

1 **HiTea: a computational pipeline to identify non-reference** 2 **transposable element insertions in Hi-C data**

3 Dhawal Jain¹, Chong Chu¹, Burak Han Alver¹, Soohyun Lee¹, Eunjung Alice Lee^{2,3} and Peter J.
4 Park¹

5 ¹Department of Biomedical Informatics, Harvard Medical School, Boston, MA 02115, USA.

6 ²Division of Genetics and Genomics, Boston Children's Hospital and Harvard Medical School,
7 Boston, MA 02115, USA.

8 ³Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA.

9

10 **Abstract**

11 Hi-C is a common technique for assessing three-dimensional chromatin conformation. Recent
12 studies have shown that long-range interaction information in Hi-C data can be used to generate
13 chromosome-length genome assemblies and identify large-scale structural variations. Here, we
14 demonstrate the use of Hi-C data in detecting mobile transposable element (TE) insertions
15 genome-wide. Our pipeline HiTea (**Hi-C based Transposable element analyzer**) capitalizes on
16 clipped Hi-C reads and is aided by a high proportion of discordant read pairs in Hi-C data to
17 detect insertions of three major families of active human TEs. Despite the uneven genome
18 coverage in Hi-C data, HiTea is competitive with the existing callers based on whole genome
19 sequencing (WGS) data and can supplement the WGS-based characterization of the TE insertion
20 landscape. We employ the pipeline to identify TE insertions from human cell-line Hi-C samples.
21 HiTea is available at <https://github.com/parklab/HiTea> and as a Docker image.

22 **Keywords**

23 retrotransposons, structural variants, Hi-C, split read analysis

24 INTRODUCTION

25 Over half of the human genome is composed of repetitive DNA sequences(de Koning *et al.*,
26 2011). The repeats belong to two major classes: (i) tandem repeats, consisting of DNA
27 sequences from few bases to few hundreds of bases that have expanded in tandem, stretching up
28 to millions of bases in the genome; and (ii) transposable elements (TEs), interspersed throughout
29 the genome and accounting for 44% of the human genome(Mills *et al.*, 2007). Unlike tandem
30 repeats, TEs are capable of transposition, in which they move from one genomic location to
31 another. The distinct self- or *trans*- encoded mechanisms used by the TEs for transposition are
32 used to group them into several families(Wicker *et al.*, 2008). Although a vast majority of the
33 TEs are inactive, a small fraction (<0.05%) still remains active in the human genome(Mills *et al.*,
34 2007), primarily SINEs (Small Interspersed Nuclear Elements), LINEs (Long Interspersed
35 Nuclear Elements), and SVAs (SINE-VNTR-*Alu*).

36 The transposition events are a major source of genomic structural variation (SV) and play an
37 important role in a multitude of human genetic diseases(Hancks and Kazazian, 2016). For
38 example, elevated levels of non-reference L1Hs (LINE) insertions are associated with epithelial
39 carcinomas(Hancks and Kazazian, 2016; Lee *et al.*, 2012; Chenais, 2015); *Alu* (SINE) insertions
40 are associated with cystic fibrosis and hemophilia(Chen *et al.*, 2008; Vidaud *et al.*, 1993); and a
41 recent case of Batten's disease that led to the development of an individualized antisense
42 oligonucleotide therapy(Kim *et al.*, 2019) was caused by an SVA insertion. The TE sequences
43 may also encode a range of regulatory features such as promoters, enhancers, transcription factor
44 binding sites, and non-coding regulatory RNA transcripts(Chuong *et al.*, 2017). Thus at the
45 molecular level, transposition can result in altered gene expression, splicing/RNA stability
46 defects, genome instability, or decreased integrity of centromere and telomeres(Bourque *et al.*,
47 2018).

48 In particular, TE sequences are a rich source of binding sites for an insulator protein CTCF,
49 which plays a key role in regulating the 3D structure of chromatin. The extended loops of the
50 DNA are maintained by binding of CTCF at the base of the loop; indeed, the Hi-C chromatin
51 maps suggest enrichment of SINE elements at the topologically associated domains (TAD)
52 boundaries(Rao *et al.*, 2014). The TE-derived CTCF binding sites are a fundamental source for
53 mammalian genome evolution at various time scales, with some highly conserved across species

54 and some species-specific expansions of CTCF sites co-occurring with species-specific
55 TADs(Schmidt *et al.*, 2012; Cournac *et al.*, 2016). Given the important regulatory role of
56 TEs(Ayarpadikannan and Kim, 2014; Garcia-Perez *et al.*, 2016; Ahmed and Liang, 2012),
57 identification of their transposition is important in understanding the disease biology, gene
58 regulation, and 3D chromatin organization.

59 Several computational tools are available for identifying non-reference (either somatic and
60 germline) TE insertions from WGS data(Rishishwar *et al.*, 2017). A key component of such
61 methods is the identification of discordant read pairs (RP), whose genome alignments display
62 unexpected between-pair distance or orientation. A discordant RP with one end mapping to the
63 consensus TE sequence and the other end mapping to the reference genome is indicative of a TE
64 insertion. Discordant RPs are typically accompanied by ‘clipped’ reads, whose partial alignment
65 can be used to obtain base-pair resolution of the breakpoints. With judicious integration of these
66 criteria and appropriate thresholds, candidate TEs insertions can be predicted across genome.

67 Besides WGS, another data type that involves a large amount of sequencing is Hi-C, an unbiased
68 genome-wide extension of the chromosome conformation capture technique. Hi-C
69 experiments(Rao *et al.*, 2014; Schmitt *et al.*, 2016) are conducted primarily to understand the
70 long-distance regulatory relationships in the genome (e.g., which enhancer interacts with which
71 promoter). In this experiment, the cross-linked DNA fragments are first digested with a suitable
72 restriction endonuclease (RE). Then, random ligation is performed in a condition that favors
73 ligation between cross-linked fragments. The resulting ligation product contains pairs of
74 fragments that were close in 3D proximity. Sequenced Hi-C reads indeed show that the effective
75 insert sizes—the distance between the mapped mates—range from few hundred to millions of
76 bases. Consequently, the proportion of discordant RPs, that are <20% in WGS, are in the excess
77 of 50-70% for Hi-C data. Furthermore, as the sequenced fragments are generated post-ligation
78 step, the proportion of reads carrying split mapping (due to encompassed RE sites) is higher in
79 the Hi-C data. These features thus limit the use of WGS-based TE detection tools on Hi-C data.

80 Here, we present a computational pipeline HiTea (**Hi-C based Transposable element analyzer**),
81 which identifies non-reference TE insertions of the LINE, SINE and SVA families using Hi-C
82 data. Our comparisons show that HiTea (run on Hi-C) performs similarly to a commonly-used
83 WGS-based tool (run on WGS at similar coverage)(Gardner *et al.*, 2017). With increasing

84 realization of 3D chromosomal structure as a regulatory component of gene regulation, large
85 scale efforts such as 4D Nucleome(Dekker *et al.*, 2017) are underway to aim to map genome
86 organization across cell-types and disease models. Our results indicate that Hi-C data can be
87 used not only to study 3D genome organization but also to characterize the non-reference TE
88 insertions.

