Supplementary Table 1. **Rationale for the selection or omission of computational polyTE detection tools for this benchmark study and their relevance to human next-generation sequencing (NGS) data.** Extensive benchmarking was done on seven tools that were selected based on the criteria adopted in this study (see "Polymorphic TE detection tools" section). Additionally, four more previously not included polyTE detection tools were tested on the low coverage dataset. Other existing polyTE detection tools that were omitted from the benchmark are also listed along with the rationale of their omission. Briefly, tools that are not specialized for polyTE detection or requires specific TSDs were not included in the benchmark.

| PolyTE detection tools selected for benchmarking | | | |
|---|---|---|---|
| Tool name | Rationale for selection | Tool's success | Relevance to human NGS data |
| MELT | All criteria | Success | High |
| Mobster | All criteria | Success | High |
| RetroSeq | All criteria | Success | High |
| Tangram | All criteria | Failure | High |
| TEMP | All criteria | Success | High |
| ITIS | Criterion #1 and #4 | Success | Medium |
| T-lex/T-lex2 | Criterion #1 and #4 | Aborted | Medium |
| DD_DETECTION | Expanded set | Failure | High |
| Jitterbug | Expanded set | Failure | High |
| TE-Locate | Expanded set | Failure | Medium |
| TE-Tracker | Expanded set | Failure | Medium |
| PolyTE detection tools omitted from the benchmarking | | | |
| Tool name | Rationale for omission | | Relevance to human NGS data |
| GRIPper | Detects non-reference gene copy insertion | | High |
| TIGRA | Breakpoint assembler, not an SV caller | | High |
| TranspoSeq | Requires paired tumor/normal WGS data | | High |
| Tea | Requires paired tumor/normal WGS data | | High |
| TraFiC | Requires paired tumor/normal WGS data | | High |
| VariationHunter | General purpose SV detection tool | | High |
| HYDRA-SV | General purpose SV detection tool | | High |
| MetaSV | General purpose SV detection tool | | High |
| ngs_te_mapper | Requires TSDs to be provided | | Medium |
| RelocaTE | Requires TSDs to be provided | | Low |

Supplementary Table 2.  **Summary of algorithmic differences between the computational polyTE detection tools benchmarked in this study**.  More detailed differences are listed in Supp. Table 3.

| Tool | Read mapping | | | | Breakpoint estimation | | | Filtering criteria | | | | Output features | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | All DP reads | Treats SR & DP independently | SR searched after DP | Mobilome Aligner | Fragment size distribution | SR dependent | Holistic | Read depth - flanking | Read depth - site | Known TEs | Mapping quality | VCF file | Predicts TSD | Predicts zygosity |
| MELT | ✓ | ✗ | ✓ | Bowtie2 | ? | ✗ | ✓ | ✓ | ✓ | ✓ | ? | ✓ | ✓ | ✓ |
| Mobster | ✓ | ✓ | ✗ | MOSAIK | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✗ | ✗ | ✗ |
| RetroSeq | ✓ | ✗ | ✓ | Exonerate | ✗ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✓ | ✗ | ✓ |
| Tangram | ✓ | ✓ | ✗ | MOSAIK | ✓ | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ | ✗ | ✗ |
| TEMP | ✓ | ✗ | ✓ | BWA | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✗ | ✗ | NA |
| ITIS | ✗ | ✓ | ✗ | BWA | ✗ | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ | ✗ | ✗ |

Supplementary Table 3. **Detailed Algorithmic differences between the computational polyTE detection tools benchmarked in this study**.

