

# Developing a standard interface for sets of $k$ -mer sets index structures

---

Andreas Rempel

# Software for Computational Pangenomics

## ❖ Data structures

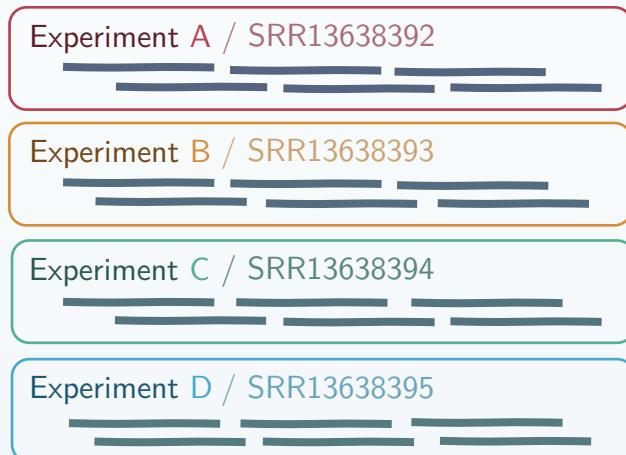


### Computational Pangenomics

- Thousands of input sequences
- Process large number of strings
  - in a reasonable time
  - within the limits of RAM
- Construct a compact index that allows fast queries

# Software for Computational Pangenomics

## ❖ Data structures

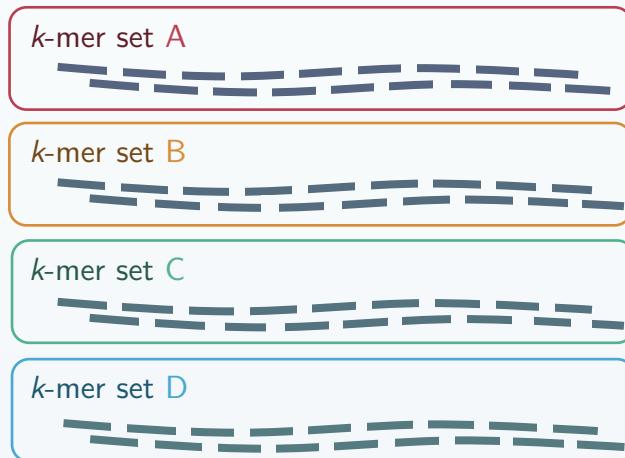


### Computational Pangenomics

- Thousands of input sequences
- Process large number of strings
  - in a reasonable time
  - within the limits of RAM
- Construct a compact index that allows fast queries

# Software for Computational Pangenomics

## ❖ Data structures



### Computational Pangenomics

- Requires advanced algorithms and data structures
- Sets of  $k$ -mer sets

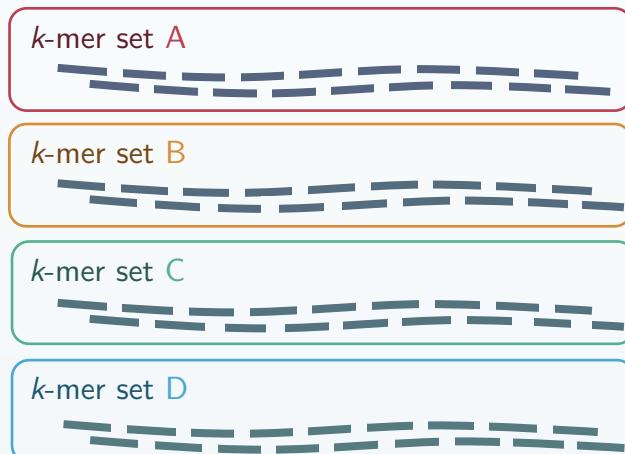
#### Idea: $k$ -mers

Split each sequence into (overlapping) substrings of length  $k$

# Software for Computational Pangenomics

## ❖ Data structures

For each input, store the set of  $k$ -mers



|                            |                                     |
|----------------------------|-------------------------------------|
| <b>AACG</b><br>Colors: {A} | <b>CAAG</b><br>Colors: {A, B}       |
| <b>AACT</b><br>Colors: {B} | <b>CACG</b><br>Colors: {C, D}       |
| <b>ACCG</b><br>Colors: {C} | <b>CAGG</b><br>Colors: {B, C, D}    |
| <b>ACCT</b><br>Colors: {D} | <b>CCGG</b><br>Colors: {A, B, C, D} |