89 **METHODS**

90 **Informative Hi-C read pairs for non-reference TE detection**

91 To understand the methodology underlying HiTea, we first describe the different types of read
92 pair (RP) mappings observed in Hi-C data (Fig.1A). Discordant RPs, defined in paired-end
93 sequencing, are RPs with unexpected distance or orientations between paired mate reads when
94 mapped to the reference genome. Due to the intrinsic design of Hi-C experiments for detecting
95 interactions between two distant genomic loci, a major proportion of RPs (typically 50-70%) in
96 Hi-C data are discordant with large (>1kb) mapping distances or atypical orientations of the
97 paired mates. A small proportion (6-30%) of RPs display WGS-like concordant read mapping
98 configuration (Fig. 1A, panel i), where both mates map close (< 500bp) to each other in
99 convergent orientation.

100 The RPs in Hi-C data can also be classified into two different categories. First, we introduce the
101 terminology *conforming* RPs to refer to those with mapping configuration explained solely by
102 the Hi-C experiment. For instance, conforming RPs with unique mapping of the entire mate
103 reads on two proximal or distant genomic loci are prevalent in Hi-C data (Fig. 1A-i,ii). Here, the
104 between-pair distance can range from WGS-like insert size (*i.e.*, ~500bp) to millions of bases
105 (Fig. 1B). A third type of conforming RPs are those in which the 5' portion of a mate maps
106 uniquely to the genome and the 3' portion maps convergent on the genomic locus of the
107 matching mate, and the two portions are connected with the RE ligation motif (Fig.1A-iii). These
108 mappings are referred to as chimeric Hi-C pairs (~10-20%) and are included in the 3D-contact
109 matrices. Second, the remaining RPs (~10-30%) do not conform to any expected configuration
110 of read mappings, and thus are discarded in standard analyses. In those *non-conforming* RPs, one
111 or both mates remain unmapped, multi-mapped, or their partial mapping does not produce
112 chimeric Hi-C pairs (Fig. 1A-iv,v,vi). To identify non-reference TE insertions, HiTea uses non-

113 conforming RPs whose partial (clipped) sequences or one entire mate read map to TE sequence
114 assemblies.

115 In Fig. 1C, we show the distribution of reads along a small genomic region. In WGS data, the
116 genomic coverage is relatively even. In Hi-C data, the coverage is more variable; however, much
117 of the region is still covered with at least some reads, thus allowing for the possibility that most
118 TE insertions can be captured. The proportion of discordant RPs (non-gray colors) is very high in
119 Hi-C data.

120 **Identification of TE insertion breakpoints**

121 HiTea starts by identifying non-conforming RPs using Pairtools
122 (<https://github.com/mirnylab/pairtools>). In the discovery step, the clipped reads without
123 legitimate RE-ligation motif are then mapped (using BWA-MEM(Li and Durbin, 2010) with '-a
124 -k 13 -T 20') to family-wise TE consensus assemblies published earlier(Gardner *et al.*, 2017) for
125 Alu (SINE), L1Hs (LINE) and SVA (<https://melt.igs.umaryland.edu/downloads.php>).
126 Additionally, it uses a separate 200 base long PolyA sequence to improve detection sensitivity of
127 TEs, especially those with long PolyA tails. For the alignment, we note that many polymorphic
128 insertions may have sequences distinct from the family-based consensus. To accommodate such
129 cases, HiTea offers an option to remap clipped reads that initially fail to map to a TE family
130 consensus, to a user-provided set of polymorphic sequences for a TE-family or sequences of the
131 members of its subfamily (e.g., from Repbase(Bao *et al.*, 2015)). HiTea, in principle, can also
132 detect insertions of other template-based transposons such as an active human endogenous
133 retrovirus (HERV-K), as long as adequate TE-consensus sequences are provided.

134 The clipped sequences are derived from non-conforming Hi-C RPs, where minimum clip length
135 (default: -s 20) can be defined by the users. Using a two base-pair leeway, a breakpoint on the
136 reference genome is determined as the location with the maximum number of clipped reads at a
137 locus (Supl.Fig.1). HiTea simultaneously records all non-conforming RPs in which a read maps
138 to the reference genome and its 'anchor' mate maps to the TE-consensus assembly (using default
139 BWA-MEM settings). We refer to these as **Repeat-Anchored non-conforming Hi-C Mates**
140 (RAMs) pairs (Fig. 1D), following the terminology introduced earlier(Lee *et al.*, 2012). All
141 breakpoints supported by at least two clipped reads with partial mapping to a TE-consensus are
142 further interrogated for enrichment of available TE supporting clipped reads and RAMs using a

143 negative binomial model (Supl.Fig.1). The candidate sites where the numbers of clipped reads
144 and RAMs are less than 5% and 2.5%, respectively, of the total Hi-C coverage at the locus are
145 omitted as unreliable.

146 Unlike WGS, where the RAM pairs are clustered around the sites of TE-insertion, Hi-C data
147 exhibits wider mapping area. Though, both WGS and Hi-C data are biased by GC-content or
148 overall mappability, the coverage in Hi-C is additionally clearly biased by the density of RE sites
149 at the locus. Hence, HiTea uses a negative binomial model to assess the enrichment of TE-
150 insertion supporting reads (i.e., RAM pairs and clipped-reads) at the locus. To model the biases,
151 HiTea uses randomly selected loci in the genome that have similar coverage of the non-
152 conforming RPs as the site under investigation. Then, the count of TE-supporting reads at a locus
153 is assessed against negative binomial model built from the random set.

