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

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

July 1, 2021

Johannes Hausmann, Luis Kress

# Background
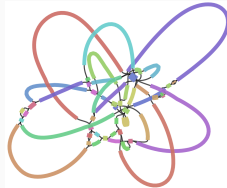
- Assembly: reconstruct target sequence from the reads

- Assembly: reconstruct target sequence from the reads
- Different assemblers, different graph structures

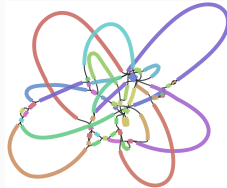- Assembly: reconstruct target sequence from the reads

- Different assemblers, different graph structures

- Repeats $\rightarrow$ assembly fragmentation



**Figure 1:** Tangled assembly graph[1]

# Background

- Assembly: reconstruct target sequence from the reads
- Different assemblers, different graph structures
- Repeats $\rightarrow$ assembly fragmentation



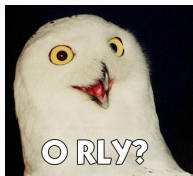**Figure 1:** Tangled assembly graph[1]

- Small differences between repeat copies $\rightarrow$ hard to resolve with error-prone reads

- Most assemblers spend much time on correct contig assembly

## Disjointigs

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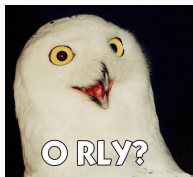- Flye uses a different approach [1]:

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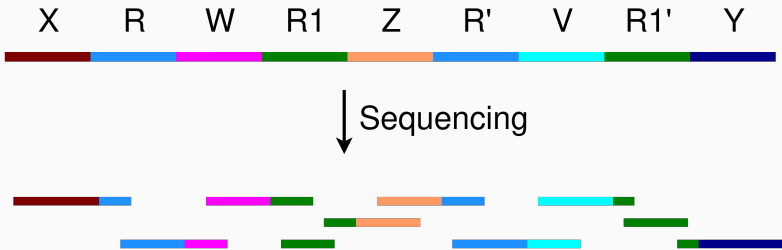  - we don't care

# Disjointigs

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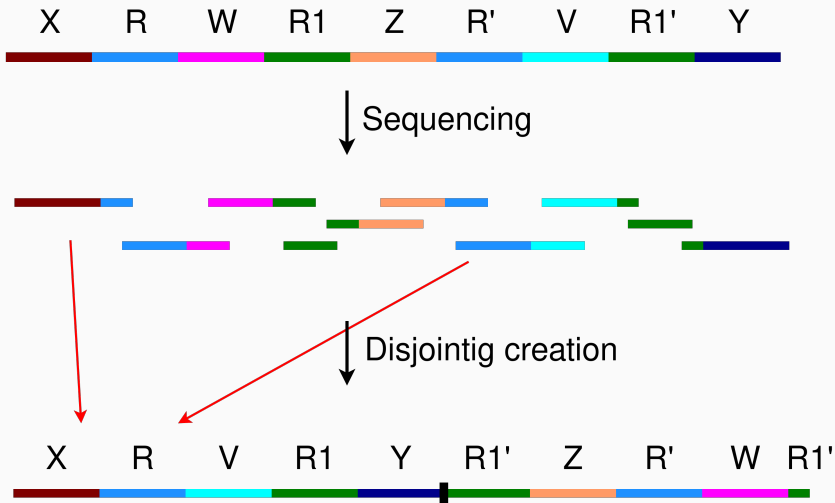  - we don't care (at least at the initial stage)



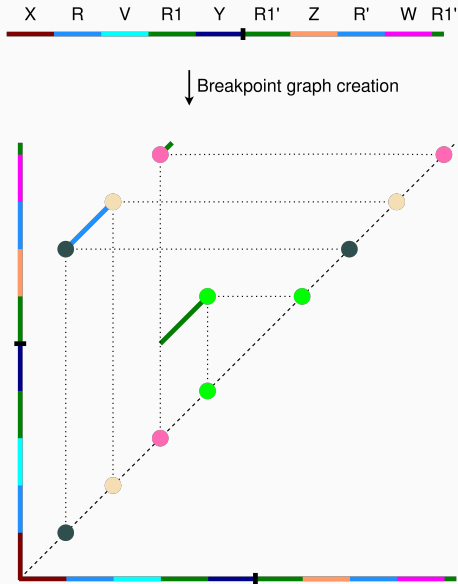- Generate paths from overlapping reads without checking for correct repeat resolution → Disjointigs
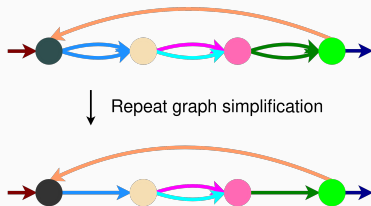
2

# Repeat Graph Creation



X  R  W  R1  Z  R'  V  R1'  Y

Sequencing

3

# Repeat Graph Creation

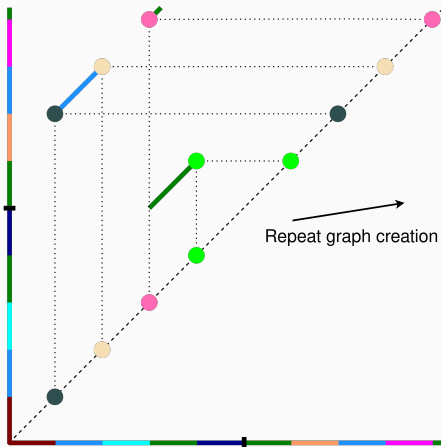X   R   V   R1   Y   R1'   Z   R'   W   R1'

