# Cdbgtricks: strategies to update a compacted de Bruijn graph

PSC 2024

Presented by:
**Khodor Hannoush** (Genscale, Inria, Rennes university, France)
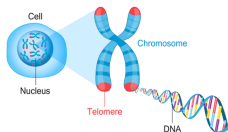Co-authors:
**Camille Marchet** (BONSAI, CNRS, Lille university, France)
**Pierre Peterlongo** (Genscale, Inria, Rennes university, France)

Aug 26-27, 2024

# Reading the DNA



Figure: Accessing the DNA

# DNA as a string

### The alphabet of the genetic language

Computer scientists treat DNA as a language of alphabet $\delta = \{A, C, G, T\}$.

# DNA as a string

### The alphabet of the genetic language

Computer scientists treat DNA as a language of alphabet
$\delta = \{A, C, G, T\}$.

### Grammar of the DNA

A DNA sequence is a string over $\delta$. For example ACGCCGTAA.

# DNA as a string

### The alphabet of the genetic language

Computer scientists treat DNA as a language of alphabet
$\delta = \{A, C, G, T\}$.

### Grammar of the DNA

A DNA sequence is a string over $\delta$. For example ACGCCGTAA.

### $k$-mer

A $k$-mer is a string of $k$ characters over $\delta$. ACGGT is a 5-mer.

# Storage of DNA sequences



**Genomes**

ACCGAGAGTCC
ACCGAGTCC

**Sequenced Reads**

ACCGAGAGTCC
CCGAGAGTCC
GAGAGTCC
ACCGAGTCC
CCGAGTCC
GAGTCC

**de Bruijn Graph (k=5)**

*k*-mers
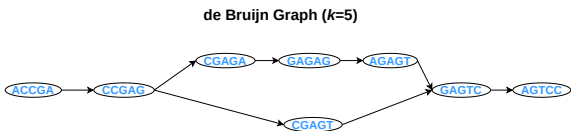
ACCGA
CCGAG
CGAGA
GAGAG
AGAGT
CGAGT
GAGTC
AGTCC

Figure: Storing sequences in a de Bruijn graph
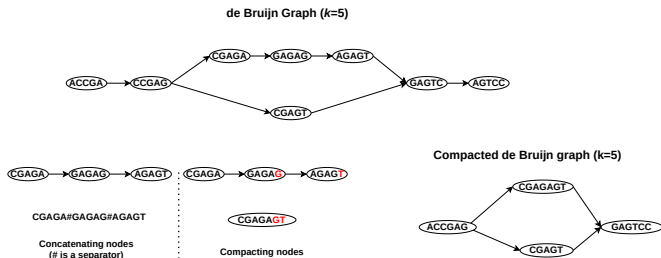
# Compacted de Bruijn graph



Figure: Compacting the de Bruijn graph

### Unitig

A unitig is a maximal non-branching path $p = \{f, v_1, v_2, ..., v_n, l\}$ such that every node $s \in p$ has only one in-coming and one out-going edges except for $f$ and $l$.

# Objective: Augmenting a cDBG

Given a compacted de Bruijn graph $G$ and a new genome sequence $S$ not in $G$, we need to add $S$ to $G$.

# Objective: Augmenting a cDBG

Given a compacted de Bruijn graph $G$ and a new genome sequence $S$ not in $G$, we need to add $S$ to $G$.

The update operation should be time-efficient as we need to support large datasets.

# Objective: Augmenting a cDBG

Given a compacted de Bruijn graph $G$ and a new genome sequence $S$ not in $G$, we need to add $S$ to $G$.

The update operation should be time-efficient as we need to support large datasets.

