# Assembly of long, error-pront reads using repeat graphs

Mikhail Kolmogorov, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner

05.07.2021

# Genome assembly in general

▶ reconstruct target sequence from the reads

# Genome assembly in general

- ▶ reconstruct target sequence from the reads
- ▶ different graph structures (De-Bruijn, Overlap-layout, String)
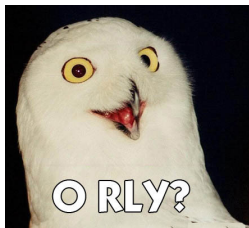
# Genome assembly in general

- ▶ reconstruct target sequence from the reads
- ▶ different graph structures (De-Bruijn, Overlap-layout, String)
- ▶ repeats $\rightarrow$ assembly fragmentation

# Genome assembly in general

- ▶ reconstruct target sequence from the reads
- ▶ different graph structures (De-Bruijn, Overlap-layout, String)
- ▶ repeats $\rightarrow$ assembly fragmentation
- ▶ error rate long read <-> short read

# Genome assembly in general

▶ reconstruct target sequence from the reads

▶ different graph structures (De-Bruijn, Overlap-layout, String)

▶ repeats → assembly fragmentation

▶ error rate long read <-> short read

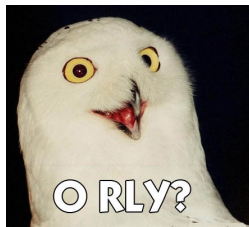▶ Flye should resolve these repeats correctly

# Disjointigs

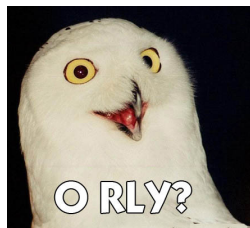▶ most assemblers spent much time on correct contig assembly

# Disjointigs

- ▶ most assemblers spent much time on correct contig assembly
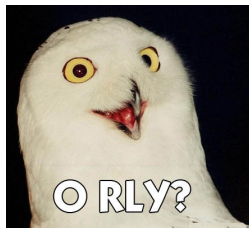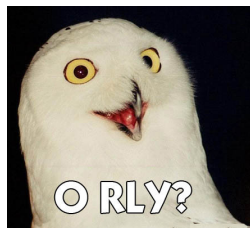- ▶ Flye uses a different approach:

# Disjointigs

- ▶ most assemblers spent much time on correct contig assembly
- ▶ Flye uses a different approach:
- ▶ we don't care (at least at the initial stage)

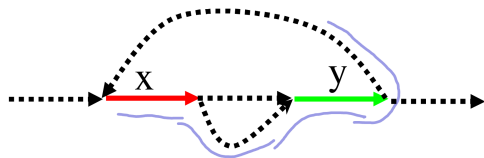# Disjointigs

▶ most assemblers spent much time on correct contig assembly

▶ Flye uses a different approach:
▶ we don't care (at least at the initial stage)
▶ correct assembly graph

# Disjointigs

- most assemblers spent much time on correct contig assembly

- Flye uses a different approach:
- we don't care (at least at the initial stage)
- correct assembly graph
- generate paths from overlapping reads without checking for correct assembly -> disjointigs
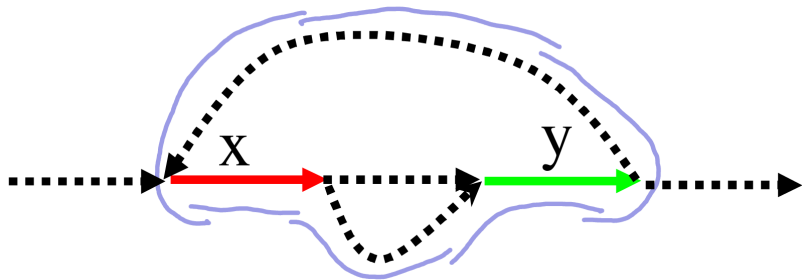
# Disjointigs

- ▶ current assemblers use much time on correct contig assembly
- ▶ Flye uses a different approach:
- ▶ generate paths from overlapping reads without checking for correct assembly -> disjointigs
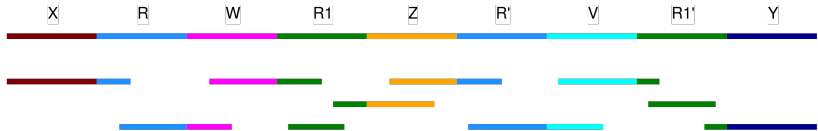
# Disjointigs

