



Abstract

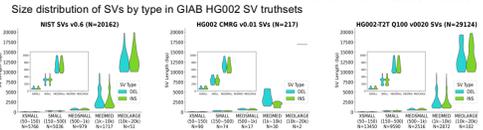
Structural variants (SVs) are an important class of genomic variation with significant functional and clinical impact, yet their discovery remains an open problem. The performance of SV callers can vary as a function of variant type, size, genomic context, coverage, and sequencing platform. However, SV truthsets are scarce and typically contain a restricted set of SV types with skewed size and context distributions, hindering our ability to uncover systematic performance bias. Here, we introduce VARium: an extensive suite of synthetic genomes designed to assess the performance of SV discovery tools as a function of key domain parameters and confounders. We evaluated 22 SV callers, including alignment-based and assembly-based tools, on each genome from VARium using multiple sequencing platforms and varying read depth. As expected, we found that both recall and precision of most methods varied significantly across different benchmark settings, demonstrating that synthetic benchmarks can serve as a key tool for revealing systematic biases and limitations of SV methods, while providing upper-bound performance estimates for real datasets.

Motivation

Real SV benchmarks are scarce and include only a limited set of SV types – primarily deletions (DEL) and insertions (INS) – with skewed size distributions and uneven genomic context coverage.

	Number of variants	Number of variants in repeats	Number of variants in SEGDPs
v0.6	20,162	16,080	195
CMRG	217	101	3
Q100	29,124	16,071	2,360

Number of variants stratified by genomic context in GIAB HG002 SV truthsets



VARium v0 genome collection

Tier 1 (79 genomes)

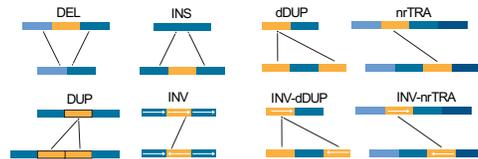
Designed to stress-test simple SV recall as a function of type, size, and context

- 4 SV types: DEL, INS, DUP, INV
- 5 size ranges: 50-150bp, 150-500bp, 500-1kbp, 1k-10kbp, 10k-20kbp, 20k-100kb, 100k-1Mb
- 6 genomic contexts: UNIQUE, SEGDP, NONUNIQUE, Aiu*, L1HS*, TR* (*DEL only)
- 3 platforms: Illumina, PacBio, ONT
- 5 coverages: 30x, 15x, 10x, 5x, 0.5x

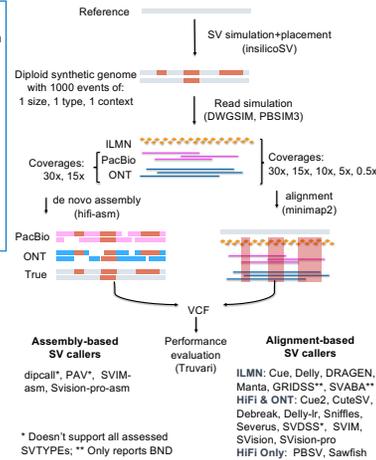
Tier 2 (40 genomes)

Designed to stress-test simple SV precision (in the presence of complex variants)

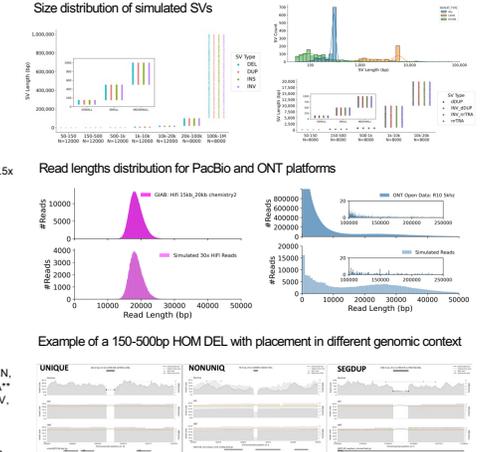
- 4 SV types: dDUP, nrTRA, INV-dDUP, INV-nrTRA
- 5 size ranges: 50-150bp, 150-500bp, 500-1kbp, 1k-10kbp, 10k-20kbp
- 2 dispersion distance regimes: 10k-50kbp, 1M-100Mbp
- 3 platforms: Illumina, PacBio, ONT
- 1 coverage: 30x



