

1 Population-wide Sampling of Retrotransposon 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, Xiaoyun Huang², Nannan Li², Xinlan Zhou², 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 houyong@genomics.cn; Bo Li: libo@genomics.cn; Shiping Liu: liushiping@genomics.cn.

1 23 **Abstract**

2
3 24 **Background:** Active retrotransposons play important roles during evolution and continue to
4
5 25 shape our genomes today, especially in genetic polymorphisms underlying a diverse set of
6
7 26 diseases. However, studies of human retrotransposon insertion polymorphisms (RIPs) based
8
9 27 on whole-genome deep sequencing at the population level have not been sufficiently
10
11 28 undertaken, despite the obvious need for a thorough characterization of RIPs in the general
12
13 29 population.

14
15 30 **Findings:** Herein, we present a novel and efficient computational tool named Specific
16
17 31 Insertions Detector (SID) for the detection of non-reference RIPs. We demonstrate that SID is
18
19 32 suitable for high depth whole-genome sequencing (WGS) data using paired-end reads
20
21 33 obtained from simulated and real datasets. We construct a comprehensive RIP database
22
23 34 using a large population of 90 Han Chinese individuals with a mean 68x depth per individual.
24
25 35 In total, we identify 9342 recent RIPs, and 8433 of these RIPs are novel compared with dbRIP,
26
27 36 including 5826 Alu, 2169 long interspersed nuclear element 1 (L1), 383 SVA, and 55 long
28
29 37 terminal repeats (LTR). Among the 9342 RIPs, 4828 were located in gene regions and five
30
31 38 were located in protein-coding regions. We demonstrate that RIPs can, in principle, be an
32
33 39 informative resource to perform population evolution and phylogenetic analyses. Taking the
34
35 40 demographic effects into account, we identify a weak negative selection on SVA and L1 but
36
37 41 approximately neutral selection for Alu elements based on the frequency spectrum of RIPs.

38
39 42 **Conclusions:** SID is a powerful open-source program for the detection of non-reference RIPs.
40
41 43 We built a non-reference RIP dataset that greatly enhanced the diversity of RIPs detected in
42
43 44 the general population and should be invaluable to researchers interested in many aspects of
44
45 45 human evolution, genetics, and disease. As a proof-of-concept, we demonstrate that the RIPs
46
47 46 can be used as biomarkers in a similar way as single nucleotide polymorphisms (SNPs).

48
49 47 **Keywords:** Transposable element, retrotransposon insertion polymorphism, next-generation
50
51 48 sequencing, whole-genome sequencing

52
53
54
55 49

56
57 50

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

80 Chinese individuals (YH90) acquired at a mean depth of 68x.

81 **Materials and methods**

82 **Samples and whole genome sequencing**

83 We obtained B-lymphocyte cell lines from 90 Han Chinese individuals at the Coriell Institute
84 (Camden, New Jersey, USA). These individuals were selected from Beijing, Hunan province
85 and Fujian province, respectively. We broadly separated the samples into “Northern group” (45
86 samples) and “Southern group” (45 samples). DNA was extracted from the B-lymphocyte cells
87 of each individual, and libraries were then constructed following the manufacturer’s
88 instructions. High-coverage paired-end 100 bp WGS libraries were sequenced on the Illumina
89 HiSeq 2000 Platform. In addition, we also used a Chinese sample [16] for which the data were
90 previously released in the European Nucleotide Archive (ENA) repository (Additional file 1:
91 Table S1). The Institutional Review Board on Bioethics and Biosafety of BGI (BGI-IRB)
92 approved the study.

93 **Processing of the WGS data**

94 Reads were aligned to the human genome reference (HG19, Build37) using *BWA* [17].
95 Duplications were removed using Picard tools, and the quality values of each reads were
96 recalibrated using the Genome Analysis Toolkit (GATK) [18]. The resulting Binary
97 Alignment/Map (BAM) files were used as input for SID (Additional file 2: Text S1).