- ▶ current assemblers use much time on correct contig assembly

- ▶ Flye uses a different approach:

- ▶ generate paths from overlapping reads without checking for correct assembly -> disjointigs
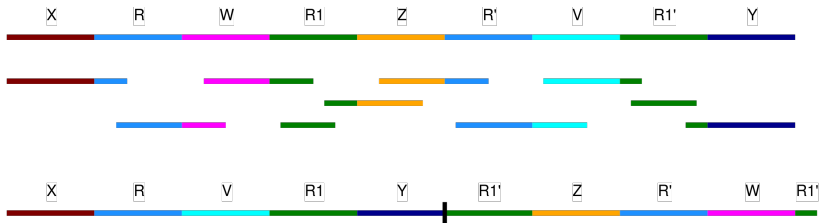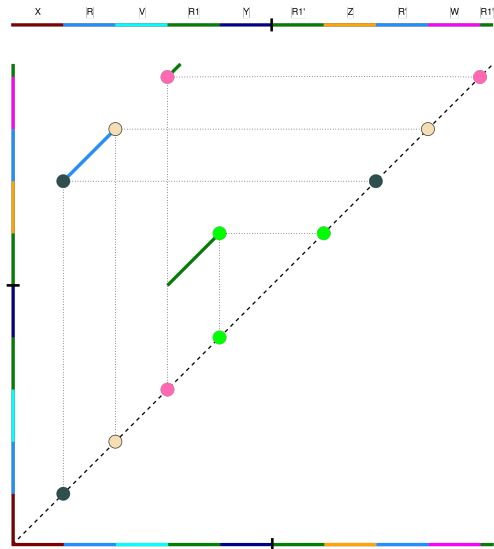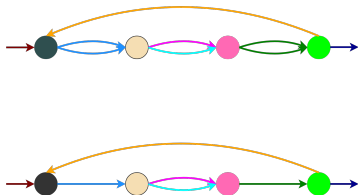
# Repeat graph creation



X    R    W    R1    Z    R'    V    R1'    Y

# Repeat graph creation

# Repeat graph creation

# Repeat graph creation

# Repeat graph creation

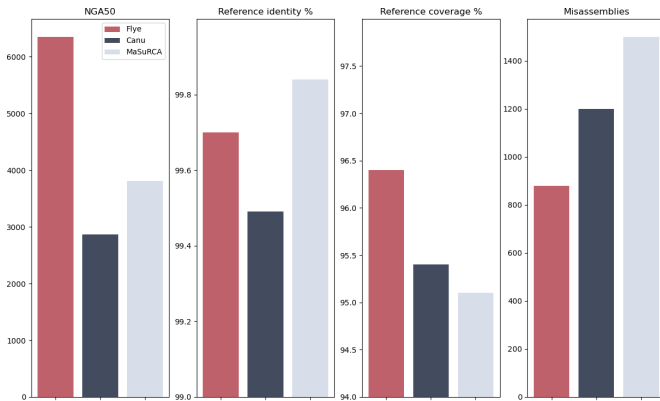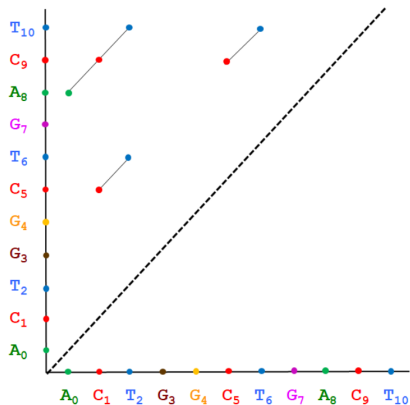# Results

# References

# Appendix

# Dot plot creation

# Repeat graphs

▶ generalization of de bruijn graphs

```
>- from disjointigs = random walk of reads on the repeat gr
>- means the repeat graph hasn't to be known
```

# Repeat graphs

▶ generalization of de bruijn graphs

▶ structure

```
>- from disjointigs = random walk of reads on the repeat gr
>- means the repeat graph hasn't to be known
```

# Repeat graphs

- generalization of de bruijn graphs

- structure

- creation

```
>- from disjointigs = random walk of reads on the repeat gr
>- means the repeat graph hasn't to be known
```

