

ORIGINAL ARTICLE

Benchmarking Oxford Nanopore read alignment-based insertion and deletion detection in crop plant genomes

Gözde Yildiz  | Silvia F. Zanini  | Nazanin P. Afsharyan  | Christian Obermeier | Rod J. Snowdon  | Agnieszka A. Golicz 

Department of Plant Breeding, Justus Liebig University Giessen, Giessen, Germany

Correspondence

Agnieszka A. Golicz, Department of Plant Breeding, Justus Liebig University Giessen, Giessen, Germany.

Email:

Agnieszka.Golicz@agrar.uni-giessen.de

Assigned to Associate Editor Hon-Ming Lam.

Funding information

Alexander von Humboldt Foundation; German Research Foundation, Grant/Award Number: 458716530

Abstract

Structural variations (SVs) are larger polymorphisms (> 50 bp in length), which consist of insertions, deletions, inversions, duplications, and translocations. They can have a strong impact on agronomical traits and play an important role in environmental adaptation. The development of long-read sequencing technologies, including Oxford Nanopore, allows for comprehensive SV discovery and characterization even in complex polyploid crop genomes. However, many of the SV discovery pipeline benchmarks do not include complex plant genome datasets. In this study, we benchmarked insertion and deletion detection by popular long-read alignment-based SV detection tools for crop plant genomes. We used real and simulated Oxford Nanopore reads for two crops, allotetraploid *Brassica napus* (oilseed rape) and diploid *Solanum lycopersicum* (tomato), and evaluated several read aligners and SV callers across 5×, 10×, and 20× coverages typically used in re-sequencing studies. We further validated our findings using maize and soybean datasets. Our benchmarks provide a useful guide for designing Oxford Nanopore re-sequencing projects and SV discovery pipelines for crop plants.

1 | INTRODUCTION

Structural variations (SVs) are a major type of polymorphisms, which consist of insertions, deletions, inversions, duplications, and translocations. SVs are larger polymorphisms (> 50 bp) compared with single nucleotide polymorphisms (SNPs) and small indels (insertions and deletions). Copy number variations (CNVs) and presence/absence variations (PAVs) occur due to these genomic polymorphisms (Alkan et al., 2011; Sedlazeck et al., 2018a). Insertions and deletions are the most abundant type of SV (Alonge et al.,

2020; Fuentes et al., 2019; Goel et al., 2019), can have a strong effect on crop traits, and have been shown to play a role in domestication and environmental adaptation (Gill et al., 2021; Tao et al., 2019; Yildiz et al., 2022; Zanini et al., 2022). Until recently, the lack of high-quality reference assemblies and the complex nature of often large, polyploid genomes made comprehensive SV exploration challenging in crop genomic research (Meyers & Levin, 2006; Yuan et al., 2021).

Development of long-read sequencing technologies such as Oxford Nanopore Technologies (ONT) (Jain et al., 2016) and Pacific Biosciences (PacBio) (Roberts et al., 2013) provided new opportunities for comprehensive SV discovery in crop plants. The sequencing accuracy of these technologies is continuously improving. Currently, PacBio HiFi

Abbreviations: CNV, copy number variant; ONT, Oxford Nanopore Technologies; PacBio, Pacific Biosciences; PAV, presence/absence variant; SNP, single nucleotide polymorphism; SV, structural variant.

This is an open access article under the terms of the [Creative Commons Attribution-NonCommercial-NoDerivs](https://creativecommons.org/licenses/by-nc-nd/4.0/) License, which permits use and distribution in any medium, provided the original work is properly cited, the use is non-commercial and no modifications or adaptations are made.

© 2023 The Authors. *The Plant Genome* published by Wiley Periodicals LLC on behalf of Crop Science Society of America.

consensus reads exceed 99% accuracy (Wenger et al., 2019) while ONT R10.3 raw reads accuracy exceeds 95% (Delahaye & Nicolas, 2021). The reduction in error rates facilitates downstream applications, including the production of high-quality genome assemblies, and SV detection. ONT sequencing in particular is being adopted in crop plant research for large scale re-sequencing projects of tens to hundreds of individuals (Alonge et al., 2020; Chawla et al., 2021; Lemay et al., 2022; Vollrath et al., 2021; Zhang et al., 2022). Despite the constant decrease in sequencing error rate, long-read technologies require specialized computational approaches to take advantage of them efficiently.

The two main approaches for SV discovery are *de novo* assembly-based and read alignment-based. *De novo* assembly-based approaches assemble reads into longer contigs and identify SVs by aligning assemblies (Wenger et al., 2019). Read alignment-based approaches directly align reads to reference genomes to discover SVs. *De novo* assembly-based methods perform better at finding larger variants (tens to hundreds of kbp long; exceeding the length of individual reads) but require sufficient amount of data to produce high-quality assemblies, which leads to substantial increase in cost of the experiments for larger crop genomes. However, read alignment-based approaches can perform well even at modest sequencing depths of 5× to 10× and use less computational resources, but the discovered SVs are limited to differences with the reference genome which makes this approach more suitable for larger re-sequencing projects (Coster et al., 2021). Several algorithms were developed for SV discovery from long-reads including Sniffles (Sedlazeck et al., 2018b), NanoVar (Tham et al., 2019), SVIM (Heller & Vingron, 2019), cuteSV (Jiang et al., 2020), and dysgu (Cleal & Baird, 2022), which have been comprehensively reviewed recently (Mahmoud et al., 2019; Yuan et al., 2021). Additionally, several long-read aligners are available such as minimap2 (Li, 2018), NGMLR (Sedlazeck et al., 2018a), Vulcan (Fu et al., 2021), and Ira (Ren & Chaisson, 2021). Considering the continued development and improvement in read-alignment and SV detection algorithms and multitude of their possible combinations, their combined performances in SV detection demand realistic and up-to-date benchmarks to guide the selection of SV discovery tools.

In this study, we hypothesized that certain combination(s) of read aligners and SV discovery software will have superior performance in datasets representing complex crop genomes. We used real and simulated ONT reads for two crop plant genomes and evaluated several mappers and SV callers across coverages including 5×, 10×, and 20× typically utilized in re-sequencing studies. We chose to perform benchmarking on allotetraploid *Brassica napus* (oilseed rape) and diploid *Solanum lycopersicum* (tomato) as these two species represent different ploidy, have different SV profiles, and were already studied using Oxford Nanopore Technology. We further val-

Core Ideas

- Structural variants (SVs) have strong impact on crop traits and play an important role in environmental adaptation.
- Long read based SV discovery tools have not been comprehensively evaluated in crops.
- We benchmarked popular SV discovery tools using real and simulated data for two contrasting crop genomes.
- Our benchmarks provide a guide for choosing insertion and deletion discovery tools for low to medium sequencing coverage experiments.

idated our findings using maize and soybean datasets. Our benchmarks provide a guide for choosing insertion and deletion discovery tools for low to medium coverage sequencing projects.

2 | MATERIALS AND METHODS

2.1 | Read aligners, SV callers, and benchmarking datasets

The SV callers included in the study were selected using several criteria: (1) citation count (adjusted by number of years since publication and used as a proxy for popularity in the research community); (2) publication date and maintenance status (excluding older tools that were no longer maintained); (3) ability to detect both insertion and deletion SVs from ONT data. The benchmarking approach involved four long-read aligners, including minimap2 (Li, 2018), NGMLR (Sedlazeck et al., 2018a), Ira (Ren & Chaisson, 2021), and Vulcan (Fu et al., 2021) as well as five SV calling software namely Sniffles (v2) (Sedlazeck et al., 2018b), NanoVar (Tham et al., 2019), SVIM (Heller & Vingron, 2019), cuteSV (Jiang et al., 2020), and dysgu (Cleal & Baird, 2022). All aligners and SV caller versions are provided in detail in (Table S1). Three simulated datasets (Sim_ONT_Bn1, Sim_ONT_Bn2, and Sim_ONT_Sl) and publicly available data, for *B. napus* and *S. lycopersicum* genomes, were used. The real-world datasets for whole genome Nanopore sequencing of *B. napus* cv. King 10 (accession number: SRR15731030) (Vollrath et al., 2021), *S. lycopersicum* cv. M82 (accession number: SRR16966224) (Alonge et al., 2021), *Zea mays* cv. Mo17 (accession number: SRR15447413), and *Glycine max* cv. Maple Isle (accession number: SRR15342671 and SRR15342672) were downloaded from NCBI Sequencing Read Archive. All but soybean datasets were randomly subsampled to 5×, 10×, and 20×

coverages using Rasusa (Hall, 2022) to test the effect of sequencing depth on SV discovery.

2.2 | Simulated dataset generation

For three simulated datasets (workflow for all simulations is presented in (Figure S1), new haplotypes including SVs were generated, and synthetic ONT reads were simulated using VISOR v1.1 (Bolognini et al., 2020). For simulation one (Sim_ONT_Bn1), 20,000 genomic intervals (mean: 750 bp, SD: 500 bp) were randomly drawn from the *B. napus* genome (Express 617 v1). A subset of 10,000 was denoted as deletions. For the remaining 10,000, denoted as insertions, the genomic start coordinate was retained, while the sequences corresponding to the genomic intervals were extracted, randomly re-assigned to the coordinates, and served as insertion sequences at those coordinates (Figure S1).