We need to be able to identify the regions in the graph where the update will take place.
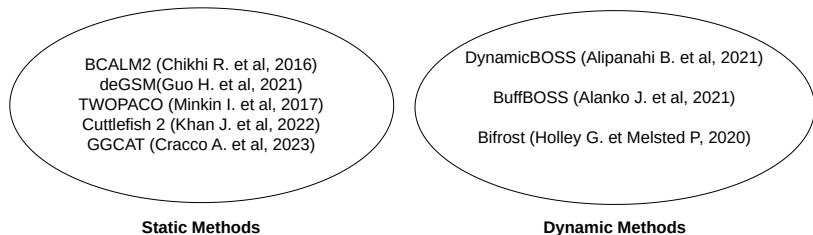
# Existing Methods



BCALM2 (Chikhi R. et al, 2016)
deGSM(Guo H. et al, 2021)
TWOPACO (Minkin I. et al, 2017)
Cuttlefish 2 (Khan J. et al, 2022)
GGCAT (Cracco A. et al, 2023)

**Static Methods**

DynamicBOSS (Alipanahi B. et al, 2021)

BuffBOSS (Alanko J. et al, 2021)

Bifrost (Holley G. et Melsted P, 2020)

**Dynamic Methods**

Figure: Methods to construct and update a de Bruijn graph

# Graph and genomes as sets of *k*-mers
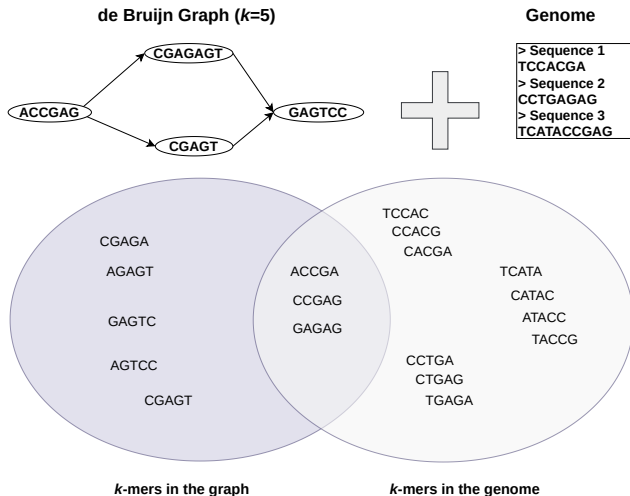


Figure: Venn diagram of the graph and the new genome
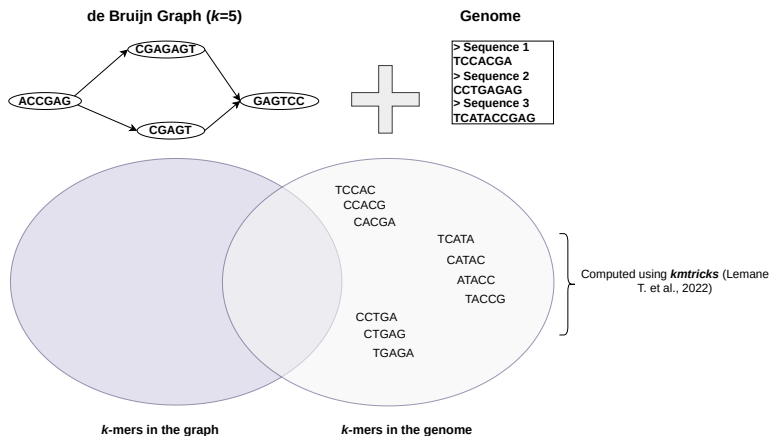
# Set of new *k*-mers



Figure: Set of new *k*-mers

## Possible cases



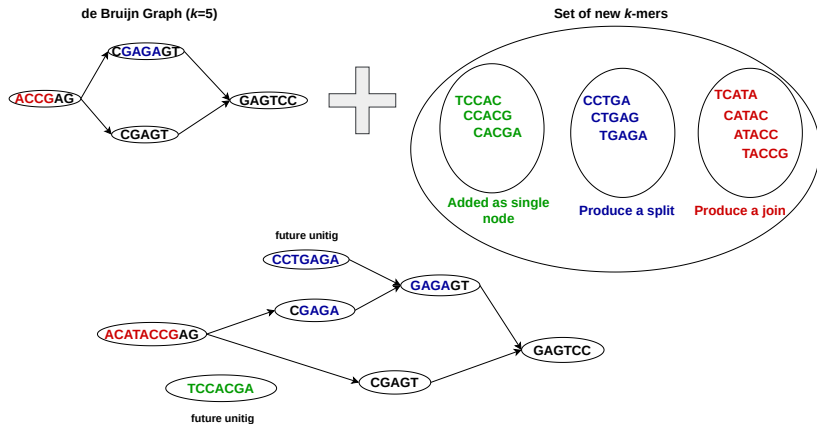Figure: Possible cases while compacting the set of new *k*-mers

