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

Mikhail Kolmogorov, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner
June 28, 2021

Johannes Hausmann, Luis Kress

- reconstruct target sequence from the reads

# Genome assembly in general

- reconstruct target sequence from the reads
- different graph structures (De-Bruijn, Overlap-layout, String)

- reconstruct target sequence from the reads
- different graph structures (De-Bruijn, Overlap-layout, String)
- repeats → 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 $\leftrightarrow$ short read

## 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 $\leftrightarrow$ short read
- Flye should resolve these repeats correctly

- most assemblers spent much time on correct contig assembly

## Disjointigs

- most assemblers spent much time on correct contig assembly
- Flye uses a different approach [1]:

## Disjointigs

- most assemblers spent much time on correct contig assembly
- Flye uses a different approach [1]:
  - we don't care (at least at the initial stage)

- most assemblers spent much time on correct contig assembly
- Flye uses a different approach [1]:
  - 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 [1]:
  - we don't care (at least at the initial stage)
  - correct assembly graph
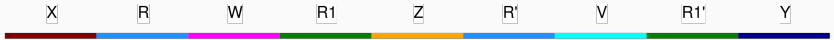- generate paths from overlapping reads without checking for correct assembly → disjointigs

**Figure 1:** Example Genome

**Figure 2:** Example Genome and Reads

**Figure 3:** Example Genome, Reads and Disjointigs

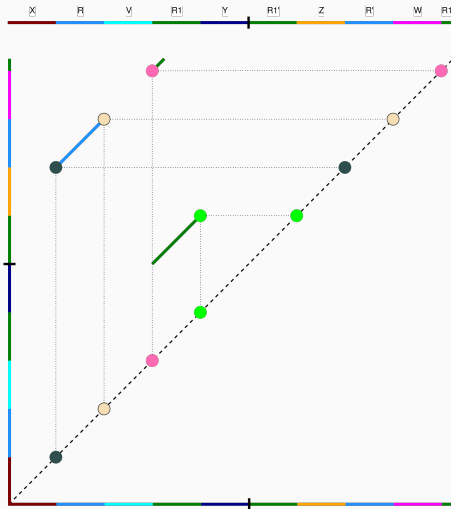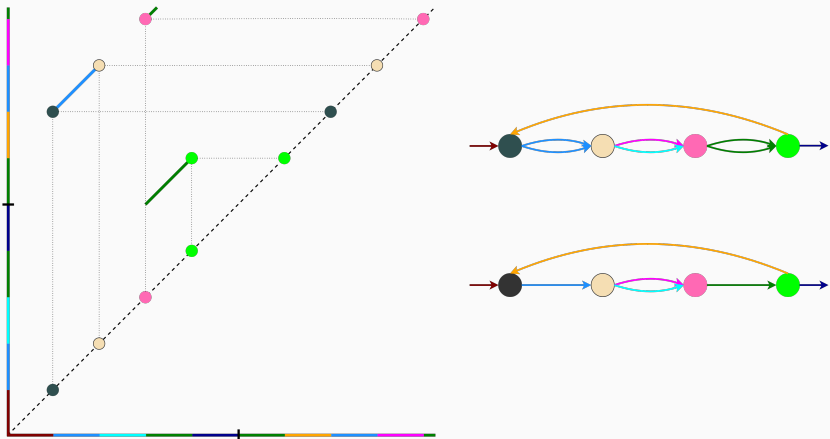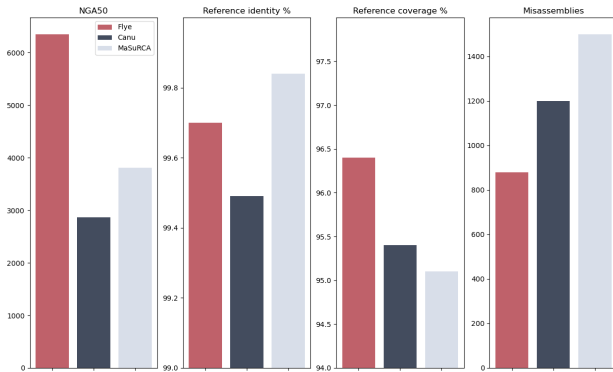**Figure 4:** Breakpoint Graph

**Figure 5:** Repeat Graph

**Figure 6:** Results for HUMAN testset

M. Kolmogorov, J. Yuan, Y. Lin, and P. A. Pevzner.
**Assembly of long, error-prone reads using repeat graphs.**
*Nature Biotechnology*, 37(5):540–546, May 2019.

# Git (presentation)



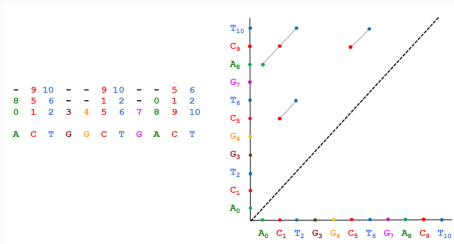**Figure 7:** Link to our git repo

# Appendix

# Dot plot creation



**Figure 8:** Dot plot creation

- generalization of de bruijn graphs

- generalization of de bruijn graphs
- structure

- generalization of de bruijn graphs
- structure
- creation

- generalization of de bruijn graphs

- structure

- creation
  - from disjointigs = random walk of reads on the repeat graph

- generalization of de bruijn graphs

- structure

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

- 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
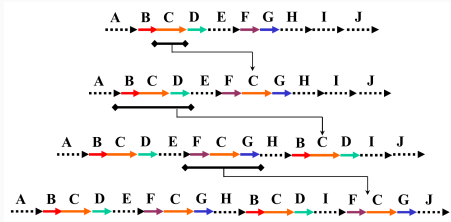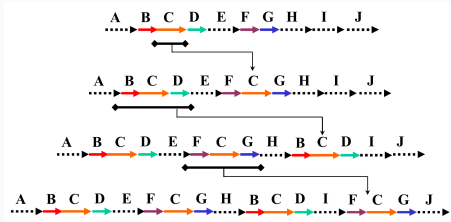- otherwise tangled graph

**Figure 9:** Segmental Duplications

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

**Figure 9:** 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.

**Figure 10:** Contigity improvements
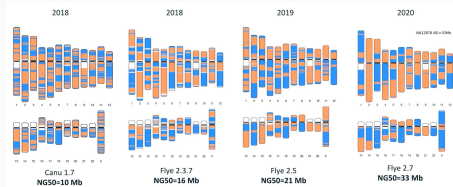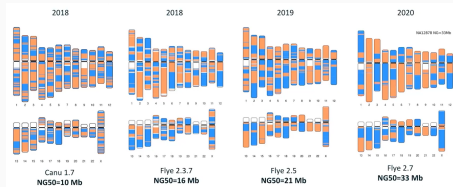
- colors are contigs

**Figure 10:** Contigity improvements

- colors are contigs
- colors are contigs

# Contigity improvement



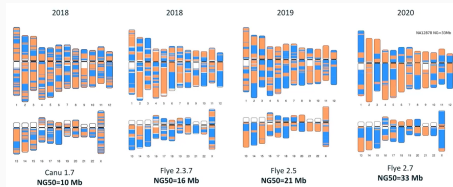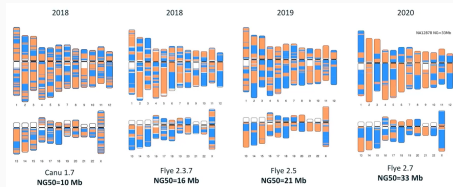**Figure 11:** Contigity improvements

- colors are contigs

**Figure 11:** Contigity improvements

- colors are contigs
- colors are contigs