Simulations two and three, denoted Sim_ONT_Bn2 and Sim_ONT_SI, were designed to reflect SVs found in real-world datasets. For Sim_ONT_Bn2, the assembled *B. napus* genomes Express 617 v1 (Lee et al., 2020) and Westar (Song et al., 2020) were aligned using minimap2 v2.24. SVs were detected using SVIM-asm v1.0.2 (Heller & Vingron, 2020). To reduce the effect of using minimap2 for benchmarking dataset generation, the SV locations were shifted by a randomly selected number in the (−5000, 5000) interval. This changed the exact SV site while maintaining the realistic distribution of SV sizes and locations along the genome. A random subset of 10,000 insertions and 10,000 deletions was drawn from all SVs to create the benchmarking dataset. SNPs discovered from short reads using bcftools v1.15.1 were also included. The SVs and SNPs were provided to VISOR to generate a new haplotype, which in turn was used for Oxford Nanopore read simulation. Sim_ONT_SI was generated using the same strategy as for Sim_ONT_Bn2 but designed to reflect SVs of the *S. lycopersicum* genome. Heinz 1706 (*Slycopersicum_691_SL4.0*) and M82 (Alonge et al., 2021) assemblies were used for whole genome alignments. Due to smaller number of SVs, a random subset of 2500 insertions and 2500 deletions were drawn from all SVs. For maize, we used Zmays_493_APGv4 (B73) and ZmaysB84_681 (B84) (Bornowski et al., 2021).

To test the effect of sequencing depth on SV discovery, the datasets were simulated at 5×, 10×, and 20× coverage. The simulations provided the objective truth sets, which could be used to calculate SV precision, recall, and combined F1-scores. Precision describes the proportion of correct positive predictions among all positive predictions. It is calculated by dividing the true positives by overall positives. Recall describes the proportion of positive predictions made out of all positive elements in the dataset. It is calculated by dividing true positives by total number of relevant elements. F1-score

combines precision and recall by taking their harmonic mean. Its value ranges from 0 to 1. F1-score close to 1 indicates high precision and recall. Using two different strategies for generating simulated datasets will make it possible to minimize analytical bias. If the same combination of tools performed best on all simulated datasets, this will likely reflect true superior performance.

2.3 | Comparative analyses

Express 617 v1 for *B. napus* (Lee et al., 2020) and *Slycopersicum_691_SL4.0* for *S. lycopersicum* (Hosmani et al., 2019) were used as reference sequences. Simulated datasets and real subsampled reads at each coverage depth were aligned to respective reference genomes. The SV call sets were filtered using the following criteria: (1) number of minimum supporting reads: 5×: 3, 10×: 5, and 20×: 8; (2) SV type: INS or DEL (the most abundant SVs supported by all the benchmarked tools); (3) minimum SV length: 50 bp; (4) SV quality: SVs flagged as “PASS”; (5) genotype: homozygous genotype for alternative allele (‘1/1’). For simulated data, precision, recall, and F1-scores of the SVs were computed for each combination of coverage depth, read aligner, and SV caller using Truvari v3.0.0 (English et al., 2022). Comparisons between results from the same tool combination across different coverages and different tool combinations across the same coverages were performed using surpyvor v0.8.1 (Jef-fares et al., 2017). For real datasets, where no truth sets were available, we focused on within-dataset comparisons and how those compared to the results from simulated data. All the relevant commands for simulated data generation and SV discovery are available in the Supporting Information. To ensure that the datasets were comparable, soybean SV calls were filtered using the same criteria as described in Lemay et al. (2022).

3 | RESULTS

3.1 | Selecting the benchmarking datasets

We chose to focus on two crop plant species *B. napus* (oilseed rape; genome size ~1.1 Gbp) and *S. lycopersicum* (tomato; genome size ~900 Mbp) because they are both important crops and their structural variation was previously studied using Oxford Nanopore Technologies (Alonge et al., 2020; Chawla et al., 2021). Whole Genome Alignment (WGA)-based SV discovery also suggested that they have quite different SV profiles with 38,666 SVs (Real_WGA_Bn, mean size: 2068 bp, median size: 593 bp, 19,450 insertions and 19,216 deletions) discovered for *B. napus* and 7108 SVs (Real_WGA_SI, mean size: 3029 bp, median size:

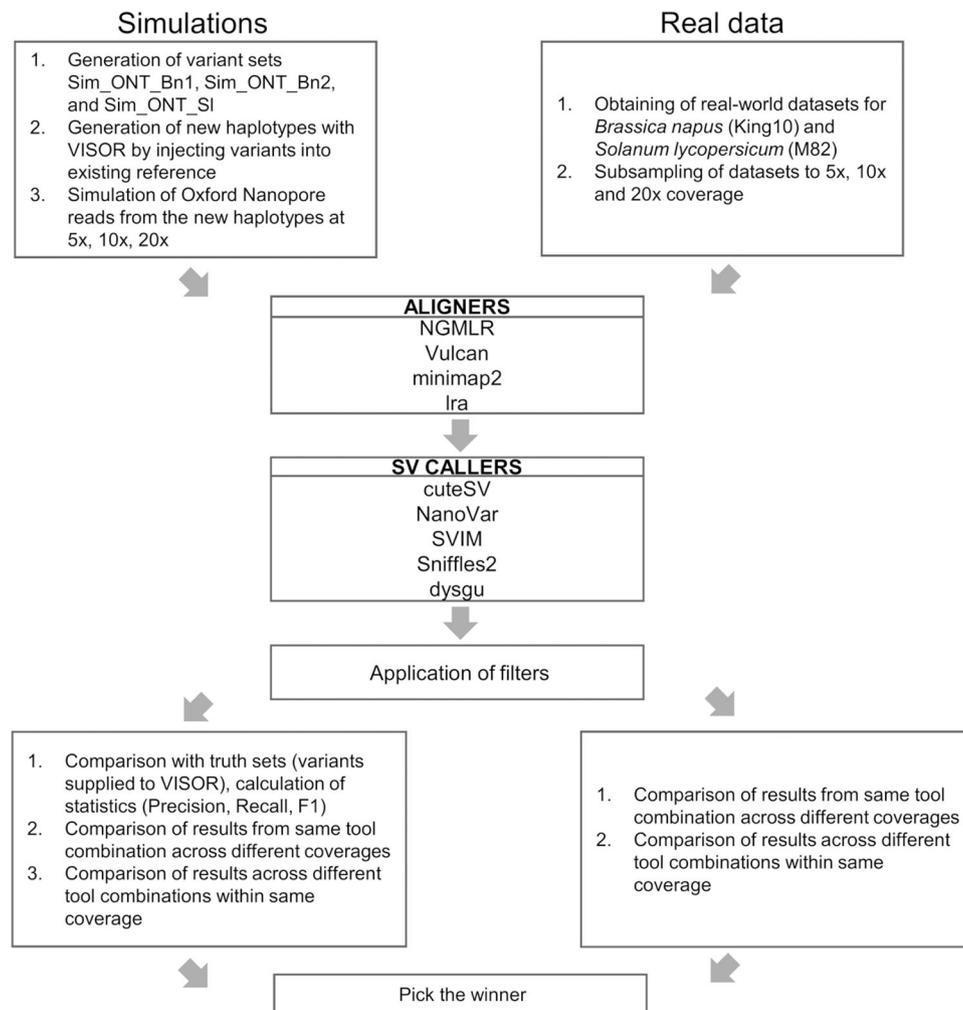


FIGURE 1 Graphical overview of the benchmarking workflow. SV, structural variant.

178 bp, 4159 insertions and 2949 deletions) discovered for *S. lycopersicum*.

Two simulated *B. napus* haplotypes (Sim_ONT_Bn1 and Sim_ONT_Bn2) and one simulated *S. lycopersicum* haplotype (Sim_ONT_Sl) were used to generate Oxford Nanopore reads at 5x, 10x, and 20x to test the effect of sequencing depth on SV discovery. The two publicly available real-world datasets, from *B. napus* (38x) and *S. lycopersicum* (68x), were subsampled with the same logic (Real_ONT_Bn, Real_ONT_Sl). The available graphical representation of a workflow for simulation and real data are shown in Figure 1.

