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

Mikhail Kolmogorov, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner
June 29, 2021

Johannes Hausmann, Luis Kress

- Assembly: reconstruct target sequence from the reads

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

# Background

- Assembly: reconstruct target sequence from the reads
- Different assemblers, different graph structures (De-Bruijn, Overlap-layout, String)
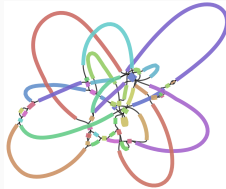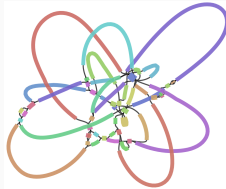- Repeats → assembly fragmentation



**Figure 1:** Tangled assembly graph

- Assembly: reconstruct target sequence from the reads
- Different assemblers, different graph structures (De-Bruijn, Overlap-layout, String)
- Repeats $\rightarrow$ assembly fragmentation



**Figure 1:** Tangled assembly graph

- Error rate long read $\leftrightarrow$ short read

- Assembly: reconstruct target sequence from the reads
- Different assemblers, different graph structures (De-Bruijn, Overlap-layout, String)
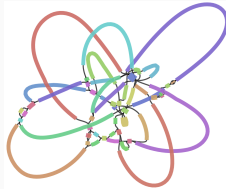- Repeats → assembly fragmentation



**Figure 1:** Tangled assembly graph

- Error rate long read ↔ short read
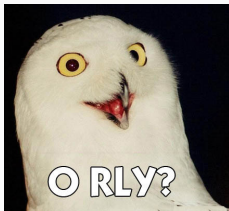- Flye → resolve these repeats correctly

- Most assemblers spend much time on correct contig assembly

# Disjointigs

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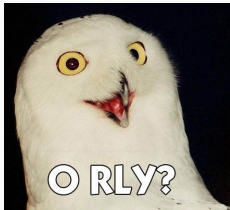- Flye uses a different approach [1]:

# Disjointigs

- Most assemblers spend much time on correct contig assembly
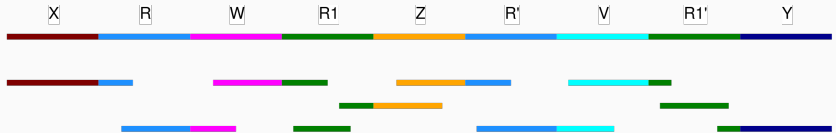- Flye uses a different approach [1]:
  - we don't care

# Disjointigs

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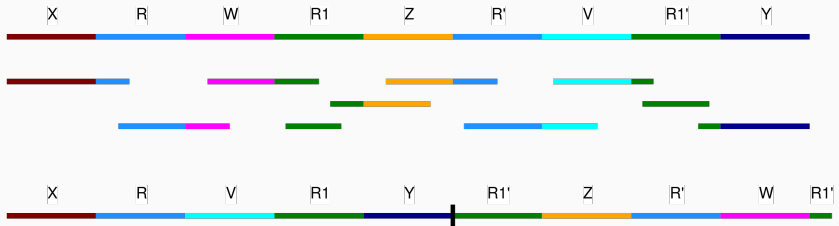  - we don't care (at least at the initial stage)



- Generate paths from overlapping reads without checking for correct repeat resolution → Disjointigs
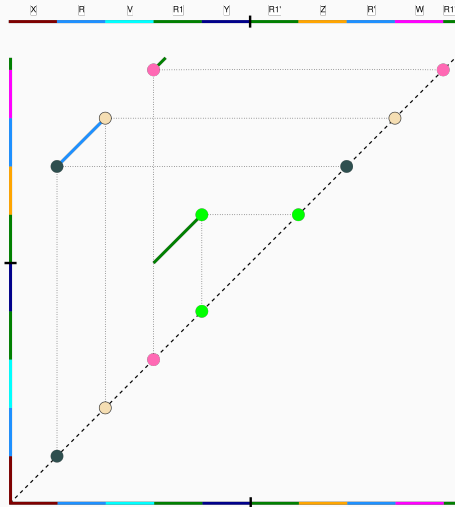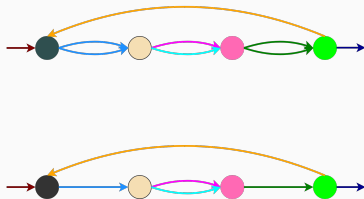
# Repeat Graph Creation