For each  $k$ -mer, indicate the input data in which the  $k$ -mer is present

# Software for Computational Pangenomics

## ❖ Available tools

**SeqOthello**

**Mantis**

**SBT/SSBT**

**AllSomeSBT**

**HowDeSBT**

**BFT**

**Bifrost**

**Cortex**

**McCortex**

**BIGSI**

**COBS**

**RAMBO**

**Raptor**

**Themisto**

**kmtricks**

**kcollections**

**VARI**

**VARI-merge**

**Rainbowfish**

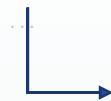
**REINDEER**

**GATB-core**

# Software for Computational Pangenomics

## ❖ Motivation

```
#my_bioinformatics_pipeline  
./path/to/programA index ...
```

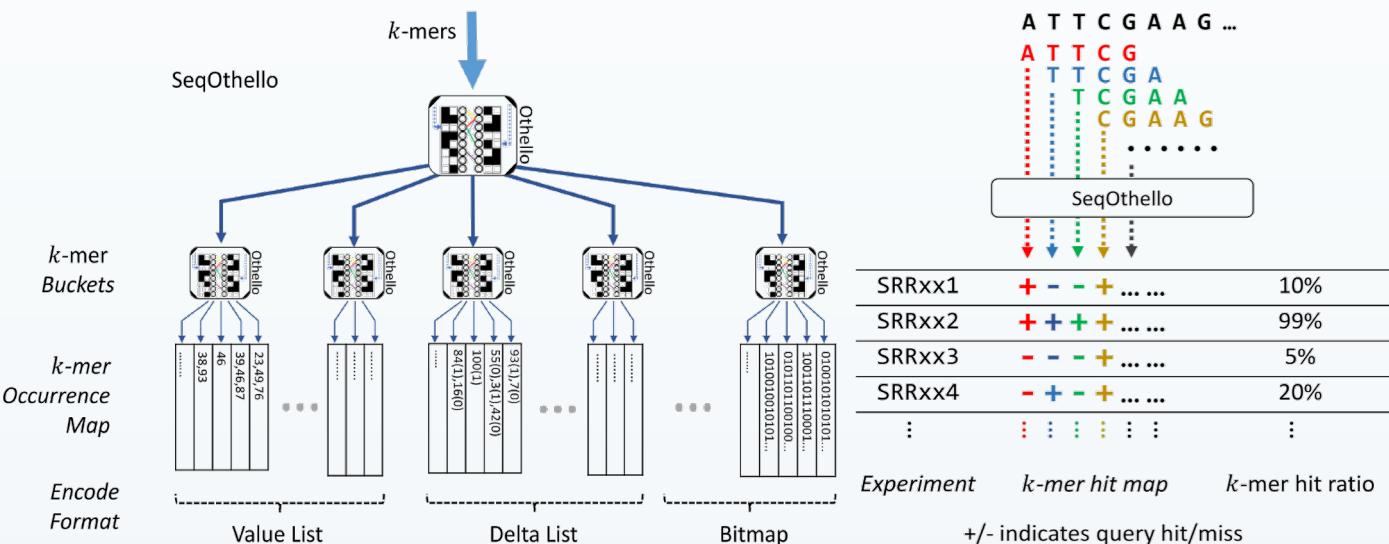


```
#my_bioinformatics_pipeline  
./path/to/programB index ...  
...  
./path/to/programB query ...  
...
```