98 **The specific insertion detector pipeline**

99 SID is compiled in Perl and includes the following two steps: discordant reads detection and
100 reads clustering. Generally, the first step collects informative reads and generates other
101 necessary files, whereas the second step discovers the specific insertion sites and exports the
102 final results into plain text.

103 *Detection of discordant reads.* The “discordant reads” were extracted for the subsequent
104 clustering step. Paired-end reads were determined as “discordant reads” if they met one of the

1 105 following criteria: a. one read mapped to HG19 uniquely and the other read mapped to the
2
3 106 retrotransposon library (multi-mapped or unmapped to HG19); b. one read mapped to HG19
4
5
6 107 uniquely and the other soft-clipped read mapped to HG19, and the clipped sequence could be
7
8
9 108 mapped to the retrotransposon library; c. one soft-clipped read mapped to HG19, and the
10
11
12 109 clipped sequence could be mapped to the retrotransposon library. The other read mapped to
13
14 110 the retrotransposon library (multi-mapped or unmapped to HG19). The retrotransposon library
15
16
17 111 includes objective TE classes, such as L1, Alu, and SVA. In this study, the TE reference
18
19
20 112 database contains known TE sequences collected from RepBase version 17.07 [19], dbRIP
21
22
23 113 [20] and Hot L1s [4]. To reduce the long processing time due to large volumes of WGS data,
24
25 114 we implemented a parallel approach to process each bam files of samples simultaneously in
26
27
28 115 the discordant reads detection step.
29
30
31 116 *Reads clustering and detection of breakpoints.* First, the “discordant reads” were scanned and
32
33
34 117 clustered into blocks that supported potential RIPs based on the Maximal Valid Clusters
35
36
37 118 algorithm. Second, we extracted all reads located within the cluster regions and determined
38
39
40 119 the breakpoints. Although high-depth, data-enabled RIP detection with high sensitivity was
41
42
43 120 possible given that more soft-clipped reads neighboring target site duplication (TSD) could be
44
45
46 121 detected, alignments neighboring the TSDs had apparently lower depth compared with the
47
48
49 122 mean sequencing depth of the whole genome due to occasional sequencing and system
50
51
52 123 errors. This feature made breakpoint detection difficult and increased the false discovery rate
53
54
55 124 (FDR). Thus, we added the recalibration process of clipped points to determine breakpoints.
56
57
58 125 Each read located within the cluster regions flanking potential breakpoints was used to confirm
59
60
61 126 the precise location of the breakpoints. Small deletions were extracted to perform breakpoint

1 127 recalibration, and the mismatched bases were removed from the deletion sequences.

2
3 128 The clipped sequences were realigned to local regions on HG19 to determine the actual
4
5
6 129 breakpoints. Breakpoints were assigned as “clips” if greater than half of the new clipped
7
8
9 130 sequences were discordant with the reference sequence and the length of gap within the new
10
11 131 clipped sequence was less than 30%. The point would not be a candidate unless it was a “clip”
12
13
14 132 and the mismatch was less than 5 bp or contained poly-A/T.

15
16
17 133 Some terminals of reads containing mismatched bases may be the clipped parts because
18
19
20 134 these bases were treated as mismatches rather than clips. The breakpoints candidates were
21
22
23 135 re-estimated by SID if mismatches accounted for greater than half of the read terminals.

24
25 136 Notably, we implemented the “Asynchronism Scanning” algorithm. Using this algorithm,
26
27
28 137 once the program clustered one possible insertion region by scanning unique reads, the
29
30
31 138 process of breakpoint detection in this region was immediately performed, rendering it possible
32
33
34 139 to detect TE insertions in one chromosome in only a few minutes. The detailed algorithm for
35
36 140 RIP candidate determination is provided in Additional file 2: Text S2.