154 **Filtering and annotation of non-reference TE insertions**

155 A substantial fraction of clipped reads in Hi-C data display chimeric mapping (Fig. 1A, panel iii)
156 carrying a ligation motif at the clip position. To avoid calling such canonical Hi-C interactions as
157 TE insertions, HiTea filters out insertion candidates whose predicted breakpoints on either the
158 reference genome or TE-consensus are within 3-bases (user-defined) of the ligation motif
159 (Fig.1D, clip reads at RE site; Supl.Fig.1 for detailed filtering steps). It also filters out candidates
160 when multiple breakpoints are predicted around a putative breakpoint, as it is likely to be a
161 complex variant other than a TE insertion. At the sites of insertion, clipped mapping positions of
162 the reads indicate a breakpoint where reads mapping to the reference genome cluster (Fig. 1D).
163 HiTea expects that the genuine breakpoint should also show reciprocal cluster of the clipped
164 sequences when mapped to the TE-consensus. Insertions defying this expectation are removed as
165 invalid. Furthermore, insertions where clipped reads mapping only to the PolyA sequences are
166 omitted as potential simple repeat expansions. The genuine breakpoints are expected to have
167 clip-sequences mapping to PolyA sequence or presence of a degenerate polyA sequence (here we
168 look for a stretch of 7 As or Ts in the proximal 10 bases at the breakpoint on clipped sequences).
169 Subfamily annotation of the insertion is done by mapping the longest clipped sequence to the
170 subfamily consensus sequence derived from Repbase(Bao *et al.*, 2015). HiTea further detects
171 target site duplication, strand information, and estimates the size of insertion from the observed
172 mapping of the clipped sequences on the TE-consensus. HiTea is written in PERL and R. It uses

173 GNU-parallel(Tange, 2011) for parallelization over available cores. The insertions are reported
174 in bed format, with following status. Status-3 insertions are supported by right- and left-hand
175 side mapping of the clipped reads (Fig. 1D), whereas status-2 insertions represent a subset of
176 status-3 cases that overlap the reference copy of the same TE family. If the insertion is supported
177 by clipped reads at one side but have unmapped reads on the other site with polyA stretches (as
178 defined earlier), such instances are flagged with status-1.

179 **RESULTS**

180 **HiTea shows performance comparable to that of a WGS-based method**

181 To assess the performance of HiTea, we utilized Hi-C data generated from the HapMap cell line
182 GM12878(Rao *et al.*, 2014). This cell line has been extensively characterized using a wide range
183 of technologies and sequencing platforms. To generate the gold standard for comparison, we
184 used an improved version of our algorithm Tea(Lee *et al.*, 2012) on PacBio HiFi long
185 reads(Zook *et al.*, 2016) with extensive manual curation (hereafter referred to as the PacBio
186 reference). For WGS, we employed Mobile Element Locator Tool (MELT)(Gardner *et al.*,
187 2017), a popular software package with reportedly superior performance at moderate sequencing
188 depth(Rishishwar *et al.*, 2017). The full datasets consisted of ~5B RPs for Hi-C(Rao *et al.*, 2014)
189 (MboI-digested dataset; downloaded from 4DN data portal) and ~1.4B RPs for WGS
190 (downloaded from the 1000 Genomes project). Sequencing depths have considerable impact on
191 the precision and recall(Rishishwar *et al.*, 2017), thus we randomly down-sampled Hi-C data
192 to 1.4B RPs (~80X coverage) to provide a fair comparison between platforms. At this coverage,
193 79% of the genome in WGS and 57% in Hi-C data are covered with at least 60X coverage (Fig.
194 2A). The coverage was calculated by counting reads with mapping quality of at least 10 (MAPQ
195 ≥ 10).

196 The candidate insertions predicted by HiTea (ran on Hi-C data) and MELT (ran on WGS data)
197 were compared against the PacBio reference set (Fig. 2B). We used two sets of insertions
198 reported by MELT for GM12878: (i) the stringent “PASS” set (1122 insertions, referred as
199 MELT-PASS) and (ii) a more lenient set that includes the PASS variants and others for which
200 genotype could still be inferred (1443 insertions, referred as MELT-GT) in the comparisons. A
201 total of 1251 insertion were identified by HiTea while the PacBio reference set consisted of 1747
202 insertions.

203 Overall, HiTea correctly identified 1085 insertions (Fig 2B). The precision (fraction of the true
204 positives among all identified insertions) was 0.87; recall (fraction of true positives among all
205 positives) was 0.62 with F1 score of 0.72. MELT-PASS and MELT-GT correctly recovered 925
206 (precision 0.82, recall 0.53, F1 0.64) and 1115 (precision 0.77, recall 0.64, F1 0.7) insertions,
207 respectively.

208 (i) *Alu*. Most of the insertions were Alu, as expected. Among the 1493 Alu insertions from
209 our reference set, HiTea correctly identified 1000 (precision 0.89, recall 0.67, F1 0.76)
210 insertions from the Hi-C data. Whereas, MELT-PASS correctly identified 825 (precision
211 0.87, recall 0.55, F1 0.68) and MELT-GT recovered 986 (precision 0.83, recall 0.66, F1
212 0.74) Alu insertions from the WGS data. These results suggest that HiTea (ran on Hi-C)
213 has considerably better performance at detecting Alu compared to MELT (ran on WGS)
214 (Fig. 2B). Notably, HiTea can detect Alu insertions with competitive precision and recall
215 from Hi-C samples with lower coverages (Fig. 2C). For instance, at 600M RPs (~40X
216 sample; recommended sequencing depth by the 4DN consortium) and 300M RPs (~20X
217 coverage), the precisions are nearly uniform (i.e. 0.89 for 1.4B, 0.89 for 600M and 0.90
218 for 300M) and the recalls decrease only slightly, from 0.67 (1.4B, F1 0.76) to 0.65 (600M,
219 F1 0.75) and 0.59 (300M, F1 0.71) (Fig. 2C). We compared the proportions of the clipped
220 reads, which are the starting point of TE insertion identification in HiTea, and RAM reads
221 that map to Alu consensus between Hi-C (identified by HiTea) and WGS (identified by
222 MELT) at the equal sequencing depth of 1.4B. Although the proportions of Alu-mapping
223 clipped reads (44% in Hi-C and 53% in WGS) were higher, we observed that the
224 proportion of RAMs pairs mapping to the Alu consensus is much higher for Hi-C (43% of
225 total RAMs) than WGS (13% of total RAMs). Taken together, better proportions of
226 mapping of clipped and RAM reads in Hi-C is likely associated with better performance
227 of HiTea on Alu.

228 (ii) *L1Hs*. Our PacBio reference set contained 194 high-confidence L1Hs insertions. HiTea
229 correctly identified 67 (precision 0.64, recall 0.35, F1 0.45), whereas MELT-PASS and
230 MELT-GT detected 73 (precision 0.61, recall 0.38, F1 0.47) and 91 (precision 0.52, recall
231 0.47, F1 0.49), respectively (Fig. 2B). With respect to sequencing depths, recall increased
232 as the depth increased, from 0.11 for 300M (F1 0.19) to 0.4 for 5B RPs (F1 0.48), while
233 the precision remained in a similar range (0.61 to 0.71) (Fig.2C). Interestingly, the

234 proportions of both clipped and RAM reads mapping to the L1Hs consensus were
235 substantially higher in the WGS data (39.5% and 84.2% respectively) compared to the Hi-
236 C data (27.5% and 52.7% respectively). Transposed copies of L1Hs are frequently
237 associated with 5' truncation and/or inversion. Moreover, during target-primed reverse
238 transcription (TPRT), L1 RNA often accommodates sequences from the downstream
239 genomic region (Pickeral *et al.*, 2000). These additional features may lower the
240 performance of HiTea for L1Hs compared to Alu.

