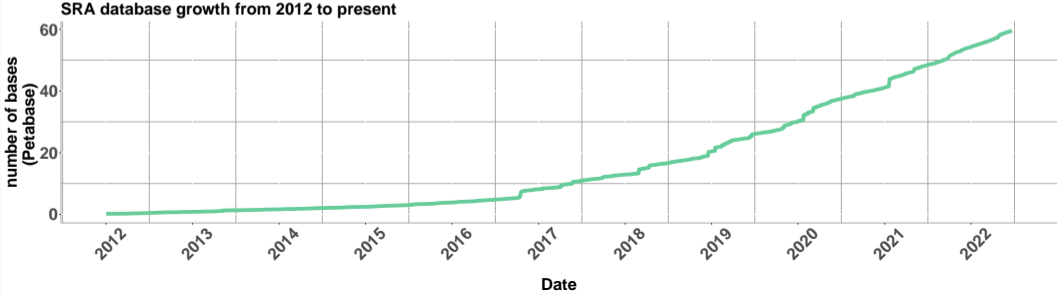
→ CTCGAGGATT...

DNA SEQUENCING

DNA samples →



→ CTCGAGGATT...



k -mer: word of size k

we typically index the k -mers of a sequence
instead of the sequence itself

```
CTGAAATG...  
CTGAA  
TGAAA  
GAAAT  
AAATG
```

most existing space-efficient data structures for storing k -mers are **static**
(e.g. spectral BWT [Alanko et al. 22], SShash [Pibiri 22])

[Conway & Bromage 11]

- we can see k -mers as integers in $\llbracket 4^k \rrbracket$
- since they're usually very sparse, we can use a **sparse bitvector** to store them

A \rightarrow 00 C \rightarrow 01 G \rightarrow 10 T \rightarrow 11

Limitations

- the data structure is static
 - it's not cache-efficient
 - $\text{index}(\text{ATGTC}) = 237$
 - $\text{index}(\text{ TGTCG}) = 950$
- average distance of $4^k/3$

[Conway & Bromage 11]

- we can see k -mers as integers in $\llbracket 4^k \rrbracket$
- since they're usually very sparse, we can use a **sparse bitvector** to store them

A \rightarrow 00 C \rightarrow 01 G \rightarrow 10 T \rightarrow 11

Limitations

- the data structure is static
 - it's not cache-efficient
 - $\text{index}(\text{ATGTC}) = 237$
 - $\text{index}(\text{ TGTCG}) = 950$
- average distance of $4^k/3$

Can we improve this approach?

THE QUEST FOR AN IDEAL DATA STRUCTURE

- **space-efficient**: close to the theoretical lower bound
- **dynamic**: support insertion and deletion after construction
- **efficient queries**:
 - membership
 - enumeration
 - insertion
 - deletion
- **locality-preserving**: reduce cache misses when querying consecutive k -mers (we often perform batch queries on many overlapping k -mers)

A COMPRESSED REPRESENTATION OF SPARSE INTEGER SETS

UNDER THE HOOD: ELIAS-FANO ENCODING

$$S = \{2, 3, 251, 403, 406, 407, 995, 999\} \quad n = 8 \quad u = 1000 \quad l = \lceil \lg \frac{u}{n} \rceil = 7 \text{ bits}$$

h_i	l_i
000	0000010
000	0000011
001	1111011
011	0010011
011	0010110
011	0010111
111	1100011
111	1100111

UNDER THE HOOD: ELIAS-FANO ENCODING

$$S = \{2, 3, 251, 403, 406, 407, 995, 999\} \quad n = 8 \quad u = 1000 \quad l = \lceil \lg \frac{u}{n} \rceil = 7 \text{ bits}$$

h_i	l_i
000	0000010
000	0000011
001	1111011
011	0010011
011	0010110
011	0010111
111	1100011
111	1100111

$n \times l$ bits

UNDER THE HOOD: ELIAS-FANO ENCODING

$$S = \{2, 3, 251, 403, 406, 407, 995, 999\} \quad n = 8 \quad u = 1000 \quad l = \lceil \lg \frac{u}{n} \rceil = 7 \text{ bits}$$

