

1 Population-wide Sampling of Retroposon Insertion

2 Polymorphisms Using Deep Sequencing and Efficient

3 Detection

4
5 Qichao Yu^{1,2,†}, Wei Zhang^{1,2,†}, Xiaolong Zhang², Yongli Zeng², Yeming Wang², Yanhui Wang²,
6 Liqin Xu², Nannan Li², Xinlan Zhou², Xiaoyun Huang², Jie Lu³, Xiaosen Guo², Guibo Li^{2,4}, Yong
7 Hou^{2,4}, Shiping Liu^{2,5,*} and Bo Li^{2,6,*}

8
9 ¹ BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,
10 China

11 * Correspondence: libo@genomics.cn; liushiping@genomics.cn

12 † Equal contributors

13 Full list of author information is available at the end of the article.

14

15 **Emails of all authors:**

16 Qichao Yu: yuqichao@genomics.cn; Wei Zhang: zhangwei7@genomics.cn; Xiaolong Zhang:
17 13528497060@163.com; Yongli Zeng: zeoly100@163.com; Yeming Wang: 1398738509@qq.com;
18 Yanhui Wang: 839584901@qq.com; Liqin Xu: xuliqin@genomics.cn; Nannan Li: linannan@genomics.cn;
19 Xinlan Zhou: zhouxinlan@genomics.cn; Xiaoyun Huang: huangxiaoyun@genomics.cn; Jie Lu:
20 lujie1@genomics.cn; Xiaosen Guo: guoxs@genomics.cn; Guibo Li: liguibo@genomics.cn; Yong Hou:
21 huyong@genomics.cn; Bo Li: libo@genomics.cn; Shiping Liu: liushiping@genomics.cn.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51 and those not (non-LTR). The majority of human TEs result from the present and past activity
52 of non-LTR retrotransposons, including the L1 (long interspersed nuclear element 1), Alu and
53 SVA elements, which collectively account for approximately one-third of the human genome
54 [1]. While most retrotransposons are inactive remnants prevalent among the human population,
55 younger retrotransposons account for much of the structural variation among individual
56 genomes [3]. There exist only a small proportion of total L1s which are highly active [4]. The
57 current rate of retrotransposition in humans has been estimated as approximately 1 for every
58 20 births for Alu, approximately 1 for every 200 births for L1 and approximately 1 for every 900
59 births for SVA [5, 6].

60 Retrotransposon insertion is known as a disease causing mechanism [7], and the
61 Next-generation sequencing (NGS) technology has been widely used to explore the
62 association between retrotransposon insertions and disease, such as cancer [8-10]. In this
63 respect, a comprehensive RIPs dataset of healthy population is necessary to serve as a
64 reference to identify disease related RIPs. Based on the database of the 1000 Genomes
65 Project (1000GP), researchers were able to carry out RIPs detection on an unprecedented
66 scale through whole-genome sequencing and detect thousands of novel RIPs [11-13].
67 However, the 1000GP relied mainly on pooled low-coverage sequencing data (1~3x per
68 individual) from many individuals for RIPs analysis. Because an insertion allele present in
69 multiple individuals would effectively receive high coverage across the pooled dataset, this
70 approach was biased towards common insertions. According to previous calculation, to detect
71 heterozygous RIPs with high sensitivity using whole-genome sequencing, at least 30x
72 coverage of sequencing was needed [14].

73 In addition, the current post-sequencing bioinformatic methods such as RetroSeq [11, 15],
74 TEA [10], PTEMD [16], Jitterbug [17], T-lex2 [18], Mobster [19] are challenged to deal with deep
75 whole genome sequencing data especially at the population level because of time
76 consumption. A fast and efficient method is required to detect RIPs in order to satisfy the
77 increasing amount of whole genome sequencing (WGS) data.

78 Here we developed a new computer program named Specific Insertions Detector (SID) to
79 detect RIPs, which has much higher detection efficiency but comparable detection accuracy
80 and sensitivity compared with TEA and Retroseq, two of the most-cited algorithms of RIPs

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

81 calling. We next presented a non-reference TEs insertion polymorphism database by
82 employing SID to analyze whole genome sequences of 90 Han Chinese individuals (YH90),
83 acquired at a mean depth of 68x.

84 **Materials and methods**

85 **Samples and whole genome sequencing**

86 We obtained B-lymphocyte cell lines of 90 Han Chinese individuals from Coriell institute
87 (Camden, New Jersey, USA). These samples were selected in Beijing, Hunan province and
88 Fujian province respectively, and we broadly separated them into ‘Northern group’ (45 samples)
89 and ‘Southern group’ (45 samples). DNA was extracted from the B-lymphocyte cell line of each
90 individual and libraries were then constructed following the manufacturer’s instructions, and
91 high-coverage paired-end 100 bp WGS libraries were sequenced on the Illumina HiSeq 2000
92 Platform. In addition, we also used the Chinese sample [20] whose data has already been
93 released in the European Nucleotide Archive (ENA) repository (for details see Additional file 1:
94 Table S1). The study was approved by the Institutional Review Board on Bioethics and
95 Biosafety of BGI (BGI-IRB).

96 **Processing of the WGS data**

97 Reads were aligned to human genome reference (HG19, Build37) using *BWA-0.6.1* [21] with
98 parameters ‘-n 3 -o 1 -e 50’. Duplications were removed and the quality values of each reads
99 were recalibrated using the Genome Analysis Toolkit (GATK) [22] and Samtools [23]. The
100 resulting Binary Alignment/Map (BAM) files were required by SID.

101 **The specific insertions detector pipeline**

102 SID is compiled in Perl and includes two steps, discordant reads detection and reads
103 clustering. Generally, the first step collects informative reads and generates other necessary
104 files, while the second step discovers the specific insertion site and outputs the final results
105 into a plain text.

1 106 *Discordant reads detection.* The ‘discordant reads’ are extracted for the subsequent clustering
2
3 107 step. Paired-end reads were determined as ‘discordant reads’ if they met one of the following
4
5
6 108 criteria: a. one read mapped to HG19 uniquely and the other read mapped to retrotransposons
7
8
9 109 library (multi-mapped or unmapped to HG19); b. one read mapped to HG19 uniquely and the
10
11 110 other soft-clipped read mapped to HG19, and the clipped sequence could be mapped to
12
13 111 retrotransposons library; c. one soft-clipped read mapped to HG19, and the clipped sequence
14
15 112 could be mapped to retrotransposons library, while the other read mapped to retrotransposons
16
17 113 library (multi-mapped or unmapped to HG19) (Fig. 1a). This retrotransposons library includes
18
19 114 the objective TE classes, such as L1, Alu, SVA, etc. In this study, the TE reference database
20
21 115 contains known TE sequences collected from RepBase 17.07 [24], dbRIP [25] and Hot L1s [4].
22
23 116 In order to reduce the long processing time due to the large whole genome sequencing data,
24
25 117 we implemented a parallel approach to process each bam files of samples simultaneously in
26
27 118 the discordant reads detection step.
28
29
30
31
32
33
34
35
36 119 *Reads clustering and breakpoints detection.* First, the ‘discordant reads’ would be scanned and
37
38 120 clustered into blocks which support potential RIPs based on Maximal Valid Clusters algorithm.
39
40
41 121 Second, we exacted all reads located within the cluster regions and determined the
42
43 122 breakpoints. Although high depth data enabled RIPs detection with high sensitivity because
44
45 123 more soft-clipped reads neighboring target site duplication (TSD) could be detected,
46
47 124 alignments neighboring the TSDs had apparently lower depth than the mean sequencing
48
49 125 depth of whole genome due to some sequencing and system errors. This made the
50
51 126 breakpoints detection difficult and increased the false positive rate inevitably. Thus, we added
52
53 127 the recalibration process of clipped points to determine breakpoints. For each read that
54
55
56
57
58
59
60
61
62
63
64
65