241 (iii) SVAs. HiTea has relatively poor sensitivity towards SVAs. Of 60 SVAs in the PacBio
242 reference set, HiTea correctly identified 18 (precision 0.75, recall 0.3, F1 0.43), whereas
243 MELT-PASS and MELT-GT respectively detected 27 (precision 0.51, recall 0.45, F1
244 0.48) and 38 (precision 0.48, recall 0.63, F1 0.55) instances. Although the proportions of
245 RAMs mapping on the SVA-consensus were comparable (2.7% for Hi-C vs 2.5% for
246 WGS), the proportions of SVA mapping clipped reads were substantially different (4.6%
247 in Hi-C vs 7.1% in WGS). SVAs comprise of frequently expanded hexameric repeats at
248 the 5', variable number of tandem repeats (VNTR) in the middle, and Alu-like sequences
249 at the 3'. This complex structure may lead to the relatively poor mapping of SVA-
250 originating reads to the SVA consensus (*e.g.*, some SVA reads map to the Alu consensus
251 instead), and thus affect the performance of HiTea for SVAs. Nonetheless, the precision of
252 detecting SVAs was strikingly high for HiTea (0.73 to 0.75) as compared to the MELT
253 calls (<0.51) (Fig. 2B,C). The impact of sequencing depth for SVAs was similar to that for
254 L1Hs.

255 Of the 1251 HiTea insertions (at 1.4B), ~13% (166) did not overlap with the PacBio reference
256 set. Hence, we interrogated them against a collection of 1000 Genome TE insertion set (at a
257 population allele frequency $\geq 10\%$; results were similar for $AF \geq 0.01\%$ and $AF \geq 0.1\%$),
258 identified on the low coverage WGS data by MELT (Gardner *et al.*, 2017). Our comparison
259 suggested that 117/166 (~71%) HiTea-specific insertions overlap with the population-based TE-
260 insertion set, suggesting that these are true insertions missed by the PacBio reference set. This
261 also suggests that the precision and recall measures above represent lower bounds.

262 HiTea missed ~38% (662/1747) of the insertions from the PacBio reference set. Of the 662, 197
263 insertions overlapped with 1000G set. We assessed the 5' end coverage of RAMs whose mates

264 or clipped sequences map to the TE consensus in Fig. 2D. This coverage plot(Gu *et al.*, 2018)
265 shows that the missed events by HiTea do not have a sufficient number of clipped reads (lower
266 right panels in Fig. 2D; 381/465 and 95/197 have less than two non-Hi-C chimeric clipped reads
267 mapping to the TE-consensus at the locus.

268 Since the 1.4B-RPs datasets used above are larger than typical datasets, we repeated the above
269 analysis with down-sampled datasets with ~600M RPs (~35-40X). Our comparison suggests that
270 HiTea (run on Hi-C) shows consistently higher precision in detecting Alu, SVA and L1Hs
271 compared to MELT (run on WGS data) (Supl.Fig.2A). A total of 1016/1152 HiTea insertions
272 (precision 0.88, recall 0.58, F1 0.70) and 908/1134 MELT-PASS insertions (precision 0.80,
273 recall 0.52, F1 0.63) overlapped with the PacBio reference set (Supl.Fig.2A). The insertions
274 missed by HiTea did not seem to show clip-read coverage at the respective loci (Supl.Fig.2B).

275 We tested HiTea on a range of human Hi-C datasets generated using different REs. A 4-cutter
276 RE (MboI, DpnII) is expected to cut the DNA at every 256 bases whereas a 6-cutter (HindIII,
277 NcoI) will digest the DNA at 4096bp on average. The infrequent cuts by a 6-cutter are expected
278 to provide low spatial resolution of the Hi-C (Supl.Fig.3A), resulting in a smaller number of
279 clipped reads along the genome. Indeed, when Hi-C datasets generated using different REs for
280 GM12878 cell line(Rao *et al.*, 2014) were compared, the overall recall dropped from 0.62 (MboI
281 digested Hi-C, 1.4B RPs, F-score 0.72) to 0.41 (1.8B RPs, HindIII digested Hi-C, F-score 0.56).
282 For comparison, the overall recalls for WGS sample were 0.53 (MELT-PASS, F-score 0.64) and
283 0.64 (MELT-GT, F-score 0.7) at 1.4B RPs (Supl.Fig.3B). Nonetheless, HiTea showed a high
284 precision (0.88) compared to MELT-PASS (0.82) and MELT-GT (0.77). Besides 17 unique,
285 remaining 794 (98%) insertions from the HiTea run either overlapped with PacBio reference set
286 or the 1000G set, whereas about 79% (811/1032) of the missed insertions displayed poor
287 coverage of clipped reads (Supl.Fig.3C, D). With the decreasing sequencing cost, many studies
288 are now using either a 4-cutter or a mix of 4-cutter enzymes, and these high-resolution Hi-C
289 datasets will be suitable for HiTea analysis.

290 Next, we assessed the performance of HiTea on another widely-characterized cell line, K562.
291 We obtained WGS and Hi-C (MboI digested Hi-C) data from Cancer Cell Line Encyclopedia
292 project(Barretina *et al.*, 2012) and a published study(Rao *et al.*, 2014), respectively. As a PacBio
293 reference set was not available for this cell line, we resorted to comparing the TE-insertions

294 called by HiTea (on Hi-C) to those from MELT (on WGS). At comparable sequencing depth of
295 1.2B RPs between Hi-C and WGS data for K562 cells, a substantial fraction (769/958, ~80%) of
296 HiTea insertions overlapped with either MELT-derived (i.e. MELT-GT) insertions or 1000G set.
297 In comparison, previously analyzed GM12878 (MboI digested, 1.4B RPs) exhibited similar
298 (1101/1251, ~88%) degree of overlap (Supl.Fig.4).

299 **HiTea aids in the characterization of the non-reference TE insertions**

300 To assess whether HiTea can correctly identify insertions otherwise missed by MELT, we
301 compared MELT-GT (better recall compared to MELT-PASS) and HiTea insertions (both at
302 1.4B RPs sequencing depth) using the PacBio reference set. Our analysis suggests that a
303 substantial number of insertions overlapping with reference-genome copy of the same TE family
304 are missed by MELT (Fig. 3A, B). TE detection along the reference TE copy of the same family
305 can be challenging due to multiple reasons, such as poor mappability of the reads and structural
306 variation within the reference-copies of the TE family. Therefore, several WGS-based tools filter
307 out these insertions to limit the number of false positives (Ewing, 2015). However, when
308 supporting reads are available and their mappings on both TE-consensus and reference genome
309 provide sufficient confidence for the insertion, HiTea reports these events. Our reference set
310 included 436 TE-insertions overlapping with the reference copies of the same TE-family. HiTea
311 correctly identified 70 insertions reported in the PacBio reference, outperforming MELT (5 and 8
312 by MELT PASS and GT) (Fig. 3B).

