

Minimizer-space de Bruijn graphs for pangenomics

Rayan Chikhi

Institut Pasteur & CNRS

DSB 2022

Maybe you've seen this talk before..

- **Compression+Computation '22**
- **Pangenome Bio Hacking, Dec '21**
- **RECOMB '21**

Today's diff: Even more pangenomics :)

Bacterial Pangenomics: representing and searching in 100,000s bacterial genomes

- k -mer indexes (VARI, Bifrost, MetaGraph, Reindeer, SShash ..)¹
- MinHash sketches (sourmash)
- ¿Pangenome graphs?

¹review: Chikhi, Holub, Medvedev 2019, Marchet *et al* 2020

Bacterial Pangenomics: representing and searching in 100,000s bacterial genomes

- *k*-mer indexes (VARI, Bifrost, MetaGraph, Reindeer, SShash ..)¹
- MinHash sketches (sourmash)
- ¿Pangenome graphs?

Graph challenges:

- terabyte-sized input
- construction
- visualization

¹review: Chikhi, Holub, Medvedev 2019, Marchet *et al* 2020

Bacterial Pangenomics: representing and searching in 100,000s bacterial genomes

- k -mer indexes (VARI, Bifrost, MetaGraph, Reindeer, SShash ..)¹
- MinHash sketches (sourmash)
- ¿Pangenome graphs?

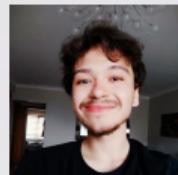
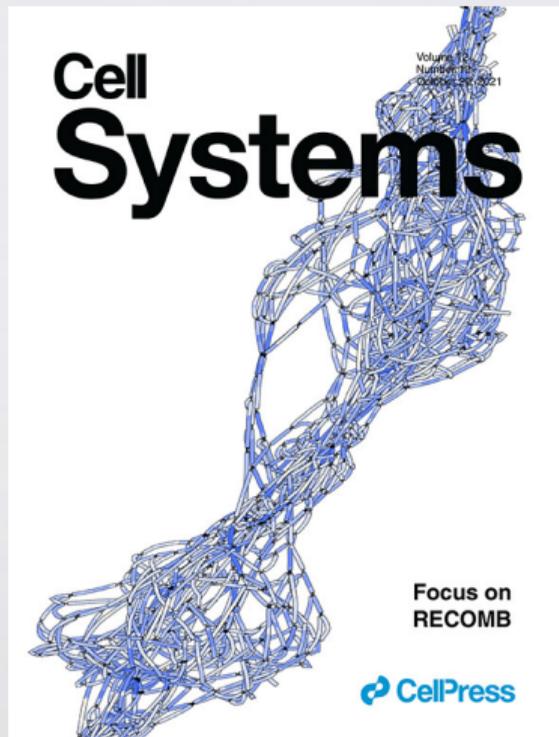
Graph challenges:

- terabyte-sized input
- construction
- visualization

In this talk: 100x-1000x cheaper pangenome graphs through controlled information loss

¹review: Chikhi, Holub, Medvedev 2019, Marchet *et al* 2020

highly scalable dBGs: Minimizer-space de Bruijn graphs



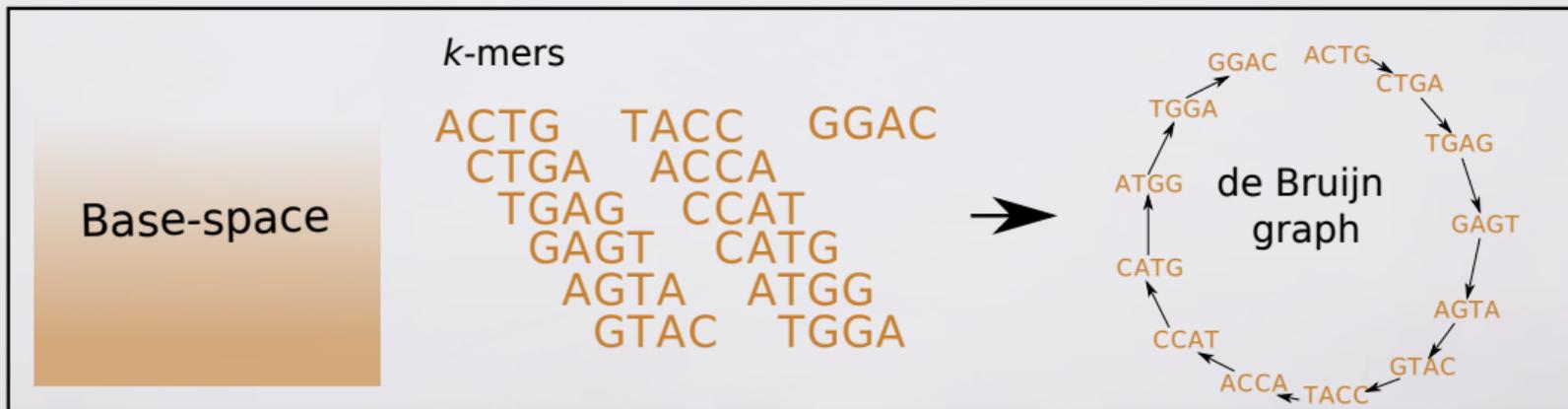
Barış Ekim



Bonnie Berger

Preliminaries k -mers, de Bruijn graph (dBG)

Reference genome ACTGAGTACCATGGAC
ACTGAGTAC
Reads CTGAGTACCAT
GAGTACCATGGAC



Preliminaries: Minimizers

Two kinds:

- **window**. Local: “smallest” l -mer in a window

AATGACATGATCATGA

AA

AC

AC

- **universe**. Global: set of l -mers with low hash values

Fixed set of
universe minimizers

AATGACATGATCATGA

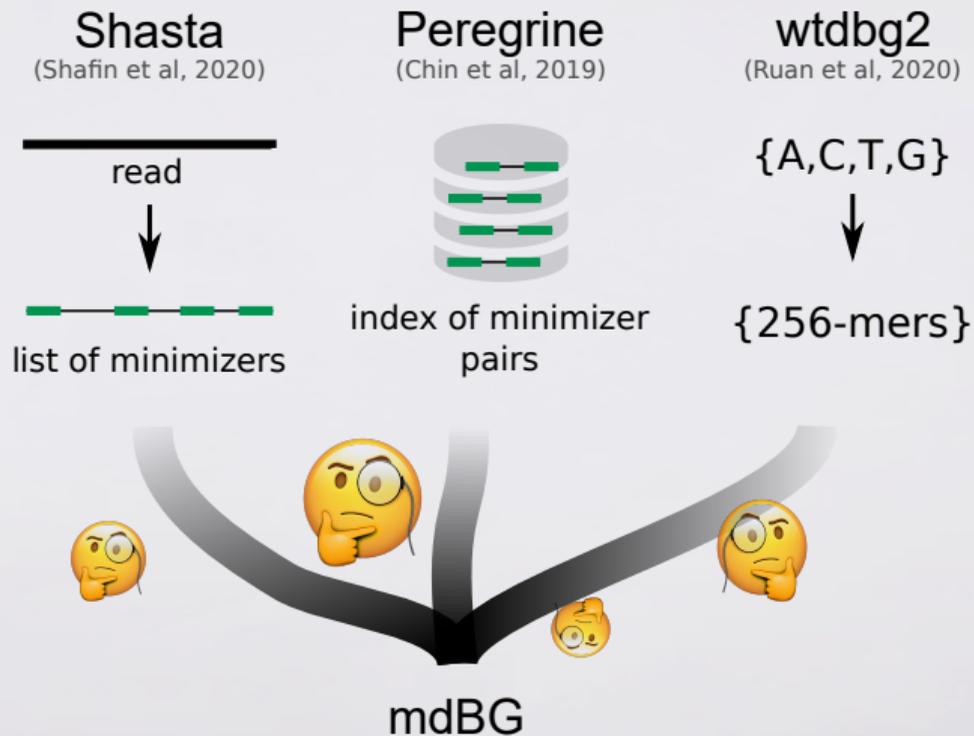
GA

TC

GA CC
TC

From now on: **universe**. (Also called Scaled MinHash)

This work: stems from three ideas



Our approach: Minimizers as *tokens* of the alphabet

Classical alphabet: $\Sigma_{DNA} = \{A, C, T, G\}$

A k -mer with $k = 3$: AGT

Our approach: Minimizers as *tokens* of the alphabet

Classical alphabet: $\Sigma_{DNA} = \{A, C, T, G\}$

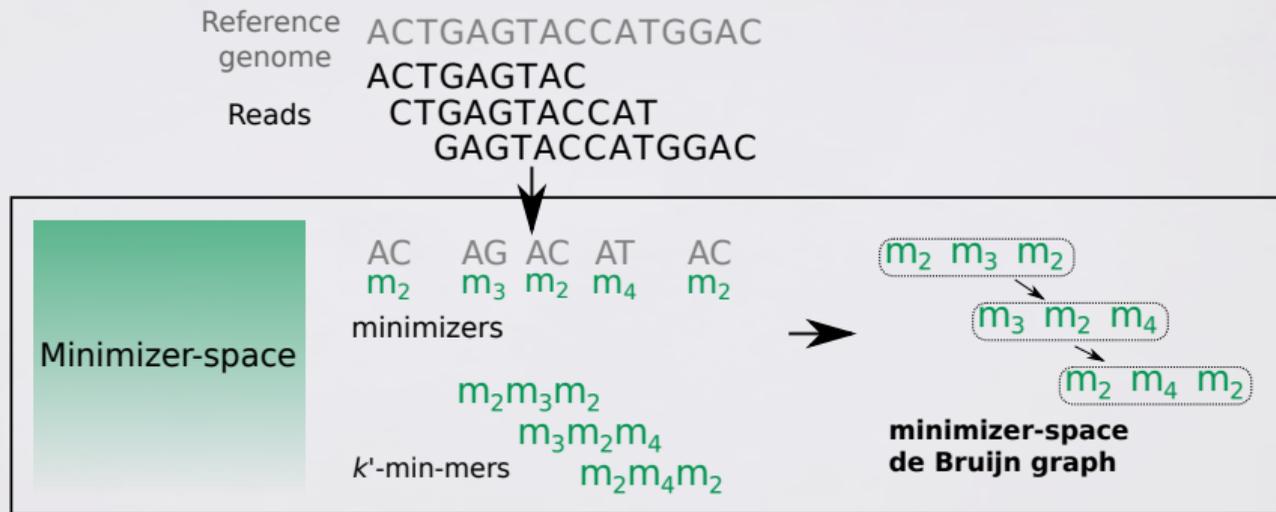
A k -mer with $k = 3$: AGT

Minimizer alphabet: $\Sigma^\ell = \{\text{all minimizers of length } \ell\} = \{m_1, m_2, m_3, \dots\}$

where e.g. $\ell = 2$, $m_1 = AA$, $m_2 = AC$, $m_3 = AG$, $m_4 = AT$

A k -mer over Σ^ℓ (a k -min-mer): $m_1 m_3 m_2$

Minimizer-space de Bruijn graph



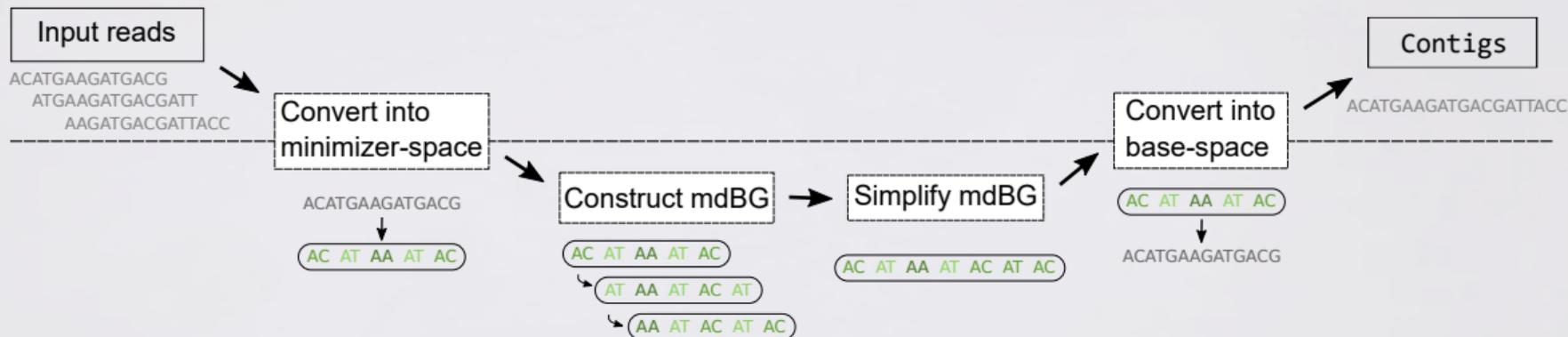
A **minimizer-space de Bruijn graph** is a **de Bruijn graph** over the **minimizer alphabet**.