37 38 39 141 **Annotation of TE insertions**

40
41
42 142 *Orientation annotation for the TE insertions.* We annotated the orientation of TE insertions
43
44
45 143 based on the BLAST results [21]. First, we extracted the discordant repeat anchored mate
46
47
48 144 (RAM) reads and clipped reads that supported the TE insertion and made the reads’
49
50
51 145 orientations the same as HG19. Then, we realigned the supporting reads against the
52
53
54 146 consensus sequences of known active retrotransposons to identify the mapped orientation in
55
56
57 147 known active retrotransposons. The orientations of TE insertions were judged by the reads’
58
59 148 orientation (for details see Additional file 2: Text S3). The accuracy of orientation annotation

1 149 was assessed by comparing 396 matched insertions from dbRIP and 21 fully sequenced
2
3 150 insertions from PCR validation experiments (Additional file 1: Table S2). In total, 326 insertions
4
5
6 151 were verified, and the FDR of orientation annotation was 21.82%.

7
8
9 152 *Subfamily annotation for RIPs.* The subfamily annotation of RIPs was performed according to
10
11 153 known active retrotransposons. We first constructed a comprehensive retrotransposons
12
13
14 154 sequence library. Alu subfamily consensus sequences were acquired from RepBase 17.07
15
16
17 155 [19]. L1 subfamily consensus sequences were acquired from Eunjung Lee [10]. SVA and LTR
18
19
20 156 consensus sequences were acquired from Baillie [22]. Next, we performed multiple subfamily
21
22
23 157 sequence alignment for each type of retrotransposon and discovered the diagnostic nucleotide
24
25
26 158 for each subfamily (for details see Additional file 1: Table S3-5). Specially, we discovered the
27
28
29 159 diagnostic nucleotide of L1 from previous studies [23-26]. We then assembled the “discordant
30
31 160 reads” of each RIP into contigs using CAP3 [27] and realigned them against all of the
32
33
34 161 subfamily sequences using BLAST [28] (Additional file 2: Text S3-4).

35
36 162 *Length annotation for RIPs.* During mapping the contigs to subfamily sequences, we identified
37
38
39 163 the first mapped site of the 5' and 3' ends of the subfamily sequence and accordingly counted
40
41
42 164 the lengths from the initial site (L_{min} and L_{max}). The length of inserted retrotransposon (L_{retro})
43
44
45 165 was calculated as the difference between the maximum and the minimum length of the aligned
46
47
48 166 sequence, as follows:

$$L_{retro} = L_{max} - L_{min} + 1.$$

167 **Simulation of RIP data**

168
169 In total, 761 TEs were randomly selected from our reference TE database (see Materials and
170 methods: Annotation of TE insertions) and inserted into HG19 autosomes randomly to

1 171 generate a new human genome (for details see Additional file 1: Table S6). The pIRS [29]
2
3 172 software was used to generate approximately 60x paired-end 100 bp reads; then, we mapped
4
5
6 173 these reads to the HG19 genome by BWA. Then, we used SID to detect these RIPs in the
7
8
9 174 simulated genome. By repeating this process, we obtained results from simulated data with
10
11
12 175 different depths to assess the sensitivity and specificity of RIP detection in sequence data with
13
14
15 176 distinct depth using SID.

17 177 **Reference RIP detection**

18
19
20 178 The reference RIPs were detected as a subset of deletions of the samples relative to the HG19
21
22
23 179 reference (Additional file 2: Figure S1). These deletions were selected from the results of
24
25
26 180 structural variation (SV) detection of YH90 samples, and the RIPs were annotated based on
27
28
29 181 matched deletion coordinates to HG19 annotation of RepeatMasker (greater than 90% of them
30
31
32 182 overlap with each other) [30].

33
34 183 The reference RIPs should be absent in the chimpanzee genome. The alignments of
35
36
37 184 chimpanzee mapped to the human genome were downloaded from UCSC
38
39
40 185 (<http://hgdownload.cse.ucsc.edu>). One reference RIP candidate should correspond to a
41
42
43 186 gap with an overlap of greater than 90% to each other, and no gaps were present in the
44
45
46 187 chimpanzee genome at this locus. The RIP candidates were filtered if no polymorphisms were
47
48
49 188 present in the YH90 samples (i.e., the allele frequency was equal to 180).