# Difference repeat graph de Bruijn graph

▶ A-Bruijn graph (alignments) generalizes the de Bruijn graph

# Difference repeat graph de Bruijn graph

▶ A-Bruijn graph (alignments) generalizes the de Bruijn graph

▶ We thus argue that the time has come to explain that the breakpoint graphs and the de Bruijn graphs are two identical data structures (if one ignores a cosmetic difference between them) as they both represent specific instances of a general notion of the A-Bruijn graph introduced in [13]. The A-Bruijn graphs are based on representing genomes as sets of labeled paths and further gluing identically labeled edges (breakpoint graphs) or vertices (de Bruijn graphs) in the resulting paths.
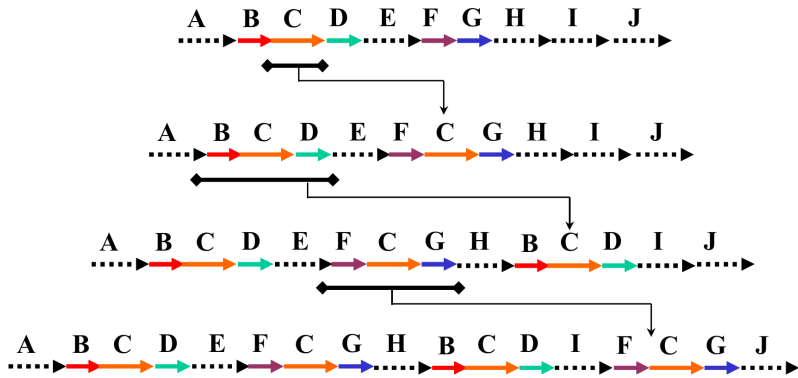
# Difference repeat graph de Bruijn graph

- ▶ A-Bruijn graph (alignments) generalizes the de Bruijn graph

- ▶ We thus argue that the time has come to explain that the breakpoint graphs and the de Bruijn graphs are two identical data structures (if one ignores a cosmetic difference between them) as they both represent specific instances of a general notion of the A-Bruijn graph introduced in [13]. The A-Bruijn graphs are based on representing genomes as sets of labeled paths and further gluing identically labeled edges (breakpoint graphs) or vertices (de Bruijn graphs) in the resulting paths.

- ▶ de Bruijn graphs need correct bases
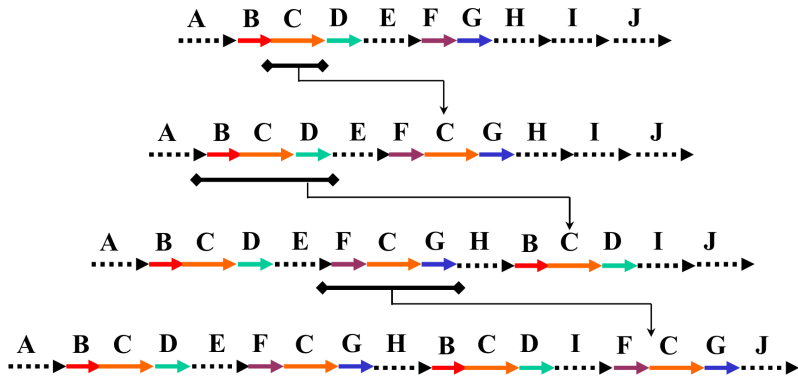
# Difference repeat graph de Bruijn graph

▶ A-Bruijn graph (alignments) generalizes the de Bruijn graph

▶ We thus argue that the time has come to explain that the breakpoint graphs and the de Bruijn graphs are two identical data structures (if one ignores a cosmetic difference between them) as they both represent specific instances of a general notion of the A-Bruijn graph introduced in [13]. The A-Bruijn graphs are based on representing genomes as sets of labeled paths and further gluing identically labeled edges (breakpoint graphs) or vertices (de Bruijn graphs) in the resulting paths.

▶ de Bruijn graphs need correct bases

▶ otherwise tangled graph
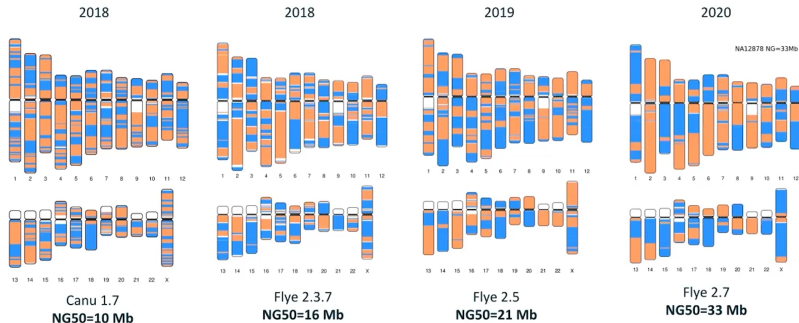
# Segmental duplications