1 128 located within the cluster regions around potential breakpoints was used to confirm the precise
2
3 129 location of the breakpoints. Small deletions were extracted to perform breakpoints
4
5
6 130 recalibration, and the mismatched bases were removed from the deletion sequences.
7
8

9 131 The clipped sequences were realigned to local regions on HG19 to gain the actual
10
11 132 breakpoints. Breakpoints were taken as 'clips' if more than half of new clipped sequences were
12
13
14 133 discordant to the reference sequence and the length of gap within new clipped sequence was
15
16
17 134 less than 30%. The point would not be candidate unless it's a 'clip' and the mismatch is less
18
19
20 135 than 5 bp or contains polyA/T.
21

22 136 Some terminals of reads (Fig. 1b, c) that contain mismatched bases may be the clipped
23
24
25 137 parts because the alignment software usually treats these bases as mismatches rather than
26
27
28 138 clips. SID re-estimates the breakpoints candidates if mismatches were more than half of the
29
30
31 139 read terminals.
32

33 140 Of note, we implemented 'Asynchronism Scanning' algorithm. Using this algorithm, once
34
35
36 141 the program clustered one possible insertion region by scanning unique reads, the process of
37
38
39 142 breakpoints detection in this region was carried out immediately, rendering it possible to detect
40
41
42 143 RIPs in one chromosome in just few minutes. The detailed algorithm for RIP candidate
43
44
45 144 determination is provided in Additional file 2: Supplementary Methods.
46

47 145 **RIPs data simulation**

48
49 146 In total 761 TEs were randomly selected from a retrotransposon database (homebrew from
50
51
52 147 dbRIP) and inserted into HG19 autosomes randomly to generate a new human genome. The
53
54
55 148 pIRS [26] software was used to generate about 60x paired-end reads and then we mapped
56
57
58 149 these reads to HG19 genome by BWA. After that we used SID to detect these RIPs in the
59

1 150 simulated genome. By repeating this process, we got results in different depth simulated data
2
3 151 to assess the sensitivity and specificity RIPs detection in sequence data with distinct depth
4
5
6 152 using SID.

8 153 **Subfamily, length and orientation annotation of the inserted retrotransposons**

10 154 Subfamily annotation of RIPs was performed according to known active retrotransposons. We
11
12
13 155 first constructed a comprehensive retrotransposons sequence library: Alu subfamily
14
15
16 156 consensus sequences were acquired from RepBase 17.07 [24], L1 subfamily consensus
17
18
19 157 sequences were acquired from Eunjung Lee [10], SVA and LTR consensus sequences were
20
21
22 158 acquired from Baillie [27]. Next, we did the multiple subfamily sequence alignment of each type
23
24
25 159 of retrotransposon and discovered the diagnosis position of each subfamily (for details see
26
27
28 160 Additional file 1: Table S2-4). To ensure the complete diagnosis positions in each subfamily
29
30
31 161 sequence, we used Ns to fill the gaps of each subfamily sequence that did not harbor the
32
33
34 162 diagnosis positions sites. Specially, we discovered the diagnosis position of L1 from previous
35
36 163 studies [28-31].

37
38 164 We then assembled the 'discordant reads' of each RIP into contigs using CAP3 [32] and
39
40
41 165 realigned them against all of the subfamily sequences by BLAST [33]. We determined the
42
43
44 166 maximum similarity score (S_{MS}) for each subfamily based on a simple penalty algorithm as
45
46
47 167 following:

$$168 \quad S_{MS} = \sum S_i$$

51
52 169 where S_i indicates the score of the specific diagnosis position i . $S_i = 1$ when the query
53
54
55 170 genotype was same as the diagnosis position i of this subfamily; $S_i = -0.5$ when the
56
57
58 171 position of query contigs were a gap while the diagnosis position i was not 'N', or the query

1 172 contigs were mismatch against the diagnosis position i .

2
3 173 We also determined divergence value (V_D) for each subfamily as following:

4
5
6 174
$$V_D = \frac{N_{mis} + N_{gap}}{L_{map}}$$

7
8
9 175 where N_{mis} and N_{gap} indicated mismatched base number and gaps number of query contigs,
10
11 176 respectively. The L_{map} stood for the mapped length of the certain subfamily.

12
13
14 177 Subfamily with the maximum similarity based on the genotype of diagnosis position would
15
16
17 178 be reported. If two or more subfamilies harbor the same maximum similarity, the subfamily with
18
19
20 179 the smallest divergence value would be reported.

21
22 180 We treated the retrotransposon subfamily classification in dbRIP as 'golden control', and
23
24
25 181 compared the classification result of 909 overlapped RIPs of our result and golden control, to
26
27
28 182 evaluate the accuracy of the subfamily classification.

29
30
31 183 During the contigs mapping to subfamily sequences, we identified the first mapped site of
32
33
34 184 5' and 3' ends of the subfamily sequence, and accordingly counted the lengths from the initial
35
36
37 185 site (L_{min} and L_{max}). The length of inserted retrotransposon (L_{retro}) was calculated as the
38
39 186 difference between the maximum and the minimum length of aligned sequence:

40
41
42 187
$$L_{retro} = L_{max} - L_{min} + 1$$

43
44
45 188 The assembled contigs of both 5' and 3' ends of insertions had the same orientation of HG19
46
47
48 189 sequence, which we defined as 'positive orientation'. If the mapping orientations of the contigs
49
50
51 190 were different, the orientation of RIPs was judged as the mapping orientation which most
52
53 191 contigs supported. Also, the poly-A tail of retrotransposon would be annotated if the RIPs is
54
55
56 192 'positive' and there were more than four 'A' bases in the first 6 bases at 3' end of the contigs.

57
58 193 And the poly-T tail the insertion orientation is 'negative' and there were more than four 'T' bases

1 194 in the first 6 bases at 5' end of the contigs.

2
3 195 **Reference RIPs detection**

4
5 196 The reference RIPs can be detected as a subset of deletions of the samples relative to the
6
7
8 197 reference of HG19 (for details see Additional file 2: Figure S1). These deletions were selected
9
10
11 198 from the structural variation (SV) detection result of YH90 samples (data not shown) and the
12
13
14 199 RIPs can be annotated based on matching deletion coordinates to HG19 annotation of
15
16 200 RepeatMasker (more than 90% overlap with each other) [34].

17
18
19 201 The reference RIPs should be absent in the chimpanzee genome. The alignments of
20
21
22 202 chimpanzee mapped to human genome was downloaded from UCSC
23
24 203 (<http://hgdownload.cse.ucsc.edu>). One reference RIP candidate should correspond to a
25
26
27 204 gap with an overlap of more than 90% to each other, and there would be no gaps in the
28
29
30 205 chimpanzee genome on this locus. The RIPs candidates would be filtered if there was no
31
32
33 206 polymorphism in YH90 samples (allele frequency equal to 180).

34
35
36 207 **Results**

37
38 208 **Non-reference retrotransposon insertions calling**

39
40 209 To detect non-reference RIPs from WGS data accurately and time-efficiently, we developed a
41
42
43 210 computer program called Specific Insertions Detector (SID). Through discordant reads
44
45
46 211 detection and reads clustering, it could detect non-reference RIPs easily and quickly (see
47
48
49 212 Materials and Methods).