Breakpoint graph creation

Repeat graph creation

Repeat graph simplification
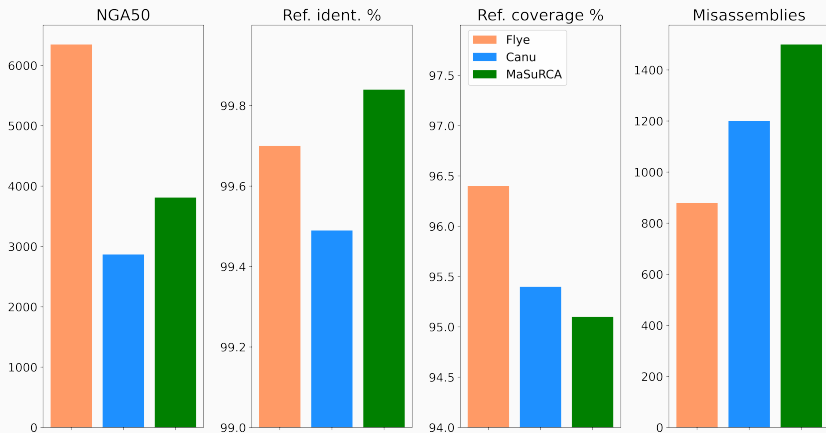
# Results



**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.

Y. Lin, S. Nurk, and P. A. Pevzner.
**What is the difference between the breakpoint graph and the de Bruijn graph?**
*BMC genomics*, 15 Suppl 6:S6, 2014.

P. A. Pevzner, P. A. Pevzner, H. Tang, and G. Tesler.
**De novo repeat classification and fragment assembly.**
*Genome Research*, 14(9):1786–1796, Sept. 2004.

B. J. Walker, T. Abeel, T. Shea, M. Priest, A. Abouelliel, S. Sakthikumar, C. A. Cuomo, Q. Zeng, J. Wortman, S. K. Young, and A. M. Earl.
**Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement.**
*PloS One*, 9(11):e112963, 2014.

GOETHE
UNIVERSITÄT
FRANKFURT AM MAIN

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

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

July 1, 2021
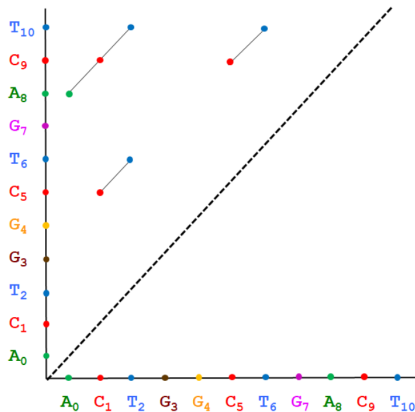
Johannes Hausmann, Luis Kress

# Appendix

# Dot plot creation



Figure 3: Dot plot creation
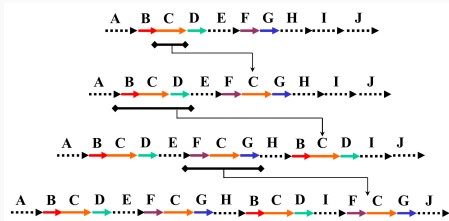
- generalization of de bruijn graphs
- structure
- creation
  - from disjointigs = random walk of reads on the repeat graph
  - means the repeat graph hasn't to be known

## 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 [2]. 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." [3]
- de Bruijn graphs need correct bases
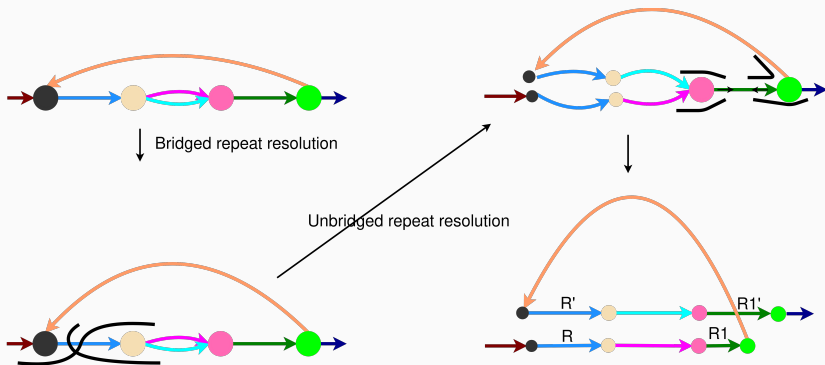- otherwise tangled graph

# Segmental duplications



**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.
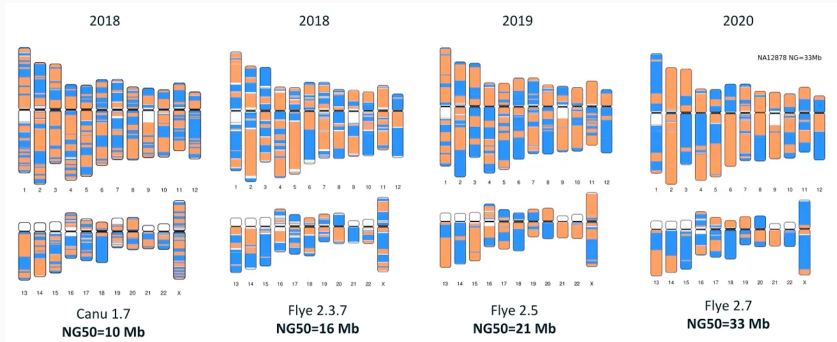
**Figure 5:** Repeat resolution of the given example

## Contigity improvement



**Figure 6:** Contigity improvements

- colors are contigs → change in color means fragmentation