

| | CUTESV | | | SNIFFILES | | | Notes |
|----------------------------|----------------------------|---------|--|----------------------------|--------|--|--|
| Common parameters | | | | -i | | | For single-sample calling: a coordinate-sorted and indexed .bam/cram. For multi-sample calling: multiple .snf generated before by running Sniffles2 for individuals samples with -snf |
| | | | | -v | | | VCF output filename. If filename ends with .gz, it will be automatically gzipped. |
| | | | | -snf | | | |
| | | | | -reference | .fasta | | |
| | | | | -tandem-repeat | .bed | | Bed containing tandem repeat annotation for the reference genome. |
| | | | | -non-germline | FALSE | | |
| | | | | --phase | FALSE | | Require input to be phased |
| | -t, -threads | | | -t, -threads | | | |
| | -b, -batches | 1000000 | Batch of genome segmentation interval | | | | window size for binning sequencing depth |
| | -S, -sample | | Sample name/id. | | | | |
| | -retain_work_dir | FALSE | Enable to retain temporary folders and files | | | | |
| | -report_readid | FALSE | Enable to report supporting read ids for each SV | | | | |
| SV filtering parameters | -p, -max_split_parts | 7 | Maximum number of split segments a read may be aligned before it is ignored. All splits are considered of 1. Recommended for assembly-based alignment. | --max-splits-kb | 0.1 | | Additional number of splits per kilobase read sequence allowed before reads are ignored |
| | | | | --max-splits-base | 3 | | Base number of splits allowed before reads are ignored (in addition to --max-splits-kb) |
| | | | | --phase-conflict-threshold | 0.1 | | Maximum fraction of conflicting reads permitted for SV phase information to be labelled as PASS (only for --phase) (default: 0.1) |
| | | | | --detect-large-ins | TRUE | | Infer insertions that are longer than most reads and therefore are spanned by few alignments only. |
| | -q, -min_mapq | 20 | | -mapq | 25 | | Minimum mapping quality value to be taken in account |
| | | | | --no-qc | FALSE | | Output all SV candidates, disregarding quality control steps. |
| | | | | --qc-sddev | TRUE | | Apply filtering based on SV start position and length standard deviation |
| | | | | --qc-sddev-abs-max | 500 | | Maximum standard deviation for SV length and size (in bp) |
| | | | | --qc-strand | FALSE | | Apply filtering based on strand support of SV calls |
| | -r, -min_read_len | 500 | Ignores reads that only reports alignment with not longer than [value] bp | --min-alignment-length | 1000 | | Reads with alignments shorter than this length (in bp) will be ignored |
| | -md, -merge-del-threshold | | Maximum distance of deletion signals to be merged. In their paper, they used -md 500 to process HG002 real human sample data.[0] | | | | |
| | -mi, -merge_ins_thresho ld | | Maximum distance of insertion signals to be merged. They also used -mi 500. | | | | |
| | | | | --qc-coverage | 1 | | Minimum surrounding region coverage of SV calls |
| Generation of SV clusters: | -s, -min_support | 10 | | -minsupport | Auto | | Minimum number of reads that support a SV to be reported |
| | | | | | | | with the default setting of --min_size (-min_size = 30), cuteSV achieved the best yields when --min_support was configured from s=4 to 10. And there is an obvious trade-off between precision and recall, that is, setting a smaller --min_support value might result in higher sensitivity but lower precision, and vice versa. Should be between 1/3 and 1/4 of the median coverage. |
| | -l, -min_size | 30 | minimal size of SV signature considered in clustering | -minsvlen | 25 | | In cuteSV setting --min_size with smaller numbers might result in higher sensitivity but lower accuracy, and vice versa. It is also worth noting that, although the trade-off exists, for each coverage, the F1 scores of cuteSV with various settings are quite close to each other (the difference is less than 1%). |
| | -L, -max_size | 10000 | All SVs are reported if -l | | | | |
| | -sl, -min_siglength | 10 | Minimum length of SV signal to be extracted. | --minsvlen-screen-ratio | 0.9 | | Minimum length for SV candidates (as fraction of --minsvlen) |
| | diff_ratio_mergin g_INS | 0.3 | Do not merge breakpoints with basepair identity more than [value] for insertion. | | | | reads spanning the same SV usually have heterogeneous breakpoints in their alignments, which also cause false-positive SV calls. In step 2 of cuteSV, the spe- cifically designed clustering-and-refinement approach enables to adaptively cluster alignment breakpoints mapped to relatively large local regions but potentially belonging to identical SVs, so that heterogeneous breakpoints can be merged more effectively and more false positives can be prevented. See this foldid page ; for breakpoint identification and sequence identity. As breakpoint occurs in regions of low sequence identity, breakpoints identification are often challenging. To overcome that, breakpoint with low sequence identity in local regions are merged together as one unique breakpoint. |
| | -diff_ratio_mergin g_DEL | 0.3 | Do not merge breakpoints with basepair identity more than [0.5] for deletion. | | | | |
| | -diff_ratio_filter ing_TRA | 0.6 | Filter breakpoints with basepair identity less than [0.6] for translocation. | | | | |
| | | | | --long-ins-length | 2500 | | Insertion SVs longer than this value are considered as hard to detect based on the aligner and read length and subjected to more sensitive filtering. |