50
51 213 To investigate the influence of sequencing depth on RIPs detection sensitivity and
52
53
54 214 accuracy, we simulated sequence data at different depth. It was observed that the detection
55
56
57 215 sensitivity dramatically increased along with rising sequencing depth, and achieved 95% when

1 216 the sequencing depth was more than 30x. By contrast, the accuracy of detected RIPs had
2
3 217 slight changes along with increasing sequencing depth (Fig. 2a).
4
5

6 218 We next estimated the RIPs detection sensitivity using two real sequencing datasets: One
7
8 219 dataset was CEU trio data, which was deep-sequenced (>75x) Illumina HiSeq data generated
9
10 220 at the Broad Institute (father NA12891, mother NA12892 and the female offspring NA12878)
11
12 221 from the 1000GP. We first used SID to detect RIPs of each individual in CEU dataset (for
13
14 222 details see Additional file 1: Table S5), and evaluated the sensitivity by comparing the
15
16 223 detection results with the PCR-validated datasets from Stewart et al. [12]. For Alu, the mean
17
18 224 sensitivity reached 96.3% among individuals. We also obtained mean sensitivity of 80.3% and
19
20 225 83.3% for L1 and SVA, respectively.
21
22
23
24
25
26
27

28 226 The other dataset including the NA18571, NA18572 and NA18537 were also recruited in
29
30 227 1000GP. The RIPs datasets of these three individuals detected by SID were much larger and
31
32 228 covered 89.59% of the same sample's results in 1000GP on average (for details see Additional
33
34 229 file 2: Figure S2).
35
36
37
38

39 230 We estimated the RIPs detection accuracy using the sequencing data of Asian individual
40
41 231 lymphocytic cell line (YH_CL, ~52x) that was the first Asian diploid genome dataset, and
42
43 232 performed the PCR validation straightly. We randomly selected 103 detected RIPs and 93/96
44
45 233 (7 loci were removed because of the poor primer specificity) loci were successfully validated,
46
47 234 indicating that SID had an accuracy of 90.29% - 96.88% (for details see Additional file 1: Table
48
49 235 S6). We also used the PCR validation result to assess our genotyping accuracy. It was about
50
51 236 93.55% (87/93, Fig. 2b, for details see Additional file 2: Supplementary Methods).
52
53
54
55
56
57
58
59

60 237

1 238

2
3 239 Table 1. Run time for three different RIPs-detection programs.
4
5

Sample	Data size (GB)	Run time (h)		
		SID	TEA	RetroSeq
YH_CL	85	16.4	36.9	31.2
NA18571	127	19.9	90.6	110.3
NA18537	124	21.9	95.2	94.8
NA18572	117	31.9	82	99.8

15 240
16
17

18 241 We next compared the RIPs detection efficiency of different methods (SID, RetroSeq [11]
19
20
21 242 and TEA [35]). In addition to YH_CL (Fig. 2c), we also selected three samples (NA18571,
22
23
24 243 NA18572 and NA18537) from YH90, which were sequenced at an average depth of 67.91x.
25
26
27 244 The run time of SID was about 3 times shorter than the other two methods, showing that the
28
29
30 245 SID was the most time-saving method of these three (Table 1). SID and TEA had comparable
31
32
33 246 sensitivity that was higher than RetroSeq, and the majority of SID detected RIPs (66.33% in
34
35
36 247 average) existed in TEA's result, and an average of 16.87% SID detected RIPs could be
37
38
39 248 generally detected by all three methods (Fig. 2c and Additional file 2: Figure S2). We also
40
41
42 249 validated the uniquely detected RIPs by PCR, and gained an accuracy of 75.86% and 77.78%
43
44
45 250 for Alu and L1, respectively, revealing a higher RIPs detection accuracy (Alu: 42.10% and
46
47
48 251 82.61%, L1: 66.67% and 66.67%, for RetroSeq and TEA, respectively).

252 **A comprehensive RIPs landscape of human population**

253 We then performed RIPs detection on a much larger scale. We sequenced 90 Han Chinese
254 individuals and generated Illumina paired-end sequence data at an average depth of 68x
255 each sample (for details see Additional file 1: Table S1). The dataset included two groups in
256 different regions of China, 45 samples from Northern China and 45 samples from Southern

1 257 China. Using the SID, the high depth of the dataset (much more than 30x) allowed us to build
2
3 258 a comprehensive non-reference RIPs landscape with high confidence.
4
5

6 259 In total we identified 9342 non-reference RIPs in autosome regions, including 6483 Alu
7
8 260 elements, 2398 L1s, 61 LTRs and 400 SVAs (Fig. 3a and for details see Additional file 1: Table
9
10 261 S7). Of this dataset, 8433 RIPs including 5826 Alu elements, 2169 L1s, 383 SVAs, 55 LTRs
11
12 262 were novel compared with dbRIP (Fig. 3b). The average number of non-reference RIPs per
13
14 263 individual was 1394 (ranging from 1304 to 1493, Fig. 3c), including 1110.80 Alu elements,
15
16 264 231.34 L1s, 43.14 SVAs and 9.01 LTRs, respectively, and each type of RIPs had similar
17
18 265 proportion ($P = 0.6364$, $P = 0.2711$, $P = 0.2128$, $P = 0.5582$, respectively, Wilcoxon
19
20 266 signed-rank test). We compared pair-wise individuals of all 90 samples, and the average
21
22 267 specific loci number was 672.79, almost a half (48.25%) of non-reference RIPs of one
23
24 268 individual.
25
26
27
28
29
30
31
32

33 269 The specific inserted location information enabled us to investigate genome-wide
34
35 270 sequence patterns of these non-reference RIPs. We observed that the non-reference RIPs
36
37 271 varied between chromosomes (Fig. 3d, e). Of note, we found that the different two
38
39 272 subpopulations (from Southern and Northern China respectively) had a similar pattern of RIPs
40
41 273 distribution ($r = 0.782$, Fig. 3e and for details see Additional file 2: Figure S3). However, we did
42
43 274 not find obvious correlation between the distribution of non-reference RIPs and GC content,
44
45 275 fixed RIPs, as well as single nucleotide polymorphisms (SNPs) of the same sample within 10M
46
47 276 non-N bins (Additional file 2: Figure S4).
48
49
50
51
52
53
54

55 277 To further investigate the distribution of non-reference RIPs in functional region, we
56
57 278 annotated all of the inserted loci (Fig. 3f). More than half of RIPs (4828/9324) were located in
58
59
60
61
62
63
64
65

1 279 gene regions, and the majority of these in introns. Only 5/9324 RIPs were located in protein
2
3 280 coding regions, including three genes C1orf66 (Alu inserted), SNX31 (Alu inserted) and
4
5
6 281 APH1B (SVA inserted) with rare frequency (1/90), and two genes ADORA3 (Alu inserted) and
7
8
9 282 Slco1b3 (L1 inserted) with higher frequency (44/90 and 12/90, respectively). Compared with all
10
11 283 known genes, we noticed that RIPs inserted genes generally had a relatively lower expression
12
13
14 284 level, later replication time in the cell cycle, lower GC content and lower conservation (higher
15
16
17 285 Ka/Ks) than the average level of all genes, respectively ($P < 0.001$, Wilcoxon Test; Fig. 4).
18
19
20 286 Besides the gene regions , we also found that at average 9.78% and 4.93% RIPs were located
21
22
23 287 in enhancer regions and promoter regions per sample, respectively (Fig. 3f).

24
25 288 Furthermore, we annotated the subfamily, the orientation and the sequence length of all
26
27
28 289 detected inserted retrotransposons based on regional sequence assembly and remapping to
29
30
31 290 the retrotransposon library. AluY sub-family constituted essentially all non-reference Alu
32
33
34 291 insertions, in which AluYa5 and AluYb8 were mostly active (for details see Additional file 1:
35
36
37 292 Table S7), supporting conclusions from previous studies[28, 36, 37]. L1 insertions were
38
39
40 293 dominated by the sub-family of L1-Pre, which was also in line with previous report [28].

41
42 294 The orientation of one RIP is judged from the mapping orientation of contigs to
43
44
45 295 retrotransposon reference and the existing of poly-A or poly-T tails of inserted sequence (for
46
47
48 296 details see Additional file 1: Table S7). Previous studies reported that the gene-inserted RIP
49
50
51 297 had a greater influence on gene expression if it was inserted on the same orientation with the
52
53
54 298 target gene [2, 38]. However, we detected a comparable number of direct and reverse events
55
56
57 299 (0.475 and 0.525, respectively), arguing against an obvious natural selection on the RIPs with
58
59
60 300 consistent orientation with the inserted gene.