# Indexing the $(k-1)$-mers

For each added $k$-mer we need to query its $(k-1)$-mers prefix and suffix.

# Indexing the $(k-1)$-mers

For each added $k$-mer we need to query its $(k-1)$-mers prefix and suffix.

The query should be performed in constant time.

# Indexing the $(k-1)$-mers

For each added $k$-mer we need to query its $(k-1)$-mers prefix and suffix.

The query should be performed in constant time.

Solution: index the $(k-1)$-mer of the graph.
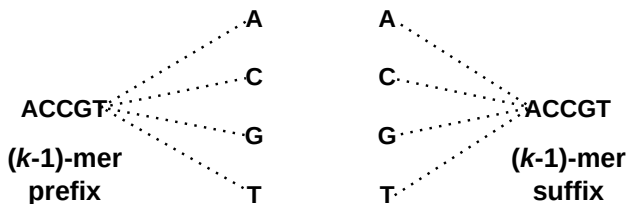
# Indexing the $(k-1)$-mers: **Drawback**



Figure: Maximum number of occurrences of a $(k-1)$-mer
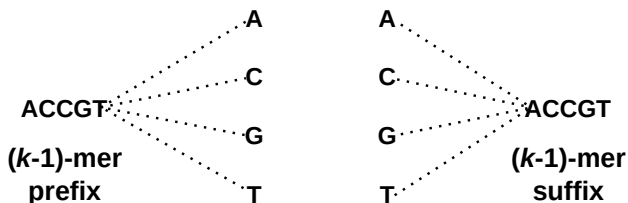
# Indexing the $(k-1)$-mers: **Drawback**



Figure: Maximum number of occurrences of a $(k-1)$-mer

A $(k-1)$-mer $x$ may have up to 8 occurrences in the graph.
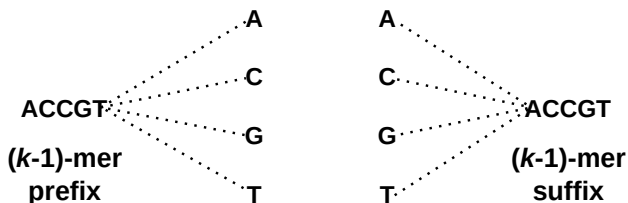
# Indexing the $(k-1)$-mers: **Drawback**



Figure: Maximum number of occurrences of a $(k-1)$-mer

A $(k-1)$-mer $x$ may have up to 8 occurrences in the graph.

It is not efficient to allocate 8 slots for each $(k-1)$-mer in the graph especially for graphs that contains billions of $(k-1)$-mers.

# Indexing $k$-mers to query $(k-1)$-mers

The $k$-mers in a compacted de Bruijn graph are unique, hence it is more memory efficient to index $k$-mers instead of $(k-1)$-mers.

# Indexing $k$-mers to query $(k-1)$-mers

The $k$-mers in a compacted de Bruijn graph are unique, hence it is more memory efficient to index $k$-mers instead of $(k-1)$-mers.

Given a $(k-1)$-mer $x$, we can query the 8 possible $k$-mers for which $x$ is either a suffix or a prefix.

# Indexing the $k$-mers

## First possible solution

- Index the $k$-mers in a hash table.

# Indexing the $k$-mers

## First possible solution

- Index the $k$-mers in a hash table.
- **Drawback**: We need to store the $k$-mers and their positions in the graph.