3.2 | Characteristics of structural variant truth sets

The SVs supplied to VISOR to generate Sim_ONT_Bn1, Sim_ONT_Bn2, and Sim_ONT_Sl haplotypes served as three truth sets for our comparisons. The truth sets included deletions and insertions. The length distribution of truth set SVs

is presented in Figure 2. Sim_ONT_Bn1 is unbiased in terms of the bioinformatics tools used, as the regions representing SVs were entirely randomly drawn from the *B. napus* genome. For any simulated dataset to reflect realistic SV distribution, SVs have to be discovered first and provided to the simulation software. Any relationship between tools used for SV identification for long-read dataset simulation and tools used for SV detection from these simulated reads (for example use of similar/same mapping algorithm) can result in inflated performance and biased results. However, Sim_ONT_Bn1 does not reflect realistic SV length and genomic distribution. To mitigate that, Sim_ONT_Bn2 and Sim_ONT_Sl were created using SVs derived from real-world datasets. The two simulation strategies are complementary and should allow both unbiased and realistic assessment of SV calls. The median (mean) sizes (bp) for insertions and deletions were 800 (834) and 795 (825) for Sim_ONT_Bn1, 629 (1959) and 594 (1904) for Sim_ONT_Bn2 and 162 (3178) and 165 (2477) for Sim_ONT_Sl. Overall, the Sim_ONT_Bn2 and Sim_ONT_Sl truth sets had a wider range of insertion and deletion sizes.

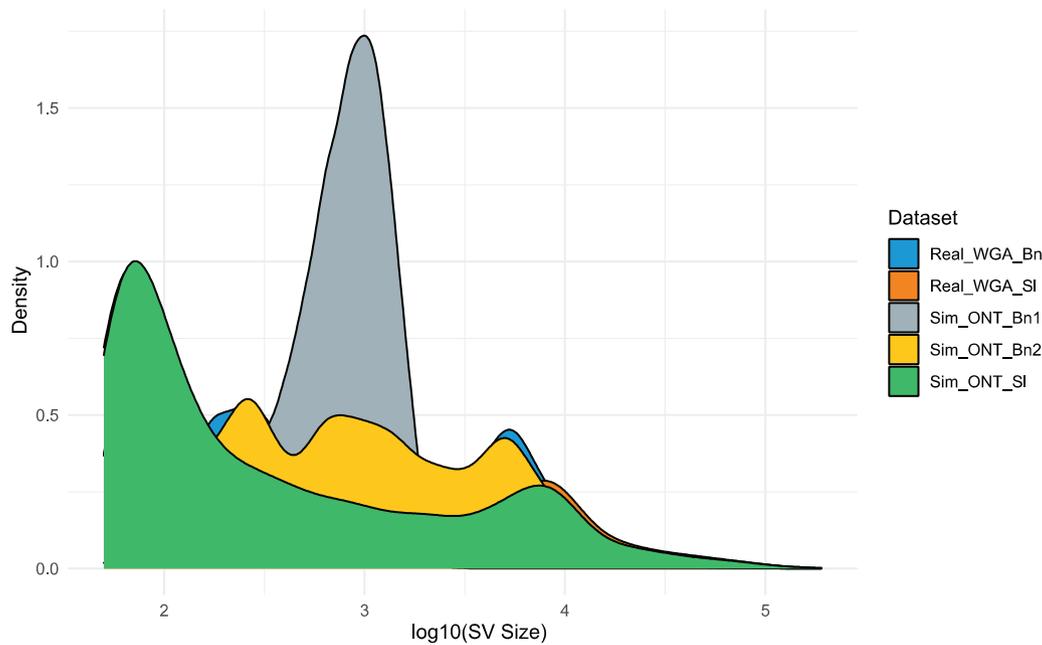


FIGURE 2 Size distribution of the real-world structural variants (SVs) and SVs from three benchmarking datasets.

They were more reflective of true biological variation, making them more realistic than the Sim_ONT_Bn1 truth set.

3.3 | Performance of long read aligners

Subsampled *S. lycopersicum*, *B. napus*, and simulated reads were aligned using Ira, minimap2, Vulcan, and NGMLR to the *Slycopersicum_691_SL4.0*, and Express 617 v1 reference genomes. Mapping statistics and run times of alignment against relevant reference genomes with different coverages of Sim_ONT_Bn1, Sim_ONT_Bn2, Sim_ONT_SI, *B. napus* (Real_ONT_Bn), and *S. lycopersicum* (Real_ONT_SI) real-world datasets are given in Table S2. Minimap2 had the shortest run time across all coverages. Conversely, NGMLR had the longest run time and also the lowest mapping rate. Figure 3 shows mapping runtime (h:mm:ss or m:ss) for both simulation and real-world datasets with eight CPUs. Real_ONT_Bn dataset with 20× coverage was aligned ~220 h by NGMLR and ~119 h by Vulcan, compared to ~4 h by minimap2 and ~5 h by Ira. Therefore, minimap2 and Ira provided a greater speed advantage than NGMLR and Vulcan. The run times increased with the higher coverages (Figure 3). Processing of real data took substantially longer than processing of simulated data. Moreover, Vulcan and minimap2 produced the highest proportion of mapped reads in Real_ONT_Bn (> 96%), Real_ONT_SI (96%–98%), and all simulated data (> 98%) (Table S2). NGMLR reported the lowest proportion of mapped reads for Real_ONT_Bn (~81%) and Real_ONT_SI (~76%), while Ira and NGMLR resulted in similar statistics (96%–97%) for Sim_ONT_Bn1,

Sim_ONT_Bn2, and Sim_ONT_SI at each coverage. The combination of fast run time, good mapping rate, and the SV calling results presented below suggest that minimap2 is the top-performing aligner for simulated and real reads.

3.4 | Performance of SV callers on simulated data

3.4.1 | Performance using Sim_ONT_Bn1 as benchmark

We calculated the precision, recall, and F1-score of the SVs generated using different mapper and SV caller combinations using the Sim_ONT_Bn1 truth set. Table S3 shows comparison of the precision, recall, and F1-scores for all mapper/SV caller combinations at the 5×, 10×, and 20× coverages. Each aligner/SV caller combination was evaluated with respect to total SVs, deletions, and insertions. Figure 4 presents the corresponding F1-scores at 5× to 20× coverages. CuteSV after minimap2 alignment reached the highest F1-scores 5×:~0.90, 10×:~0.97, and 20×:~0.99 for total SVs, 5×:~0.91, 10×:~0.97, and 20×:~0.99 for deletions, and 5×:~0.89, 10×:~0.96, and 20×:~0.99 for insertions. At the lower end of coverage (5×), the combination of minimap2/cuteSV provided a better advantage when compared to other mapper/SV caller combinations, especially in capturing insertions. Minimap2/Sniffles2 had second-best F1-scores (Figure 4). SVs detection by NanoVar was obtained directly from reads as NanoVar has its own internal mapping

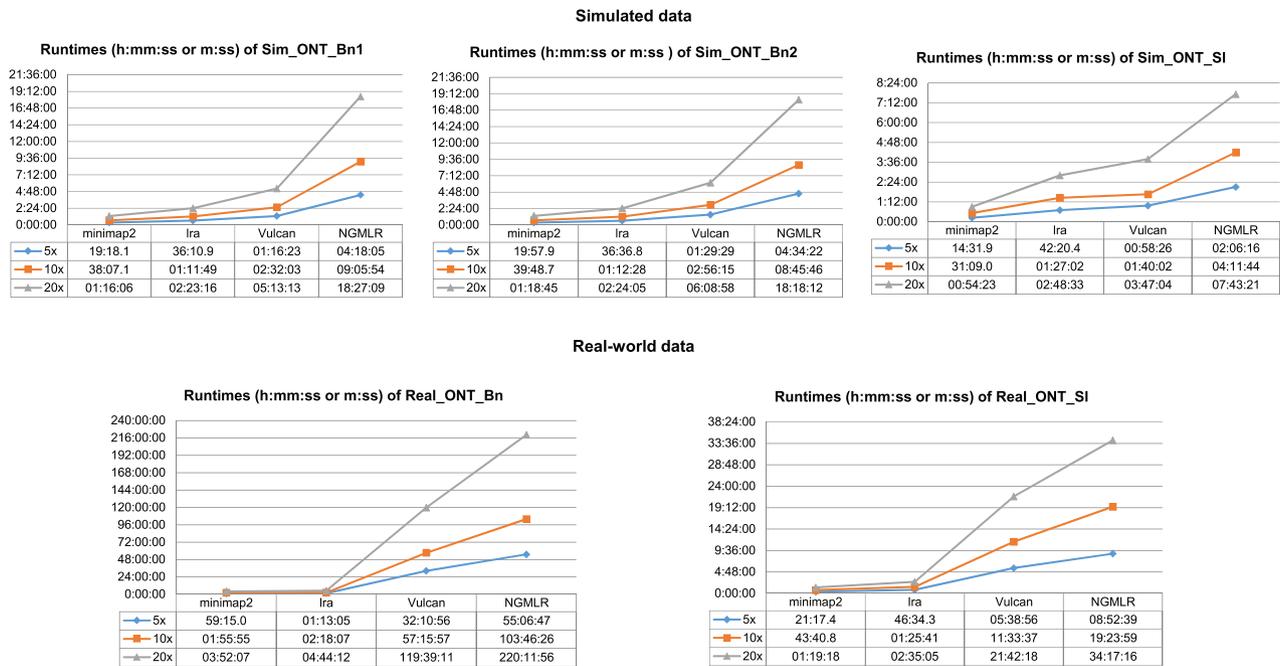


FIGURE 3 Read aligner run time (h:mm:ss or m:ss) for both simulation and real-world datasets with 5×, 10×, and 20× coverages (8 CPU). The reads were simulated with a mean length of 15,000 bp. Unplaced contigs were not included in simulations, which may reduce run time for simulated reads. Read-world reads had higher N50 (~29 Kbp for *B. napus* and ~42 Kbp for *S. lycopersicum*) compared to simulated data (~22 Kbp). In addition, *B. napus* real world data could contain non-reciprocal homeologous exchanges (HEs) uncounted for in simulations. Higher N50 and presence of HEs could increase run time for real-world data.

algorithm; therefore, the precision, recall, and F1-scores for different aligners are not included.