1 301 Along with subfamily and orientation annotation, we also calculated the length of each
2
3 302 insertion sequence. We found that different types of RIP had different length distributions
4
5
6 303 (Additional file 2: Figure S5). More than half of Alu elements (~70%) were full-length while the
7
8
9 304 length of the L1 distributed more discretely. Most of L1s (> 80%) were fractured during the
10
11
12 305 process of retrotransposon, which verified previous study [13].

13 306 **RIPs reference of healthy people**

14
15
16 307 The pure and comprehensive dataset of RIPs can be used as a baseline of healthy people for
17
18
19 308 other disease-related research, especially in single-gene disease. The candidate
20
21
22 309 disease-related retrotransposon insertions that were found in this dataset would be filtered.
23
24
25 310 We explicitly measured the overlap between our dataset and the disease-related
26
27
28 311 retrotransposon insertions data in dbRIP (<http://dbrip.org>) [39]. None of the insertion sites
29
30
31 312 existed in our dataset, indicating the accuracy of the database. We also tested some data of
32
33
34 313 cancer research. We tested the dataset of candidate cancer related somatic retrotransposon
35
36
37 314 insertions which was strictly generated from 11 tumor types data of The Cancer Genome Atlas
38
39
40
41 315 (TCGA) Pan-Cancer Project. None of overlapped RIPs were detected, whereas 43.36%
42
43
44 316 germline retrotransposons were detected. According to the comparison of colon cancer
45
46
47 317 specific data [9], we found two L1 insertions consistent with our dataset with frequency of
48
49
50 318 51/90 and 50/90. These two L1 insertions were germline retrotransposon insertions that were
51
52
53 319 further validated by PCR validation in Solyom's research. We also tested the candidate of
54
55
56 320 Hepatocellular Carcinoma specific insertions [8] and found one L1 insertion was also present
57
58
59 321 in our dataset with frequency of 9/90. This site was finally validated as a germline insertion by
60
61
62 322 PCR in that research. All of these indicated that our data provided a reference panel to wipe off
63
64
65

1 323 false positive insertions related to cancer.
2
3

4 324 **Conclusions**
5
6
7

8 325 In this paper, we developed a computer program SID to detect non-reference RIPs of 90
9
10 326 healthy Han Chinese individuals through high depth whole-genome sequencing. Compared
11
12 327 with TEA and RetroSeq, the SID has the fastest detection speed as well as high sensitivity and
13
14 328 accuracy. We described the landscape of RIPs distribution on population genomes, and
15
16 329 annotated the subfamily, orientation, and length of RIPs. We demonstrated that the RIPs could
17
18 330 be used as a normal baseline for retrotransposon related disease research.
19
20
21
22
23

24 331 To our knowledge, this dataset is the largest dataset for human by now. Compared with
25
26 332 1000GP result of the same samples, the majority (mean 69.68%) of RIPs in our dataset has
27
28 333 not been previously observed, suggesting that our deep-sequenced data had much higher
29
30 334 detection sensitivity than the low coverage ones. For example, it was reported that the serum
31
32 335 ACE level was determined by the Alu insertion/deletion (I/D) polymorphism in the following
33
34 336 order: DD > ID > II [40], and the D allele of ACE gene was found to be associated with
35
36 337 essential hypertension in different populations [41-44]. We found that ACE gene harbored Alu
37
38 338 insertion in the 15th intron with a frequency of 81/90 in our 90 Chinese genomes, compared
39
40 339 with the much lower frequency (7/63) in CEPH individuals [12], which was supported by
41
42 340 previous study [45]. To our surprise, during the analysis of retrotransposon insertions of ACE,
43
44 341 we found that there was no RIPs of ACE in Han Chinese samples of 1000GP dataset, which
45
46 342 was a high-frequency inserted gene in our RIPs data. ACE specific PCR validation (for details
47
48 343 see Additional file 2: Figure S6) and previous study of ACE [46] indicated our result was in line
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 344 with the real situation. It can be seen that the enough depth of sequencing is very important to
2
3 345 investigate RIPs frequency and our data is able to present a better result in line with the actual
4
5
6 346 situation. The highly sensitive and accurate RIPs dataset gave us a perfect opportunity to
7
8
9 347 perform fitness analysis of RIPs.

10
11 348 This dataset can be used to compare with others to give guidance to research the
12
13
14 349 disease-causing mechanisms in particular population and successfully coalesced the insert
15
16
17 350 time of a specific locus. This dataset can also be used as a standard to other RIPs research
18
19
20 351 and can be a baseline to filter meaningless RIPs in the disease-causing retrotransposon
21
22
23 352 research. Genome-wide Association Studies (GWAS) have proven their utility in identifying
24
25
26 353 genomic variants associated with risk for many diseases. Unlike SNPs and copy number
27
28
29 354 variations (CNVs) that were widely used in GWAS, RIPs, the major contributor to human
30
31
32 355 variation, have always been overlooked. It is significant that this dataset provides a valuable
33
34
35 356 source to do GWAS and collects more markers related to complex diseases.

36 357 The high cost of whole-genome sequencing at high depth is still the main limitation,
37
38
39 358 preventing it from being widely used in TEs research. Furthermore, the large amount of data
40
41
42 359 yielded by high depth WGS makes it difficult to undertake bioinformatic analysis. With the
43
44
45 360 development of biotechnology (BT) and IT, this situation would be changed soon. However, it
46
47
48 361 may not be solved in a short time.

49
50 362 The next step is to research RIPs on the transcriptome level. The impact of RIPs on gene
51
52
53 363 expression is still unclear. Combining the genome and transcriptome would give us a
54
55
56 364 comprehensive picture of the regulation of RIPs. In this way, we can further expound the
57
58
59 365 position of the retrotransposon in the course of human evolution.

1 366

2
3 367 **Availability and requirements**

4
5
6
7 368 ● Project name: SID

8
9
10 369 ● Project home page: <https://github.com/Jonathanyu2014/SID>

11
12 370 ● Operating system(s): Linux

13
14
15 371 ● Programming language: Perl

16
17
18 372 ● Other requirements: Perl 5.14 or higher

19
20
21 373 ● License: Apache License 2.0

22
23
24 374 ● Any restrictions to use by non-academics: None

25
26
27 375 **Additional files**

28
29
30 376 Additional file 1: Supporting data description and the results of RIPs calling. (XLSX 1797 kb)

31
32
33 377 Additional file 2: The method of RIP candidate determination and all the supplementary figures.

34
35
36 378 (PDF 1095 kb)

37
38
39
40 379 **Abbreviations**

41
42
43 380 RIP, retrotransposon insertion polymorphism; TE, transposable element; LTR, long terminal

44
45
46 381 repeat; L1, long interspersed nuclear element 1; WGS, whole genome sequencing; NGS,

47
48
49 382 next-generation sequencing; SID, specific insertions detector; TSD, target site

50
51
52 383 duplication;1000GP, 1000 Genomes Project; CNV, copy number variation; SNP, single

53
54
55 384 nucleotide polymorphism; ENA, European Nucleotide Archive; GWAS, genome-wide

56
57
58 385 association study.

