


Genome analysis

GAVISUNK: genome assembly validation via inter-SUNK distances in Oxford Nanopore reads

Philip C. Dishuck ¹, Allison N. Rozanski¹, Glennis A. Logsdon¹, David Porubsky¹ and Evan E. Eichler^{1,2,*}

¹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA 98195, USA and ²Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195, USA

*To whom correspondence should be addressed.

Associate Editor: Pier Luigi Martelli

Received on July 11, 2022; revised on September 26, 2022; editorial decision on October 19, 2022; accepted on October 31, 2022

Abstract

Motivation: Highly contiguous *de novo* phased diploid genome assemblies are now feasible for large numbers of species and individuals. Methods are needed to validate assembly accuracy and detect misassemblies with orthologous sequencing data to allow for confident downstream analyses.

Results: We developed GAVISUNK, an open-source pipeline that detects misassemblies and produces a set of reliable regions genome-wide by assessing concordance of distances between unique *k*-mers in Pacific Biosciences high-fidelity assemblies and raw Oxford Nanopore Technologies reads.

Availability and implementation: GAVISUNK is available at <https://github.com/pdishuck/GAVISUNK>.

Contact: eee@gs.washington.edu

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

Highly accurate and contiguous phased diploid *de novo* assemblies of long-read sequencing data have made reference-grade assemblies feasible for many species and individuals (Ebert *et al.*, 2021; Rhie *et al.*, 2021). Pacific Biosciences high-fidelity (HiFi) sequencing in particular has facilitated some of the first complete assembly of centromeres, acrocentric regions, as well as other complex segmental duplications (Logsdon *et al.*, 2021; Miga *et al.*, 2020; Nurk *et al.*, 2022). These assemblies make possible comprehensive, whole-genome evaluations of sequence variation, including some of the most difficult-to-assemble regions, unbiased by reference alignments for the first time. Phased genome assemblies, however, are still subject to the collapse of repetitive sequences, incorrect orientations and misassemblies. Thus, any discoveries based on these automated shotgun sequence assemblies raise the question of whether the assembled sequence is, in fact, valid.

Here, we present genome assembly validation via inter-singly unique nucleotide *k*-mer (SUNK) distances in Oxford Nanopore Technologies (ONT) reads, known hereafter as GAVISUNK. GAVISUNK is a method of validating phased diploid HiFi-driven assemblies with orthogonal ONT sequence. It specifically assesses the contiguity of regions, flagging potential haplotype switches or misassemblies. Although the ONT platform has a significantly higher error rate than that of HiFi (Logsdon *et al.*, 2020), such reads are typically much longer, making it a powerful orthogonal approach for assessing both contiguity and read depth across regions

of interest. Whereas previous genome blacklists or masks of inaccessible regions, such as those used by the ENCODE Consortium, are determined based on annotation of a reference genome (Amemiya *et al.*, 2019), GAVISUNK may be applied to any region or genome assembly to identify misassemblies and potential collapses and is, thus, particularly valuable for validating the integrity of regions with large and highly identical repeats that are more prone to assembly error. This method can be applied genome-wide or at fine scale to closely examine regions of interest across multiple haplotype assemblies.

2 Materials and methods

This assembly validation method relies on identifying SUNKs, *k*-mers that occur just a single time within the HiFi-based assembly (Sudmant *et al.*, 2010) and confirming these SUNKs within long ONT sequencing reads. Because of the relatively lower accuracy of ONT data, false SUNK overlaps may occur at an appreciable frequency between paralogous regions of the genome or haplotypes. Therefore, we leverage not only the presence or absence of SUNKs, but also the intervening distance between pairs of SUNKs, referred to hereafter as inter-SUNK distances. The approach then compares the expected inter-SUNK distance in the assembly and observed distance within ONT reads to recruit reads to their corresponding genomic location. The failure of ONT reads to span between SUNKs in the assembly defines a misjoin, while an excess of reads flags a potential collapse.

We apply Jellyfish (Marçais and Kingsford, 2011) to identify unique k -mers based on all HiFi contigs within an assembly, generating a set of SUNKs for validation. ONT reads are haplotypes phased using parental Illumina whole-genome sequencing data via Canu (Koren et al., 2017), or, in the absence of parental data, with Hi-C and HapCUT2 (Edge et al., 2016). The position of all SUNKs within each ONT read are identified and used for downstream analysis. Each ONT read is assigned to its best-matching HiFi-assembled contig and orientation by comparing the locations of read SUNKs to assembly SUNKs within a diagonal band centered on the median SUNK location for that read. SUNKs observed in ONT reads at an implausibly high or low frequency (i.e. greater than four standard deviations above the mean or fewer than twice) for either haplotype are excluded from consideration.