Nodes = k -min-mers,

Edges = exact overlaps between $k-1$ minimizers

Applied to whole-genome *de novo* assembly

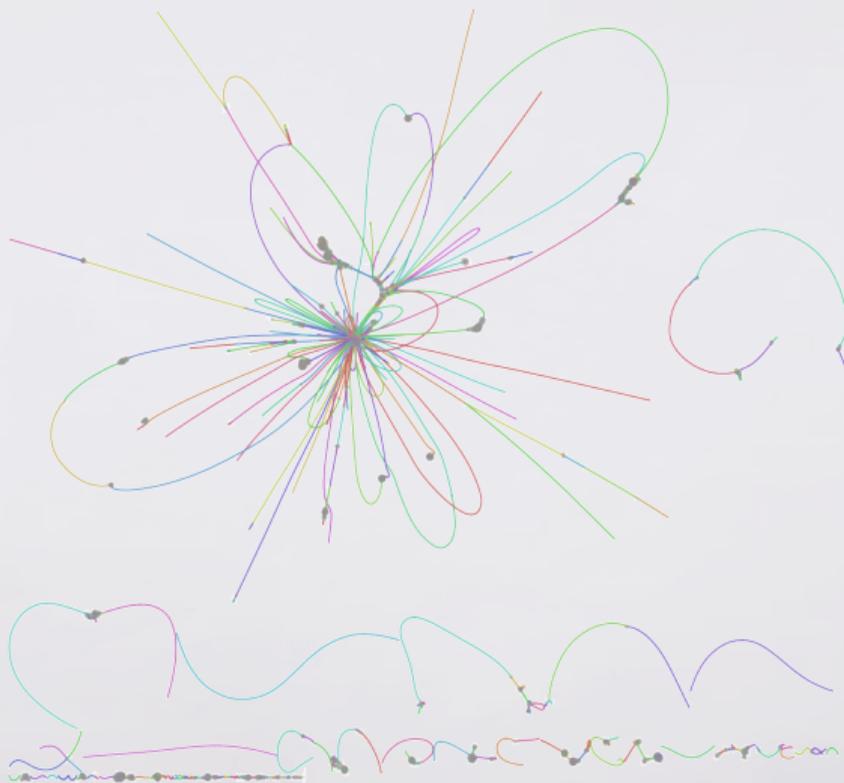
From accurate HiFi (< 1% error-rate) reads



Whole human PacBio HiFi (HG002) 50x coverage:

Tool name	Peregrine	hifiasm	rust-mdbg
Wall-clock time	14h8m	58h41m	10m23s
Memory usage	188 GB	195 GB	10 GB
# contigs	8109	431	805
NG50 (Mbp)	18.2	88.0	16.1
Genome fraction	97.0%	94.2%	95.5%

Human HiFi mdBG



Assembly implementation details

- `gfatools` ([H. Li, unpublished](#)) for graph simplifications
- Automatic parameters (suboptimal):

$$\ell = 12$$

$$\text{density} = 0.003$$

$$k = \frac{3}{4} \cdot \text{avg_readlen} \cdot \text{density}$$

- Multi-k script (à la IDBA/SPAdes).
- Code available at github.com/ekimb/rust-mdbg/

Minimizer considerations

- Universe minimizers, NtHash (rolling).
- Tested: Syncmers, Locally Consistent Parsing
- Untested: Strobemers, any other exotic minimizer scheme
- Barış' insight: we may need Minimal Confidently Alignable Substrings (winnomap2)

Results: Pangenome graph of 661,405 bacterial genomes

Data from Blackwell et al, 2021:

2.9T 661k_assemblies.fa

1.6T 661k_assemblies.fa.lz4

```
rust-mdbg -k 10 -l 12 --density 0.001 --minabund 1 661k_assemblies.fa.lz4
```

Largest 5
connected
components:



Taxons in component

18

22

4

22

10

Dominant species

*Mycobacterium
tuberculosis*

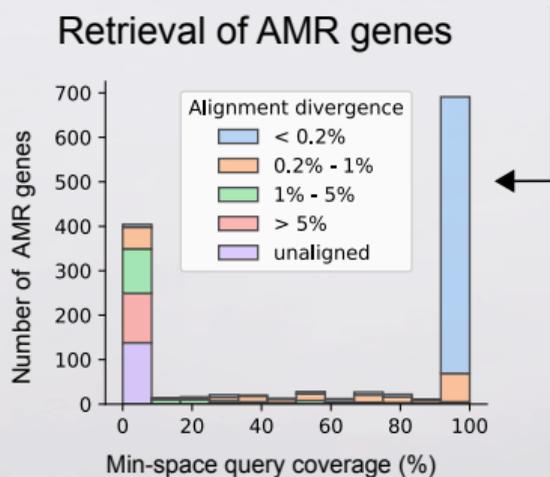
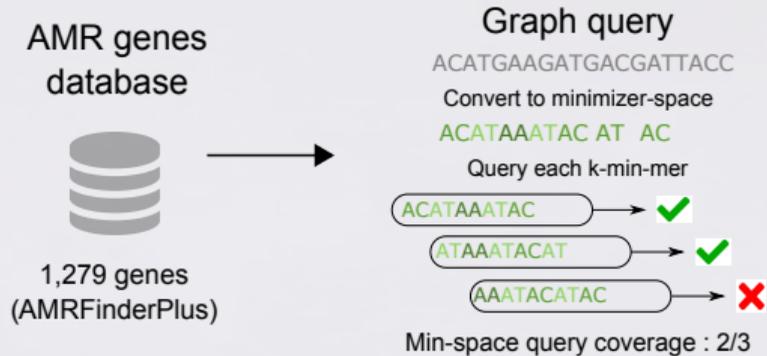
*Salmonella
enterica*

*Burkholderia
gladioli*

*Pseudomonas
protegens*

*Cupriavidus
alkaliphilus*

Biological results: Querying AMR genes



Behind the scenes

- rust-mdbg tool: from reads to raw mdBG
- set of scripts
(github.com/ekimb/rust-mdbg/tree/master/experiments/661k_genomes)

In particular:

- grep for k-min-mers search (10mins)
- lz4 k-min-mer compression
- pangenome .gfa.gz just the topology: 2-20GB
- "Resolution": 10-100kbp (kminmer span)

Behind the scenes

- rust-mdbg tool: from reads to raw mdBG
- set of scripts
(github.com/ekimb/rust-mdbg/tree/master/experiments/661k_genomes)

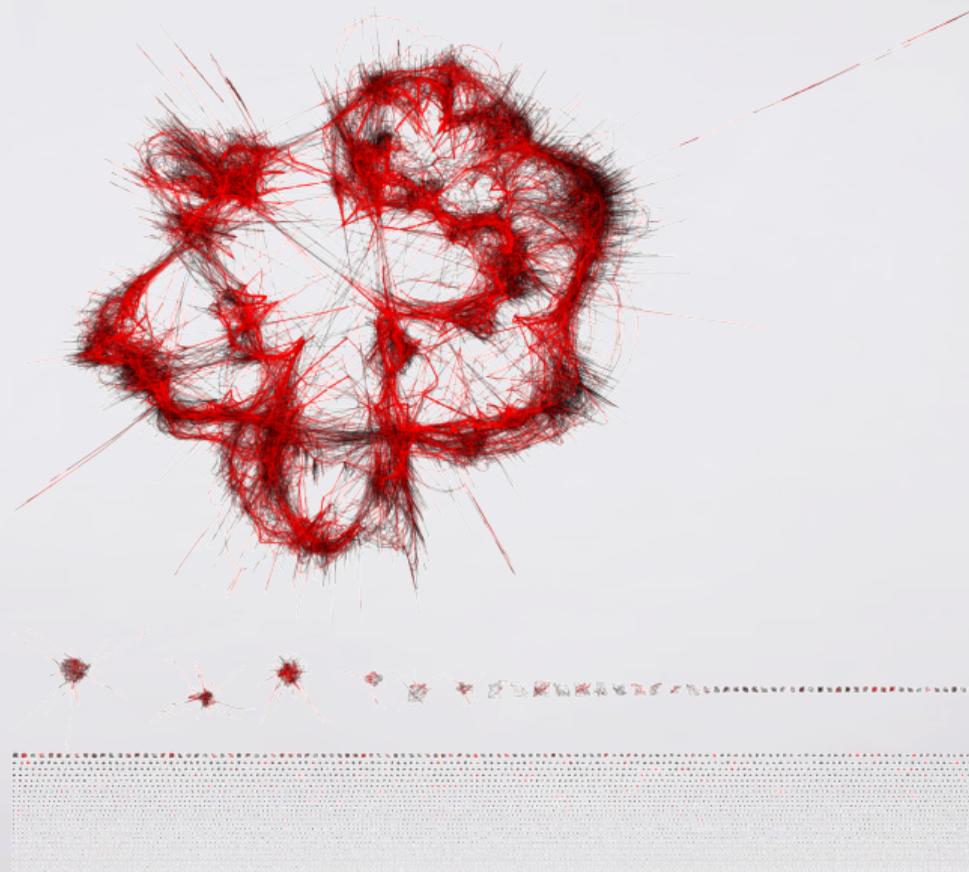
In particular:

- grep for k-min-mers search (10mins)
- lz4 k-min-mer compression
- pangenome .gfa.gz just the topology: 2-20GB
- "Resolution": 10-100kbp (kminmer span)

What we *don't* have:

- Succinct (colored) mdBGs
- O(1) k-min-mer sequence search
- visualization of pangenome mdBGs (100k-1M nodes)

Results: Pangenome graph of 160,000 *E. coli* genomes



Results: Pangenome graph of 160,000 E. coli genomes

Exploring 167,000 E. colis

Graph drawing

Scope:

Node(s):

Match: Exact Partial

Distance:

Style: Single Double

Graph display

Zoom:

Node width:

Colour by depth

Node labels



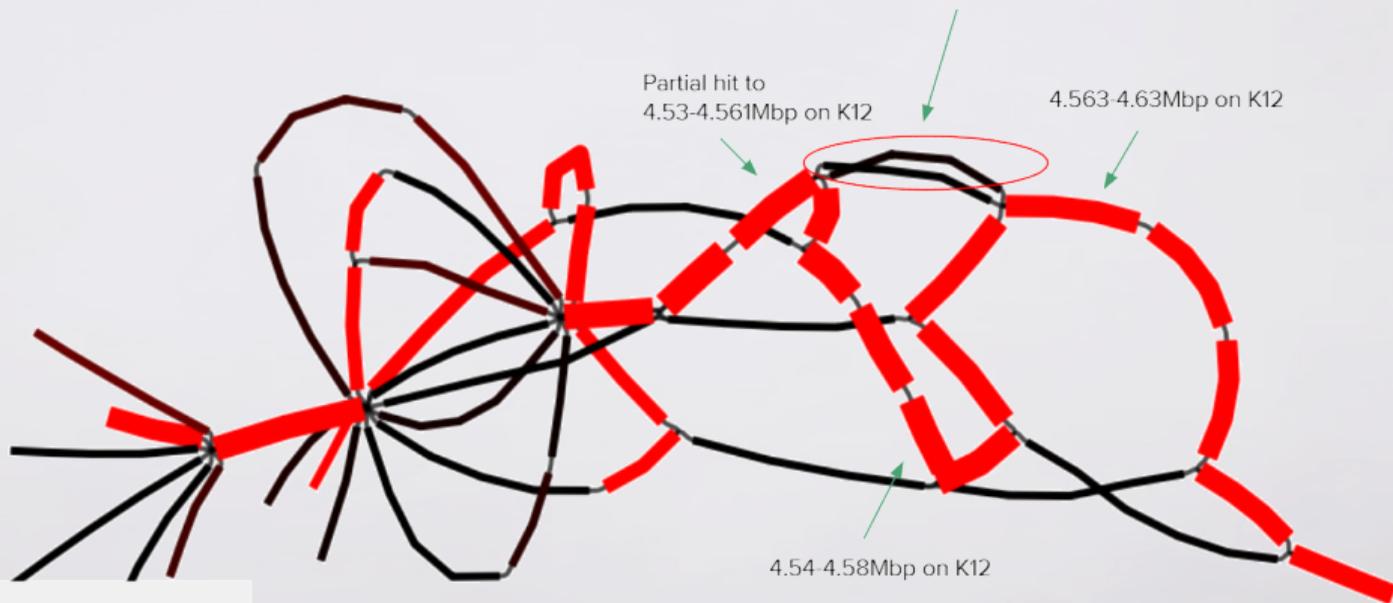
Results: Pangenome graph of 160,000 E. coli genomes

Graph comparative genomics

Split hit on K12 -> exchange?

4.53-4.57Mbp (start -> 34kbp)

4.58-4.62Mbp (44kbp -> end)