1 386 **Acknowledgments**

2
3
4 387 We are grateful for Zengli Yan, Nan Li, Na Li and Runze Jiang for optimizing and testing the
5
6
7 388 SID program. We thank Haoxiang Lin and Wenjuan Zhu for providing technical assistance to
8
9
10 389 us. We acknowledge the support by the 1000 Genomes Project Consortium. This work was
11
12
13 390 supported by the Shenzhen Municipal Government of China [JSGG20140702161347218] and
14
15 391 [KQCX20150330171652450].
16
17
18

19 392 **Availability of data and materials**

20
21
22 393 The source code of SID is available in GitHub repository[47]. The human (Homo sapiens)
23
24
25 394 reference genome sequence (HG19) and its annotation files were downloaded from UCSC
26
27
28 395 Genome Bioinformatics (<http://genome.ucsc.edu/>). The raw sequence data of YH_CL from
29
30
31 396 previous reports is available in ENA repository (accession number ERA000005) [48]. All the
32
33
34 397 YH90 raw sequences have been released in ENA repository (accession number ERA496654).
35
36

37 398 **Authors' contribution**

38
39
40
41 399 BL, SL, YH initiated this project and reviewed the manuscript. QY, XZ, YZ drafted the
42
43
44 400 manuscript. XH, JL polished up the manuscript. QY, WZ, XZ, YW performed the data analysis
45
46
47 401 and drew the pictures. YZ, YW designed and developed the SID program. NL, XZ, GL
48
49
50 402 conducted the experiment for sequencing. LX designed the primers and did the PCR validation.
51
52
53 403 YH, BL, SL, XZ, XG contributed with fruitful discussions.
54

55 404 **Competing interests**

56
57
58
59 405 The authors declare that they have no competing interests.
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

406 **Author details**

407 ¹ BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,
408 China. ² BGI-Shenzhen, Shenzhen 518083, China. ³ BGI College, Shenzhen 518083, China.
409 ⁴ Department of Biology, University of Copenhagen, Copenhagen 1599, Denmark. ⁵ School of
410 Life Sciences, Sun Yat-sen University, Guangzhou 510006, China. ⁶ BGI-Forensics,
411 Shenzhen 518083, China.

412 **Ethics, consent and permissions**

413 This study was approved by BGI-IRB.

414 **Consent to publish**

415 Both BGI-IRB and participants involved consented to publish this research.

416

417 **References**

418 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle
419 M, FitzHugh W *et al*: **Initial sequencing and analysis of the human genome**. *Nature* 2001,
420 **409**(6822):860-921.
421 2. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution**.
422 *Nature reviews Genetics* 2009, **10**(10):691-703.
423 3. Kidd JM, Graves T, Newman TL, Fulton R, Hayden HS, Malig M, Kallicki J, Kaul R, Wilson
424 RK, Eichler EE: **A human genome structural variation sequencing resource reveals**

1 425 **insights into mutational mechanisms. *Cell* 2010, **143**(5):837-847.**

2

3 426 4. Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran JV, Kazazian HH, Jr.: **Hot**

4

5

6 427 **L1s account for the bulk of retrotransposition in the human population. *Proceedings of the***

7

8

9 428 *National Academy of Sciences of the United States of America* 2003, **100**(9):5280-5285.

10

11 429 5. Xing J, Zhang Y, Han K, Salem AH, Sen SK, Huff CD, Zhou Q, Kirkness EF, Levy S, Batzer

12

13 430 **MA et al: Mobile elements create structural variation: analysis of a complete human**

14

15

16 431 **genome. *Genome research* 2009, **19**(9):1516-1526.**

17

18

19 432 6. Cordaux R, Hedges DJ, Herke SW, Batzer MA: **Estimating the retrotransposition rate of**

20

21

22 433 **human Alu elements. *Gene* 2006, **373**:134-137.**

23

24

25 434 7. Hancks DC, Kazazian HH, Jr.: **Active human retrotransposons: variation and disease. *Curr***

26

27 435 *Opin Genet Dev* 2012, **22**(3):191-203.

28

29

30 436 8. Shukla R, Upton KR, Munoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, Brennan PM,

31

32 437 Baillie JK, Collino A, Ghisletti S et al: **Endogenous retrotransposition activates oncogenic**

33

34

35 438 **pathways in hepatocellular carcinoma. *Cell* 2013, **153**(1):101-111.**

36

37

38 439 9. Solyom S, Ewing AD, Rahrmann EP, Doucet T, Nelson HH, Burns MB, Harris RS, Sigmon DF,

39

40 440 Casella A, Erlanger B et al: **Extensive somatic L1 retrotransposition in colorectal tumors.**

41

42 441 *Genome research* 2012, **22**(12):2328-2338.

43

44

45 442 10. Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ 3rd, Lohr JG, Harris CC, Ding L,

46

47 443 Wilson RK et al: **Landscape of somatic retrotransposition in human cancers. *Science* 2012,**

48

49 444 **337**(6097):967-971.

50

51

52 445 11. Keane TM, Wong K, Adams DJ: **RetroSeq: transposable element discovery from**

53

54 446 **next-generation sequencing data. *Bioinformatics* 2013, **29**(3):389-390.**

55

56

57

58

59

60

61

62

63

64

65

- 1 447 12. Stewart C, Kural D, Stromberg MP, Walker JA, Konkel MK, Stutz AM, Urban AE, Grubert F,
2
3 448 Lam HY, Lee WP *et al*: **A comprehensive map of mobile element insertion polymorphisms**
4
5
6 449 **in humans**. *PLoS Genet* 2011, **7**(8):e1002236.
7
8
9 450 13. Ewing AD, Kazazian HH, Jr.: **Whole-genome resequencing allows detection of many rare**
10
11 451 **LINE-1 insertion alleles in humans**. *Genome research* 2011, **21**(6):985-990.
12
13
14 452 14. Xing J, Witherspoon DJ, Jorde LB: **Mobile element biology: new possibilities with**
15
16 453 **high-throughput sequencing**. *Trends in genetics : TIG* 2013, **29**(5):280-289.
17
18
19 454 15. Wong K, Keane TM, Stalker J, Adams DJ: **Enhanced structural variant and breakpoint**
20
21 455 **detection using SVMerge by integration of multiple detection methods and local assembly**.
22
23 456 *Genome biology* 2010, **11**(12):R128.
24
25
26 457 16. Kang H, Zhu D, Lin R, Opiyo SO, Jiang N, Shiu SH, Wang GL: **A novel method for**
27
28 458 **identifying polymorphic transposable elements via scanning of high-throughput short**
29
30 459 **reads**. *DNA Res* 2016, **23**(3):241-251.
31
32
33
34
35
36 460 17. Henaff E, Zapata L, Casacuberta JM, Ossowski S: **Jitterbug: somatic and germline**
37
38 461 **transposon insertion detection at single-nucleotide resolution**. *BMC genomics* 2015,
39
40 462 **16**:768.
41
42
43
44 463 18. Fiston-Lavier AS, Barron MG, Petrov DA, Gonzalez J: **T-lex2: genotyping, frequency**
45
46 464 **estimation and re-annotation of transposable elements using single or pooled**
47
48 465 **next-generation sequencing data**. *Nucleic acids research* 2015, **43**(4):e22.
49
50
51
52
53 466 19. Thung DT, de Ligt J, Vissers LE, Steehouwer M, Kroon M, de Vries P, Slagboom EP, Ye K,
54
55 467 Veltman JA, Hehir-Kwa JY: **Mobster: accurate detection of mobile element insertions in**
56
57 468 **next generation sequencing data**. *Genome biology* 2014, **15**(10):488.
58
59
60
61
62
63
64
65

1 469 20. Wang J, Wang W, Li R, Li Y, Tian G, Goodman L, Fan W, Zhang J, Li J, Zhang J *et al*: **The**
2
3 470 **diploid genome sequence of an Asian individual.** *Nature* 2008, **456**(7218):60-65.
4
5
6 471 21. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform.**
7
8 472 *Bioinformatics* 2009, **25**(14):1754-1760.
9
10
11 473 22. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K,
12
13 474 Altshuler D, Gabriel S, Daly M *et al*: **The Genome Analysis Toolkit: a MapReduce**
14
15 475 **framework for analyzing next-generation DNA sequencing data.** *Genome research* 2010,
16
17 476 **20**(9):1297-1303.
18
19
20 477 23. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
21
22 478 Genome Project Data Processing S: **The Sequence Alignment/Map format and SAMtools.**
23
24 479 *Bioinformatics* 2009, **25**(16):2078-2079.
25
26
27 480 24. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J: **Rebase Update,**
28
29 481 **a database of eukaryotic repetitive elements.** *Cytogenet Genome Res* 2005,
30
31 482 **110**(1-4):462-467.
32
33
34 483 25. Wang J, Song L, Grover D, Azrak S, Batzer MA, P L: **dbRIP: a highly integrated database of**
35
36 484 **retrotransposon insertion polymorphisms in humans.** *Hum Mutat* 2006, **27**(4):323-329.
37
38
39 485 26. Hu X, Yuan J, Shi Y, Lu J, Liu B, Li Z, Chen Y, Mu D, Zhang H, Li N: **pIRS: Profile-based**
40
41 486 **Illumina pair-end reads simulator.** *Bioinformatics* 2012, **28**(11):1533-1535.
42
43
44 487 27. Baillie JK, Barnett MW, Upton KR, Gerhardt DJ, Richmond TA, De Sapio F, Brennan PM,
45
46 488 Rizzu P, Smith S, Fell M *et al*: **Somatic retrotransposition alters the genetic landscape of the**
47
48 489 **human brain.** *Nature* 2011, **479**(7374):534-537.
49
50
51 490 28. Boissinot S, Chevret P, AV F: **L1 (LINE-1) retrotransposon evolution and amplification in**
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 491 **recent human history.** *Mol Biol Evol* 2000, **17**(6):915-928.