We also compared the total number of SVs, insertions, and deletions for all tested aligner/SV caller combinations. Table S4 summarizes the number of SVs found at 5×, 10×, and 20× coverages. There were more discovered deletions than insertions regardless of coverage. The combinations of minimap2/cuteSV and minimap2/Sniffles2 detected the highest number of SVs at each coverage. We also analyzed how many of the SVs overlapped across different coverages while using the same tool combination and how many of the SVs overlapped across different tool combinations within the same coverage. Data S1 shows the number of overlapping and unique SVs across coverages. Minimap2/cuteSV combination had the highest number of overlapping SVs. It also resulted in the highest proportion of overlapping SVs; 76.99% for all SVs, 79.19% for deletions, and 74.79% for insertions, while the minimap2/Sniffles2 combination (second best according to F1-scores) had the second highest percentage overlap; 75.35% for all SVs, 78.35% for deletions, and 72.33% for insertions (Table S5). In addition, we performed comparisons across different tool combinations within the same coverage. Data S2 displays the overlap, including the intersection sizes between SV calls and the Sim_ONT_Bn1 truth set. The highest number of overlapping SVs was found at 20x coverage, following minimap2 aligner. Our Sim_ONT_Bn1 results suggest that the combination of cuteSV and Sniffles2 with

minimap2 alignment gave the best results achieving high F1-scores and capturing the highest number of overlapping SVs across coverages.

3.4.2 | Performance using Sim_ONT_Bn2 as benchmark

While Sim_ONT_Bn1 represents relatively short SVs randomly distributed along the genome, Sim_ONT_Bn2 reflects true biological variation in *B. napus*. Table S6 presents comparison of the precision, recall, and F1-scores for all mapper/SV caller combinations at the 5×, 10×, and 20× coverages. Figure 5 presents the F1-scores of SVs (total, insertions, and deletions) obtained using different combinations of aligners and variant callers across coverages. CuteSV following minimap2 alignment again was the top performing combination with the highest overall F1-score values 5×:~0.87, 10×:~0.93, and 20×:~0.96 for total SVs, 5×:~0.90, 10×:~0.96, and 20×:~0.98 for deletions, and 5×:~0.83, 10×:~0.90, and 20×:~0.94 for insertions. Especially, at low 5× coverage, this combination performed better than others. Minimap2/Sniffles2 had the second highest F1-scores at 20× coverage as in Sim_ONT_Bn1. However, minimap2/dysgu F1-score for insertions at 5× and 10× was higher than Sniffles2 after the minimap2 alignment.

	Total			Deletions			Insertions		
	minimap2			minimap2			minimap2		
	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1
cuteSV	0.9003	0.9676	0.9955	0.9074	0.9733	0.9961	0.8931	0.9620	0.9948
Sniffles2	0.8928	0.9635	0.9948	0.9037	0.9724	0.9963	0.8818	0.9544	0.9933
SVIM	0.7825	0.9645	0.9869	0.7970	0.9715	0.9957	0.7676	0.9574	0.9778
dysgu	0.8618	0.9417	0.9776	0.9057	0.9721	0.9952	0.8140	0.9092	0.9593
	Ira			Ira			Ira		
cuteSV	0.8665	0.9417	0.9829	0.8836	0.9562	0.9860	0.8488	0.9267	0.9798
Sniffles2	0.8578	0.9352	0.9801	0.8821	0.9557	0.9865	0.8324	0.9138	0.9736
SVIM	0.7291	0.9354	0.9793	0.7696	0.9563	0.9857	0.6857	0.9135	0.9728
dysgu	0.7593	0.8718	0.9148	0.8783	0.9552	0.9852	0.6112	0.7735	0.8336
	Vulcan			Vulcan			Vulcan		
cuteSV	0.8495	0.9256	0.9751	0.8707	0.9441	0.9823	0.8275	0.9065	0.9678
Sniffles2	0.8000	0.8787	0.9463	0.8544	0.9323	0.9780	0.7401	0.8191	0.9124
SVIM	0.6864	0.9024	0.9345	0.7325	0.9389	0.9809	0.6367	0.8632	0.8834
dysgu	0.7441	0.8253	0.8639	0.8695	0.9484	0.9814	0.5866	0.6689	0.7150
	NGMLR			NGMLR			NGMLR		
cuteSV	0.8001	0.8691	0.9220	0.8490	0.9152	0.9496	0.7465	0.8187	0.8927
Sniffles2	0.6524	0.7174	0.7689	0.8198	0.8924	0.9338	0.4282	0.4753	0.5424
SVIM	0.6295	0.8496	0.8980	0.7116	0.9110	0.9477	0.5358	0.7805	0.8431
dysgu	0.6275	0.7120	0.7415	0.8428	0.9260	0.9534	0.3116	0.3894	0.4193
	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1
NanoVar	0.8950	0.9593	0.9848	0.9012	0.9676	0.9913	0.8886	0.9509	0.9784

FIGURE 4 F1-scores of Sim_ONT_Bn1 including total structural variants (SVs), deletions, and insertions at 5x, 10x, and 20x coverages for different combinations of read aligners and SV callers.

	Total			Deletions			Insertions		
	minimap2			minimap2			minimap2		
	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1
cuteSV	0.8709	0.9301	0.9628	0.9060	0.9609	0.9825	0.8335	0.8973	0.9422
Sniffles2	0.8589	0.9182	0.9545	0.9011	0.9580	0.9827	0.8132	0.8752	0.9248
SVIM	0.7481	0.9195	0.9527	0.7942	0.9549	0.9791	0.6984	0.8816	0.9250
dysgu	0.8602	0.9316	0.9528	0.8968	0.9576	0.9756	0.8214	0.9045	0.9292
	Ira			Ira			Ira		
cuteSV	0.8059	0.8732	0.9203	0.8592	0.9254	0.9648	0.7474	0.8155	0.8715
Sniffles2	0.8032	0.8768	0.9237	0.8473	0.9189	0.9578	0.7556	0.8313	0.8874
SVIM	0.6726	0.8616	0.9068	0.7334	0.9178	0.9553	0.6056	0.7992	0.8535
dysgu	0.7045	0.8295	0.8686	0.8125	0.9085	0.9466	0.5752	0.7382	0.7784
	Vulcan			Vulcan			Vulcan		
cuteSV	0.8000	0.8635	0.9122	0.8469	0.9101	0.9524	0.7490	0.8126	0.8686
Sniffles2	0.7553	0.8240	0.8759	0.8136	0.8832	0.9300	0.6910	0.7581	0.8160
SVIM	0.6448	0.8382	0.8870	0.7005	0.8919	0.9358	0.5841	0.7790	0.8336
dysgu	0.7240	0.8391	0.8770	0.7788	0.8914	0.9310	0.6642	0.7819	0.8175
	NGMLR			NGMLR			NGMLR		
cuteSV	0.7762	0.8408	0.8885	0.8272	0.8887	0.9302	0.7201	0.7882	0.8429
Sniffles2	0.7219	0.7895	0.8442	0.7857	0.8531	0.9029	0.6509	0.7182	0.7785
SVIM	0.6137	0.8095	0.8608	0.6825	0.8731	0.9204	0.5372	0.7382	0.7943
dysgu	0.6703	0.7998	0.8436	0.7541	0.8749	0.9160	0.5743	0.7143	0.7612
NanoVar	0.7987	0.8583	0.8964	0.8399	0.9030	0.9432	0.7550	0.8108	0.8471

FIGURE 5 F1-scores of Sim_ONT_Bn2 including total structural variants (SVs), deletions, and insertions at 5x, 10x, and 20x coverages for different combinations of read aligners and SV callers.