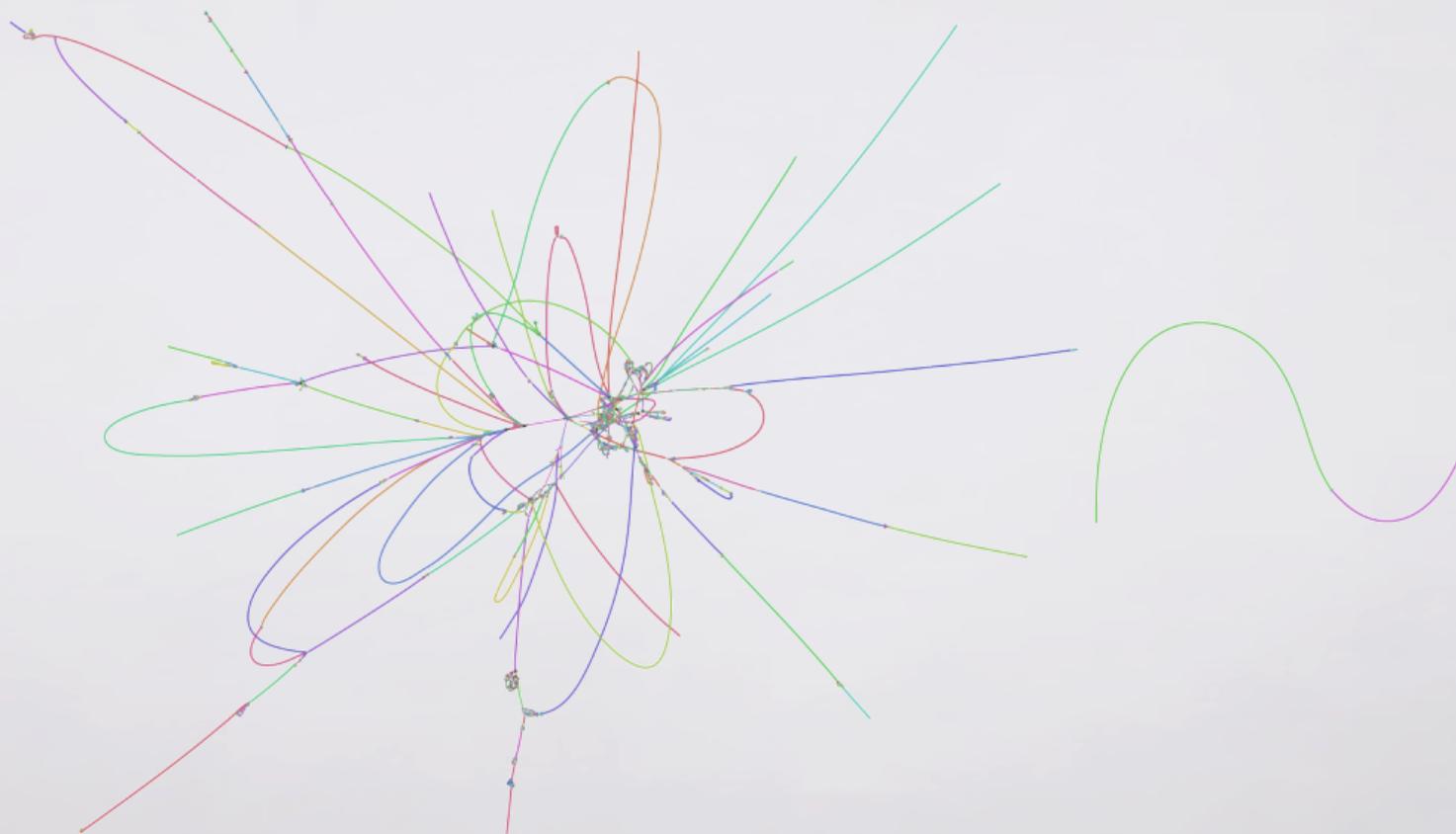
Bandage

BLAST

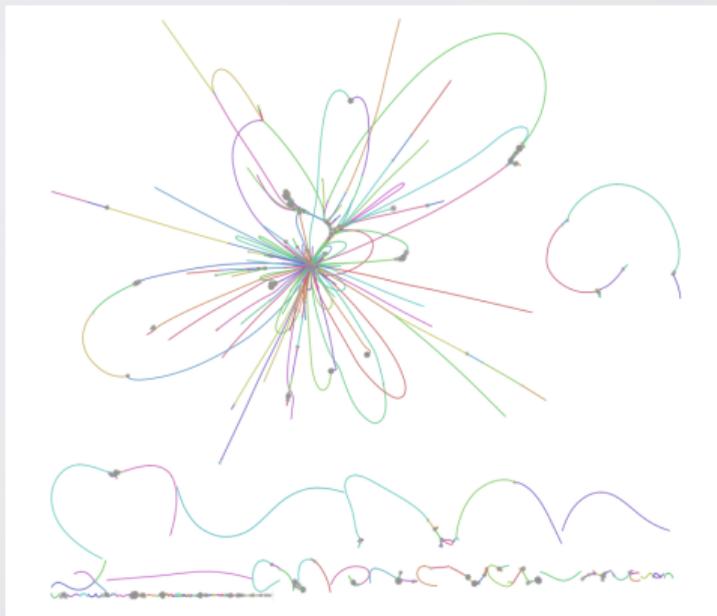
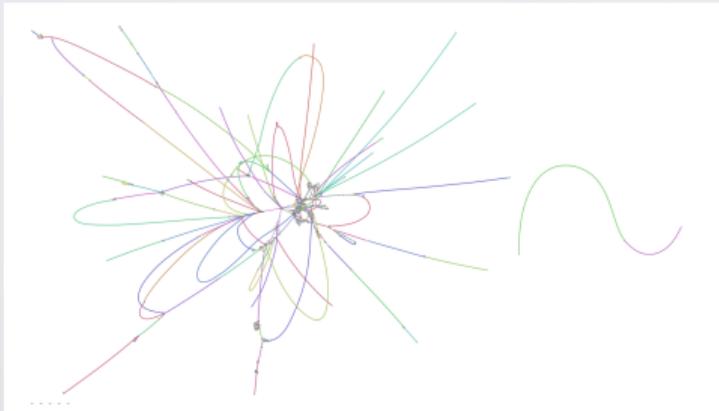
Create/view BLAST search

Query: none

Results: mDBG of CHM13 assembly

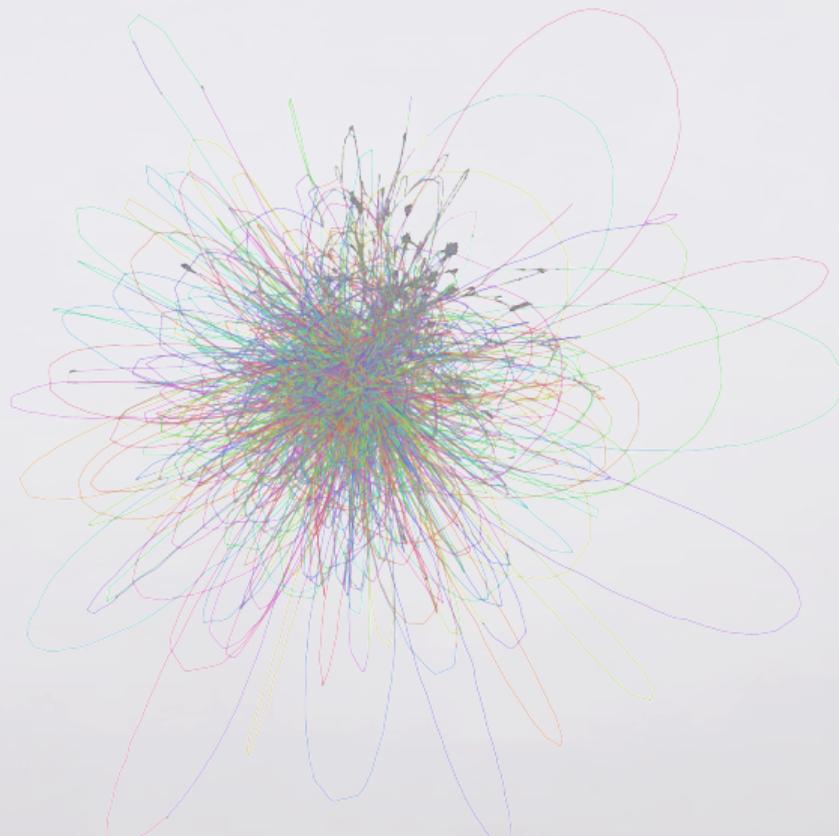


To be compared with the graph of raw reads (right)