2

3 492 29. Dombroski BA, Mathias SL, Nanthakumar E, Scott AF, Jr KH: **Isolation of an active human**

4

5

6 493 **transposable element.** *Science* 1991, **254**(5039):1805-1808.

7

8

9 494 30. Ovchinnikov I, Rubin A, GD S: **Tracing the LINEs of human evolution.** *Proceedings of the*

10

11 495 *National Academy of Sciences of the United States of America* 2002, **99**(16):10522-10527.

12

13

14 496 31. Ovchinnikov I, Troxel AB, GD S: **Genomic characterization of recent human LINE-1**

15

16

17 497 **insertions: evidence supporting random insertion.** *Genome research* 2001,

18

19

20 498 **11**(12):2050-2058.

21

22

23 499 32. Huang X, Madan A: **CAP3: A DNA sequence assembly program.** *Genome research* 1999,

24

25 500 **9**(9):868-877.

26

27

28 501 33. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.**

29

30

31 502 *Journal of molecular biology* 1990, **215**(3):403-410.

32

33

34 503 34. Tarailo-Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in**

35

36 504 **genomic sequences.** *Curr Protoc Bioinformatics* 2009, **Chapter 4**:Unit 4 10.

37

38

39 505 35. Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ, 3rd, Lohr JG, Harris CC, Ding L,

40

41

42 506 Wilson RK *et al*: **Landscape of somatic retrotransposition in human cancers.** *Science* 2012,

43

44 507 **337**(6097):967-971.

45

46

47 508 36. Hormozdiari F, Alkan C, Ventura M, Hajirasouliha I, Malig M, Hach F, Yorukoglu D, Dao P,

48

49

50 509 Bakhshi M, Sahinalp SC *et al*: **Alu repeat discovery and characterization within human**

51

52 510 **genomes.** *Genome research* 2011, **21**(6):840-849.

53

54

55 511 37. Batzer MA, Deininger PL: **Alu repeats and human genomic diversity.** *Nature reviews*

56

57

58 512 *Genetics* 2002, **3**(5):370-379.

59

60

61

62

63

64

65

1 513 38. Burns KH, Boeke JD: **Human transposon tectonics**. *Cell* 2012, **149**(4):740-752.

2

3 514 39. Wang J, Song L, Grover D, Azrak S, Batzer MA, Liang P: **dbRIP: a highly integrated**

4

5

6 515 **database of retrotransposon insertion polymorphisms in humans**. *Hum Mutat* 2006,

7

8

9 516 **27**(4):323-329.

10

11 517 40. Rigat B, Hubert C, Alhenc-Gelas F, Cambien F, Corvol P, F S: **An insertion/deletion**

12

13

14 518 **polymorphism in the angiotensin I-converting enzyme gene accounting for half the**

15

16

17 519 **variance of serum enzyme levels**. *J Clin Invest* 1990, **86**(4):1343-1346.

18

19

20 520 41. Jeng JR, Harn HJ, Jeng CY, Yueh KC, SM S: **Angiotensin I converting enzyme gene**

21

22

23 521 **polymorphism in Chinese patients with hypertension**. *Am J Hypertens* 1997,

24

25 522 **10**(5Pt1):558-561.

26

27

28 523 42. Zee RY, Lou YK, Griffiths LR, BJ M: **Association of a polymorphism of the angiotensin**

29

30

31 524 **I-converting enzyme gene with essential hypertension**. *Biochem Biophys Res Commun* 1992,

32

33

34 525 **184**(1):9-15.

35

36 526 43. Asamoah A, Yanamandra K, Thurmon TF, Richter R, Green R, Lakin T, C M: **A deletion in the**

37

38

39 527 **angiotensin converting enzyme (ACE) gene is common among African Americans with**

40

41

42 528 **essential hypertension**. *Clin Chim Acta* 1996, **254**(1):41-46.

43

44

45 529 44. Duru K, Farrow S, Wang JM, Lockette W, T K: **Frequency of a deletion polymorphism in the**

46

47

48 530 **gene for angiotensin converting enzyme is increased in African-Americans with**

49

50

51 531 **hypertension**. *Am J Hypertens* 1994, **7**(8):759-762.

52

53 532 45. Anand SS, Yusuf S, Vuksan V, Devanesen S, Teo KK, Montague PA, Kelemen L, Yi C, Lonn E,

54

55

56 533 Gerstein H *et al*: **Differences in risk factors, atherosclerosis, and cardiovascular disease**

57

58

59 534 **between ethnic groups in Canada: the Study of Health Assessment and Risk in Ethnic**

535 **groups (SHARE). *Lancet* 2000, 356(9226):279-284.**

536 46. Batzer MA, Stoneking M, Alegria-Hartman M, Bazan H, Kass DH, Shaikh TH, Novick GE,
537 Ioannou PA, Scheer WD, RJ H: **African origin of human-specific polymorphic Alu**
538 **insertions. *Proceedings of the National Academy of Sciences of the United States of America***
539 1994, **91**(25):12288-12292.

540 47. Yu Q: **Specific Insertions Detector**. In: *Zenodo*. 2016.

541 48. Zong C, Lu S, Chapman AR, Xie XS: **Genome-wide detection of single-nucleotide and**
542 **copy-number variations of a single human cell. *Science* 2012, 338(6114):1622-1626.**

544 **Figure legends**

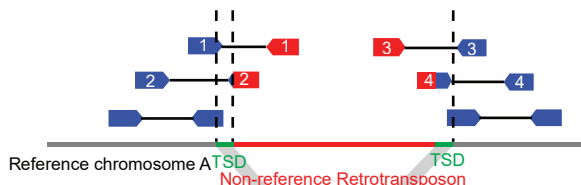
545 Fig. 1 The principle of retrotransposon insertions detection. (a) The SID schematic diagram for
546 RIPs detection in genome. TSD: target site duplication. SID: Specific Insertions Detector. (b) A
547 reads mapping example of predicted homozygous insertions. (c) A reads mapping example of
548 predicted heterozygous insertions. In (b) and (c), the red bases indicate the mismatches and
549 the sequences with orange background stand for the clipped part of the reads. The clipped
550 reads come from one allele with inserted retrotransposons and the normal reads come from
551 the other allele that same with the reference. The three reads with asterisk show no clipped
552 part but with terminal mismatches, which also can support the breakpoint and have
553 consistency with the clipped reads.