In addition, the total number of SVs, the total number of insertions, and deletions for all combinations of tested aligners and SV callers were compared. Table S7 summarizes the total number of SVs detected at 5x, 10x, and 20x coverages. Minimap2/cuteSV found the highest number of SVs at each coverage like in Sim_ONT_Bn1. Again, more dele-

tions than insertions were found for all aligner and SV caller combinations across different coverages. We also analyzed how many of the SVs overlapped across different coverages while using the same tool combination and how many of the SVs overlapped across different tool combinations within the same coverage. Data S3 lists the number of overlapping SVs

	Total minimap2			Deletions minimap2			Insertions minimap2		
	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1	5x-F1	10x-F1	20x-F1
cuteSV	0.8467	0.9167	0.9375	0.8831	0.9477	0.9654	0.8073	0.8833	0.9077
Sniffles2	0.8432	0.9174	0.9394	0.8795	0.9492	0.9671	0.8041	0.8835	0.9099
SVIM	0.7520	0.9184	0.9377	0.7944	0.9488	0.9647	0.7060	0.8858	0.9089
dysgu	0.8371	0.9008	0.9043	0.8594	0.9194	0.9226	0.8134	0.8814	0.8852
	Ira			Ira			Ira		
cuteSV	0.8158	0.8936	0.9278	0.8547	0.9308	0.9628	0.7736	0.8530	0.8897
Sniffles2	0.8334	0.9170	0.9515	0.8519	0.9329	0.9652	0.8141	0.9007	0.9376
SVIM	0.7158	0.8955	0.9315	0.7591	0.9322	0.9639	0.6688	0.8558	0.8966
dysgu	0.7682	0.8884	0.9149	0.8040	0.9215	0.9462	0.7295	0.8524	0.8808
	Vulcan			Vulcan			Vulcan		
cuteSV	0.8128	0.8881	0.9161	0.8490	0.9255	0.9525	0.7736	0.8472	0.8763
Sniffles2	0.8012	0.8785	0.9120	0.8324	0.9123	0.9448	0.7678	0.8419	0.8768
SVIM	0.6994	0.8804	0.9140	0.7396	0.9188	0.9476	0.6562	0.8387	0.8780
dysgu	0.7852	0.8759	0.8992	0.8005	0.9074	0.9241	0.7694	0.8421	0.8730
	NGMLR			NGMLR			NGMLR		
cuteSV	0.8002	0.8693	0.9001	0.8370	0.9114	0.9380	0.7601	0.8226	0.8581
Sniffles2	0.7889	0.8615	0.9015	0.8178	0.8924	0.9257	0.7580	0.8282	0.8758
SVIM	0.6832	0.8662	0.9050	0.7279	0.9106	0.9412	0.6347	0.8174	0.8658
dysgu	0.7636	0.8626	0.8881	0.7908	0.9002	0.9194	0.7348	0.8214	0.8541
NanoVar	0.7504	0.8098	0.8103	0.8488	0.9093	0.9232	0.6282	0.6841	0.6608

FIGURE 6 F1-scores of Sim_ONT_SI including total structural variants (SVs), deletions, and insertions at 5×, 10×, and 20× coverages for different combinations of read aligners and SV callers.

across different coverages using the same tool combination. Minimap2/cuteSV combination had the highest number of overlapping SVs. It also had the highest proportion of overlapping SVs; 73.95% for all SVs, 80.05% for deletions, and 67.44% for insertions. The minimap2/dysgu combination was second best detecting 73.23% for all SVs, and 67.28% for insertions. Minimap2/Sniffles2 combination was the second best for deletions with 79.14% overlap (Table S5). Data S4 displays overlap between results from different SV callers within the same coverage after each aligner, including the intersection with the Sim_ONT_Bn2 truth set. The highest number of overlapping SVs was found at 20x coverage, following minimap2 aligner. Overall, in Sim_ONT_Bn2, the combination of cuteSV after minimap2 alignment gave the best results both in terms of F1-Scores and concordance across coverages.

3.4.3 | Performance using Sim_ONT_SI as benchmark

Sim_ONT_SI represents the true biological variation of *S. lycopersicum*. Table S8 presents comparison of the precision, recall, and F1-scores for all mapper/SV caller combinations at the 5×, 10×, and 20× coverages. Figure 6 shows the F1-score of SVs (total, insertions, and deletions) identified using combinations of the different aligners and variant callers. CuteSV and Sniffles2 with minimap2 alignment were top performers with the highest F1-score values (5×:~0.85, 10×:~0.92, and 20×:~0.94) for total SVs, (5×:~0.88, 10×:~0.95, and

20×:~0.97) for deletions, and (5×:~0.81, 10×:~0.88, and 20×:~0.91) for insertions. Ira/Sniffles2 combination had the best F1-score for insertions for each coverage.

In addition, the total number of SVs, the total number of insertions, and deletions for all tested aligner/SV caller combinations were compared. Table S9 summarizes the total number of SVs at 5×, 10×, and 20× coverages. Again, more deletions than insertions were found for all aligner and SV caller combinations across coverages like in the previous simulated datasets. The number of SVs overlapping across coverages while using the same tool combination and the number of SVs overlapping across different tool combinations but within the same coverage were also calculated. Data S5 shows the number of overlapping SVs across different coverages using the same tool combination. Minimap2/dysgu combination had the highest number of overlapping SVs. However, minimap2/cuteSV combination found the highest proportion of overlap; 73.49% for all SVs, 77.52% for deletions, and 68.98% for insertions, while the minimap2/Sniffles2 combination was second best detecting 72.73% for all SVs, 76.32% for deletions, and 68.72% for insertions (Figure 7 and Table S5). Although minimap2/dysgu found the highest number of SVs at each coverage in Sim_ONT_SI, the proportion of overlapped SVs was reported as 68.82%. Data S6 displays overlap between results from different SV callers within the same coverage after each aligner, including the intersection with Sim_ONT_SI truth set. The highest number of overlapping SVs was found at 20x coverage, following minimap2 aligner. Overall, in Sim_ONT_SI, the combination of cuteSV and Sniffles2 after minimap2 alignment gave the best

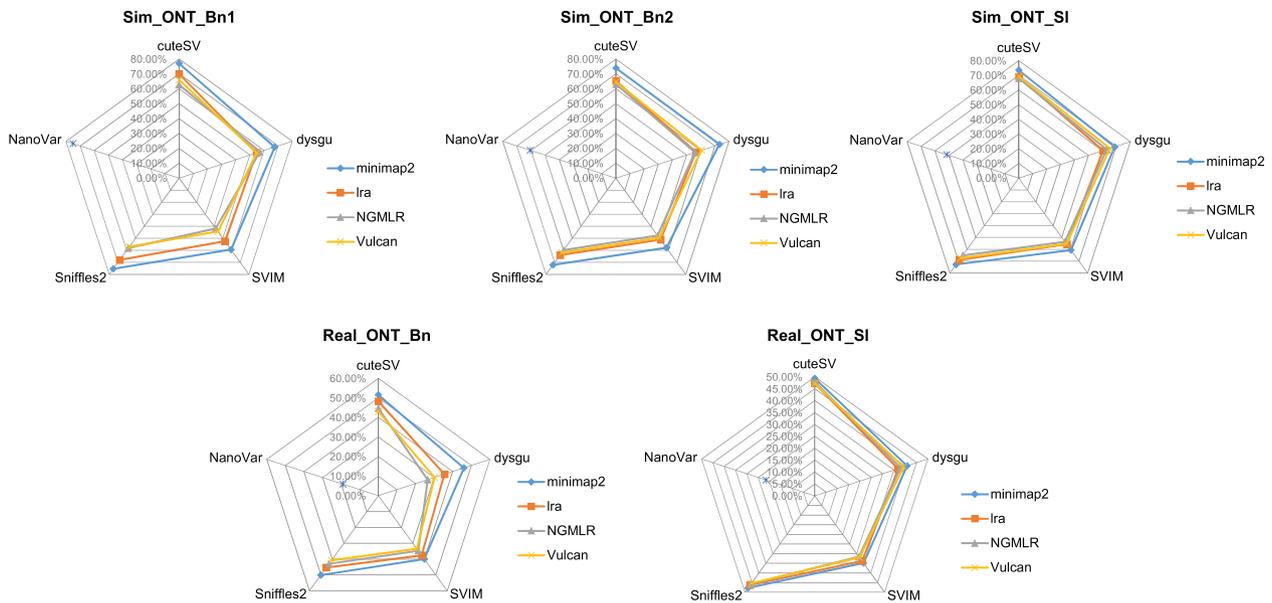


FIGURE 7 Proportion of overlapped structural variant (SVs) (%), across 5 \times , 10 \times , and 20 \times coverages for simulated and real-world datasets.

results both in terms of F1-Scores and concordance across coverages.

3.5 | Performance of SV callers on real-world data

While tool performance on simulated data provides a useful guide, real-world datasets usually provide additional unaccounted-for complexity and challenges. After finding the best combinations in simulated data, we investigated whether the pattern would be similar in real-world datasets. Since for the real-world data we do not have an objective truth set, they were only evaluated from two perspectives which are the congruence of results when using the same tool combination across different coverages and when using different tool combinations within the same coverage.