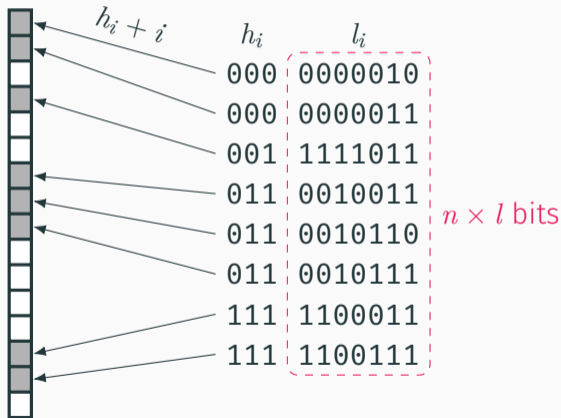


h_i	l_i
000	0000010
000	0000011
001	1111011
011	0010011
011	0010110
011	0010111
111	1100011
111	1100111

$n \times l$ bits

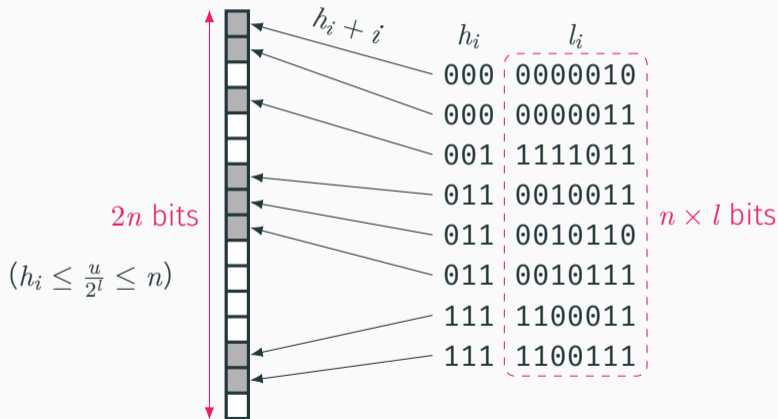
UNDER THE HOOD: ELIAS-FANO ENCODING

$$S = \{2, 3, 251, 403, 406, 407, 995, 999\} \quad n = 8 \quad u = 1000 \quad l = \lceil \lg \frac{u}{n} \rceil = 7 \text{ bits}$$



UNDER THE HOOD: ELIAS-FANO ENCODING

$$S = \{2, 3, 251, 403, 406, 407, 995, 999\} \quad n = 8 \quad u = 1000 \quad l = \lceil \lg \frac{u}{n} \rceil = 7 \text{ bits}$$



Space usage of Elias-Fano

$$EF(n, u) = 2n + n \left\lceil \lg \frac{u}{n} \right\rceil$$

Information theoretic lower bound

$$\begin{aligned} \lg \binom{u}{n} &\approx n \lg e + n \lg \frac{u}{n} \\ &\approx 1.44n + n \lg \frac{u}{n} \end{aligned}$$

Note that the bound can get lower if we have additional knowledge about the distribution.

PARTITIONING THE SET

PARTITIONING THE SET [OTTAVIANO & VENTURINI 14]



PARTITIONING THE SET [OTTAVIANO & VENTURINI 14]



Main idea: split the sequence into smaller blocks,

PARTITIONING THE SET [OTTAVIANO & VENTURINI 14]



Main idea: split the sequence into smaller blocks,
choose the best encoding depending on the density:

- for **sparse** blocks: Elias-Fano ; $2n + n \lceil \lg \frac{u}{n} \rceil$ bits

PARTITIONING THE SET [OTTAVIANO & VENTURINI 14]



Main idea: split the sequence into smaller blocks,
choose the best encoding depending on the density:

- for **sparse** blocks: Elias-Fano ; $2n + n \lceil \lg \frac{u}{n} \rceil$ bits
- for **dense** blocks: plain bitset ; u bits



Main idea: split the sequence into smaller blocks,
choose the best encoding depending on the density:

- for **sparse** blocks: Elias-Fano ; $2n + n \lceil \lg \frac{u}{n} \rceil$ bits
- for **dense** blocks: plain bitset ; u bits
- for **full** blocks: lower bound + size is enough

PARTITIONING THE SET [OTTAVIANO & VENTURINI 14]



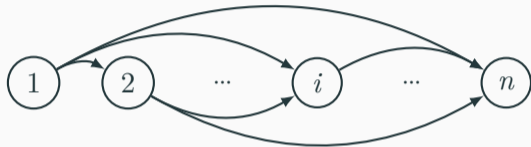
Main idea: split the sequence into smaller blocks,
choose the best encoding depending on the density:

- for **sparse** blocks: Elias-Fano ; $2n + n \lceil \lg \frac{u}{n} \rceil$ bits
- for **dense** blocks: plain bitset ; u bits
- for **full** blocks: lower bound + size is enough

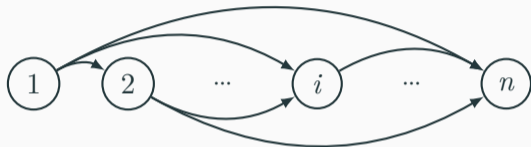
What is the optimal partition cost?