554 Fig. 2 Assessing the results of SID. (a) Detecting accuracy and sensitivity estimation along
555 cumulating sequencing depth of simulated data. (b) RIPs genotyping of YH_CL. The validation

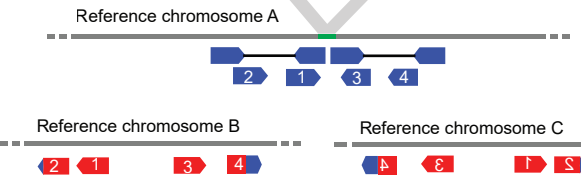
1 556 results by PCR were marked. HEE: estimated heterozygous site. HOE: estimated
2
3 557 homozygous site. HEV: validated heterozygous site. HOV: validated homozygous site. The
4
5
6 558 dash line shows the estimated boundary between heterozygous and heterozygous site. Note
7
8
9 559 that some of the validated RIPs stand in the same locus in the plot figure (for details see
10
11 560 Additional file 1: Table S6). (c) RIPs detection results of YH_CL by three different programs.
12
13
14 561 Adjacent 100bp regions of RIPs were taken into consideration.
15
16
17
18 562 Fig. 3 Comprehensive landscape of non-reference RIPs of YH90. (a) Proportions of novel
19
20
21 563 insertions found for each kind of retrotransposon. (b) Comparison of YH90 non-reference RIPs
22
23 564 results with dbRIP. Adjacent 100bp regions of RIPs were taken into consideration. (c) TE
24
25
26 565 distribution of each YH sample. (d) Box plots of non-reference RIPs distribution among
27
28
29 566 autosomes. (e) TE frequency distribution among YH90 samples. Rings from outer to inner
30
31
32 567 stand for Alu insertions frequency, L1 insertion frequency, SVA insertion frequency, LTR
33
34
35 568 insertion frequency and cytobands structure, respectively. The inside frequency of rings stands
36
37
38 569 for northern people's insertion frequency and the outside ones stand for southern people's. (f)
39
40 570 The RIPs distribution in different functional region of genome.
41
42
43
44 571 Fig. 4 Nature of non-reference TE inserted genes. The retrotransposons inserted genes were
45
46
47 572 compared with all annotated genes of UCSC in aspects of gene expression (a), replication
48
49
50 573 time of cell cycle (b), GC content (c) and conservation (d). The conservation of genes is
51
52
53 574 represented by Ka/Ks ratio, and the replication time ranges approximately from 100 (very early)
54
55 575 to 1000 (very late).
56
57
58
59
60
61
62
63
64
65

Figure 1

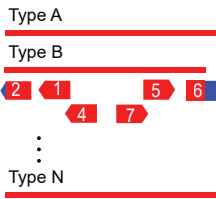
Sample genome



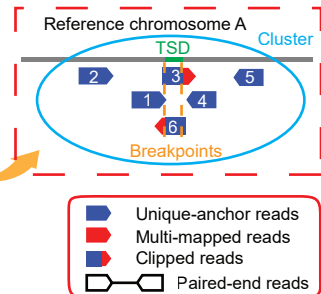
Hg19 genome



Retrotransposon library



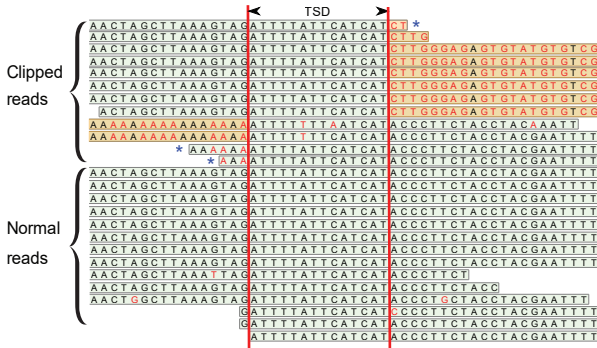
Detection result



B

[Click here to download Figure Fig. 1.pdf](#)

Reference AACTAGCTTAAAGTAGATTTTATT CAT CATA CCCTTCTACCTACGAAATTT



C

Reference GGTTTACCAATTAGTCTCCCTTAAAGGGCACTGTTTATGATCATCACCACAAATGGATGCA

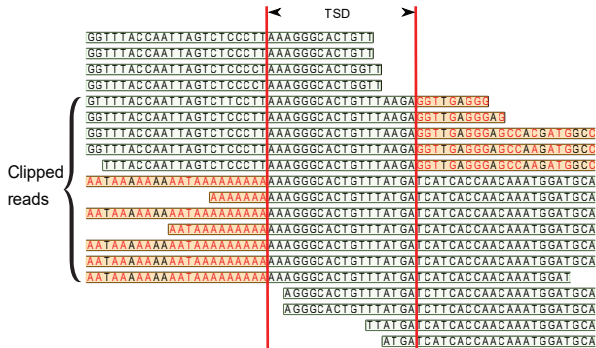
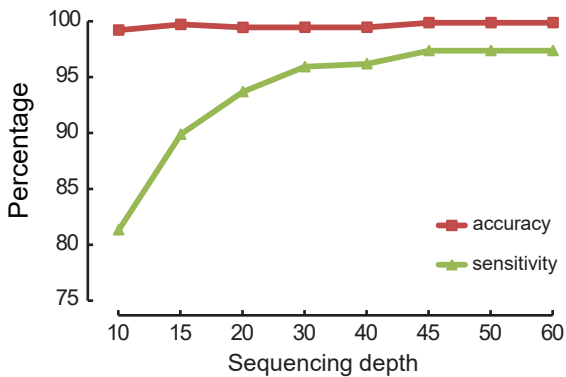


Figure 2

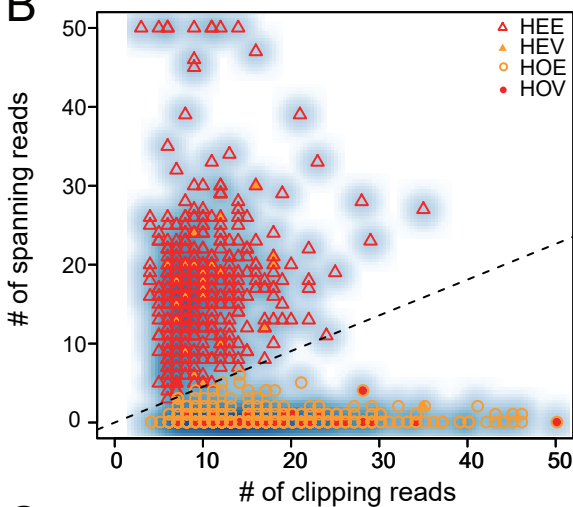
[Click here to download Figure Fig.](#)



