

GOLDRUSH: A *de novo* assembler for long reads with linear time complexity

Johnathan Wong, Vladimir Nikolic, Lauren Coombe, Emily Zhang, René L. Warren, and Inanc Birol

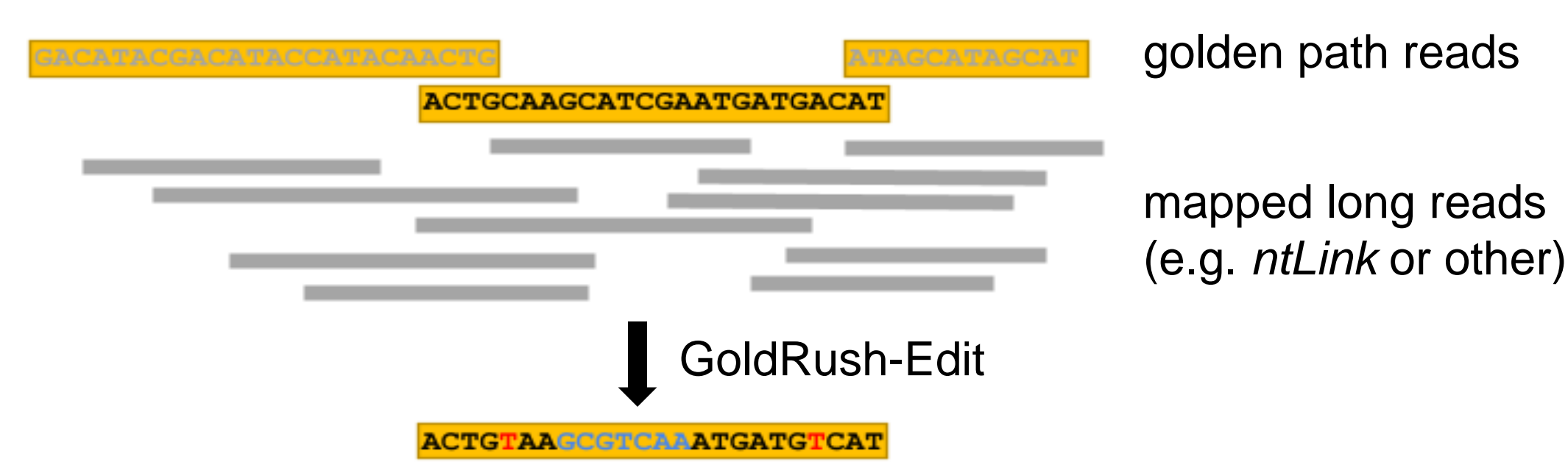
Introduction

State-of-the-art *de novo* long read assemblers

- Overlap-Layout-Consensus paradigm
- High RAM cost

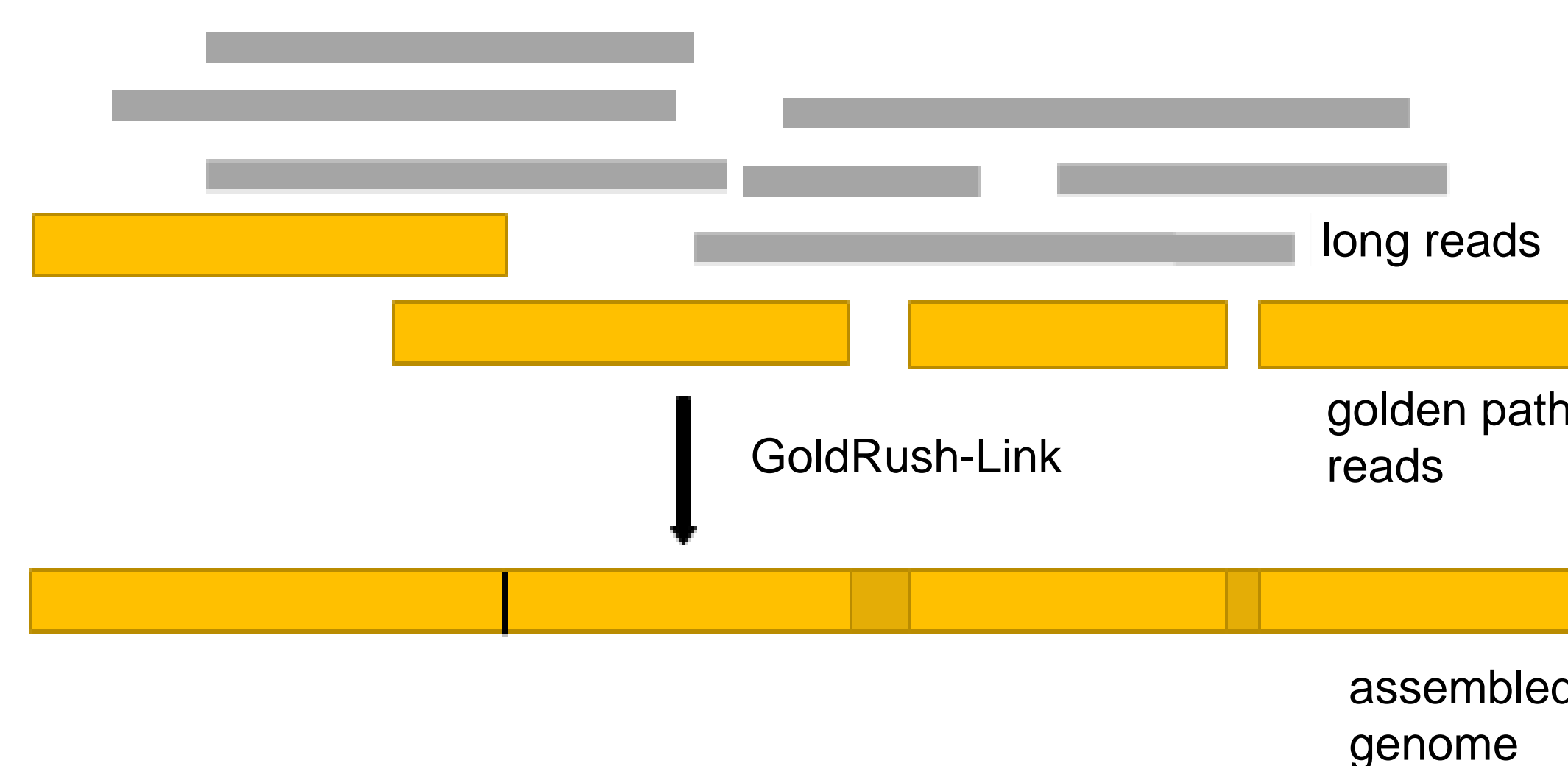