313 In total, HiTea identified 160 PacBio reference insertions missed by MELT-GT. Conversely,
314 MELT-GT identified 180 insertions missed by HiTea from the reference set (Supl.Fig.5A, B).
315 When assessed for the features that led to disqualification of these true-positive insertions by
316 either MELT or HiTea, we observed that indeed insertions within a reference-genome copy of
317 the same TE family were preferentially missed by MELT (66/160, ~41%; Supl.Fig.5C). As the
318 exact features used by MELT are unavailable (the code is not open source), we could not further
319 investigate the instances missed by MELT. Over half of the insertions (124/180, ~69%) missed
320 by HiTea were due to poor coverage of clipped reads, proximity to the RE motif, coverage
321 thresholds, and absence of clipped reads supporting polyA tails (Supl.Fig.5D).

322 Coverage in the Hi-C experiment is significantly higher around the RE sites in the genome.
323 Thus, insertions proximal to the RE sites tend to have higher coverage of supporting reads even

324 at relatively low overall sequencing depth. In the example shown in Fig.3C, the read coverage at
325 an Alu insertion site missed by MELT-GT on chromosome 20 is much higher in Hi-C than in
326 WGS, although the overall sequencing depth is the same (both bam files were subsampled to
327 10% of total reads for better visualization). To assess whether the same phenomenon is observed
328 at many sites, we counted total 5' end coverage in a 1kb window centered at the 925 insertions
329 identified by both MELT-GT and HiTea and the 160 insertions identified only by HiTea. As
330 expected, the insertions identified by both methods tend to have similar coverages, whereas those
331 missed by MELT-GT tend to have relatively lower coverage overall in WGS compared to Hi-C
332 (Fig.3D).

333 A total of 49/1251 (~4%) insertions detected by HiTea were not explained by either the PacBio
334 reference set or the 1000G set (Fig.2D, second panel from the top). Of these, 4 and 6 were
335 reported by MELT-PASS and MELT-GT, respectively. These HiTea-specific insertions exhibit
336 clear presence of TE-mapping clipped reads from Hi-C data (Fig.2D). Representative examples
337 of two Alu insertions suggest that the HiTea-unique insertions have the support of both clipped
338 and discordant reads at the insertion locus in the WGS data (Fig.3E). We suspect that many of
339 these cases may be true positives that were missed by MELT due to its stringent filtering criteria.

340 **Installation and usage**

341 HiTea is available at Github (<https://github.com/parklab/HiTea>) and as a Docker image
342 (4dndc/hitea:v1 on Docker Hub). TE (Alu, L1Hs and SVA) family-wise consensus sequences
343 and the genomic locations of the TE-family members required for running HiTea are provided
344 for hg38 and hg19 human genome references, with a description on how to generate them for
345 other types of TEs on the GitHub page. HiTea dependencies are PERL (\geq v5.24), R (\geq v3.2),
346 bedtools (\geq v2.26)(Quinlan and Hall, 2010), samtools (\geq v1.7), GNU-parallel(Tange, 2011) and
347 Pairtools (<https://github.com/mirnylab/pairtools>). Additionally, there are mandatory
348 (GenomicRanges, data.table, MASS) and optional (rmarkdown, knitr, EnrichedHeatmap(Gu *et*
349 *al.*, 2018), circlize) R packages used for computation and HTML-report generation steps
350 respectively. Users can start the analysis with a single command by providing a name-sorted bam
351 file, restriction enzyme used for the Hi-C assay and the genome build used to map the Hi-C data.
352 HiTea auto detects if the read class information is present in the bam file (e.g. files obtained from
353 4DN data portal <https://data.4dnucleome.org/> carry this information). If not, it automatically

354 employs Pairtools to generate read class information. User-defined TE-consensus or
355 polymorphic sequences and the genomic locations of the members of TE-sequences can be
356 provided using a detailed input option. A HiTea run on a typical Hi-C dataset (~600M RPs) takes
357 about 3.5-4 hrs to complete with 8 cores and 20 G memory.

358 **DISCUSSION**

359 Although used primarily for understanding three dimensional organization of the genome and its
360 regulatory role, the long-range chromatin interaction information in Hi-C data have been used to
361 assemble small scaffolds into chromosome-length assemblies(Dudchenko *et al.*, 2017; Gong *et*
362 *al.*, 2018) and to identify copy number and translocations(Chakraborty and Ay, 2018; Dixon *et*
363 *al.*, 2018; Wang *et al.*, 2020). In the present work, we have demonstrated that Hi-C can be used
364 also to identify TE insertions.

365 The strong performance of HiTea was somewhat unexpected. Given the nature of the
366 experiment, the read coverage for Hi-C is highly variable along the genome. We thus expected
367 that there would not be enough reads at some TE insertions sites, resulting in degraded
368 performance for HiTea compared to a WGS-based method. What makes HiTea competitive with
369 a WGS-based method, however, is the use of clipped reads to locate candidate TE insertions at
370 the discovery step, in contrast to the discordant RP-based candidate discovery in most WGS-
371 based methods. The higher proportion of clipped reads (carrying no RE ligation junction) in Hi-
372 C data (1.6%) than in WGS data (1.4%) is further helpful. Moreover, the proportion of RPs
373 whose one end remains unmapped or multimapped is higher in the Hi-C data (21%) compared to
374 the WGS data (14%) due to wider effective insert sizes, increasing the power of Hi-C data for
375 detecting insertions. In particular, the TE insertions in the reference genome copies of the same
376 family or those occurring in regions with comparatively lower coverage in WGS data are
377 sometimes detected by HiTea but missed by MELT.

378 The availability of PacBio HiFi data (circular consensus sequencing method, with half the reads
379 >50kb) for GM12878 made it easier to evaluate the performance of different methods. However,
380 the TE insertion map based on this one sample is obviously incomplete, as seen by the fact that
381 many HiTea candidates not present in the PacBio reference set were present in the 1000G data.
382 A small fraction (<5%) of HiTea insertions were still not explained by either PacBio reference
383 set or 1000G set. Although some of these insertion calls may be false positives, it is interesting to

384 note that both WGS and Hi-C data show presence of discordant and non-conforming RPs
385 mapping to the underlying TE consensus, respectively, along most of these loci. Additional long-
386 read data or independent experimental validations may prove useful in discerning the nature of
387 HiTea-specific calls.

388 The number of studies mapping chromatin organization in diverse organisms, cell types, and
389 disease states as well as the collective efforts to organize such data have gained
390 momentum(Dekker *et al.*, 2017). However, it is imperative to mark structural variations in the
391 genome before construing the chromatin interactions from Hi-C data as functional interactions,
392 as we have demonstrated recently(Wang *et al.*, 2020). HiTea exploits Hi-C data to identify non-
393 reference TE insertions, using reads that otherwise would be discarded. Finally, although we
394 compared call sets from Hi-C and WGS data in our analysis, the ideal scenario is to have both
395 data types for a sample of interest, so that the insertions calls can be cross-validated and
396 expanded. Continued development of more comprehensive reference TE insertions maps and
397 robust computational methods for TE identification will be important.