| Tool | DP definition | SR definition | DP and SR search paradigm | Mobilome alignment tool | Cluster definition | Merging clusters | Breakpoint estimation | Filtering criteria |
|---|---|---|---|---|---|---|---|---|
| MELT | Information not available | Information not available | DPs were used to identify potential/ candidate TE sites<br><br>SRs were used to identify breakpoints and TSDs | Bowtie2 with default parameters | Sites with at least 4 DP anchors clustered within 500bp of each other | Merges all DP and SR clusters from all BAM files (from 1KG project) | Unspecified type of model was built containing all available information for the candidate site. This model was then used to predict precise insertion site, strand, TSD, insertion sequence and length. | Based on:<br>1) minimum 4 supporting DPs<br>2) proximity to a reference TE<br>3) filter sites with depth of coverage outside 70-130% of the 100bp flanking region |
| Mobster | 1) Orientation different from expectation or<br>2) Distance between pairs significantly different or<br>3) Reads mapping to different chromosome or<br>4) One read mapped, other not<br><br>DP will have at least one uniquely mapping read referred to as the anchor read | Reads that map partially (clipped); will have one uniquely mapping anchor read and uniquely mapping unclipped part (anchor for SR) | DP and SR are searched independently<br><br>Anchor reads tagged as unmapped or by the TE family their mate/clipping maps to | MOSAIK (hash size = 9; max mismatches = 10%, min length = 20 bp) | DP clusters<br>1) Anchors map to same strand<br>2) support the same TE family<br>3) have start position in proximity to each other<br><br>SR clusters<br>1) Anchors belong to the same TE family or polyA/T stretch<br>2) same side clipping<br>3) clipped within a few bp of each other | Merge same family (or homopolymer) forward and reverse strand clusters<br><br>First merge DP and SR independently, then proceed to merge the two<br><br>Confidence assigned based on the number of clusters and orientations (5' and 3') that were merged | Breakpoints are estimated based on the inner borders of 5' and 3'clipped<br><br>If clipped reads not available, inner borders of DP clusters are used for breakpoint estimation<br><br>Else, estimated from insert size distribution and cluster length | Based on:<br>1) proximity (within 90bp) to a reference TE<br>2) user controlled read depth based filtering |
| RetroSeq | 1) SAM flag 0x0002 unset, i.e., reads that are not proper pairs, or<br>2) One mate of the pair is unmapped<br><br>Proper pairs are defined as reads whose pair maps within the expected distance | Partially mapped reads | Extracts DP in the beginning<br><br>SR are only searched for breakpoint estimation step | Exonerate (80% min identity, 36 bp min length, mapping quality 30, local alignment with affine gap penalty, report best 5 results) | Forward and reverse orientation clusters created by the start position of the anchor reads<br><br>Max gap 120bp between reads in a cluster | Uses bedtools window command to merge forward and reverse clusters | Excludes clusters with average read depth surrounding the cluster above a cutoff (def 200).<br><br>Estimates using a set of parameters: 1) read depth of DP on both strands, 2) forward to reverse reads ratio at 5' and 3' of the putative breakpoint and 3) distance between last 5' and first 3' read. | Based on:<br>1) proximity (within 100bp of an Alu or within 200bp of an L1) to a reference TE<br><br>Confidence for each genotype provided in the output VCF file |

| Tool | DP definition | SR definition | DP and SR search paradigm | Mobilome alignment tool | Cluster definition | Merging clusters | Breakpoint estimation | Filtering criteria |
|------|---------------|---------------|---------------------------|-------------------------|--------------------|------------------|-----------------------|--------------------|
| Tangram | Utilizes customized BAM format that contains both the genome and TE reference sequence alignment (No instruction provided on generating this alignment)<br><br>DP are read pairs with one read mapping uniquely to the reference genome and the other mapping to the TE reference sequence | SR have one mate mapping uniquely while the other is either soft-clipped or unmapped<br><br>The unaligned or soft-clipped reads are then realigned to both reference genome and reference TE sequences | DP and SR are searched independently | MOSAIK (Done before the process begins) | Clusters candidate read pairs using fragment center position; applies a customized nearest-neighbor algorithm for clustering<br><br>Utilizes fragment length distribution<br><br>Capable of handling multiple different libraries | | DP – identifies pair of clusters spanning on the insertion from 3' and 5'.  Leftmost position candidate insertion position<br><br>SR – Performs fast local alignment to identify the breakpoint | Based on 1) supporting reads per insertion, 2) additional filtering if only DP support and 3) proximity to a reference TE |
| TEMP | One uniquely mapping read (anchor read) and the other read that maps to multiple distant locations or is unmappable | Reads that start in genomic sequence but are interrupted by transposon or non-contiguous genomic sequence<br><br>Clipped portion which maps to the TE should be at least 7bp long | SR are looked for after DP<br><br>DP identifies insertion regions, SR helps in breakpoint estimation | BWA-aln and BWA-sampe | Defines intervals such as they contain TE "junction in the beginning (and) at the end of the anchor read, and extending into the genome by the length of the average insert size"<br><br>Reads supporting same TE type, same orientation and intervals that overlap by at least 1nt are clustered | | Extends intervals in both directions to find overlapping SRs. Estimates breakpoint based on the clipped portion.  In case of multiple locations, selects the one with highest support.<br><br>When base estimate is not available, interval midpoint is taken as insertion position | Based on 1) Read depth |
| ITIS | One end mapped to the reference genome, other mapped to the TE | At least one end covers both reference and TE | DP and SR are searched independently but SR determines the genomic location | BWA | DP and SR that are in close proximity | | Based on the SR | Based on 1) MAQ > 0 2) 2 < RD < 300 – around insertion 3) RD (DP/SR) > 2 |