switch between tools without having  
to rewrite your whole pipeline

# Software for Computational Pangenomics

## ❖ SeqOthello



# Software for Computational Pangenomics

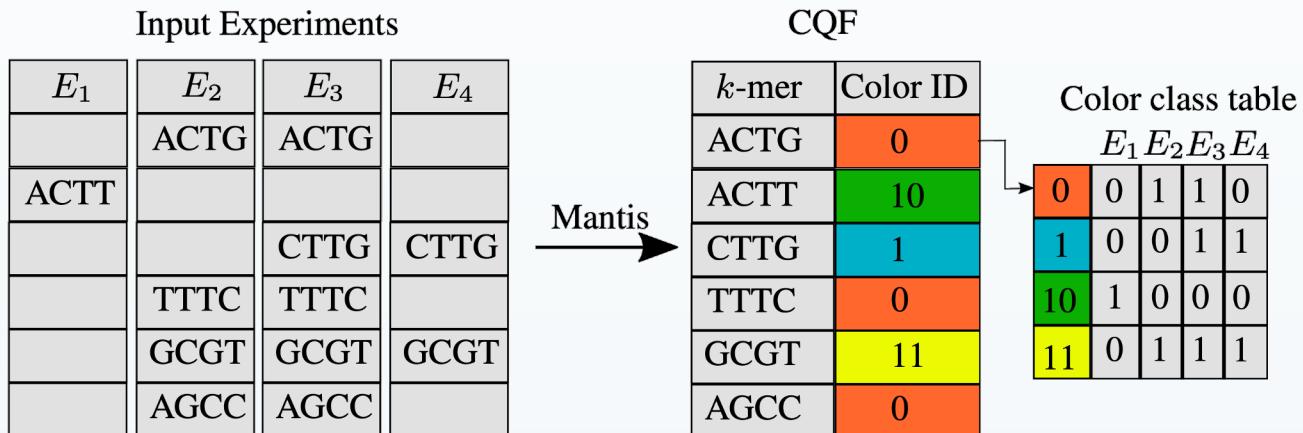
## ❖ SeqOthello



ultra-fast and memory-efficient indexing structure for RNA-seq experiments  
<https://github.com/LiuBioinfo/SeqOthello>

# Software for Computational Pangenomics

## ❖ Mantis



# Software for Computational Pangenomics

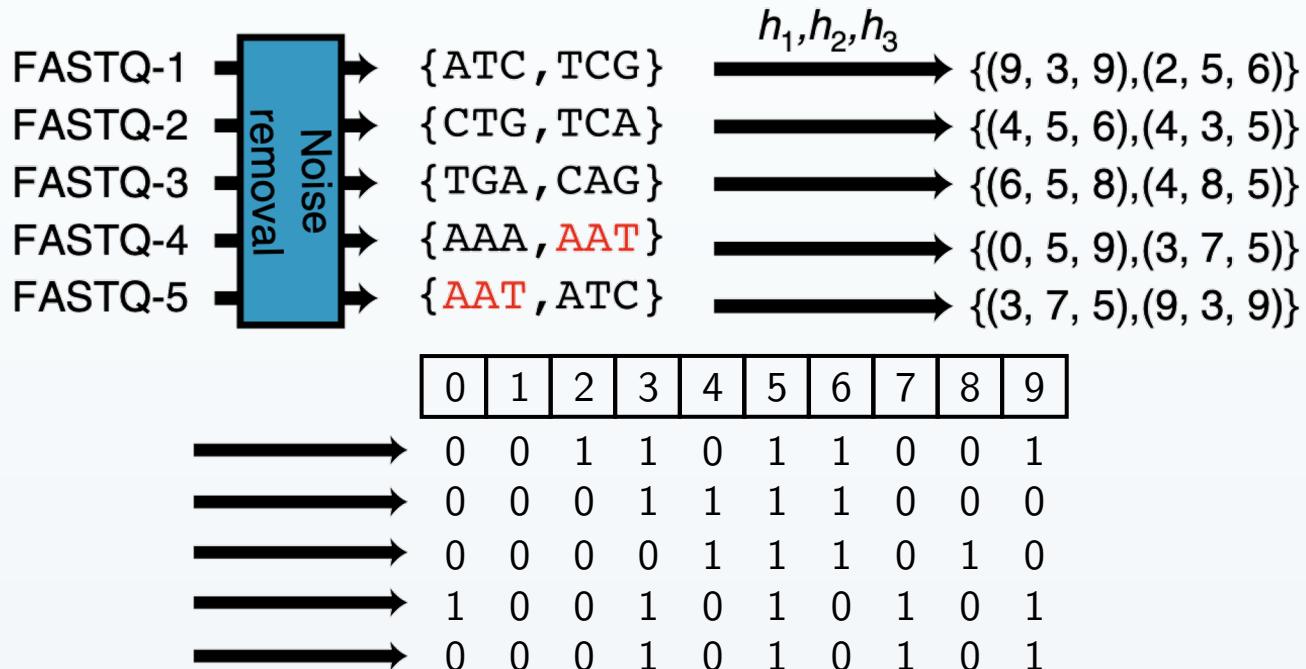
## ❖ Mantis

|                    |   |                 |
|--------------------|---|-----------------|
| squeakr count      | } | construct index |
| ./bin/mantis build |   |                 |
| ./bin/mantis mst   | → | compress index  |
| ./bin/mantis query | → | sequence query  |

fast, small & exact large-scale sequence-search index for raw-read experiments  
<https://github.com/splatlab/mantis>