GoldRush

- Linear time algorithm
- Low RAM cost



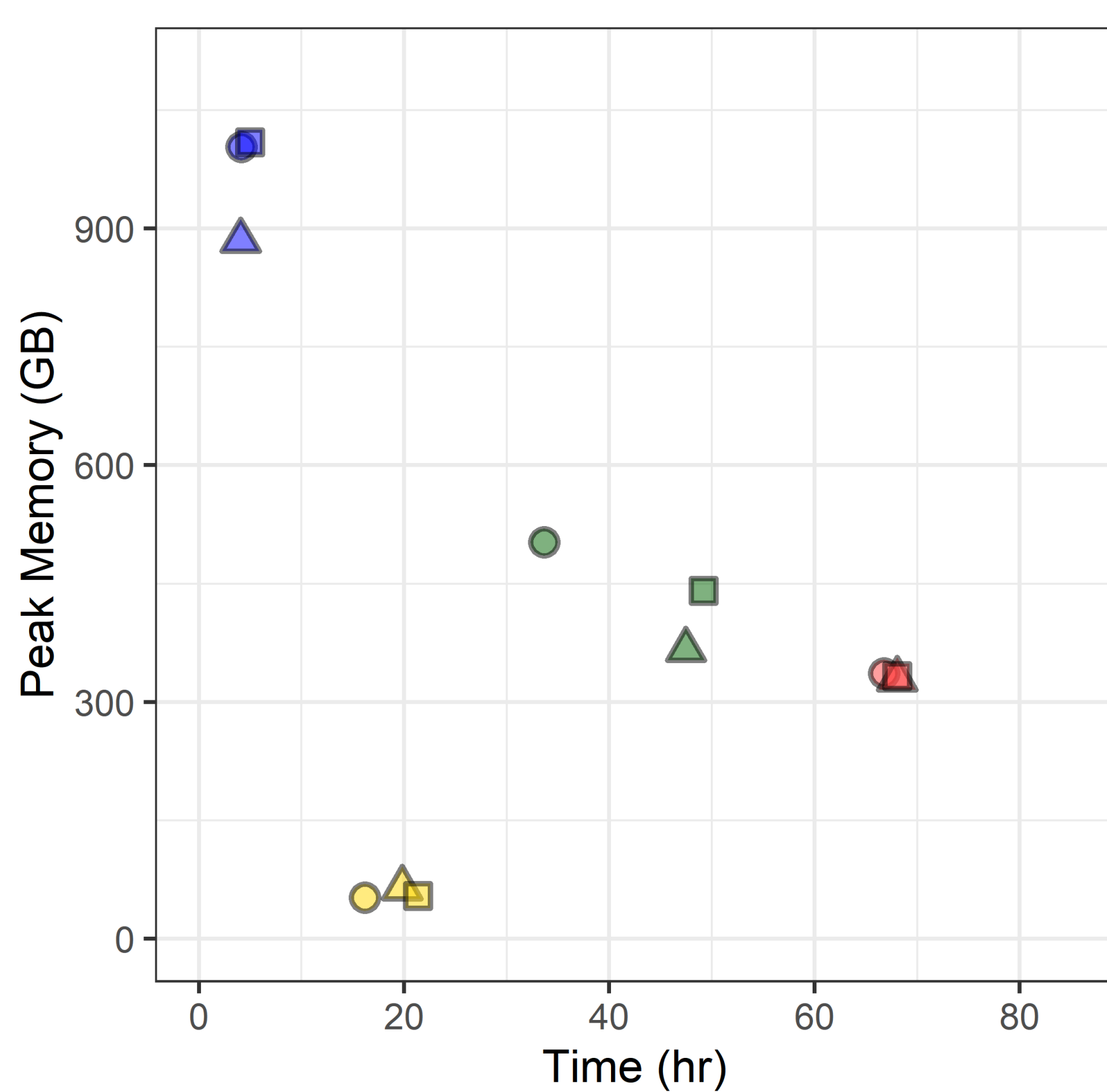
Detects and removes misassemblies

- Identify and cut chimeric golden path reads