Data generation workflow

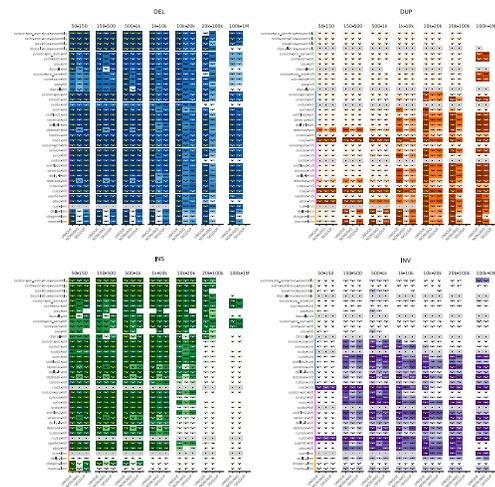


Synthetic WGS dataset characteristics

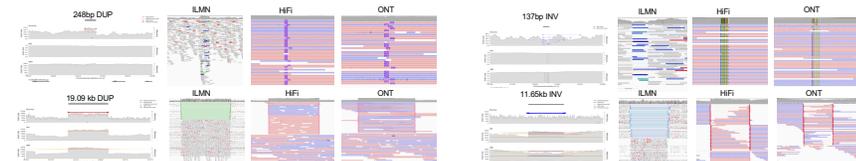


Tier 1 results

F1 scores stratified by SV type, size, context, and platform at 30x coverage

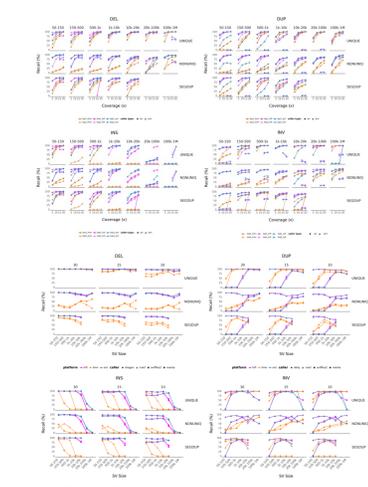


Examples of DUP and INV detection challenges across different event sizes

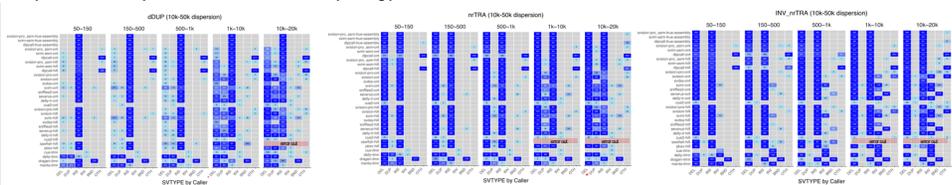


Tier 2 results

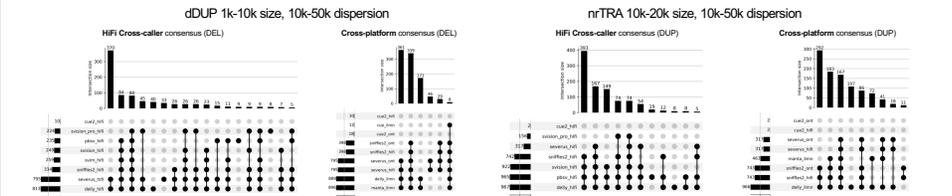
Recall stratified by SV type, size, context, platform, and coverage



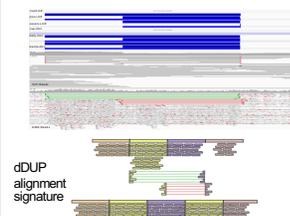
SV predictions at complex variant sites across three sequencing platforms



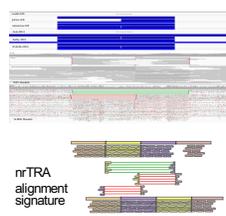
Cross-platform consensus of FP simple SV calls



Simple SV predictions at a dDUP site



Simple SV predictions at an nrTRA site



Conclusion

- 22 state-of-the-art SV callers show wide performance variability across conditions.
- Few callers achieve high recall across all type and size combinations; most perform well only in some regimes.
- Small DUPs and INVs are particularly challenging for most long-read callers due to alignment artifacts.
- Most callers incorrectly report simple events at different combinations of complex SV breakpoints. Notably, these FP calls are often shared across callers and platforms.