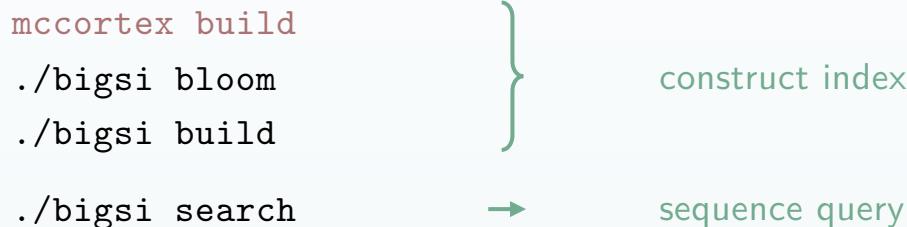
# Software for Computational Pangenomics

## ❖ BIGSI



# Software for Computational Pangenomics

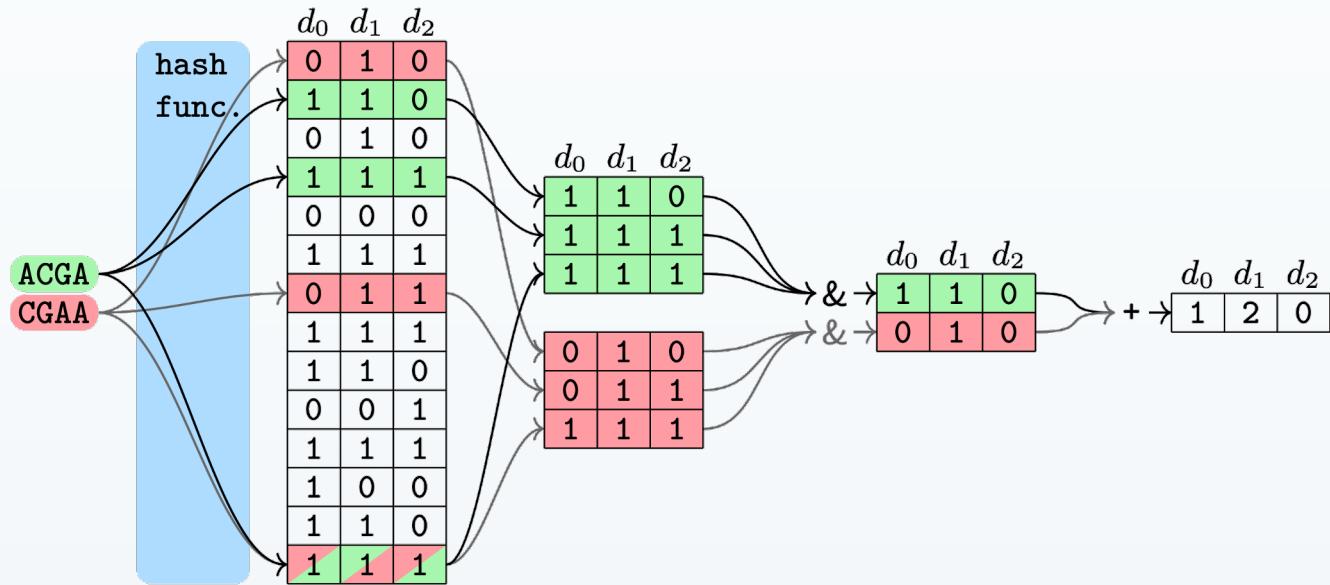
## ❖ BIGSI



bit-sliced signature index for ultra-fast search in bacterial and viral genomic data  
<https://github.com/iqbal-lab-org/bigsi>