A



B



C

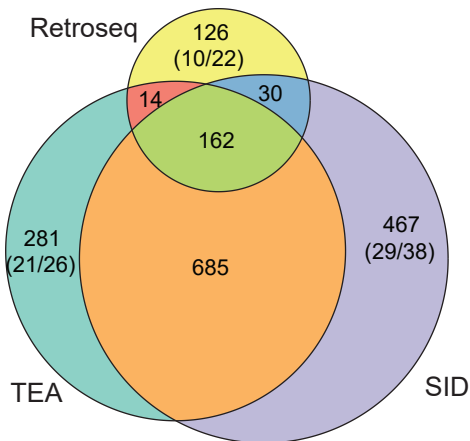


Figure 3

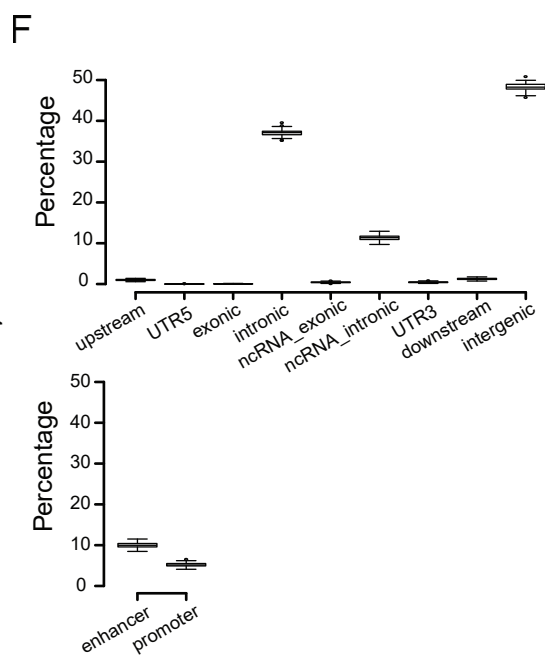
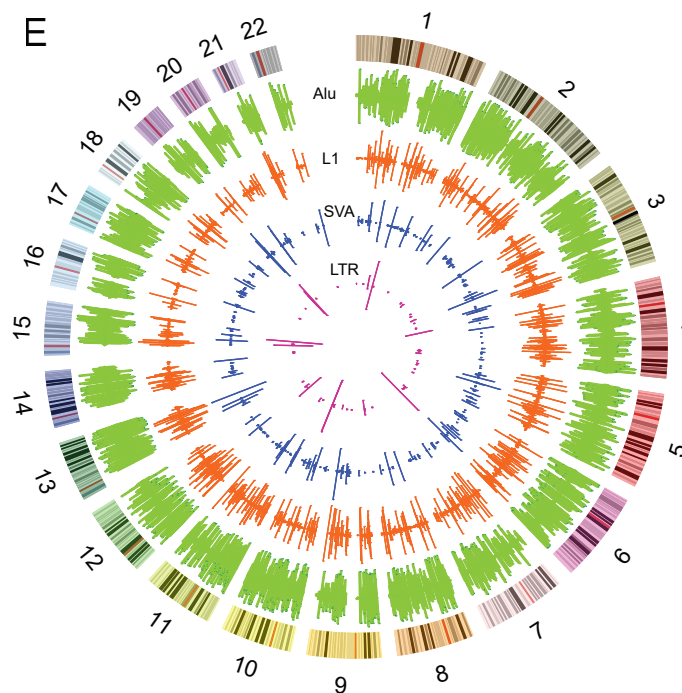
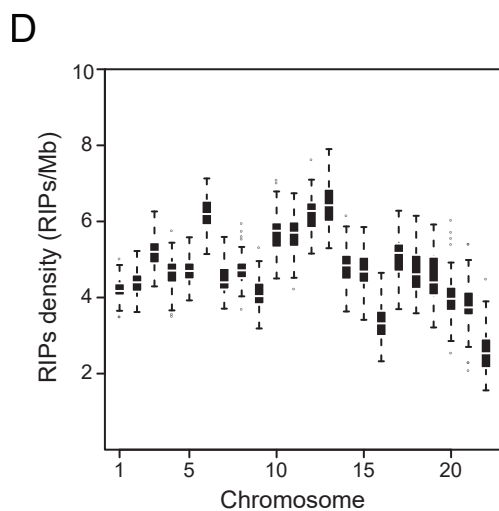
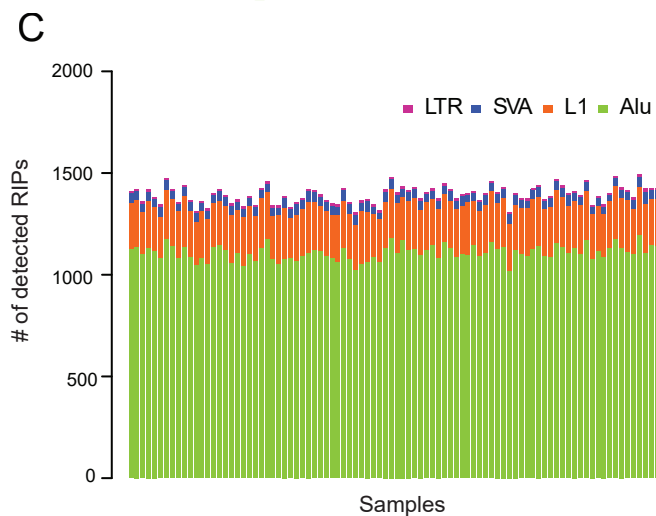
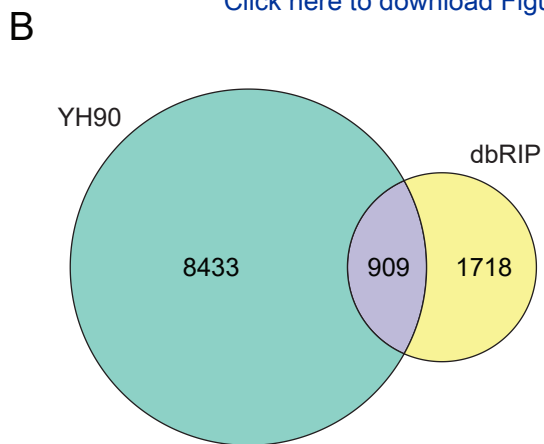
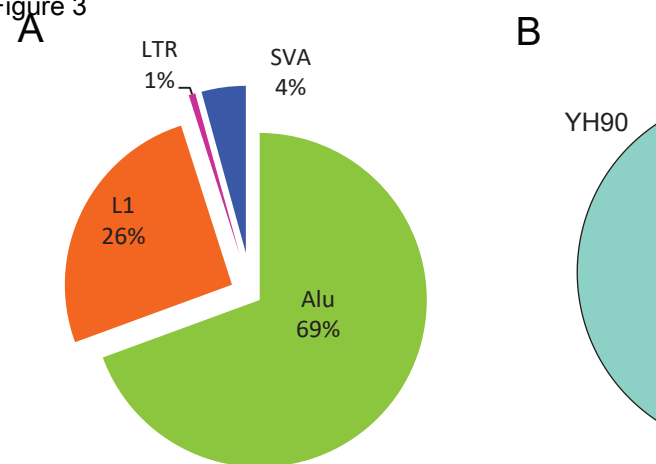
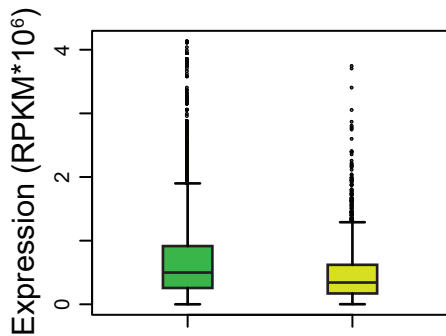
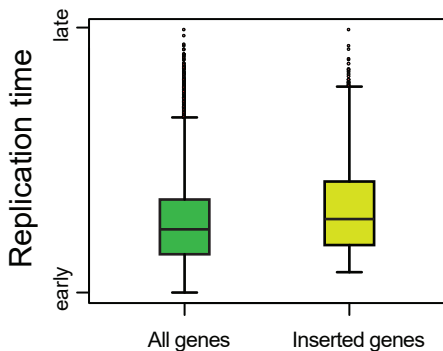
[Click here to download Figure Fig. 3.pdf](#)

Figure 4

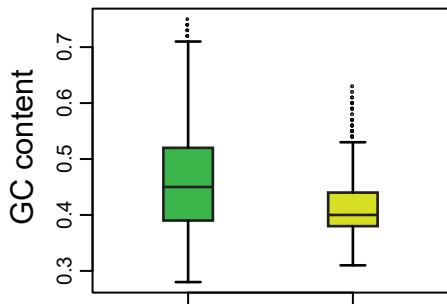
A



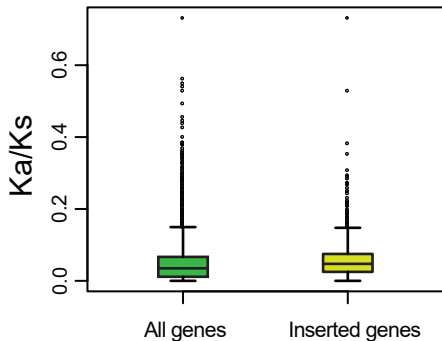
B

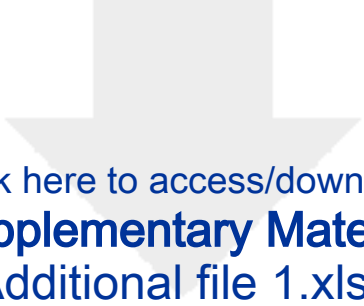
[Click here to download Figure Fig. 4.pdf](#)

C




D





Click here to access/download
Supplementary Material
Additional file 1.xlsx



Click here to access/download
Supplementary Material
Additional file 2.pdf