# Repeat Graph Creation
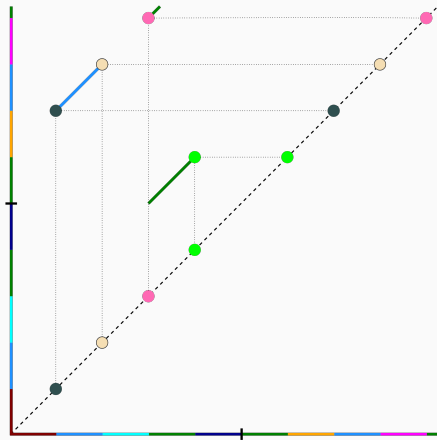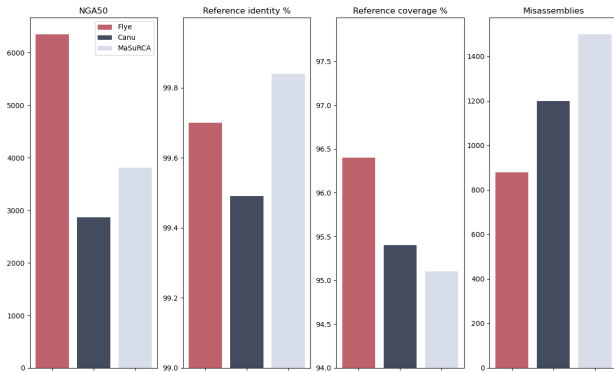
**Figure 2:** Results for HUMAN testset

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.
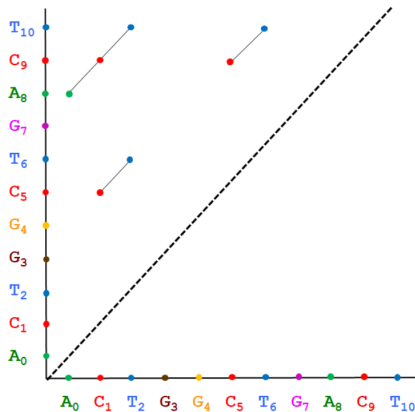
# Appendix

# Dot plot creation



**Figure 3:** Dot plot creation

## Repeat graphs

- 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

## 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
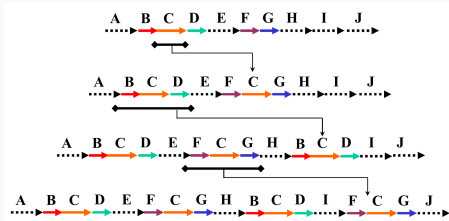- otherwise tangled graph

**Figure 4:** 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.
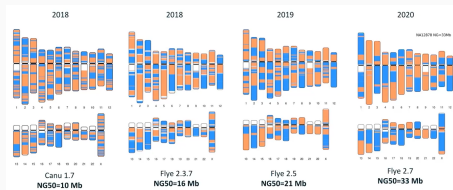
Figure 5: Contigity improvements
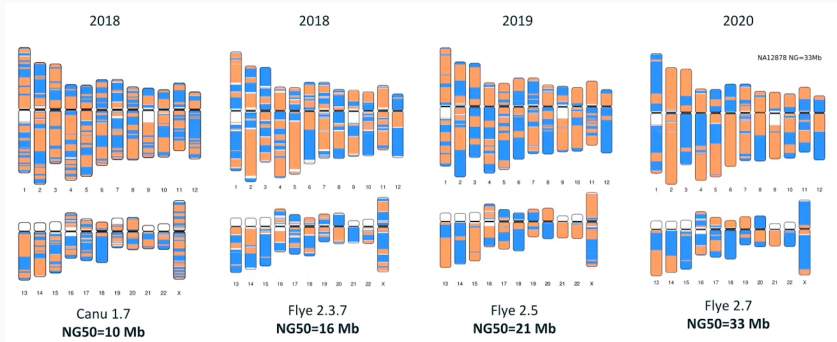
- colors are contigs

**Figure 6:** Contigity improvements

- colors are contigs → change in color means fragmentation