REDUCTION TO SHORTEST PATH [FERRAGINA ET AL. 11]

- $V = \llbracket 1, n \rrbracket$ $E =$
 $\{i < j ; i, j \in V\}$
- $w_{i,j} = \text{cost to encode } S[i, j]$



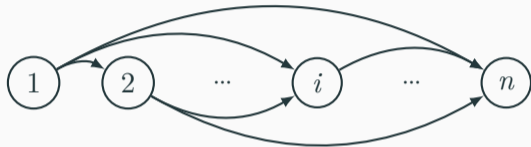
- $V = \llbracket 1, n \rrbracket$ $E = \{i < j ; i, j \in V\}$
- $w_{i,j} = \text{cost to encode } S[i, j]$



Computing the optimal partition

- **optimal solution** in $\mathcal{O}(|V| + |E|) = \mathcal{O}(n^2)$ using dynamic programming

- $V = \llbracket 1, n \rrbracket$ $E = \{i < j ; i, j \in V\}$
- $w_{i,j} = \text{cost to encode } S[i, j]$



Computing the optimal partition

- **optimal solution** in $\mathcal{O}(|V| + |E|) = \mathcal{O}(n^2)$ using dynamic programming
- **$(1 + \varepsilon)$ -approximation** in $\mathcal{O}(n \cdot \frac{1}{\varepsilon} \ln \frac{1}{\varepsilon})$ by sparsifying the graph

Main idea: **augment the partitioned data structure**

- build a B+ tree on top of the partitions
- maintain a dynamic prefix sum
- maintain dynamic successors with a y-fast trie

Good news: it only requires **$o(n)$ extra space**

Query complexity:

- membership and successor in $\mathcal{O}(\lg \lg n)$
- insertion and deletion in $\mathcal{O}(\lg n / \lg \lg n)$

BACK TO *K*-MERS

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 **bijjective** transformation
- necklaces of consecutive k -mers **share long prefixes** (a.k.a. minimizers)

RANKING NECKLACES TO IMPROVE COMPRESSION

The **number of necklaces** of size k on an alphabet with σ 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 using Sawada's algorithm

[Sawada & Williams 17]

CONCLUSION

TAKE HOME MESSAGES






- k -mer sets are ubiquitous in bioinformatics
- Elias-Fano has a near-optimal space usage
assuming we have no prior knowledge on the elements
- partitioning helps both in reducing space usage and making the structure dynamic
- a well-chosen encoding can significantly improve locality

TAKE HOME MESSAGES

- k -mer sets are ubiquitous in bioinformatics
- Elias-Fano has a near-optimal space usage
assuming we have no prior knowledge on the elements
- partitioning helps both in reducing space usage and making the structure dynamic
- a well-chosen encoding can significantly improve locality

Thank you!

REFERENCES I

-  Alanko, Jarno N, Simon J Puglisi & Jaakko Vuhtoniemi (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.
-  Elias, Peter (1974). “**Efficient storage and retrieval by content and address of static files**”. In: *Journal of the ACM (JACM)* 21.2, pp. 246–260.
-  Fano, Robert Mario (1971). ***On the number of bits required to implement an associative memory***. Massachusetts Institute of Technology, Project MAC.
-  Ferragina, Paolo, Igor Nitto & Rossano Venturini (2011). “**On optimally partitioning a text to improve its compression**”. In: *Algorithmica* 61, pp. 51–74.

REFERENCES II

-  Ottaviano, Giuseppe & Rossano Venturini (2014). **“Partitioned elias-fano indexes”**. In: *Proceedings of the 37th international ACM SIGIR conference on Research & development in information retrieval*, pp. 273–282.
-  Pibiri, Giulio Ermanno (2022). **“Sparse and skew hashing of k-mers”**. In: *Bioinformatics* 38.Supplement_1, pp. i185–i194.
-  Pibiri, Giulio Ermanno & Rossano Venturini (2017). **“Dynamic elias-fano representation”**. In: *28th Annual symposium on combinatorial pattern matching (CPM 2017)*. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
-  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.