398 **Acknowledgements:** This work was supported by the grants from the National Institutes of
399 Health Common Fund 4D Nucleome Program (U01CA200059) and National Institutes of Mental
400 Health (U01MH106883) to PJP.

401 **Conflict of Interests:** None declared

402 **References**

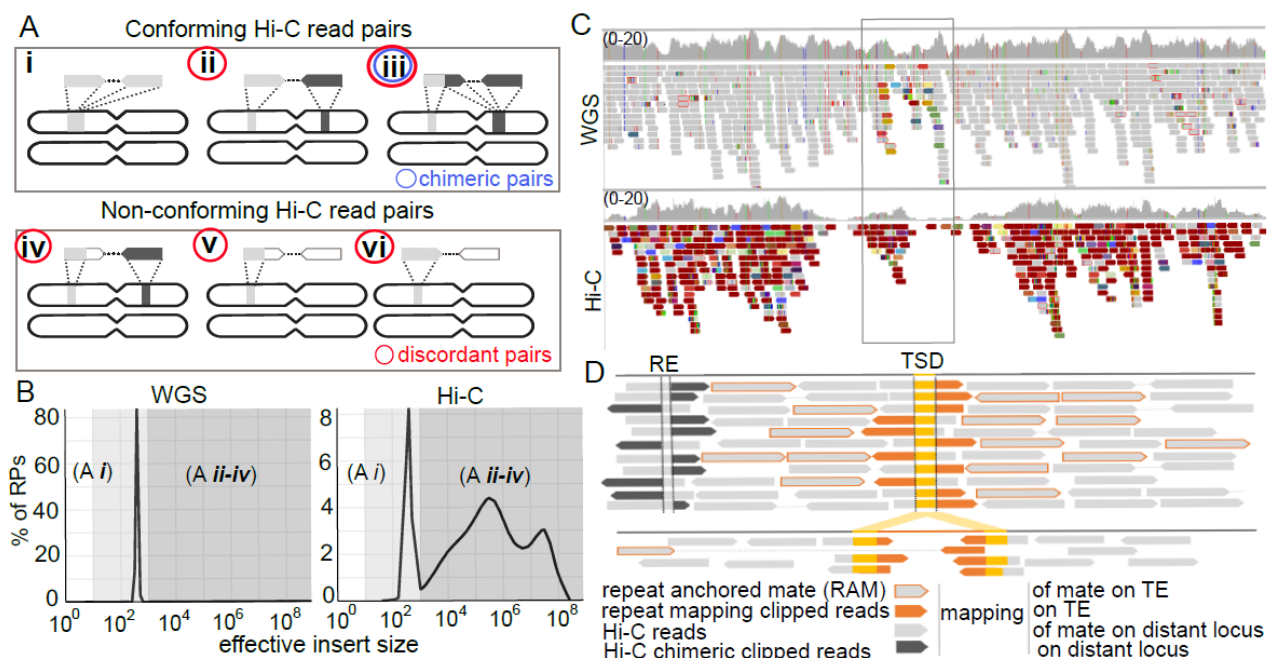
- 403 Ahmed,M. and Liang,P. (2012) Transposable Elements Are a Significant Contributor to Tandem
404 Repeats in the Human Genome. *Comp. Funct. Genomics*, **947089**.
- 405 Ayarpadikannan,S. and Kim,H.-S. (2014) The Impact of Transposable Elements in Genome
406 Evolution and Genetic Instability and Their Implications in Various Diseases. *Genomics*
407 *Inf.*, **12**, 98–104.
- 408 Bao,W. *et al.* (2015) Repbase Update, a database of repetitive elements in eukaryotic genomes.
409 *Mob. DNA*, **6:11**.
- 410 Barretina,J. *et al.* (2012) The Cancer Cell Line Encyclopedia enables predictive modelling of
411 anticancer drug sensitivity. *Nature*, **483**, 603–607.
- 412 Bourque,G. *et al.* (2018) Ten things you should know about transposable elements. *Genome*
413 *Biol.*, **19**, 199.
- 414 Chakraborty,A. and Ay,F. (2018) Identification of copy number variations and translocations in
415 cancer cells from Hi-C data. *Bioinformatics*, **34**, 338–345.
- 416 Chen,J.M. *et al.* (2008) Detection of two Alu insertions in the CFTR gene. *J. Cyst. Fibros.*, **7**,
417 37–43.

- 418 Chenais,B. (2015) Transposable elements in cancer and other human diseases. *Curr. Cancer*
419 *Drug Targets*, **15**, 227–242.
- 420 Chuong,E.B. *et al.* (2017) Regulatory activities of transposable elements: From conflicts to
421 benefits. *Nat. Rev. Genet.*, **18**, 71–86.
- 422 Cournac,A. *et al.* (2016) The 3D folding of metazoan genomes correlates with the association of
423 similar repetitive elements. **44**, 245–255.
- 424 Dekker,J. *et al.* (2017) The 4D nucleome project. *Nature*, **549**, 219–226.
- 425 Dixon,J.R. *et al.* (2018) Integrative detection and analysis of structural variation in cancer
426 genomes. *Nat. Genet.*, **50**, 1388–1398.
- 427 Dudchenko,O. *et al.* (2017) De novo assembly of the *Aedes aegypti* genome using Hi-C yields
428 chromosome-length scaffolds. *Science (80-.)*, **356**, 92–95.
- 429 Ewing,A.D. (2015) Transposable element detection from whole genome sequence data. *Mob.*
430 *DNA*, **6:24**.
- 431 Garcia-Perez,J.L. *et al.* (2016) The impact of transposable elements on mammalian development.
432 *Development*, **143**, 4101–4114.
- 433 Gardner,E.J. *et al.* (2017) The mobile element locator tool (MELT): Population-scale mobile
434 element discovery and biology. *Genome Res.*, **27**, 1916–1929.
- 435 Gong,G. *et al.* (2018) Chromosomal-level assembly of yellow catfish genome using third-
436 generation DNA sequencing and Hi-C analysis. *Gigascience*, **7**, 1–9.
- 437 Gu,Z. *et al.* (2018) EnrichedHeatmap: An R/Bioconductor package for comprehensive
438 visualization of genomic signal associations. *BMC Genomics*, **19**, 234.
- 439 Hancks,D.C. and Kazazian,H.H. (2016) Roles for retrotransposon insertions in human disease.
440 *Mob. DNA*, 7:9.
- 441 Kim,J. *et al.* (2019) Patient-customized oligonucleotide therapy for a rare genetic disease. *N.*
442 *Engl. J. Med.*, **381**, 1644–1652.
- 443 de Koning,A.P.J. *et al.* (2011) Repetitive elements may comprise over Two-Thirds of the human
444 genome. *PLoS Genet.*, **7**.
- 445 Lee,E. *et al.* (2012) Landscape of somatic retrotransposition in human cancers. *Science (80-.)*.
- 446 Li,H. and Durbin,R. (2010) Fast and accurate long-read alignment with Burrows-Wheeler
447 transform. *Bioinformatics*, **26**, 589–595.
- 448 Mills,R.E. *et al.* (2007) Which transposable elements are active in the human genome? *Trends*
449 *Genet.*, **23**, 183–191.
- 450 Pickeral,O.K. *et al.* (2000) Frequent human genomic DNA transduction driven by line-1
451 retrotransposition. *Genome Res.*, **10**, 411–415.
- 452 Quinlan,A.R. and Hall,I.M. (2010) BEDTools: A flexible suite of utilities for comparing
453 genomic features. *Bioinformatics*, **26**, 841–842.
- 454 Rao,S.S.P. *et al.* (2014) A 3D map of the human genome at kilobase resolution reveals principles
455 of chromatin looping. *Cell*, **159**, 1665–1680.
- 456 Rishishwar,L. *et al.* (2017) Benchmarking computational tools for polymorphic transposable
457 element detection. *Brief. Bioinform.*, **18**, 908–918.
- 458 Schmidt,D. *et al.* (2012) Waves of retrotransposon expansion remodel genome organization and
459 CTCF binding in multiple mammalian lineages. *Cell*, **148**, 335–348.
- 460 Schmitt,A.D. *et al.* (2016) Genome-wide mapping and analysis of chromosome architecture. *Nat.*
461 *Rev. Mol. Cell Biol.*
- 462 Tange,O. (2011) GNU Parallel: The Command-Line Power Tool. *USENIX Mag.*, **36**, 42–47.
- 463 Vidaud,D. *et al.* (1993) Haemophilia B due to a de novo insertion of a human-specific Alu