# Indexing the $k$-mers

## First possible solution

- Index the $k$-mers in a hash table.
- **Drawback**: We need to store the $k$-mers and their positions in the graph.

## Second possible solution

- Compute one minimal perfect hash function (MPHF) $f$ from the $k$-mers of the graph. An MPHF is a hash function that bijectively maps a set of $N$ keys to the set $\{i \mid 0 \leq i < N\}$.

# Indexing the $k$-mers

## First possible solution

- Index the $k$-mers in a hash table.
- **Drawback**: We need to store the $k$-mers and their positions in the graph.

## Second possible solution

- Compute one minimal perfect hash function (MPHF) $f$ from the $k$-mers of the graph. An MPHF is a hash function that bijectively maps a set of $N$ keys to the set $\{i \mid 0 \leq i < N\}$.
- **Drawback**: The MPHF is static, for every addition of $k$-mers we need to re-compute $f$.

# Indexing the $k$-mers

## First possible solution

- Index the $k$-mers in a hash table.
- **Drawback**: We need to store the $k$-mers and their positions in the graph.

## Second possible solution

- Compute one minimal perfect hash function (MPHF) $f$ from the $k$-mers of the graph. An MPHF is a hash function that bijectively maps a set of $N$ keys to the set $\{i \mid 0 \leq i < N\}$.
- **Drawback**: The MPHF is static, for every addition of $k$-mers we need to re-compute $f$.

## Our solution

**Solution**: partition the $k$-mers into **buckets**, and compute one MPHF for each bucket.

# Constructing future unitigs

From the set of new $k$-mers, we construct the future unitigs (**funitigs**) that get added to the graph.

# Updating the index

The position of $k$-mers in the unitigs that went into splits or joins are changed.

# Updating the index

The position of $k$-mers in the unitigs that went into splits or joins are changed.

### Re-computing MPHFs

The MPHF of the buckets to which we added new $k$-mers are re-computed.

# Cdbgtricks: product of this work

These functionalities (indexing and updating a compacted de Bruijn graph) are available in one open source software **Cdbgtricks**.

# Cdbgtricks: product of this work

These functionalities (indexing and updating a compacted de Bruijn graph) are available in one open source software **Cdbgtricks**.

**Cdbgtricks** is available at https://github.com/khodor14/Cdbgtricks

# Test Data and Tools

- Datasets
  - 7054 *E. coli* genomes (1 *E. coli genome* $\approx 4Mb$).
  - 100 human genomes (1 *human genome* $\approx 3Gb$).

# Test Data and Tools

- Datasets
  - 7054 *E. coli* genomes (1 *E. coli genome* $\approx 4Mb$).
  - 100 human genomes (1 *human genome* $\approx 3Gb$).
- Competitor tool of Cdbgtricks
  - Bifrost update [Holley, G. and Melsted, P,2020].
  - GGCAT (only for reconstruction)[Andrea Cracco and Alexandru I. Tomescu,2022].

# Test Data and Tools

- Datasets
  - 7054 *E. coli* genomes (1 *E. coli genome* $\approx 4Mb$).
  - 100 human genomes (1 *human genome* $\approx 3Gb$).
- Competitor tool of Cdbgtricks
  - Bifrost update [Holley, G. and Melsted, P,2020].
  - GGCAT (only for reconstruction)[Andrea Cracco and Alexandru I. Tomescu,2022].
- Experimental settings
  - The value of $k$ is 31.
  - The experiments were executed using 32 threads.
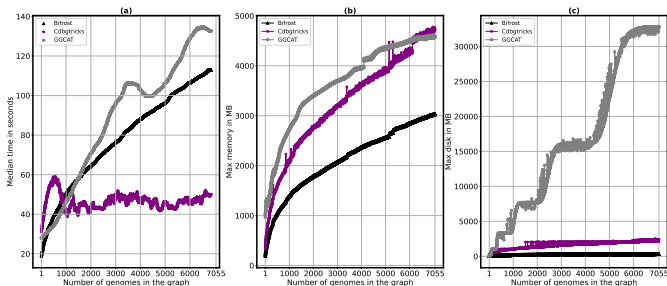
# Results of benchmark - *E.coli* genomes



Figure: Results of *E.coli* genomes dataset. Time (a), Memory (b) and Disk (c).