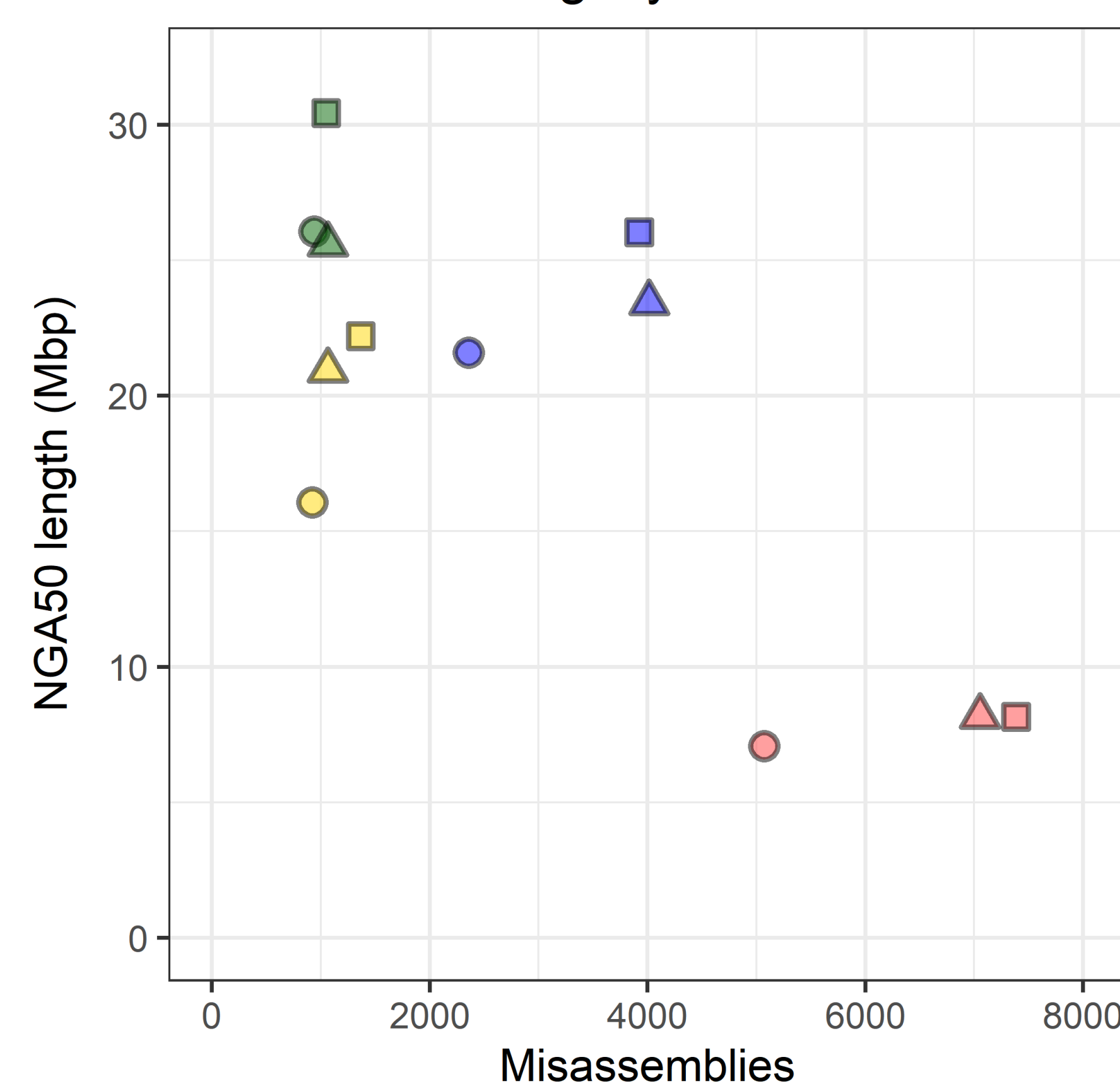


Results

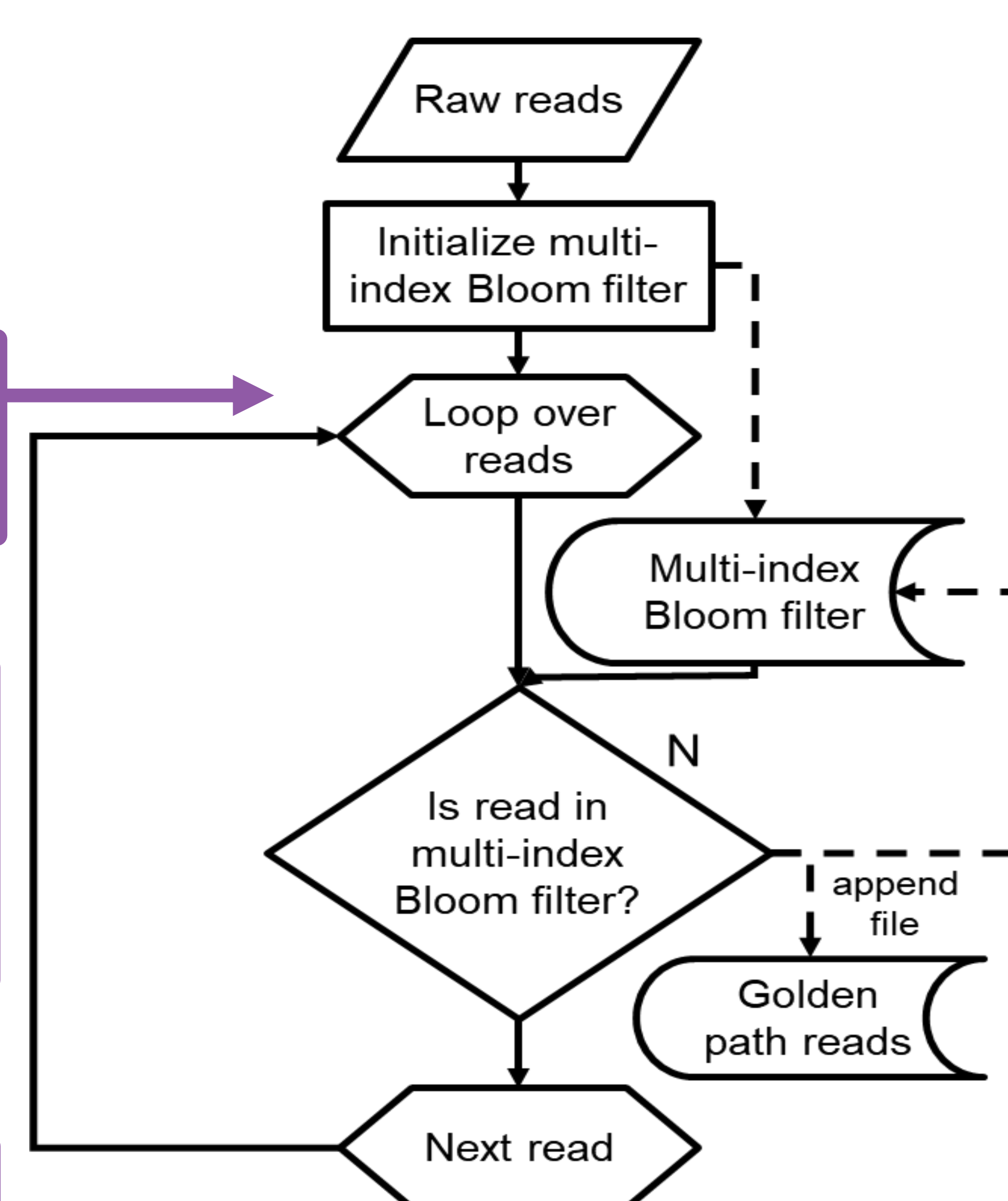
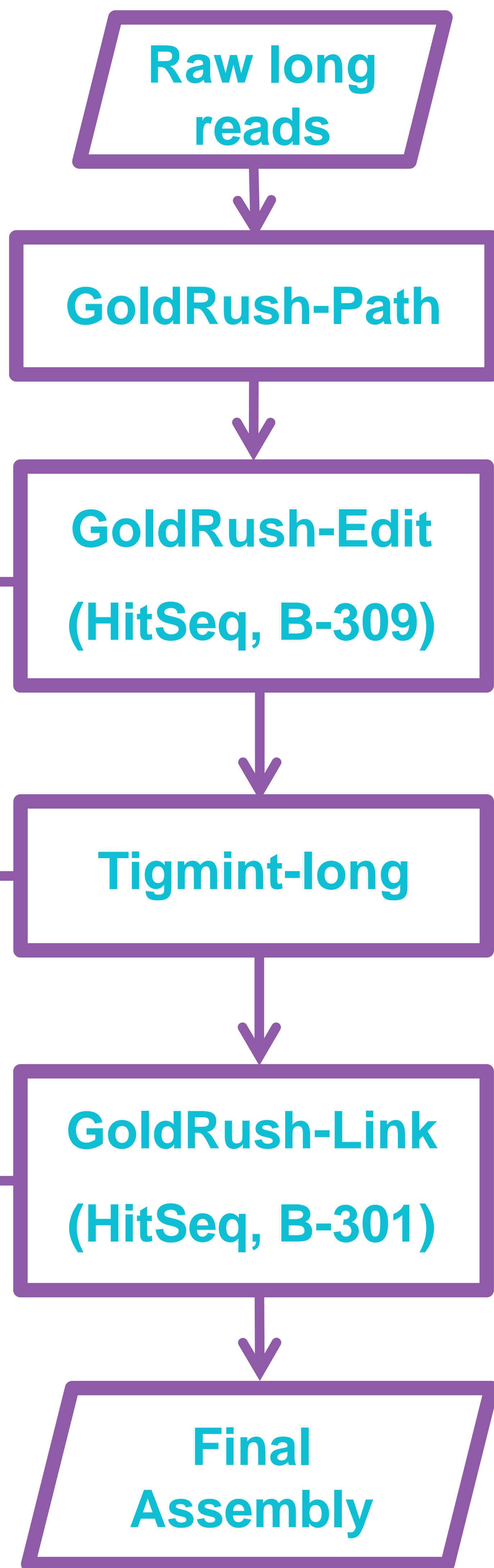
Performance Metrics