| | -max_cluster_bias_DEL | 100 | Maximum distance to cluster read together for deletion | --cluster-merge-pos | 150 | | Max. Distance to cluster reads for insertions and deletions on the same read and cluster in non-repeat regions |
| | -max_cluster_bias_INS | 100 | Maximum distance to cluster read together for insertion . | | | | In the first step of cuteSV, insertions/deletions in nearby genomic regions are combined to unbroken signatures of larger SVs. Allow to reduce the errors caused by the fragile read alignments, but also enables to produce more homogenous SV signatures from various reads, which is beneficial to the processing of later steps. |
| | | | | --long-del-length | 50000 | | Deletion SVs longer than this value are subjected to central coverage drop-based filtering |
| | | | | --long-del-coverage | 0.66 | | Long deletions with central coverage (in relation to upstream/downstream coverage) higher than this value will be filtered (Not applicable for --non-germline) |
| | --max_cluster_bias_DUP | 500 | Maximum distance to cluster read together for duplication | | | | No del with central coverage higher than [0.66], for deletion as defined in --long-del-length |
| | | | | --long-dup-length | 50000 | | Duplication SVs longer than this value are subjected to central coverage increase-based filtering (Not applicable for --non-germline) |
| | | | | --long-dup-coverage | 1.33 | | Long duplications with central coverage (in relation to upstream/downstream coverage) lower than this value will be filtered (Not applicable for --non-germline) |
| | | | | --cluster-binsize | 100 | | Initial screening bin size in bp |
| | | | | --cluster-r | 2.5 | | Multiplier for SV start position standard deviation criterion in cluster merging |
| | | | | --cluster-repeat-h | 1.5 | | Multiplier for mean SV length criterion for tandem repeat cluster merging |
| | | | | --cluster-repeat-h-max | 1000 | | merging distance based on SV length criterion for tandem repeat cluster merging |
| | | | | --cluster-merge-len | 0.33 | | Max. size difference for merging SVs as fraction of SV length |
| | | | | --cluster-merge-bnd | 1500 | | Max. merging distance for breakend SV candidates. |
| | -max_cluster_bias_INV | 500 | Maximum distance to cluster read together for inversion | | | | |
| | -max_cluster_bias_TRA | 50 | Maximum distance to cluster read together for translocation | | | | |
| Computing genotypes: | --genotype | | | --genotype-vcf | None | | Determine the genotypes for all SVs in the given input .vcf file. Re-genotyped .vcf will be written to the output file specified with --vcf. |
| | --gt_round | 500 | Maximum round of iteration for alignments searching if perform genotyping. | | | | |
| | | | | --genotype-ploidy | 2 | | Sample ploidy |
| | | | | --genotype-error | 0.05 | | Estimated false positive rate for leads |

| | CUTESV | | | SNIFFLES | | | Notes |
|--|--------|--|--|------------------------------|-------|--|-------|
| Force calling: | -lvcf | | Optional given vcf file. Enable to perform force calling | | | | |
| Multi-Sample Calling / Combine parameters: | | | | --combine-high-confidence | 0.0 | Minimum fraction of samples in which a SV needs to have individually passed QC for it to be reported in combined output (a value of zero will report all SVs that pass QC in at least one of the input samples) | |
| | | | | --combine-low-confidence | 0.2 | Minimum fraction of samples in which a SV needs to be present (failed QC) for it to be reported in combined output | |
| | | | | --combine-low-confidence-abs | 2 | Minimum absolute number of samples in which a SV needs to be present (failed QC) for it to be reported in combined output | |
| | | | | --combine-null-min-coverage | 5 | Minimum coverage for a sample genotype to be reported as 0/0 (sample genotypes with coverage below this threshold at the SV location will be output as /.) | |
| | | | | --combine-match | 250 | Multiplier for maximum deviation of multiple SV's start/end position for them to be combined across samples. Given by $\text{max_dev} = M \cdot \sqrt{\min(\text{SV_length_a}, \text{SV_length_b})}$, where M is this parameter. | |
| | | | | --combine-match-max | 1000 | Upper limit for the maximum deviation computed for --combine-match in bp | |
| | | | | --combine-separate-intra | FALSE | Disable combination of SVs within the same sample | |
| | | | | --combine-output-filtered | FALSE | Include low-confidence / putative non-germline SVs in multi-calling | |
| | | | | --output-rnames | | Output names of all supporting reads for each SV | |
| | | | | --no-consensus | | Disable consensus generation for insertion SV calls (may improve performance) | |
| | | | | --no-sort | FALSE | Do not sort output VCF by genomic coordinates (may slightly improve performance) | |
| | | | | --no-progress | FALSE | Disable progress display | |
| | | | | --quiet | FALSE | Disable all logging, except errors | |
| | | | | --max-del-seq-len | 50000 | Maximum deletion sequence length to be output. Deletion SVs longer than this value will be written to the output as symbolic SVs | |
| | | | | --symbolic | FALSE | Output all SVs as symbolic, including insertions and deletions, instead of reporting nucleotide sequences. | |
| | | | | --combine-consensus | FALSE | Output the consensus genotype of all samples | |