**Cdbtricks** is 2 to 3x faster than Bifrost and GGCAT on the *E. coli* genomes dataset.

# Results of benchmark - Human genomes



Figure: Results of human genomes dataset. Time (a), Memory (b) and Disk (c).

**Cdbtricks** is 2 to 3x faster than Bifrost and has the potential to be faster than GGCAT on a larger human genomes dataset.

# Conclusion

### Cdbgtricks

A novel method to update a compacted de Bruijn graph and its index

# Conclusion

## Cdbgtricks

A novel method to update a compacted de Bruijn graph and its index

## Indexing a compacted de Bruijn graph

**Cdbgtricks** indexes a compacted de Bruijn graph

# Conclusion

## Cdbgtricks

A novel method to update a compacted de Bruijn graph and its index

## Indexing a compacted de Bruijn graph

**Cdbgtricks** indexes a compacted de Bruijn graph

## Performance

**Cdbgtricks** outperforms the state-of-the-art tools dedicated to the creation of the update of a compacted de Bruijn graph on a data set of thousands of *E. coli* genomes and another data set of 100 *E. coli* genomes.

# Perspectives

### Optimise $k$-mer partitioning

The aim is to be able to locate consecutive $k$-mers in the same bucket.

### Colored and compacted de Bruijn graph

Store the references (**colors**) of $k$-mers in **Cdbgtricks** and implement a mechanism to update the set of colors.

# Acknowledgements

# Questions?

# Preliminaries



Figure: Double helix structure of DNA sequence.

# Preliminaries



Figure: Double helix structure of DNA sequence.

### Reverse Complement of DNA Sequence

The reverse complement $\bar{s}$ of a DNA sequence is obtained by reversing it and replacing each character by its complement (A:T, C:G,T:A,G:C). The reverse complement of ACCT is AGGT.

# Preliminaries



Figure: Double helix structure of DNA sequence.

## Reverse Complement of DNA Sequence

The reverse complement $\bar{s}$ of a DNA sequence is obtained by reversing it and replacing each character by its complement (A:T, C:G,T:A,G:C). The reverse complement of ACCT is AGGT.

## Canonical Sequence

The canonical sequence of A DNA sequence s is the smallest sequence in lexicographical order between s and its reverse complement $\bar{s}$. The canonical sequence of ACCT is ACCT.

# Bit encoding of a unitig

| A | 00 |
|---|----|
| C | 01 |
| G | 10 |
| T | 11 |

**ACCGATTATTA**

**000101100011110011100**

Figure: Encoding unitig using 2 bits per character

# Minimizer

**k=8**    **ACCGTAAT**

| | **m-mer** | **hash value** | |
|---|---|---|---|
| | ACC | 6 | |
| | CCG | 9 | |
| **m=3** | CGT | 3 | **m-mer with the smallest hash value** |
| | GTA | 11 | |
| | TAA | 5 | |
| | AAT | 8 | **lexicographically smallest m-mer** |

Figure: Hash based versus lexicographic based minimizer

# Minimal perfect hashing

**Universal**          **Perfect**          **Minimal Perfect**



Figure: Types of hashing

# Operations on adding a $k$-mer



- a. The $k$-mer $x$ forms a new unitig
- b. A unitig of the graph is extended from right
- c. A unitig of the graph is extended from left
- d. Two unitigs in the graph get merge
- e. A unitig is split into two unitigs