To validate the assembly of each HiFi contig, all reads identified in the previous step are considered. For each read, a matrix of all pairwise inter-SUNK distances within the read is generated using NumPy and compared to expected distances from the assembly, allowing $\pm 2\%$ variation in length for a given distance by default (Harris et al., 2020). Only the largest set of SUNKs with fully inter-consistent inter-SUNK distances are retained from each read for subsequent validation. In this way, chimeric reads and other spuriously connected SUNKs are separated. A graph is generated for each contig, with SUNKs as nodes and reads as edges, connecting pairs of SUNKs with consistent inter-SUNK distances, using the graph-tool library. This graph is decomposed into its connected components, each of which now corresponds to a validated region of the contig, identifying SUNKs spanned by ONT reads. These validated regions are sent as output to a BED file, and assembly SUNKs with no read support are listed as potentially artifactual.

3 Usage and examples

GAVISUNK can run on consumer hardware, research clusters, or in cloud computing environments, as its steps are automated as part of a configurable Snakemake workflow (Mölder et al., 2021). To install GAVISUNK, clone <https://github.com/pdishuck/GAVISUNK>. Upon first run, dependencies will install automatically as a conda environment.

As described in the README, two configuration files must be edited to specify the assembly and k -mer size (`config.yaml`) and input nanopore read files (`ont.tsv`). To execute the pipeline, submit

```
snakemake -use-conda -cores {thread.count}
```

By default, the pipeline produces a BED file of validated regions (`hap#.validated.bed`) and the complementary unconfirmed ‘gaps’ between validated regions (`hap#.gaps.bed`) for each haplotype of the assembly. For the region of each validation gap, the pipeline produces a visualization to show SUNK-tagged read support, in PDF, SVG and PNG formats (`contig_start_end.format`), as shown in Figure 1 and Supplementary Figure S1. Optionally, a BED file with regions of interest for visualization can be specified for each haplotype in `ont.tsv`.

To annotate the visualizations, BED files encoding regions of interest and color may be supplied. In the example shown in Figure 1, DupMasker annotations demarcate segmental duplications across the locus (Jiang et al., 2008). In addition, the sizes of unspanned inter-SUNK gaps are output for comparison with the distribution of inter-SUNK distances and expectation of spanning that distance given empirical ONT coverage at that read length (Supplementary Fig. S2). Highly identical loci may have insufficient SUNK density for validation with this method, given current ONT read lengths (Supplementary Fig. S3).

We constructed a pseudo-diploid genome assembly as a benchmark of false discovery rate for a well-curated reference assembly with matched ultra-long ONT data. We used the long-read assemblies of the CHM13 and CHM1 human cell lines, both of which are derived from complete hydatidiform moles exhibiting genome-wide uniparental disomy and are therefore guaranteed to represent a single haplotype. CHM13 was assembled by the T2T Consortium with

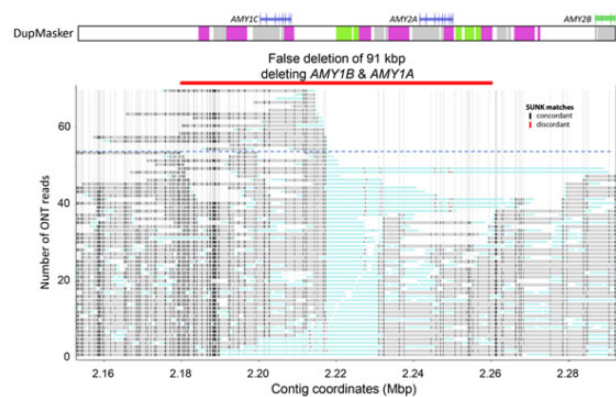


Fig. 1. Example detected misassembly (HG02723 paternal haplotype) within the amylase duplication locus. This assembly contains a false deletion of AMY1B and AMY1A, which is confirmed by the failure of ONT to anchor across the misjoin. Each ONT read is represented by a horizontal bar: the thick gray region contains valid inter-SUNK distances, while the remainder of the read is thinner cyan. Concordant SUNK matches are full-width black tick marks, while discordant matches are small red tick marks. All possible assembly SUNKs are shown as dotted vertical lines in the background, and a DupMasker track marks segmental duplications above. The horizontal dotted blue line indicates the mean genome-wide coverage of ONT sequencing data (A color version of this figure appears in the online version of this article)

a variety of methods and manual validation (Nurk et al., 2022), and CHM1 was assembled from HiFi data [Vollger et al., 2022, using hifiasm v0.12 (Cheng et al., 2021)], each with $>30\times$ coverage of ONT reads longer than 100 kbp. Results for this pseudo-diploid, along with true diploid HG02723 (hifiasm v0.14), are summarized in Supplementary Table S1, with detailed results for CHM13-T2T validation gaps in Supplementary Table S2.

Other assembly validation methods TandemTools and VerityMap (Mikheenko et al., 2020) also used rare k -mers for validation but have distinct goals to GAVISUNK. Both are intended for use on extra-long tandem repeats and developed to use HiFi to validate haploid assemblies. GAVISUNK uses ultra-long ONT data genome-wide to validate diploid assemblies, focused on long interspersed repeats such as segmental duplications.