50 189 **Results**

52 190 **Establishment of SID**

53
54
55 191 To detect non-reference RIPs from WGS data accurately and in a time-efficient manner, we
56
57
58 192 developed SID, which can detect non-reference RIPs easily and quickly through discordant
59
60
61
62
63
64
65

1 193 reads detection and reads clustering. In the first step, three types of informative discordant
2
3 194 reads were selected for further analysis (Fig. 1a). Then, the reads that had mismatched bases
4
5
6 195 at the terminals (Fig. 1b, 1c) were used for judging heterozygosity. The clipped reads were
7
8
9 196 used to confirm the sequence of TSD and the precise insertion site of certain TEs.

10 11 197 **Non-reference retrotransposon insertion calling**

12
13
14 198 To investigate the influence of sequencing depth on RIP detection sensitivity and accuracy, we
15
16
17 199 simulated sequence data at different depths. Detection sensitivity dramatically increased with
18
19
20 200 increasing sequencing depth and achieved 95% (730/761) when the sequencing depth was
21
22
23 201 greater than 30x. By contrast, detection accuracy slightly changed with increasing sequencing
24
25
26 202 depth (Fig. 2a).

27
28 203 We next estimated the RIP detection sensitivity using two real sequencing datasets. One
29
30
31 204 dataset was the CEU trio data, which was deep-sequenced (> 75x) Illumina HiSeq data
32
33
34 205 generated by the Broad Institute (father NA12891, mother NA12892 and the female offspring
35
36
37 206 NA12878) from the 1000GP. We first used SID to detect the RIPs of each individual in the CEU
38
39
40 207 dataset and evaluated the sensitivity by comparing the detection results with the
41
42
43 208 PCR-validated datasets from Stewart et al. [12]. For Alu, the mean sensitivity reached 96.3%
44
45
46 209 among individuals. We also obtained a mean sensitivity of 80.3% and 83.3% for L1 and SVA,
47
48
49 210 respectively (Additional file 1: Table S7).

50
51 211 The other dataset, including NA18571, NA18572 and NA18537, was also recruited in
52
53
54 212 1000GP. The RIP datasets of these three individuals detected by SID were larger and covered
55
56
57 213 70.08% of the same sample's results in 1000GP on average (Additional file 2: Figure S2). We
58
59
60 214 estimated RIP detection accuracy using the sequencing data from a lymphocytic cell line

1 215 (YH_CL, ~52x) obtained from an Asian individual. These data represent the first Asian diploid
2
3 216 genome dataset, and we performed PCR validation. We randomly selected 103 detected RIPs,
4
5
6 217 and 93/96 (7 loci were removed because of the poor primer specificity) loci were successfully
7
8
9 218 validated, indicating that SID had an accuracy of 90.29% - 96.88% (Additional file 1: Table S8
10
11 219 and Additional file 2: Figure S3 and Text S5). We also used the PCR validation result to access
12
13
14 220 the accuracy of genotyping, which was approximately 93.55% (87/93, Fig. 2b, Additional file 2:
15
16
17 221 Text S6).

18
19
20 222 We next compared the RIP detection efficiency of different methods (SID, RetroSeq [11]
21
22 223 and TEA [31]) using YH_CL and three samples (NA18571, NA18572 and NA18537) from
23
24
25 224 YH90 (Additional file 2: Text S7). The run time of SID was approximately 3-fold reduced
26
27
28 225 compared with the other two methods, suggesting that SID was the most time-saving method
29
30
31 226 among the three methods (Additional file 2: Table S9). SID and TEA had comparable
32
33
34 227 sensitivities that were increased compared with RetroSeq (Additional file 2: Figure S4). We
35
36
37 228 also validated the uniquely detected RIPs by PCR (Additional file 1: Table S10) with an
38
39 229 accuracy of 75.86% (22/29) and 77.78% (7/9) for Alu and L1, respectively, revealing a higher
40
41
42 230 RIP detection accuracy (Alu: 42.10% (8/19) and 82.61% (19/23) and L1: 66.67% (2/3) and
43
44
45 231 66.67% (2/3) for RetroSeq and TEA, respectively).