Contiguity Metrics



Human Individual
 △ HG01243 (63X, 9% error rate)
 □ HG02055 (71X, 11% error rate)
 ○ NA24385 (67X, 4% error rate)
 Assembler
 ● Flye
 ● GoldRush
 ● Redbean
 ● Shasta



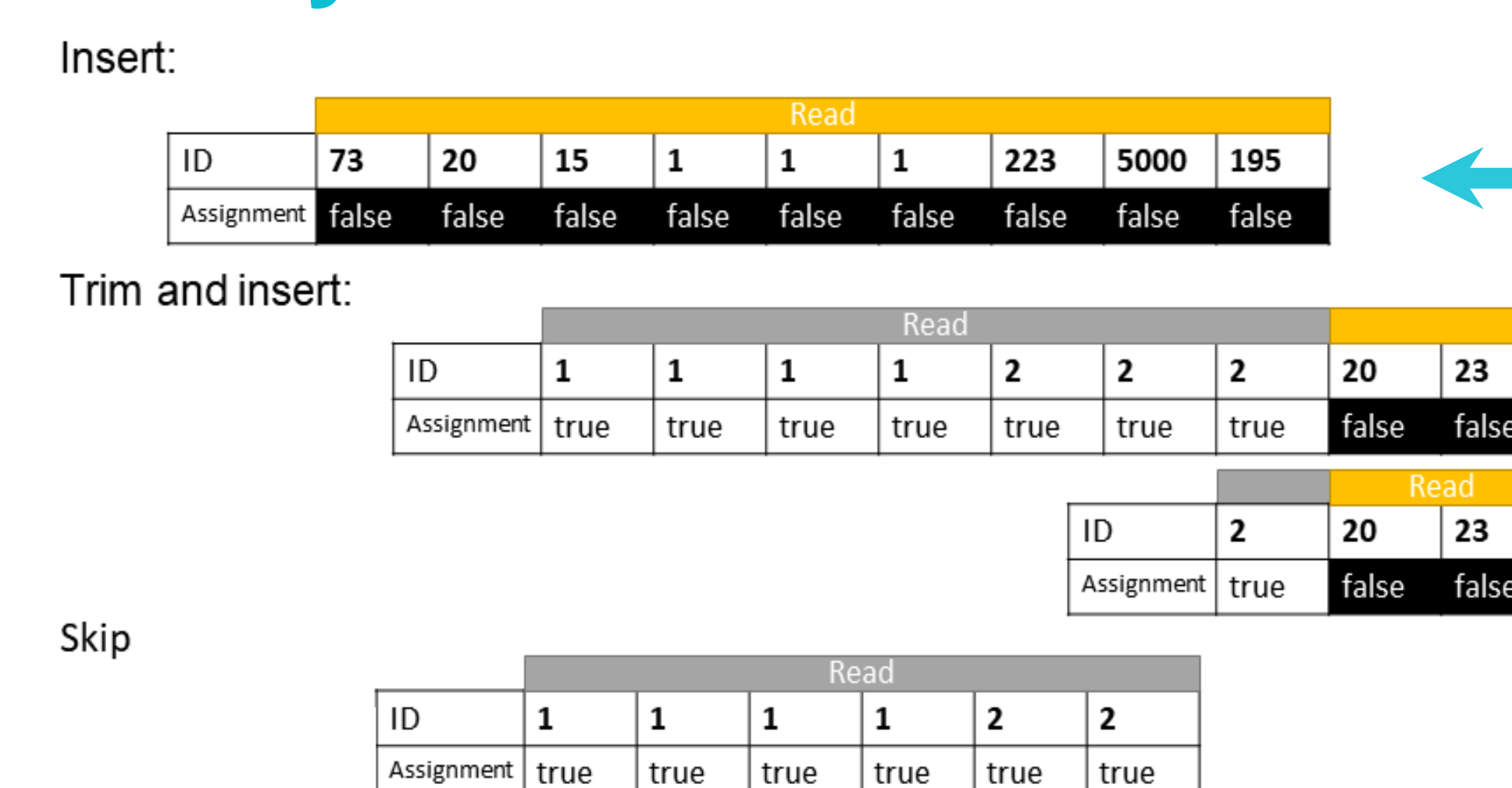
Spaced Seed

Kmer 1	CTAA C ACACG
Kmer 2	CTAA G TACACG
Spaced Seed	1111001111

Insertion

- Read is divided into blocks of length $t \times b$
 - t is the length of a tile
 - b is the number of tiles in a block
- Spaced seed inside a block are inserted with the same ID