464 subfamily member within the coding region of the factor IX gene. *Eur. J. Hum. Genet.*, **1**,
 465 30–36.
 466 Wang,S. *et al.* (2020) HiNT: a computational method for detecting copy number variations and
 467 translocations from Hi-C data. *Genome Biol.*, **22**, 73.
 468 Wicker,T. *et al.* (2008) A universal classification of eukaryotic transposable elements
 469 implemented in Repbase. *Nat. Rev. Genet.*, **9**, 414.
 470 Zook,J.M. *et al.* (2016) Extensive sequencing of seven human genomes to characterize
 471 benchmark reference materials. *Sci. Data*.

472
 473
 474

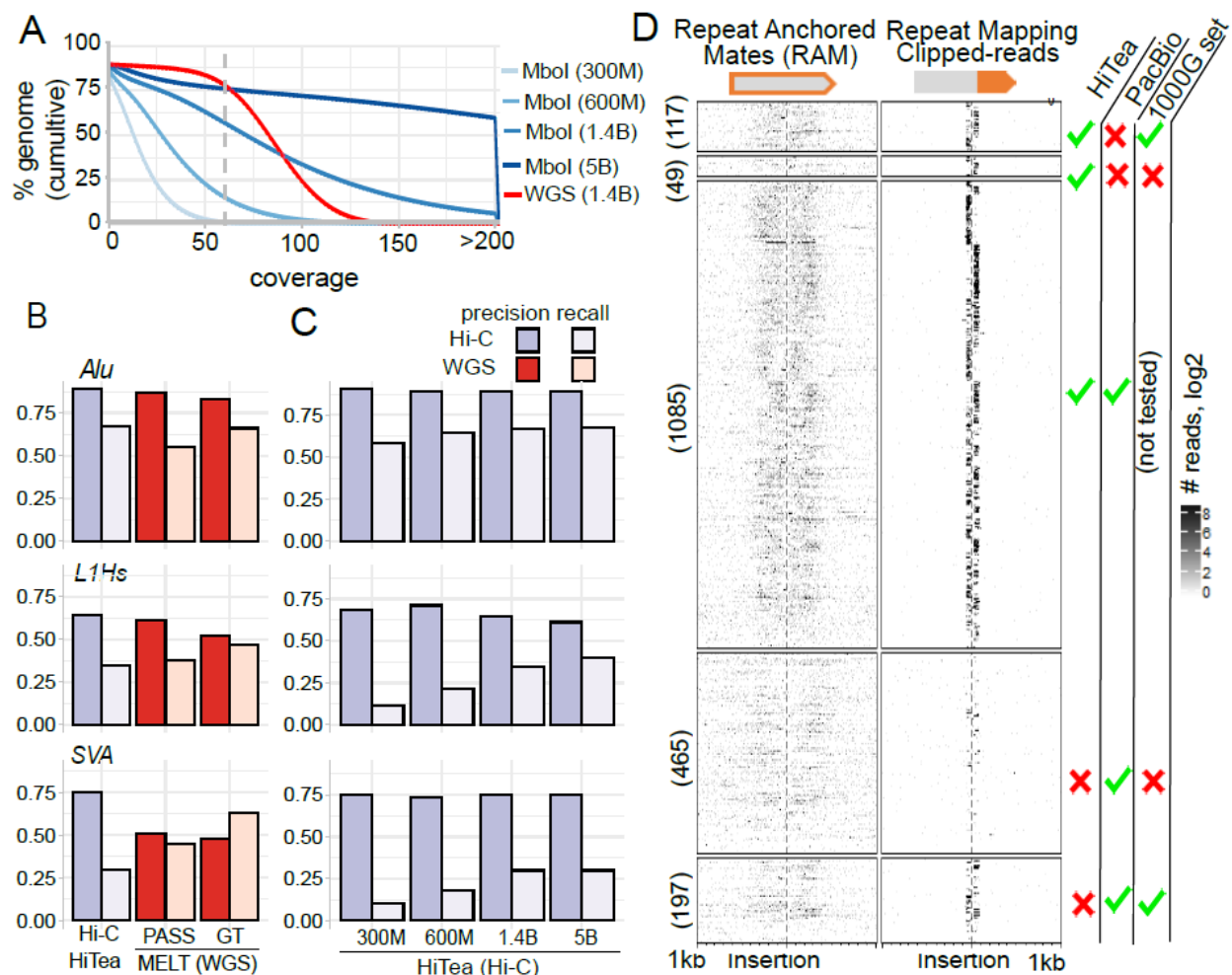
Figures and legends:



475
 476 **Figure 1: Properties of Hi-C reads supporting a TE insertion. (A)** Hi-C read pairs (RPs) can
 477 be grouped into two classes that we termed ‘conforming’ and ‘non-conforming’. Conforming RPs
 478 comprise of (i) WGS-like pairs with short insert sizes, (ii) pairs with large effective insert sizes,
 479 and (iii) chimeric RPs where the clip-sequence maps convergent to mapped locus of its paired
 480 mate. Non-conforming RPs comprise of mapping configurations where (iv) the clipped sequence
 481 does not display chimeric mapping or (v-vi) the mate remains unmapped on reference genome.
 482 **(B)** Comparison of the between-pair distances for WGS and Hi-C experiments. **(C)** A genome
 483 browser view of a true insertion event, showing both coverage and the discordant RPs (non-gray
 484 color) in WGS and Hi-C experiments. Box marks the TE-insertion site. Mapped read pairs in the
 485 display are color-coded by the insert sizes using default IGV color scheme. **(D)** A schematic of
 486 Hi-C read configuration at insertion site. Clipped reads supporting TE insertion exhibit partial

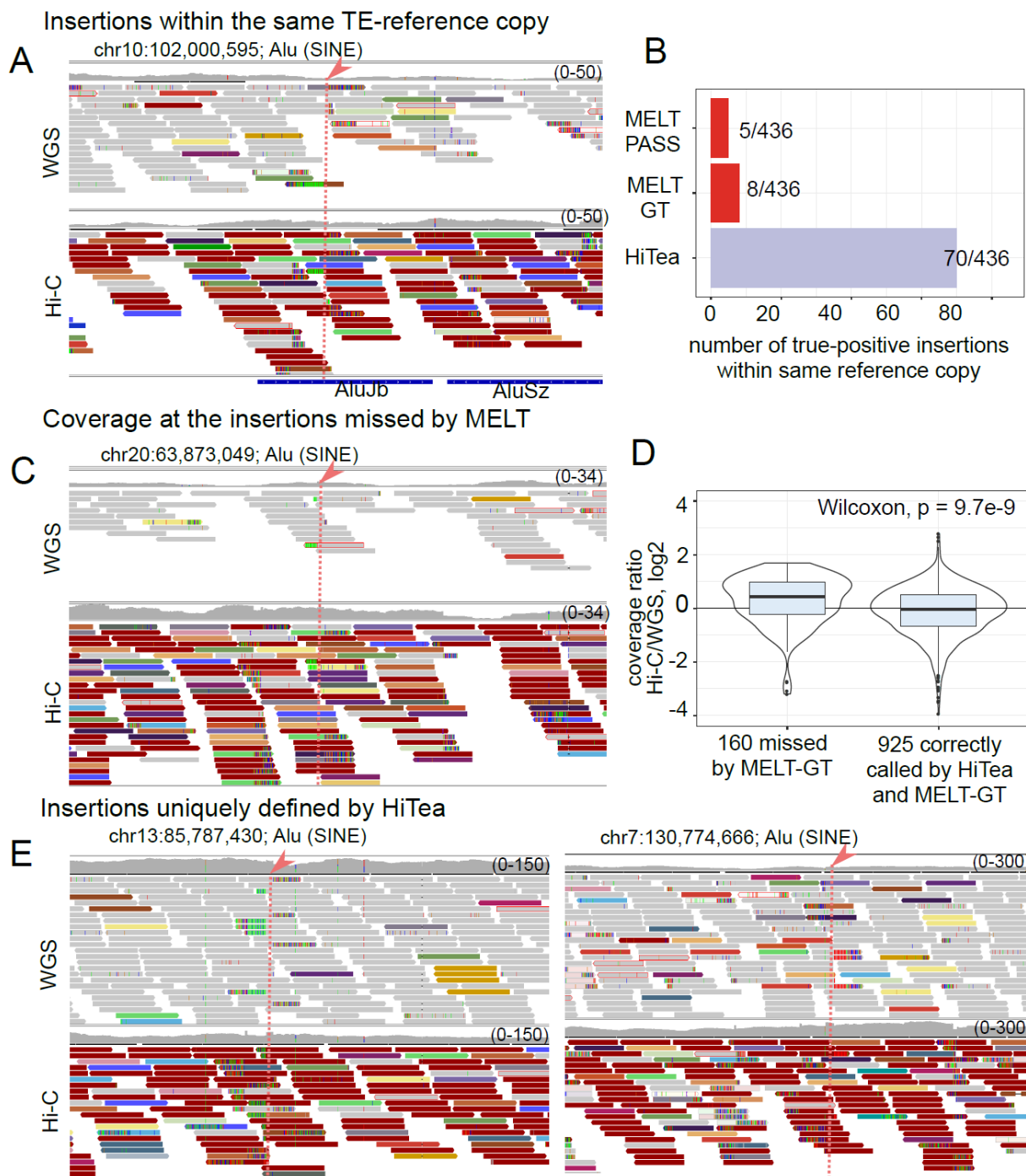
487 mapping to TE-family consensus (orange), whereas those that do not, map at distant reference
 488 locus (black). RPs with a mate mapping to the TE-family consensus are displayed with orange
 489 outline. (RE: restriction endonuclease, TSD: target site duplication)

490



491

492 **Figure 2: Performance of HiTea.** (A) Cumulative distribution of the coverage for different
 493 datasets. Gray dotted line marks 60X coverage. (B) Precision and recall for detecting insertions
 494 of Alu, L1Hs and SVA families using HiTea (on Hi-C) and MELT (on WGS) at 1.4B sequencing
 495 depth. PASS and GT refer to the more and less stringent call sets, respectively, in MELT. (C)
 496 Precision and recall comparison at different sequencing depths of Hi-C experiment. (D) 5' end
 497 coverage for the RAMs whose mates map to the TE consensus (left) or reads whose clipped-
 498 sequences map to the TE consensus (right). The insertions are grouped according to the criteria
 499 shown on the right. PacBio is the reference set constructed using PacBio HiFi reads; 1000G set
 500 refers to insertions detected in the 1000 Genome data by MELT.



501

502 **Figure 3: Examples of TE insertions detected in Hi-C but missed in WGS.**

503 **(A)** A browser view of an insertion overlapping the reference-genome copy of a TE-family. This
 504 insertion is identified by HiTea (on Hi-C) but missed by MELT (on WGS). Reads with concordant
 505 and discordant mapping configurations are displayed in gray and non-gray colors, respectively.
 506 The discordant RPs are color-coded according to their insert sizes. Dotted red line with arrowhead

507 marks the insertion site. **(B)** Summary of TE insertions detected when the insertion occurs in the
508 reference-genome copies of the TE-family. **(C)** Sequencing coverage comparison at an insertion
509 correctly called by HiTea but missed by MELT. **(D)** The boxplot for the Hi-C/WGS read coverage
510 ratios shows that Hi-C coverage is higher in cases identified by HiTea but missed by MELT-GT.
511 **(E)** More examples of insertions called only by HiTea.