46 47 232 **A comprehensive RIP landscape of the Han Chinese population**

48
49
50 233 We then performed RIP detection on a much larger scale. We sequenced 90 Han Chinese
51
52
53 234 individuals and generated Illumina paired-end sequence data at an average depth of 68x for
54
55
56 235 each sample (Additional file 1: Table S1). Using SID, the high depth of the dataset (much more
57
58
59 236 than 30x) allowed us to build a comprehensive non-reference RIP landscape with high

1 237 confidence.

2
3 238 In total, we identified 9342 non-reference RIPs in autosome regions, including 6483 Alu
4
5
6 239 elements, 2398 L1s, 61 LTRs and 400 SVAs (Fig. 3a; for details, see Additional file 1: Table
7
8
9 240 S11 and Additional file 2: Text S8). Of this dataset, 8433 RIPs, including 5826 Alu elements,
10
11 241 2169 L1s, 383 SVAs, and 55 LTRs, were novel compared with dbRIP (Fig. 3b). The average
12
13
14 242 number of non-reference RIPs per individual was 1394 (ranging from 1304 to 1493, Fig. 3c),
15
16
17 243 including 1110.80 Alu elements, 231.34 L1s, 43.14 SVAs and 9.01 LTRs, and each type of RIP
18
19
20 244 had a similar proportion ($P = 0.6364$, $P = 0.2711$, $P = 0.2128$, $P = 0.5582$, respectively,
21
22
23 245 Wilcoxon signed-rank test). We compared pair-wise individuals of all 90 samples, and the
24
25
26 246 average specific loci number was 672.79, which is approximately half (48.25%) of
27
28
29 247 non-reference RIPs of one individual.

30
31 248 We next compared our results with the 1000GP SV dataset. In total, 34.94% (3264/9342)
32
33
34 249 of RIPs in YH90 were also found in the 1000GP dataset. The Pearson correlation coefficient
35
36
37 250 was 0.7998 ($P < 2.2 \times 10^{-16}$) between YH90 and all the 26 populations in 1000GP SV dataset.
38
39
40 251 The Pearson correlation coefficient was 0.8856 between YH90 and the East Asian (EAS)
41
42
43 252 population in 1000GP, which was higher than other populations ($r = 0.7662$, $r = 0.5741$, $r =$
44
45
46 253 0.7025 and $r = 0.7627$ for American (AMR), African (AFR), European (EUR) and South Asian
47
48
49 254 (SAS) populations, respectively. Additional file 2: Text S9)[14].

50
51 255 Specific insert location information enabled us to investigate genome-wide sequence
52
53
54 256 patterns of these non-reference RIPs. We observed that the non-reference RIPs varied among
55
56
57 257 chromosomes (Fig. 3d, e). Notably, we found that the two different subpopulations (from
58
59
60 258 southern and northern China) had similar patterns of RIP distribution ($r = 0.782$, Fig. 3e and for

1 259 details see Additional file 2: Figure S5). However, the distribution of non-reference RIPs was
2
3 260 not obviously correlated with GC content, fixed RIPs, or SNPs of the same sample within 10M
4
5
6 261 non-N bins (Additional file 2: Figure S6).
7