Query Outcome



Conclusions

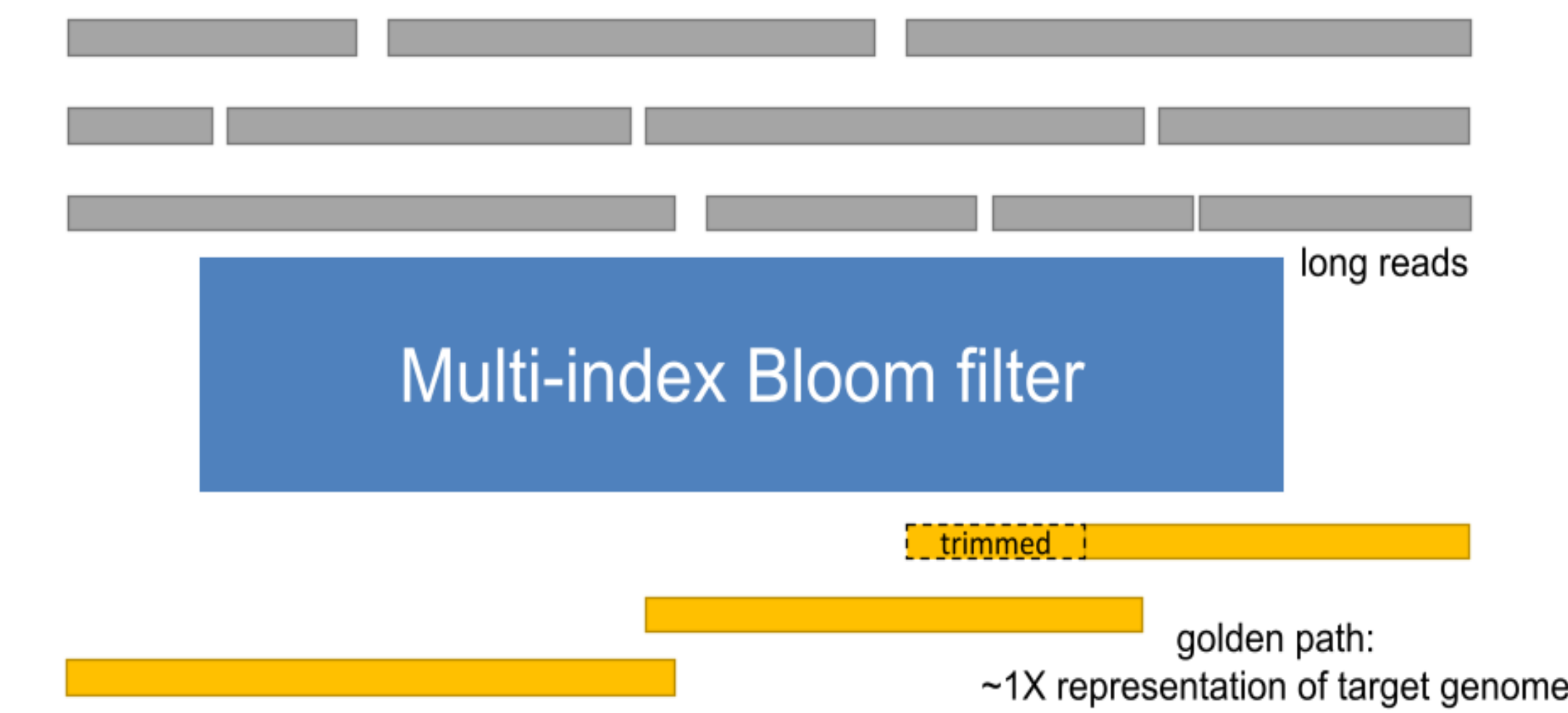
- Novel *de novo* long read assembly algorithm
- Lowest RAM cost
- High contiguity and low misassemblies

