

GoldRush-Edit : A targeted, alignment-free polishing & finishing pipeline for long read assembly, using long read k-mers

Background

- The ability to generate accurate genome sequences is cornerstone to many life sciences research projects and translational applications in clinical genomics. This is especially important today as genomics projects utilizing single molecule sequencing (SMS)/long reads exclusively, such as those from Oxford Nanopore Technologies PLC (ONT, Oxford, UK), are increasing despite the still appreciable error rates they afford.
- Most popular polishing methods for long reads, e.g. Racon [1], rely on sequence alignments. The current implementations are not scalable for large (>3Gbp) genomes, requiring large memory servers.
- We present GoldRush-Edit, a memory-efficient polishing pipeline to correct base errors in long read assemblies, using a scalable and targeted k-mer-based method.

Conclusions

- GoldRush-Edit polisher is capable of outputting draft assembly quality comparable to its competitor Racon, but at an order of magnitude less memory usage and similar run time.
- It is possible to use alignments for polishing, which slightly increases the run time, but improves the polishing quality.

Funding



National Institutes of Health

GenomeCanada



- Robert Vaser et al. "Fast and accurate de novo genome assembly from long uncorrected reads". In: Genome Research 27.5 (Jan. 2017), pp. 737–746. DOI: 10.1101/gr.214270.116. URL: https://doi.org/10.1101/gr.214270.116.
- René L Warren et al. "ntEdit: scalable genome sequence polishing". In: Bioinformatics 35.21 (May 2019). Ed. by Bonnie Berger, pp. 4430-4432. DOI: 10.1093/bioinformatics/ btz400. URL: https://doi.org/10.1093/bioinformatics/btz400.
- Daniel Paulino et al. "Sealer: a scalable gap-closing application for finishing draft genomes". In: BMC Bioinformatics 16.1 (July 2015). DOI: 10.1186/s12859-015-0663-4. URL: https: //doi.org/10.1186/s12859-015-0663-4.

regions it cannot resolve. Sealer [3] is then used to polish marked regions.

Vladimir Nikolić¹ Lauren Coombe¹ Johnathan Wong¹ Janet Li^{1,2} René Warren¹ Inanç Birol¹

¹Canada's Michael Smith Genome Sciences Centre at BC Cancer, Vancouver, BC V5Z 4S6, Canada ²Bioinformatics Graduate Program, The University of British Columbia, Vancouver, BC V6T 1Z4, Canada





Figure 2. GoldRush-Edit spawns parallel ntEdit+Sealer pipelines and uses a dedicated Bloom filter builder process to maximize concurrency.

github.com/bcgsc/goldrush-edit

Figure 3. GoldRush-Edit by default runs ntLink to obtain alignment-free read-to-contig mappings. It can also run minimap2 to obtain mappings, which comes at a cost to run time and memory usage, but improves polishing quality. GoldRush-Edit uses over an order of magnitude less memory than Racon and has a comparable run time, while coming close in terms of polishing quality.



vnikolic@bcgsc.ca