For CHM13, 1.0% of the assembled genome (103 gaps, 31.7 Mbp) is unsupported by ONT inter-SUNK distances, and 0.7% (274 gaps, 23.9 Mbp) for CHM1. For comparison, assuming random distribution of the empirical ONT read lengths compared to the distances between SUNKs for the CHM13 assembly, 24.9 Mbp is expected to be unsupported (92 gaps, Supplementary Fig. S6). Extra-long tandem repeats are difficult to validate with this method, particularly the *qh* regions of chromosomes 1, 9 and 16, with 73 of 103 CHM13 validation gaps falling in the region. An optional ‘2pass’ mode performs a second pass of analysis on putative validation gaps with more-permissive read recruitment and validates the higher-order repeats of 21/23 CHM13 centromeres (Supplementary Table S3).

4 Conclusion

We developed GAVISUNK, a method for assembly validation using inter-SUNK distances in ONT reads. Applied to HiFi genome assemblies, this tool provides orthogonal validation of regions for downstream analysis, allowing for subsequent genome analyses and annotation. GAVISUNK provides easily interoperable BED outputs and interpretable visualizations of supported and unsupported regions of interest.

Acknowledgements

The authors thank Mitchell Vollger and William Harvey for their helpful algorithmic, programming and source control advice, Tonia Brown for proof-reading, and Kendra Hoekzema for ONT sequencing. This article is subject to

HHMF's Open Access to Publications policy. HHMI lab heads have previously granted a non-exclusive CC BY 4.0 license to the public and a sublicensable license to HHMI in their research articles. Pursuant to those licenses, the author-accepted manuscript of this article can be made freely available under a CC BY 4.0 license immediately upon publication.

Funding

This work was supported, in part, by US National Institutes of Health (NIH) grants HG002385 and HG010169 to E.E.E and 1F32GM134558 to G.A.L. E.E.E. is an investigator of the Howard Hughes Medical Institute.

Conflict of Interest: E.E.E. is a scientific advisory board (SAB) member of Variant Bio, Inc.

Data availability

ONT sequencing reads and their corresponding assemblies are available from the Telomere-to-Telomere and Human Pangenome Reference consortia at <https://github.com/marbl/CHM13> and https://github.com/human-pangenomics/HPP_Year1_Data_Freeze_v1.0. ONT and HiFi sequencing of CHM1 are available on SRA as PRJNA869061 and PRJNA726974, and its hifiasm assembly is available at <https://zenodo.org/record/5502036>.

References

Amemiya, H.M. *et al.* (2019) The ENCODE blacklist: identification of problematic regions of the genome. *Sci. Rep.*, **9**, 9354.
Cheng, H. *et al.* (2021) Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods*, **18**, 170–175.

Ebert, P. *et al.* (2021) Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science*, **372**, eabf7117.
Edge, P. *et al.* (2017) HapCUT2: robust and accurate haplotype assembly for diverse sequencing technologies. *Genome Res.*, **27**, 801–812. <https://doi.org/10.1101/gr.213462.116>.
Harris, C.R. *et al.* (2020) Array programming with NumPy. *Nature*, **585**, 357–362.
Jiang, Z. *et al.* (2008) DupMasker: a tool for annotating primate segmental duplications. *Genome Res.*, **18**, 1362–1368.
Koren, S. *et al.* (2017) Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res.*, **27**, 722–736.
Logsdon, G.A. *et al.* (2020) Long-read human genome sequencing and its applications. *Nat. Rev. Genet.*, **21**, 597–614.
Logsdon, G.A. *et al.* (2021) The structure, function and evolution of a complete human chromosome 8. *Nature*, **593**, 101–107.
Marçais, G. and Kingsford, C. (2011) A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics*, **27**, 764–770.
Miga, K.H. *et al.* (2020) Telomere-to-telomere assembly of a complete human X chromosome. *Nature*, **585**, 79–84.
Mikheenko, A. *et al.* (2020) TandemTools: mapping long reads and assessing/improving assembly quality in extra-long tandem repeats. *Bioinformatics*, **36**, i75–i83.
Mölder, F. *et al.* (2021) Sustainable data analysis with snakemake. *F1000Res.*, **10**, 33.
Nurk, S. *et al.* (2022) The complete sequence of a human genome. *Science*, **376**, 44–53.
Rhie, A. *et al.* (2021) Towards complete and error-free genome assemblies of all vertebrate species. *Nature*, **592**, 737–746.
Sudmant, P.H. *et al.*; 1000 Genomes Project. (2010) Diversity of human copy number variation and multicopy genes. *Science*, **330**, 641–646.
Vollger, M.R. *et al.* (2022) Segmental duplications and their variation in a complete human genome. *Science*, **376**, eabj6965.