Software Availability

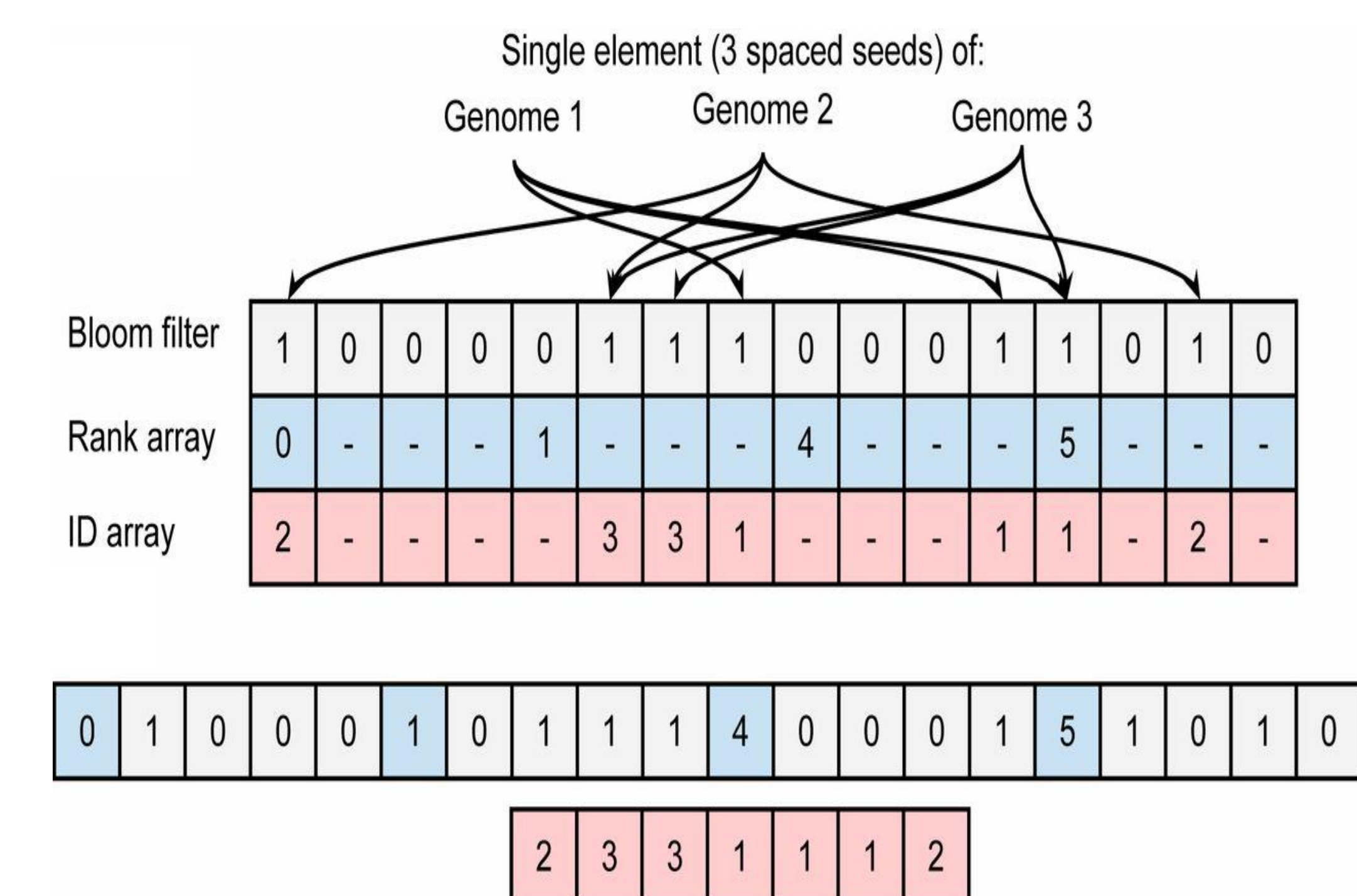
<https://github.com/bcgsc/goldrush>



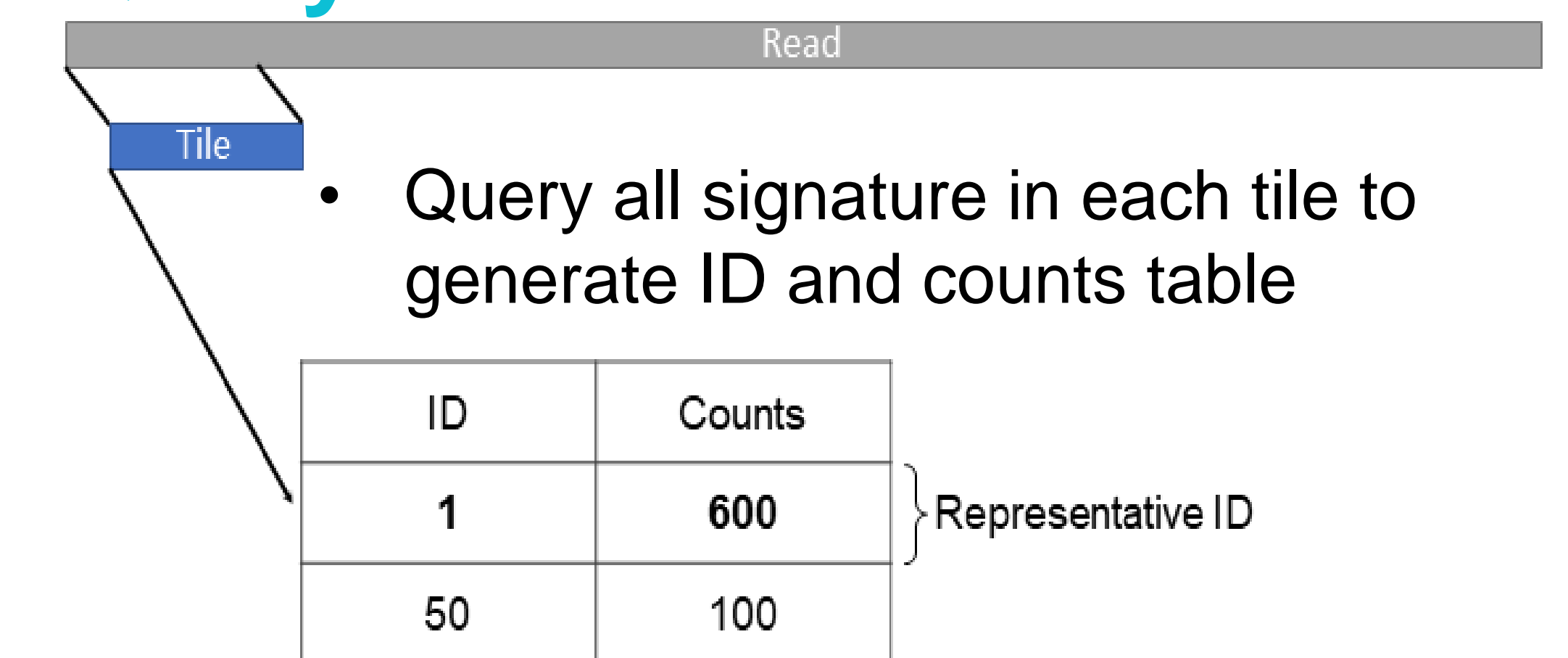
Golden Path



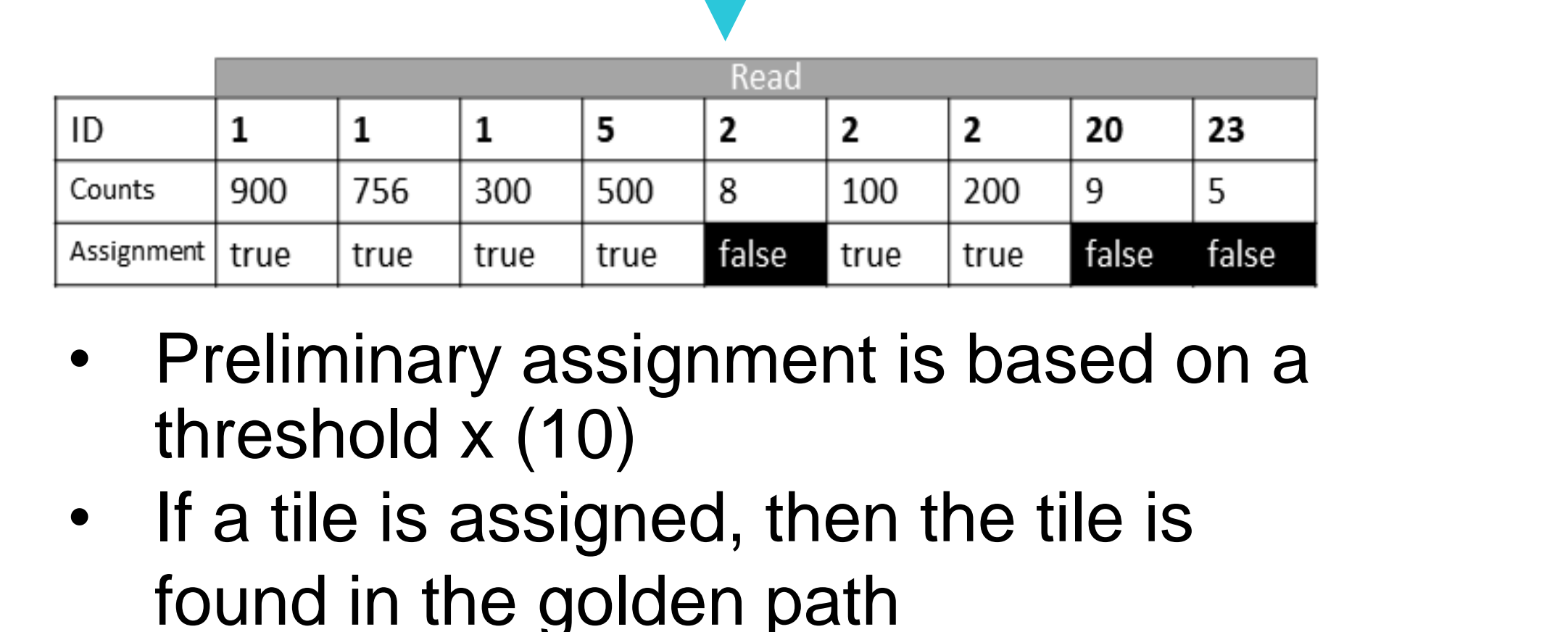
Multi-index Bloom filter



Query



Best hits



References

Coombe L, et al. 2021. LongStitch: high-quality genome assembly correction and scaffolding using long reads. *BMC Bioinformatics* 22:2021.06.17.448848. DOI: 10.1186/s12859-021-04451-7.
 Chu J, et al. 2020. Mismatch-tolerant, alignment-free sequence classification using multiple spaced seeds and multiindex Bloom filters. *Proceedings of the National Academy of Sciences* 117:2020.07.08. DOI: 10.1073/pnas.1903436117.

Contact

jowong@bcgsc.ca

Funding