3.5.1 | Performance on *B. napus* real-world ONT data

B. napus ONT real dataset (Real_ONT_Bn) was evaluated using the above-described strategy. Table S10 shows the number of SVs from all tested combinations at different coverages in *B. napus*. The minimap2/cuteSV and minimap2/dysgu combinations within all coverages captured the highest number of total SVs, deletions, and insertions. Overall, a higher number of deletions than insertions was detected for all aligner and SV caller combinations at different coverages. The number of overlapped SVs across coverages for the same SVs caller/aligner combinations was calculated

(Data S7). Minimap2/cuteSV combination found the highest proportion of overlapping SVs discovered at different coverages using the same combination of tools (51.53% of total SVs, 54.52% of deletions, and 47.91% of insertions), while the minimap2/Sniffles2 combination was second best, detecting overlap of 50.1% for all SVs, 54.56% for deletions, and 44.92% for insertions across coverages (Figure 7). Although the minimap2/dysgu combination found more SVs, the percentage of intersecting SV was low. NanoVar detected the lowest proportion of overlapping SVs across coverages (19.04% of total SVs, 25.07% of deletions, and 10.21% of insertions) and discovered more unique SVs. Surprisingly we noticed a high proportion of heterozygous genotypes (0/1) in SV calling results for Real_ONT_Bn, considering that the data represented a highly inbred elite line (Vollrath et al., 2021). Tables S11 and S12 show the number of SVs genotyped as homozygous and heterozygous in simulated and real-world data, respectively. As our SV filtering required the genotypes to be homozygous for the alternative allele (1/1), these heterozygous calls were removed prior to analysis. We also investigated the overlap in SV calls across different tool combinations within the same coverage (Data S8). We observed that a substantial proportion of deletions and insertions were shared by most SV callers, with the largest number of overlapping SVs at 20 \times , following minimap2 alignment.

3.5.2 | Performance on *S. lycopersicum* real-world ONT data

We performed a similar evaluation for the real-world dataset of *Solanum lycopersicum* (Real_ONT_SI). Table S13 shows

the number of SVs found from all tested combinations at different coverages. The minimap2/dysgu combinations at 5×, 10×, and 20× captured the most SVs. Additionally, for *S. lycopersicum* all tool combinations with the exception of NanoVar found more insertions than deletions at each coverage. We also calculated the number of overlapping SVs while using the same tool combination across different coverages (Data S9). Minimap2/cuteSV combination found the highest proportion of overlapping SVs; 49.34% for all SVs, 49.63% for deletions, and 49.16% for insertions, while the minimap2/Sniffles2 combination detected 47.80% for all SVs, 49.41% for deletions, and 46.61% for insertions. Even though the minimap2/dysgu combination found more SVs, the percentage of common SVs (40.82%) was low like Real_ONT_Bn data. NanoVar again detected the lowest proportion of overlapping SVs (21.57% for all SVs, 31.20% for deletions, and 12.16% for insertions), and it discovered more unique SVs like for the Real_ONT_Bn dataset (Table S14 and Figure 7). Again, we also tested overlaps between SV calls within the same coverage, but across different tool combinations (Data S10). The largest number of overlapping SVs was found at 20×, following minimap2 alignment.

3.5.3 | Performance of Minimap2 and cuteSV/Sniffles2 combination in other crops

To assess whether our observations are robust for other crops, we performed similar benchmarking analysis for maize and compared already published SV calls in soybean, discovered using a combination of NGMLR and Sniffles1, with our results obtained from minimap2/cuteSV and minimap2/Sniffles2 combinations (Lemay et al., 2022). For maize simulated data, we found that the combination of minimap2/cuteSV had the best performance for deletions while the combination of minimap2/dysgu had the best performance for insertions (Figures S2 and S3). However, as for *B. napus* and *S. lycopersicum*, minimap2/cuteSV combination had much higher overlap across coverages in real world data (Figure S4). For soybean, we found that minimap2/cuteSV and minimap2/Sniffles discovered over 3500 new SVs, while recovering a vast majority of existing calls (Figures S5–S9).

3.5.4 | The Unique features of real-world datasets

We found a surprisingly high proportion of heterozygous calls in the real-world datasets given the highly inbred nature of the material used for sequencing. A high proportion of those is therefore likely SV discovery/genotyping errors. More heterozygous calls were found in the *B. napus* than in the *S. lycopersicum* dataset. *B. napus* is an allotetraploid species,

which undergoes reciprocal and non-reciprocal homeologous exchanges (HEs; exchanges of large corresponding chromosome segments between subgenomes). Non-reciprocal HEs could potentially cause erroneous SV calls if there are HE present in the reference, but absent in the sample. As a result, reads will have no corresponding mapping location and may be mis-mapped. To test such a scenario, we used the Sim_ONT_Bn2 dataset (20×, minimap2 for mapping, and cuteSV for SV detection) and two versions of the modified Express 617 reference. In the first version, we replaced chromosome A01 by C01 (two C01 chromosomes and no A01). In the second version, we replaced chromosome C01 by A01 (two A01 chromosomes and no C01). In both cases, the use of the modified reference resulted in an increased number of heterozygous (162.3% for reference with A01 missing, and 237.1% for reference with C01 missing), but not homozygous calls across all chromosomes (Figure 8), suggesting the non-reciprocal HEs can contribute to produce erroneous heterozygous calls.

4 | DISCUSSION

Many of the SV detection tools are benchmarked primarily on human/animal datasets (Bolognini & Magi, 2021; Coster et al., 2019; Dierckxsens et al., 2021; Jiang et al., 2020, 2021; Zhou et al., 2019); however, the complexity and different SV profiles of crop plant genomes might bring unique challenges. Therefore, to guide the design of large-scale long-read re-sequencing studies, this study performed comprehensive benchmarking of popular SV calling tools with a focus on tool performance at lower sequencing coverage. For this purpose, we designed two data simulation strategies representing both unbiased and realistic benchmarking datasets reflecting structural variation for two major crops, oilseed rape (*B. napus*) and tomato (*S. lycopersicum*). We further validated our findings using maize and soybean datasets.

Four long-read aligners (minimap2, NGMLR, Ira, and Vulcan) and five SV callers (Sniffles2, SVIM, cuteSV, dysgu, and NanoVar) were tested to detect SVs, particularly deletions and insertions. Our analysis focused on deletions and insertions as they are by far the most abundant SV types. Alignment time varied widely between the four aligners, while differences in the proportion of mapped reads were moderate. As expected, higher sequencing coverage and reference genome size length increased the run time of the mapping algorithms. The real-world datasets required more time at the same coverage, which most likely reflected several factors: exclusion of unplaced contigs from simulations, higher N50 of real world reads, potential presence of homeologous exchanges in *B. napus* dataset, and additional complexity not captured in simulations. Overall, the results found minimap2 to be the best performing aligner for SV calling applications, which also

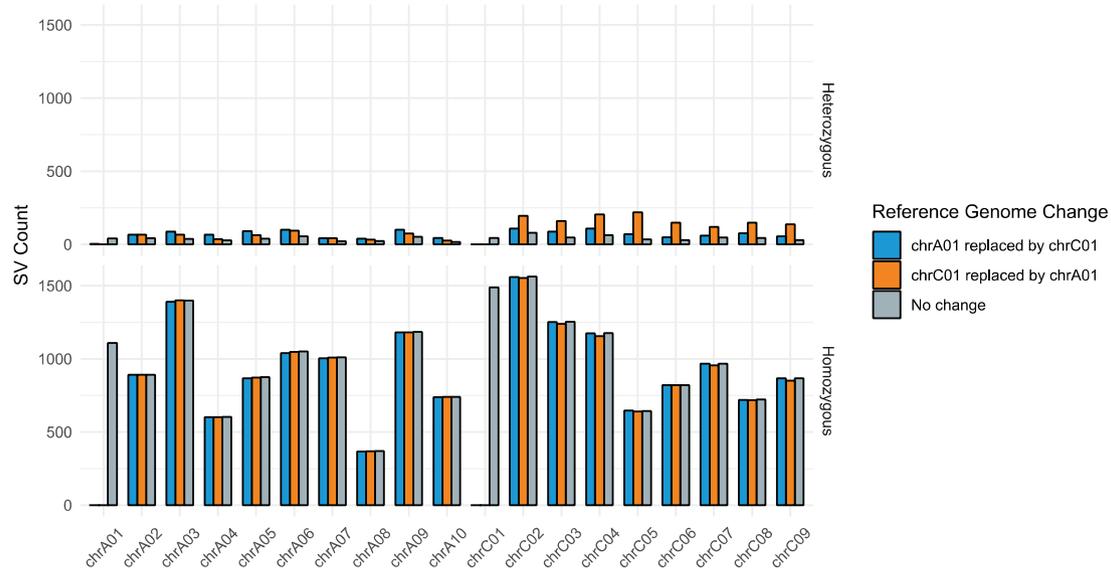


FIGURE 8 The effect of non-reciprocal homeologous exchanges on structural variant (SV) discovery. Nonreciprocal homeologous exchanges were simulated by replacing chromosome A01 by C01 and C01 by A01.

had the fastest run time and the most mapped bases. Recent benchmarking studies on human data also recommended minimap2 among tested aligners such as GraphMap, LAST, and NGMLR (Bolognini & Magi, 2021; Coster et al., 2019; Zhou et al., 2019).