# Software for Computational Pangenomics

❖ COBS



# Software for Computational Pangenomics

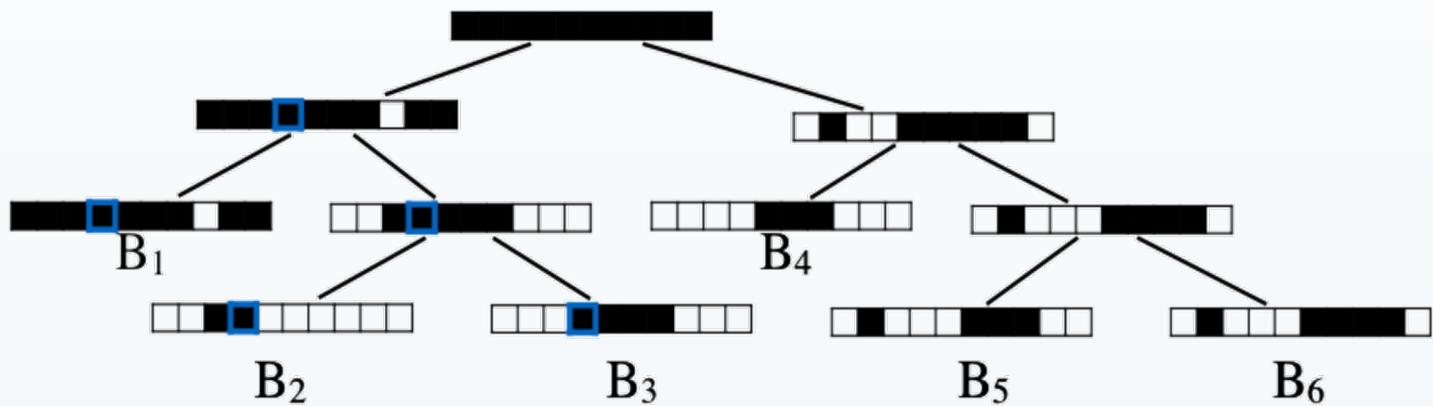
## ❖ COBS

```
./src/cobs compact_construct      construct index  
./src/cobs query                →      sequence query
```

compact bit-sliced signature index for ultra-fast search in large-scale genomic data  
<https://github.com/bingmann/cobs>

# Software for Computational Pangenomics

## ❖ Sequence Bloom Tree



# Software for Computational Pangenomics

## ❖ SBT/SSBT

|                                |   |                 |
|--------------------------------|---|-----------------|
| <code>./src/bt hashes</code>   | } | construct index |
| <code>./src/bt count</code>    |   |                 |
| <code>./src/bt build</code>    |   |                 |
| <code>./src/bt compress</code> | → | compress index  |
| <code>./src/bt query</code>    | → | sequence query  |

bloom-tree based index and search for short-read sequencing experiments  
<https://github.com/Kingsford-Group/bloomtree | splitsbt>

# Software for Computational Pangenomics

## ❖ HowDeSBT

```
./howdesbt makebf  
./howdesbt cluster  
./howdesbt build  
  
./howdesbt query
```



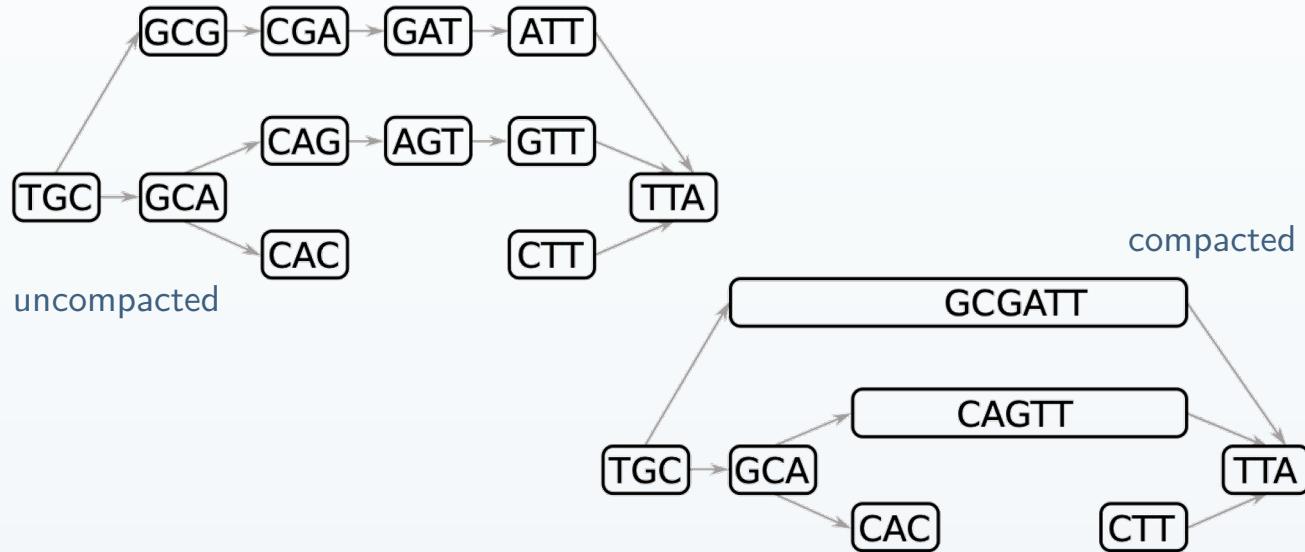
construct index

sequence query

improved index and search for short-read sequencing experiments  
<https://github.com/medvedevgroup/HowDeSBT>