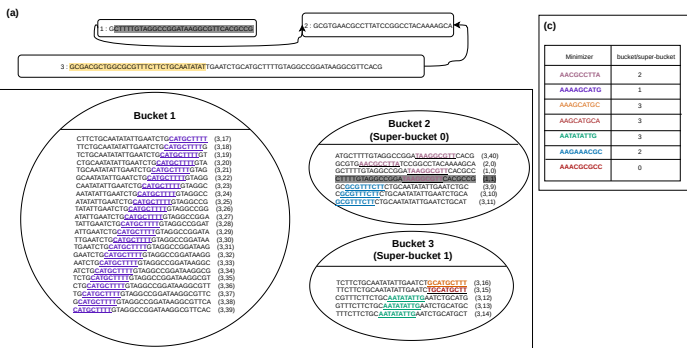# Graph index in Cdbgtricks

# Partitioning the $k$-mers



Figure: Partitioning the $k$-mers into buckets using their minimizers (the smallest $m$-mer according to some order where $0 < m < k$)

# Querying a compacted de Bruijn graph

### Presence/absence

The Cdbgtricks index helps validate whether or not the $k$-mers of a query sequence are in the graph.

### Uni-MEMs

If some $k$-mers are present, Cdbgtricks can output the unitig id and the offset in this unitig where these present $k$-mer appear.

# Test Data and Tools

- Data sets
  - 15,006 *E. coli*.
  - 10 human genomes.
- Competitor tool of Cdbgtricks
  - Bifrost [Holley, G. and Melsted, P,2020].
  - GGCAT (only for reconstruction)[Andrea Cracco and Alexandru I. Tomescu,2022].
  - SSHash (Pibiri, G., SSHash).

## Results query

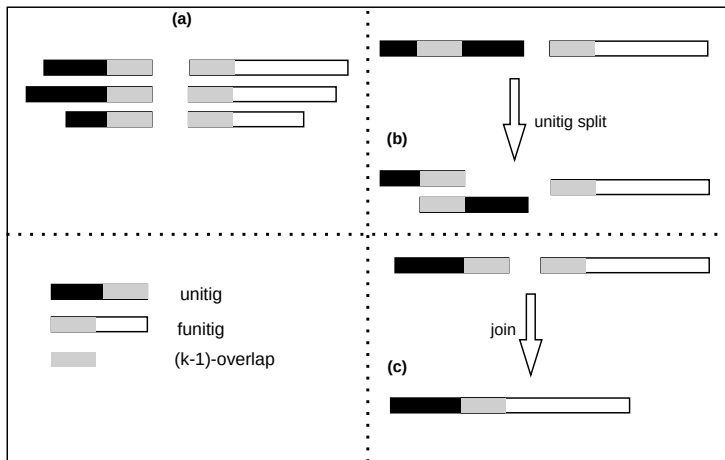| Dataset | Query type | Tool | Memory (MB) | Disk (MB) | time (mm:ss) |
|---------|-----------|------|-------------|-----------|--------------|
| *E. coli* | Negative | Cdbgtricks | 4723 | 0 | 10:02 |
| | | Bifrost | 4362 | 0 | 06:43 |
| | | SSHash | 725 | 0 | 00:07 |
| | | GGCAT | 560 | 3325 | 01:32 |
| | Positive | Cdbgtricks | 4724 | 0 | 02:15 |
| | | Bifrost | 4362 | 0 | 01:43 |
| | | SSHash | 725 | 0 | 01:10 |
| | | GGCAT | 644 | 2978 | 01:26 |
| *human* | Negative | Cdbgtricks | 25520 | 0 | 12:25 |
| | | Bifrost | 27376 | 0 | 11:37 |
| | | SSHash | 6090 | 0 | 00:07 |
| | | GGCAT | 615 | 6861 | 4:55 |
| | Positive | Cdbgtricks | 25520 | 0 | 04:23 |
| | | Bifrost | 27376 | 0 | 06:37 |
| | | SSHash | 6090 | 0 | 01:14 |
| | | GGCAT | 746 | 7053 | 05:04 |

# Performing splits and joins



Figure: **a**. Unitigs remain unchanged. **b**. The split case. **c**. The join case.