We found that similar tool combinations (especially cuteSV, followed closely by Sniffles2 and dysgu after minimap2 alignment) had superior performance across all the simulated datasets. The findings are in line with a recent study reporting that cuteSV performed better than other tested SV tools such as Sniffles1, SVIM, and pbsv for precision and recall at both SV calling and genotyping in human datasets (Bolognini & Magi, 2021). Increasing coverage improved recall and F1-scores for all tested SVs calling combinations, confirming that the probability of detecting quality SVs increases with more sequencing coverage (Jiang et al., 2021). However, even at low coverages (5×) using cuteSV, Sniffles2, and dysgu for SV detection from reads aligned by minimap2 achieved > 0.8 F1-scores on simulated datasets, suggesting that Oxford Nanopore technology might be suitable for large-scale low coverage re-sequencing projects. While the lack of objective truth sets for real-world datasets precludes similar comparisons, the results revealed that tool combinations with best performance for simulated datasets also had the most consistent outcome across the range of coverages.

The criteria for filtering SV in this study were quite stringent, including retaining only SV genotyped as homozygous for alternative allele (1/1). While in simulated datasets the number of SV genotyped as heterozygous was relatively low, the proportion was much higher for real-world datasets, especially in *B. napus*. We found that in *B. napus*, the presence of homeologous exchanges will likely contribute to the erro-

neous discovery of heterozygous SV. *B. napus* is well known to harbor wide-spread nonreciprocal homeologous chromosomal exchanges even extending to whole chromosomes, for example, for chromosomes A01 and C01 as simulated here (Udall et al., 2005). The finding underlies the importance of species-specific consideration when interpreting SV discovery results. The presence of HEs likely explains only a proportion of the observed heterozygous calls and other factors need to be considered as well, including other sources of mis-mappings, genotyping errors, and residual heterozygosity in samples.

In conclusion, we found that for homozygous/inbred genotypes often used in crop studies, a substantial proportion of SVs can be discovered/genotyped at coverages as low as 5×, making Oxford Nanopore technology a suitable option for larger-scale re-sequencing studies. At this time, following our benchmarks, we recommend using the minimap2 aligner in combination with either cuteSV or Sniffles2, as it achieves good precision and recall at insertion and deletion calling and found the highest overlap between SVs across coverages.

AUTHOR CONTRIBUTIONS

Gözde Yildiz: Formal analysis; Methodology; Writing – original draft; Writing – review & editing. **Silvia F. Zanini:** Conceptualization; Methodology; Writing – original draft; Writing – review & editing. **Nazanin P. Afsharyan:** Methodology; Writing – review & editing. **Christian Obermeier:** Methodology; Writing – review & editing. **Rod J. Snowdon:** Methodology; Writing – review & editing. **Agnieszka A. Golicz:** Conceptualization; Funding acquisition; Methodology; Project administration; Supervision; Writing – original draft; Writing – review & editing.

ACKNOWLEDGMENTS

This work was supported by the Alexander von Humboldt Foundation in the framework of Sofja Kovalevskaja Award to Agnieszka A. Golicz and the German Research Foundation (DFG) project number 458716530 to Rod J. Snowdon. This work was performed with support from Justus Liebig University Bioinformatics Core Facility (BCF).

Open Access funding enabled and organized by Projekt DEAL.

CONFLICT OF INTEREST STATEMENT

The authors declare no conflict of interest.

ORCID

Gözde Yildiz  <https://orcid.org/0000-0003-0407-1829>

Silvia F. Zanini  <https://orcid.org/0000-0002-9137-8783>

Nazanin P. Afsharyan  <https://orcid.org/0000-0003-0298-988X>

Rod J. Snowdon  <https://orcid.org/0000-0001-5577-7616>

Agnieszka A. Golicz  <https://orcid.org/0000-0002-9711-4826>

REFERENCES

- Alkan, C., Coe, B. P., & Eichler, E. E. (2011). Genome structural variation discovery and genotyping. *Nature Reviews Genetics*, *12*, 363–376.
- Alonge, M., Lebeigle, L., Kirsche, M., Aganezov, S., Wang, X., Lippman, Z. B., Schatz, M. C., & Soyk, S. (2021). Automated assembly scaffolding elevates a new tomato system for high-throughput genome editing. *bioRxiv*, 2021.11.18.469135.
- Alonge, M., Wang, X., Benoit, M., Soyk, S., Pereira, L., Zhang, L., Suresh, H., Ramakrishnan, S., Maumus, F., Ciren, D., Levy, Y., Harel, T. H., Shalev-Schlosser, G., Amsellem, Z., Razifard, H., Caicedo, A. L., Tieman, D. M., Klee, H., Kirsche, M., ... & Lippman, Z. B. (2020). Major impacts of widespread structural variation on gene expression and crop improvement in tomato. *Cell*, *182*, 145–161.e23. <https://doi.org/10.1016/j.cell.2020.05.021>
- Bolognini, D., & Magi, A. (2021). Evaluation of germline structural variant calling methods for nanopore sequencing data. *Frontiers in Genetics*, *12*. <https://doi.org/10.3389/fgene.2021.761791>
- Bolognini, D., Sanders, A., Korb, J. O., Magi, A., Benes, V., & Rausch, T. (2020). VISOR: A versatile haplotype-aware structural variant simulator for short- and long-read sequencing. *Bioinformatics*, *36*, 1267–1269. <https://doi.org/10.1093/bioinformatics/btz719>
- Bornowski, N., Michel, K. J., Hamilton, J. P., Ou, S., Seetharam, A. S., Jenkins, J., Grimwood, J., Plott, C., Shu, S., Talag, J., Kennedy, M., Hundley, H., Singan, V. R., Barry, K., Daum, C., Yoshinaga, Y., Schmutz, J., Hirsch, C. N., Hufford, M. B., ... & Buell, C. R. (2021). Genomic variation within the maize stiff-stalk heterotic germplasm pool. *Plant Genome*, *14*, e20114. <https://doi.org/10.1002/tpg2.20114>
- Chawla, H. S., Lee, H., Gabur, I., Vollrath, P., Tamilselvan-Nattar-Amutha, S., Obermeier, C., Schiessl, S. V., Song, J. -M., Liu, K., Guo, L., Parkin, I. A. P., & Snowdon, R. J. (2021). Long-read sequencing reveals widespread intragenic structural variants in a recent allopolyploid crop plant. *Plant Biotechnology Journal*, *19*, 240–250. <https://doi.org/10.1111/pbi.13456>
- Cleal, K., & Baird, D. M. (2022). Dysgu: Efficient structural variant calling using short or long reads. *Nucleic Acids Research*, *50*, e53. <https://doi.org/10.1093/nar/gkac039>
- Coster, W. d., Rijk, P. d., Roeck, A. d., Pooter, T. d., D'Hert, S., Strazisar, M., Slegers, K., & van Broeckhoven, C. (2019). Structural variants identified by Oxford Nanopore PromethION sequencing of the human genome. *Genome Research*, *29*, 1178–1187. <https://doi.org/10.1101/gr.244939.118>
- Coster, W. d., Weissensteiner, M. H., & Sedlazeck, F. J. (2021). Towards population-scale long-read sequencing. *Nature Reviews Genetics*, *22*, 572–587. <https://doi.org/10.1038/s41576-021-00367-3>
- Delahaye, C., & Nicolas, J. (2021). Sequencing DNA with nanopores: Troubles and biases. *PLoS ONE*, *16*, e0257521. <https://doi.org/10.1371/journal.pone.0257521>
- Dierckxsens, N., Li, T., Vermeesch, J. R., & Xie, Z. (2021). A benchmark of structural variation detection by long reads through a realistic simulated model. *Genome Biology*, *22*, 342. [10.1186/s13059-021-02551-4](https://doi.org/10.1186/s13059-021-02551-4)
- English, A. C., Menon, V. K., Gibbs, R., Metcalf, G. A., & Sedlazeck, F. J. (2022). Truvari: Refined structural variant comparison preserves allelic diversity. *BioRxiv*, 2022.02.21.481353.
- Fu, Y., Mahmoud, M., Muraliraman, V. V., Sedlazeck, F. J., & Treangen, T. J. (2021). Vulcan: Improved long-read mapping and structural variant calling via dual-mode alignment. *Gigascience*, *10*, giab063. <https://doi.org/10.1093/gigascience/giab063>
- Fuentes, R. R., Chebotarov, D., Duitama, J., Smith, S., La Hoz, J. F. d., Mohiyuddin, M., Wing, R. A., McNally, K. L., Tatarinova, T., Grigoriev, A., Mauleon, R., & Alexandrov, N. (2019). Structural variants in 3000 rice genomes. *Genome Research*, *29*, 870–880. <https://doi.org/10.1101/gr.241240.118>
- Gill, R. A., Scossa, F., King, G. J., Golicz, A. A., Tong, C., Snowdon, R. J., Fernie, A. R., & Liu, S. (2021). On the role of transposable elements in the regulation of gene expression and subgenomic interactions in crop genomes. *Critical Reviews in Plant Sciences*, *40*, 157–189. <https://doi.org/10.1080/07352689.2021.1920731>
- Goel, M., Sun, H., Jiao, W. -B., & Schneeberger, K. (2019). SyRI: Finding genomic rearrangements and local sequence differences from whole-genome assemblies. *Genome Biology*, *20*, 277. <https://doi.org/10.1186/s13059-019-1911-0>
- Hall, M. B. (2022). Rasusa: Randomly subsample sequencing reads to a specified coverage. *Journal of Open Source Software*, *7*, 3941. <https://doi.org/10.21105/joss.03941>
- Heller, D., & Vingron, M. (2019). SVIM: Structural variant identification using mapped long reads. *Bioinformatics*, *35*, 2907–2915. <https://doi.org/10.1093/bioinformatics/btz041>
- Heller, D., & Vingron, M. (2020). SVIM-asm: Structural variant detection from haploid and diploid genome assemblies. *Bioinformatics*, *36*, 5519–5521. <https://doi.org/10.1093/bioinformatics/btaa1034>
- Hosmani, P. S., Flores-Gonzalez, M., van de Geest, H., Maumus, F., Bakker, L. V., Schijlen, E., van Haarst, J., Cordewener, J., Sanchez-Perez, G., Peters, S., Fei, Z., Giovannoni, J. J., Mueller, L. A., & Saha, S. (2019). An improved de novo assembly and annotation of the tomato reference genome using single-molecule sequencing, Hi-C proximity ligation and optical maps. *BioRxiv*, 767764.
- Jain, M., Olsen, H. E., Paten, B., & Akeson, M. (2016). The Oxford Nanopore minion: Delivery of nanopore sequencing to the genomics community. *Genome Biology*, *17*, 239. <https://doi.org/10.1186/s13059-016-1103-0>