# Software for Computational Pangenomics

## ❖ Bifrost



# Software for Computational Pangenomics

## ❖ Bifrost

|                               |   |                               |
|-------------------------------|---|-------------------------------|
| <code>./bifrost build</code>  | → | construct index               |
| <code>./bifrost update</code> | → | add sequences<br>update index |
| <code>./bifrost query</code>  | → | sequence query                |

parallel construction, indexing & querying of colored compacted de Bruijn graphs  
<https://github.com/pmelsted/bifrost>

# Software for Computational Pangenomics

## ❖ COSMO/VARI-merge

|                            |   |                 |
|----------------------------|---|-----------------|
| <code>./cosmo-build</code> | → | construct index |
| <code>./pack-color</code>  | → | compress index  |
| <code>./vari-merge</code>  | → | merge graphs    |
| <code>./color-merge</code> |   | update index    |

fast construction and merging of succinct colored de Bruijn graphs  
[https://github.com/cosmo-team/cosmo/tree/VARI | VARI-merge](https://github.com/cosmo-team/cosmo/tree/VARI%20|VARI-merge)

# Software for Computational Pangenomics

## ❖ Cortex

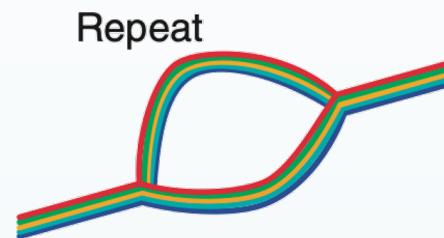


Heterozygous

Homozygous



Repeat



Repeat



Polymorphic site

# Software for Computational Pangenomics

## ❖ Cortex

```
./bin/cortex_var  
    --se_list/pe_list  
    --colour_list } construct index
```

```
./bin/cortex_var  
    --detect_bubbles  
    --output_bubbles } detect bubbles  
                      call variants
```

efficient & low-memory software for consensus assembly and variation analysis  
<https://github.com/iqbal-lab/cortex>

# Software for Computational Pangenomics

## ❖ McCortex

```
./bin/mccortex31
  - build
  - index
  - clean
  - correct
  - unitigs
  - contigs
```

construct index

clean/correct errors

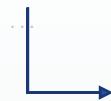
output unitigs  
assemble contigs

population de novo assembly and variant calling using linked de Bruijn graphs  
<https://github.com/mcveanlab/mccortex>

# Software for Computational Pangenomics

## ❖ Motivation

```
#my_bioinformatics_pipeline  
./path/to/programA index ...
```



```
#my_bioinformatics_pipeline  
./path/to/programB index ?? build ??  
...
```

switch between tools without having  
to rewrite your whole pipeline  
**would be nice...**

# Software for Computational Pangenomics

## ❖ Step I: Define a common standard interface

- Contact authors of existing tools for exchange of ideas
- Decide on a set of core features and consistent naming
- Update tools or write wrapper implementing the interface

|            |          |            |
|------------|----------|------------|
| AllSomeSBT | Cortex   | VARI       |
| HowDeSBT   | McCortex | VARI-merge |
| Themisto   | BIGSI    | REINDEER   |
| kmtricks   | COBS     | GATB-core  |

# Software for Computational Pangenomics

## Core interface:

- construct index
  - **input:** set of sequences<sup>1</sup>, each with a distinct color
  - **output:** index
- add | remove
  - **input:** set of sequences, index
  - **output:** index
- read/write index
  - **input:** index
  - **output:** proprietary format | gfa | fasta |  $k$ -mer matrix
- $k$ -mer query
  - **input:**  $k$ -mer
  - **output:** set of colors
- color query
  - **input:** color
  - **output:** all  $k$ -mers with this color