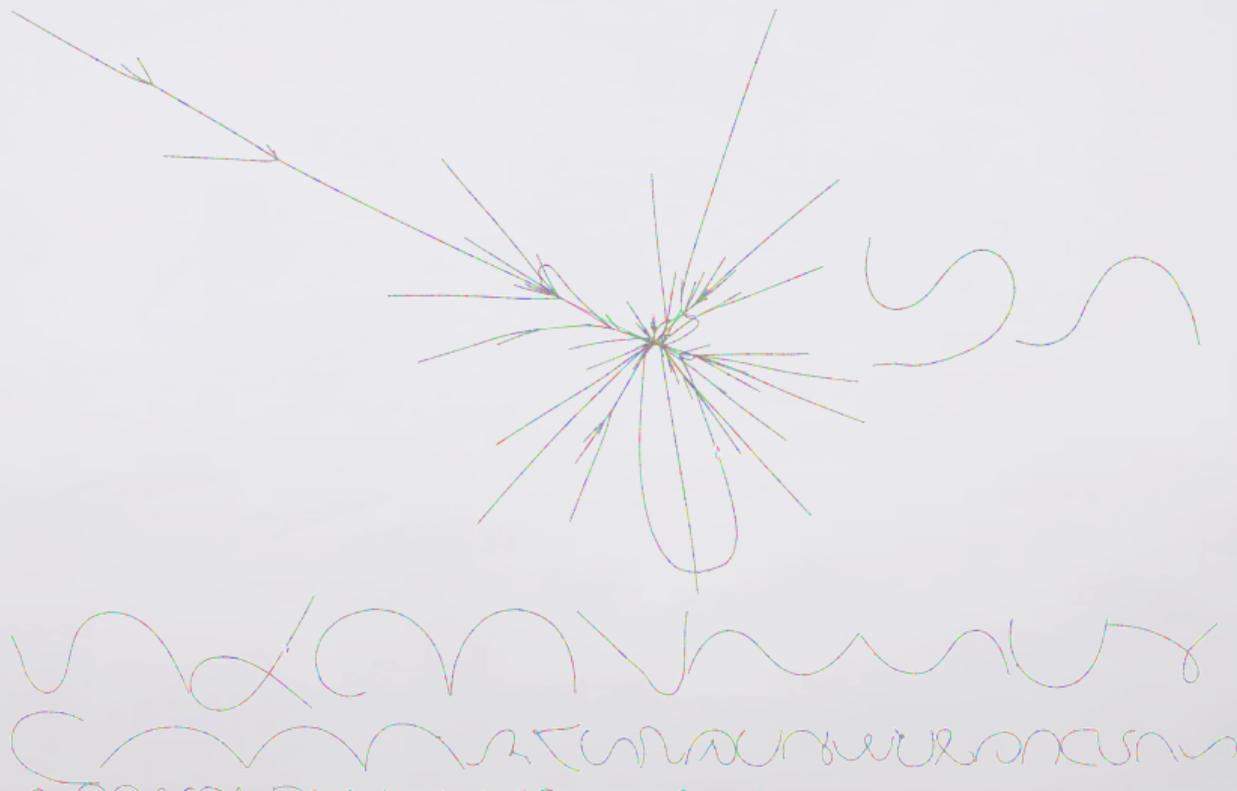


Results: mDBG of CHM13 assembly

Beware of setting parameters “wrong”! ($d=0.001$ here, $d=0.0001$ before)

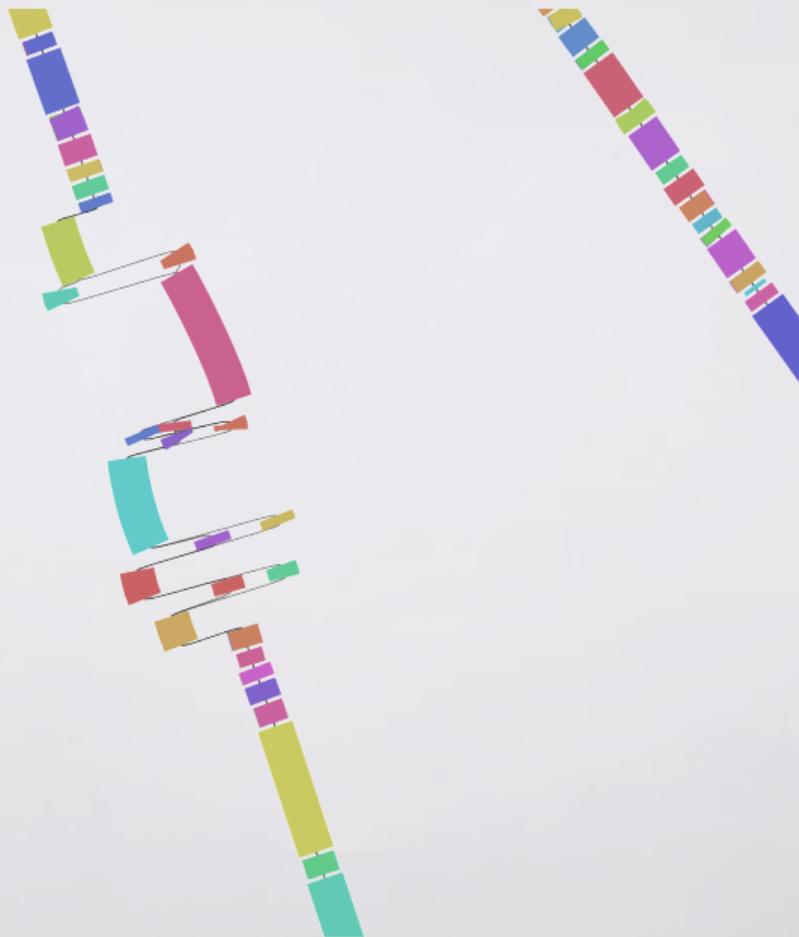


mdBG of another human genome (diploid): HG002 (19k nodes)