- Jeffares, D. C., Jolly, C., Hoti, M., Speed, D., Shaw, L., Rallis, C., Balloux, F., Dessimoz, C., Bähler, J., & Sedlazeck, F. J. (2017). Transient structural variations have strong effects on quantitative traits and reproductive isolation in fission yeast. *Nature Communications*, 8, 14061. <https://doi.org/10.1038/ncomms14061>
- Jiang, T., Liu, S., Cao, S., Liu, Y., Cui, Z., Wang, Y., & Guo, H. (2021). Long-read sequencing settings for efficient structural variation detection based on comprehensive evaluation. *BMC Bioinformatics [Electronic Resource]*, 22, 552. <https://doi.org/10.1186/s12859-021-04422-y>
- Jiang, T., Liu, Y., Jiang, Y., Li, J., Gao, Y., Cui, Z., Liu, Y., Liu, B., & Wang, Y. (2020). Long-read-based human genomic structural variation detection with cuteSV. *Genome Biology*, 21, 189. <https://doi.org/10.1186/s13059-020-02107-y>
- Lee, H., Chawla, H. S., Obermeier, C., Dreyer, F., Abadi, A., & Snowdon, R. (2020). Chromosome-Scale assembly of winter oilseed rape *Brassica napus*. *Frontiers in Plant Science*, 11. <https://doi.org/10.3389/fpls.2020.00496>
- Lemay, M. -A., Sibbesen, J. A., Torkamaneh, D., Hamel, J., Levesque, R. C., & Belzile, F. (2022). Combined use of Oxford Nanopore and illumina sequencing yields insights into soybean structural variation biology. *BMC Biology*, 20, 53. <https://doi.org/10.1186/s12915-022-01255-w>
- Li, H. (2018). Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics*, 34, 3094–3100. <https://doi.org/10.1093/bioinformatics/bty191>
- Mahmoud, M., Gobet, N., Cruz-Dávalos, D. I., Mounier, N., Dessimoz, C., & Sedlazeck, F. J. (2019). Structural variant calling: The long and the short of it. *Genome Biology*, 20, 246. <https://doi.org/10.1186/s13059-019-1828-7>
- Meyers, L. A., & Levin, D. A. (2006). On the abundance of polyploids in flowering plants. *Evolution; International Journal of Organic Evolution*, 60, 1198–1206.
- Ren, J., & Chaisson, M. J. P. (2021). Ira: A long read aligner for sequences and contigs. *PLoS Computational Biology*, 17, e1009078. <https://doi.org/10.1371/journal.pcbi.1009078>
- Roberts, R. J., Carneiro, M. O., & Schatz, M. C. (2013). The advantages of SMRT sequencing. *Genome Biology*, 14, 405. <https://doi.org/10.1186/gb-2013-14-6-405>
- Sedlazeck, F. J., Lee, H., Darby, C. A., & Schatz, M. C. (2018a). Piercing the dark matter: Bioinformatics of long-range sequencing and mapping. *Nature Reviews Genetics*, 19, 329–346. <https://doi.org/10.1038/s41576-018-0003-4>
- Sedlazeck, F. J., Rescheneder, P., Smolka, M., Fang, H., Nattestad, M., Haeseler, A. v., & Schatz, M. C. (2018b). Accurate detection of complex structural variations using single-molecule sequencing. *Nature Methods*, 15, 461–468. <https://doi.org/10.1038/s41592-018-0001-7>
- Song, J. -M., Guan, Z., Hu, J., Guo, C., Yang, Z., Wang, S., Liu, D., Wang, B., Lu, S., Zhou, R., Xie, W.-Z., Cheng, Y., Zhang, Y., Liu, K., Yang, Q.-Y., Chen, L.-L., & Guo, L. (2020). Eight high-quality genomes reveal pan-genome architecture and ecotype differentiation of *Brassica napus*. *Nature Plants*, 6, 34–45. <https://doi.org/10.1038/s41477-019-0577-7>
- Tao, Y., Zhao, X., Mace, E., Henry, R., & Jordan, D. (2019). Exploring and exploiting Pan-genomics for crop improvement. *Molecular Plant*, 12, 156–169. <https://doi.org/10.1016/j.molp.2018.12.016>
- Tham, C. Y., Tirado-Magallanes, R., Goh, Y., Fullwood, M. J., Koh, B. T., Wang, W., Ng, C. H., Chng, W. J., Thiery, A., Tenen, D. G., & Benoukraf, T. (2019). NanoVar: Accurate characterization of patients' genomic structural variants using low-depth nanopore sequencing. *BioRxiv*, 662940.
- Udall, J. A., Quijada, P. A., & Osborn, T. C. (2005). Detection of chromosomal rearrangements derived from homeologous recombination in four mapping populations of *Brassica napus* L. *Genetics*, 169, 967–979. <https://doi.org/10.1534/genetics.104.033209>
- Vollrath, P., Chawla, H. S., Schiessl, S. V., Gabur, I., Lee, H., Snowdon, R. J., & Obermeier, C. (2021). A novel deletion in FLOWERING LOCUS t modulates flowering time in winter oilseed rape. *Theoretical and Applied Genetics*, 134, 1217–1231. <https://doi.org/10.1007/s00122-021-03768-4>
- Wenger, A. M., Peluso, P., Rowell, W. J., Chang, P. -C., Hall, R. J., Concepcion, G. T., Ebler, J., Fungtammasan, A., Kolesnikov, A., Olson, N. D., Töpfer, A., Alonge, M., Mahmoud, M., Qian, Y., Chin, C. -S., Phillippy, A. M., Schatz, M. C., Myers, G., DePristo, M. A., ... & Hunkapiller, M. W. (2019). Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nature Biotechnology*, 37, 1155–1162. <https://doi.org/10.1038/s41587-019-0217-9>
- Yildiz, G., Zanini, S. F., Knight, P., & Golicz, A. A. (2022). *Pangenomics in agriculture*. *CABI Biotechnology Series*. CABI.
- Yuan, Y., Bayer, P. E., Batley, J., & Edwards, D. (2021). Current status of structural variation studies in plants. *Plant Biotechnology Journal*, 19, 2153–2163. <https://doi.org/10.1111/pbi.13646>
- Zanini, S. F., Bayer, P. E., Wells, R., Snowdon, R. J., Batley, J., Varshney, R. K., Nguyen, H. T., Edwards, D., & Golicz, A. A. (2022). Pangenomics in crop improvement—from coding structural variations to finding regulatory variants with pangenome graphs. *Plant Genome*, 15, e20177. <https://doi.org/10.1002/tpg2.20177>
- Zhang, F., Xue, H., Dong, X., Li, M., Zheng, X., Li, Z., Xu, J., Wang, W., & Wei, C. (2022). Long-read sequencing of 111 rice genomes reveals significantly larger pan-genomes. *Genome Research*, 32, 853–863.
- Zhou, A., Lin, T., & Xing, J. (2019). Evaluating nanopore sequencing data processing pipelines for structural variation identification. *Genome Biology*, 20, 237. <https://doi.org/10.1186/s13059-019-1858-1>

SUPPORTING INFORMATION

Additional supporting information can be found online in the Supporting Information section at the end of this article.

How to cite this article: Yildiz, G., Zanini, S. F., Afsharyan, N. P., Obermeier, C., Snowdon, R. J., & Golicz, A. A. (2023). Benchmarking Oxford Nanopore read alignment-based insertion and deletion detection in crop plant genomes. *The Plant Genome*, 16, e20314. <https://doi.org/10.1002/tpg2.20314>