<sup>1</sup> sequence = assembled genome | read set | set of  $k$ -mers

# Software for Computational Pangenomics

## Optional interface:

- general query
  - **input:** set of sequences, set of colors | quorum<sup>2</sup>
  - **output:** all hits matching the query
- forward/backward extension
  - **input:**  $k$ -mer / unitig
  - **output:** neighboring  $k$ -mers / unitigs in De Bruijn Graph
- abundance query
  - **input:**  $k$ -mer or longer sequence (e.g. transcript), color
  - **output:** abundance of the sequence in this color

## Additional interface:

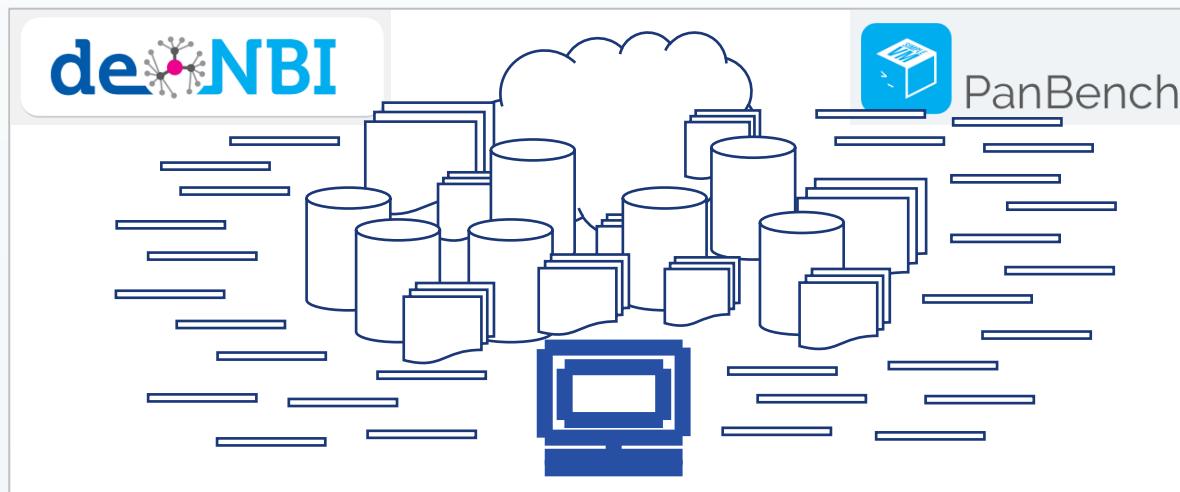
- use case specific features
  - genome assembly, e.g. filter erroneous  $k$ -mers, assemble contigs
  - phylogenomics, e.g. output list of splits for the  $k$ -mer/color sets
  - and more ...

<sup>2</sup> quorum = required min. ratio of colors that have a match

# Software for Computational Pangenomics

## ❖ Step II: Set up an automated test environment

- Select suitable test data sets and criteria for an evaluation
- Implement a benchmarking platform on the de.NBI cloud



# Software for Computational Pangenomics

## ❖ Step III: Perform the actual benchmarking study

- Contact authors to participate in a benchmark of tools
- Provide an open catalog of evaluated tools with results

**Bifrost** ★★★★☆  
Guillaume Holley

[download from github](#)  
22.371 downloads

Bifrost is a software tool for highly parallel construction and indexing of colored and compacted de Bruijn graphs.

**SeqOthello** ★★★★☆  
Ye Yu, Jinpeng Liu

[download from github](#)  
13.843 downloads

SeqOthello is an ultra-fast and memory-efficient indexing structure to support sequence query against collections...

**Mantis** ★★★★☆  
Prashant Pandey

[download from github](#)  
12.256 downloads

Mantis is a space-efficient data structure that can be used to index thousands of raw-read experiments and facilitate...

### Bifrost

Core interface: **fully implemented**  
Optional interface: **fully implemented**

Genome Assembly: **not implemented**  
Phylogenomics: **partly implemented**

### Benchmarking

Runtime:   
Memory: 

[Which tests were performed?](#)

# Software for Computational Pangenomics

## ❖ Conclusion

- Pangenomes can hold a variety of information
- Many different software tools are available
- Establishing a standard interface has the potential to enrich comparative genomics and genome research



# Software for Computational Pangenomics

## ❖ Special Thanks

Prof. Dr. Jens Stoye  
Genome Informatics group  
Graduate School “DILS”

**Thank you for  
your attention!**

Contact: [andreas.rempel  
@uni-bielefeld.de](mailto:andreas.rempel@uni-bielefeld.de)