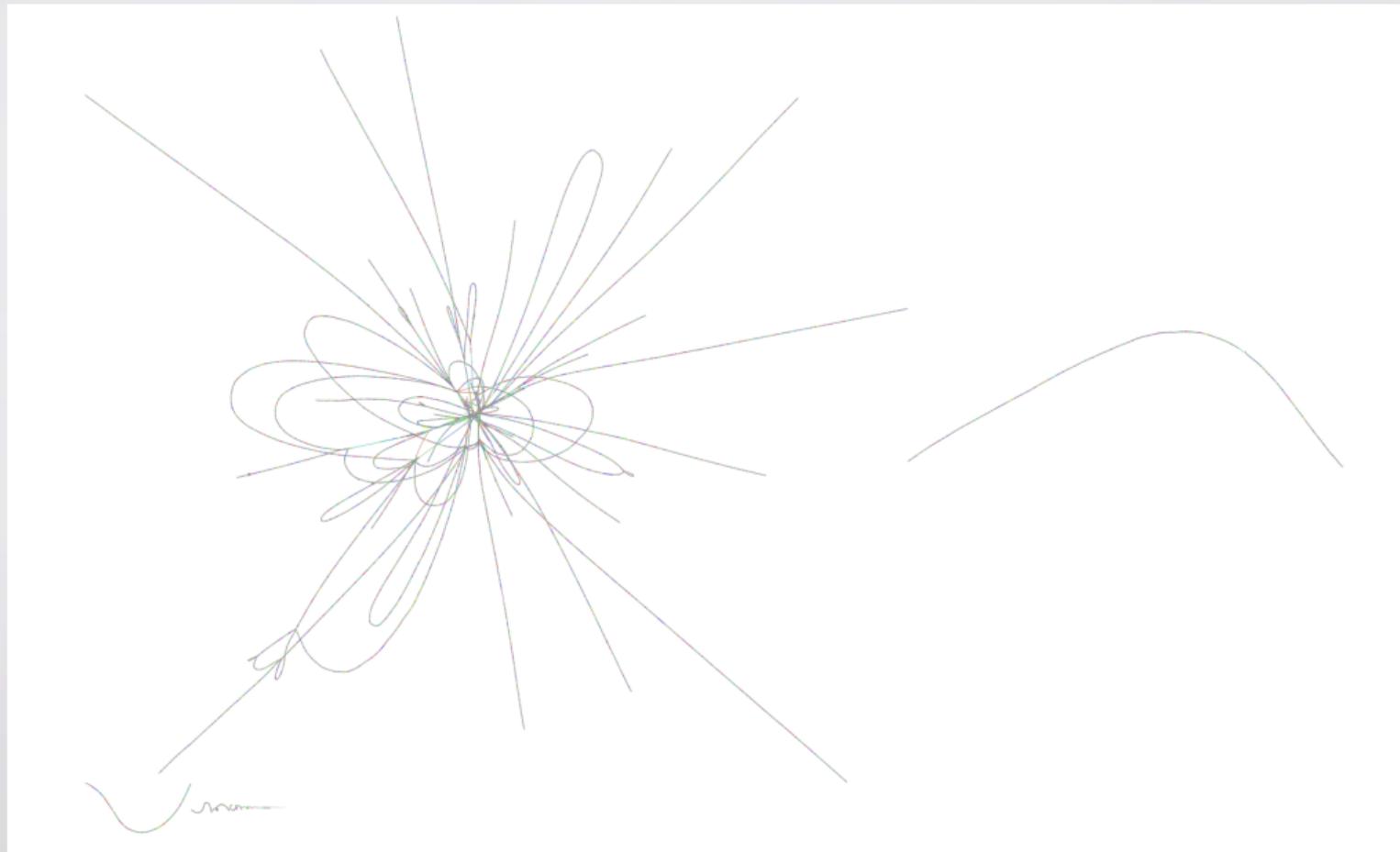


HG002 (zoomed in)

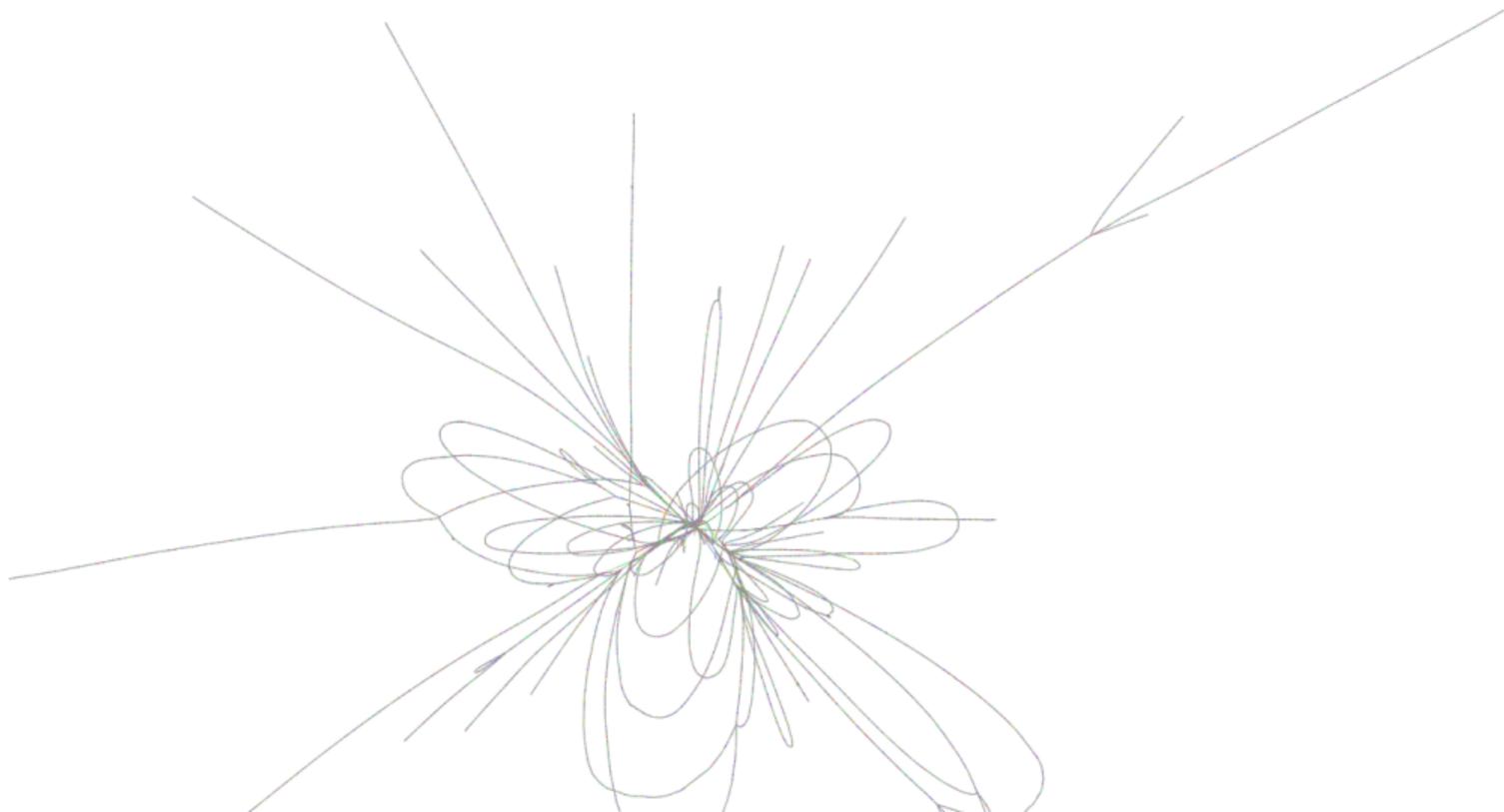
3



mdBG of 2 human genomes: HG002 + HG003 (30k nodes)



mdBG of 5 human genomes: HG002 + HG003 + HG004 +
HG006 + HG007 (63k nodes)



mdBGs of whole genomes are fast to construct

Roughly 2.5 minutes and 6 GB mem per human genome

Compare with:

- \approx 1 hour with Bifrost, Minigraph
- \approx 1 day with Pggp

FAQ on mdBGs

Short reads?

- Doesn't seem applicable

Improving assembly N50?

- Better graph simplifications

Nanopore data?

- 5% error rate is too much, but 1% sounds promising

Some open questions

1. Can one represent all life 31-mers?
2. Can one represent all life 31-mers up to 2 edit mutations?
3. Can one represent all life k-min-mers? ($k=10$, $l=12$, density to be determined)

4. Can one represent all prokaryote+viral 31-mers known to date?
5. Can one represent all human 31-mers known to date?

Conclusion

- **mdBGs** can not only perform genome assembly but also represent pangenome graphs for large collections (661k bacterial genomes) efficiently (10s of GB) at 10-100kbp resolution

Main idea: “ k -mers” over sequences of minimizers “characters”

For pangenomics:

- Higher-resolution mdBGs (1 kbp span for k -min-mers)
- Large graphs (Eukaryotes)
- Automated differential analysis on colored GFAs
- Large structural variant calling on colored GFAs

Thank you! any questions?

mdBG is joint work with Barış Ekim and Bonnie Berger.