▶ Segmental duplications are duplicated blocks of
genomic DNA typically ranging in size from 1-200 kb
(IHGSC 2001)
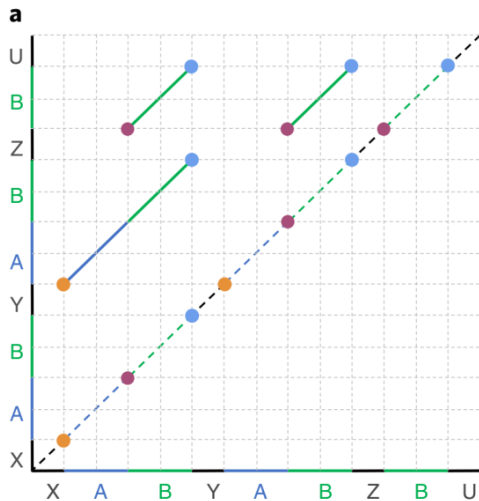
# Segmental duplications



- ▶ Segmental duplications are duplicated blocks of genomic DNA typically ranging in size from 1-200 kb (IHGSC 2001)
- ▶ They often contain sequence features such as high-copy repeats and gene sequences with intron-exon structure
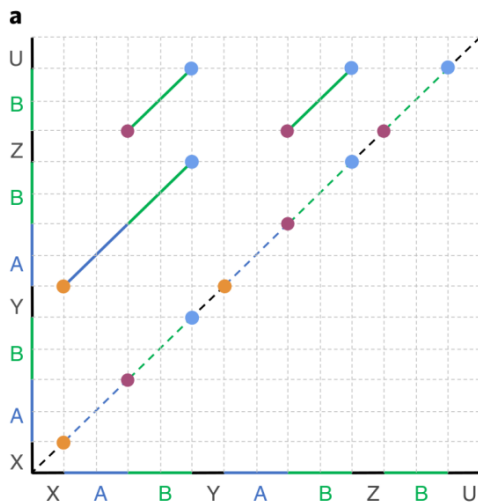
# Contigity improvement



2018 — Canu 1.7, NG50=10 Mb
2018 — Flye 2.3.7, NG50=16 Mb
2019 — Flye 2.5, NG50=21 Mb
2020 — Flye 2.7, NG50=33 Mb (NA12878 NG=33Mb)

▶ colors are contigs

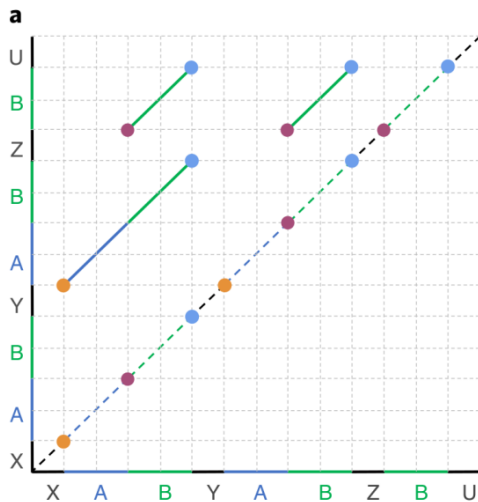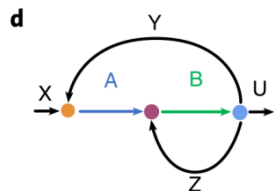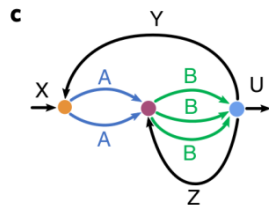▶ color changes -> fragmented
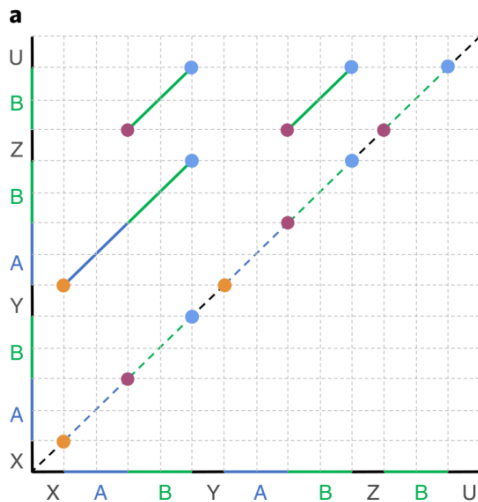
# Repeat graph creation

# Repeat graph creation

# Repeat graph creation

# Repeat graph creation

# Repeat resolution