Acknowledgements: Yoann Dufresne, Francesco Andreatta.

Funding: H2020 ITN Alpaca, H2020 RISE Pangaia, ANR

Application 1: Long read genome assembly

- Oxford Nanopore, PacBio CLR
 - ▶ 10-1,000 kbp reads, **5-12%** error rate
- PacBio HiFi
 - ▶ 10-25 kbp reads, \leq **1%** error rate



Classical *de Bruijn* graphs not applicable (no long error-free k -mers). Instead:

- Overlap graphs (Canu, miniasm, Shasta, Peregrine, hifiasm, ...)
- Fuzzy dBGs (wtDBG2)
- Sparse dBGs: A-Bruijn or minimizers (Flye, MBG, LJA)

Challenge: Approaches don't scale (high resource usage, slow assembly time)!

Sequencing errors propagate to minimizer-space



Minimizer-space insertion

TACCATAGAC

\overline{m}_2 \overline{m}_4 \overline{m}_3 \overline{m}_2

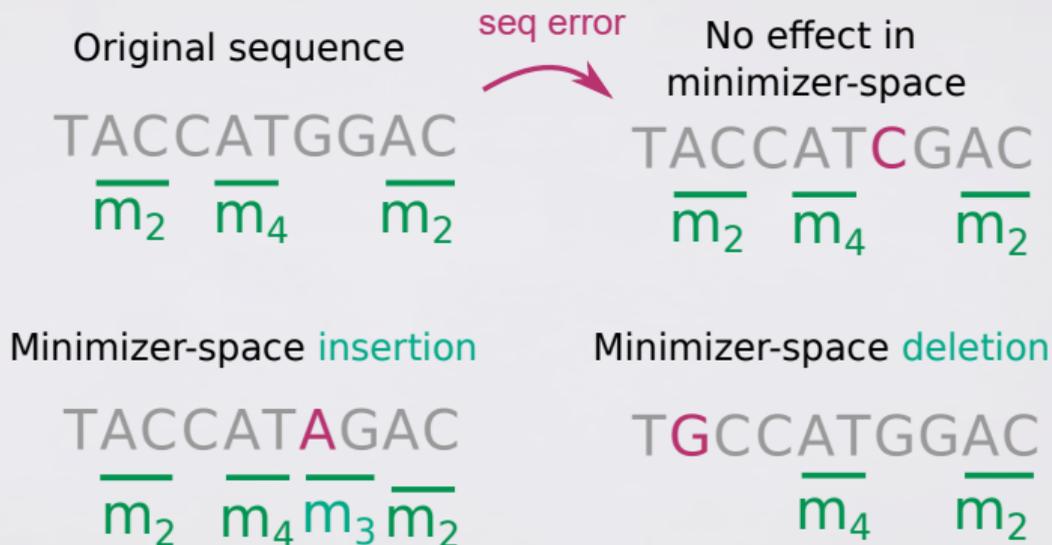
Minimizer-space deletion

TGCCATGGAC

\overline{m}_4 \overline{m}_2

Error rate: base-space \ll minimizer-space,
e.g. 5% in base-space corresponds to 50% in minimizer-space.

Sequencing errors propagate to minimizer-space



Error rate: base-space \ll minimizer-space,
e.g. 5% in base-space corresponds to 50% in minimizer-space.

Error correction: minimizer-space POA (base-space POA: Lee *et al*, '02))

Density minimizers vs syncmers in mDBG

D. mel 100x	Density minimizers	Downsampled syncmers
Best N50	3.9 Mbp	3.8 Mbp
Asm size	111 Mbp	111 Mbp
<i>k</i>	30	25
<i>l</i>	10	10
<i>s</i>	N/A	6
<i>density</i>	0.0035	0.02

Notes:

- Best result out of a coarse parameter grid search
- Both schemes use same hash function
- LCP: in journal (similar)

Results: Metagenome assembly

Zymo D6331 mock metagenome HiFi

Species	Abundance	hifiasm	rust-mdbg
<i>A. muciniphila</i>	1.36%	100.000%	100.000%
<i>B. fragilis</i>	13.13%	99.994%	99.997%
<i>B. adolescentis</i>	1.34%	100.000%	99.730%
<i>C. albican</i>	1.61%	67.832%	39.821%
<i>C. difficile</i>	1.83%	99.996%	99.978%
<i>C. perfringens</i>	0.00%	0.005%	0.005%
<i>E. faecalis</i>	0.00%	0.006%	0.006%
<i>E. coli B1109</i>	8.44%	100.000%	97.918%
<i>E. coli b2207</i>	8.32%	100.000%	98.663%
<i>E. coli B3008</i>	8.25%	100.000%	99.558%
<i>E. coli B766</i>	7.83%	96.913%	96.270%

Species	Abundance	hifiasm	rust-mdbg
<i>E. coli JM109</i>	8.37%	100.000%	97.852%
<i>F. prausnitzii</i>	14.39%	100.000%	100.000%
<i>F. nucleatum</i>	3.78%	100.000%	99.960%
<i>L. fermentum</i>	0.86%	100.000%	100.000%
<i>M. smithii</i>	0.04%	99.840%	87.175%
<i>P. corporis</i>	5.37%	99.561%	99.561%
<i>R. hominis</i>	3.88%	100.000%	100.000%
<i>S. cerevisiae</i>	0.18%	69.522%	39.556%
<i>S. enterica</i>	0.02%	6.232%	4.619%
<i>V. rogosae</i>	11.02%	100.00%	100.000%

	hifiasm	rust-mdbg
Running time	34h29m	55s
Memory usage	83 GB	0.9 GB

Towards bigger and bigger pangenomes..

Community explores many directions:

1. Sketches

- ▶ Mash, sourmash
- ▶ Low-resolution search, graph

2. All nucleotides

- ▶ Approx: BIGSI, HowDeSBT
- ▶ Exact: Cuttlefish 2, MetaGraph, vg, minigraph, ..
- ▶ High-resolution search, graph
- ▶ Expensive to store

3. “In-between”

- ▶ mdBG
- ▶ Low-resolution search, graph
- ▶ Inexpensive to store

Towards bigger and bigger pangenomes..

Community explores many directions:

1. Sketches

- ▶ Mash, sourmash
- ▶ Low-resolution search, graph

2. All nucleotides

- ▶ Approx: BIGSI, HowDeSBT
- ▶ Exact: Cuttlefish 2, MetaGraph, vg, minigraph, ..
- ▶ High-resolution search, graph
- ▶ Expensive to store

3. “In-between”

- ▶ mdBG
- ▶ Low-resolution search, graph
- ▶ Inexpensive to store

4. Gene families

- ▶ What biologists actually do
- ▶ Lowest-resolution, seq search, sometimes graph
- ▶ Inexpensive to store