8
9 262 To further investigate the distribution of non-reference RIPs in the functional region, we
10
11 263 annotated all the inserted loci (Fig. 3f). Greater than half of RIPs (4828/9342) were located in
12
13
14 264 gene regions, and the majority of these were located in introns. Only 5/9342 RIPs were located
15
16
17 265 in protein-coding regions, including three genes, C1orf66 (Alu-inserted), SNX31 (Alu-inserted)
18
19
20 266 and APH1B (SVA-inserted), with low frequency (1/90) and two genes, ADORA3 (Alu-inserted)
21
22
23 267 and Slco1b3 (L1-inserted), with higher frequency (44/90 and 12/90, respectively). In addition to
24
25
26 268 gene regions, we also found that on average 9.78% and 4.93% RIPs were located in enhancer
27
28
29 269 regions and promoter regions per sample, respectively (Fig. 3f).
30

31 270 Furthermore, we annotated the subfamily, orientation and sequence length of all detected
32
33
34 271 inserted retrotransposons based on regional sequence assembly and remapping to the
35
36
37 272 retrotransposon library. The AluY sub-family constituted essentially all non-reference Alu
38
39
40 273 insertions, in which AluYa5 and AluYb8 were mostly active (Additional file 1: Table S11),
41
42
43 274 supporting conclusions from previous studies [23, 32, 33].
44

45 275 The orientation of one RIP is determined from the mapping orientation of contigs to a
46
47
48 276 retrotransposon reference and the existence of poly-A or poly-T tails of the inserted sequence
49
50
51 277 (Additional file 1: Table S11). Previous studies have reported that the gene-inserted RIP had a
52
53
54 278 greater influence on gene expression if it was inserted on the same orientation as the target
55
56
57 279 gene [2, 34]. However, we detected a comparable number of direct and reverse events (0.475
58
59
60 280 and 0.525, respectively), arguing against an obvious natural selection on the RIPs with
61

1 281 consistent orientation with the inserted gene.

2
3 282 Along with subfamily and orientation annotation, we also calculated the length of each
4
5
6 283 insertion sequence. We found that different types of TE insertions had different length
7
8
9 284 distributions (Additional file 2: Figure S7). Greater than half of Alu elements (~70%) were
10
11 285 full-length, whereas the length of the L1 was distributed more discretely. Most L1s (> 80%)
12
13
14 286 were fractured during the process of retrotransposon, which is consistent with a previous study
15
16
17 287 [13].
18
19

20 288 **RIPs of a healthy population**

21
22
23 289 The pure and comprehensive RIP dataset can be used as a baseline of healthy people for
24
25
26 290 other disease-related research, especially single-gene diseases. The candidate
27
28
29 291 disease-related retrotransposon insertions found in this dataset were filtered. We explicitly
30
31
32 292 measured the overlap between our dataset and the disease-related retrotransposon insertion
33
34
35 293 data in dbRIP (<http://dbrip.org>) [35]. None of the insertion sites existed in our dataset,
36
37
38 294 indicating the accuracy of the database. We also tested some cancer research data. We
39
40
41 295 tested the dataset of candidate cancer-related somatic retrotransposon insertions that was
42
43
44 296 strictly generated from data of The Cancer Genome Atlas (TCGA) Pan-Cancer Project for 11
45
46
47 297 tumor types. No overlapping RIPs were detected, whereas 43.36% germline retrotransposons
48
49
50
51 298 were detected. According to the comparison of colon cancer-specific data [9], we identified two
52
53
54 299 L1 insertions consistent with our dataset with frequency of 51/90 and 50/90. These two L1
55
56
57 300 insertions were germline retrotransposon insertions that were further validated by PCR
58
59
60 301 validation in Solyom's research. We also tested the candidate hepatocellular
61
62
63 302 carcinoma-specific insertions [8] and identified one L1 insertion that was also present in our
64
65

1 303 dataset with a frequency of 9/90. This site was finally validated as a germline insertion by PCR
2
3 304 in that research. In conclusion, our data provide a reference panel to exclude false positive
4
5
6 305 insertions related to cancer.
7

8 9 306 **Population evolution analysis**

10
11 307 To perform the population evolution analysis of RIPs, we first merged the non-reference RIP
12
13 308 dataset with the “reference” retrotransposon insertions that were polymorphic in YH90
14
15
16
17 309 samples (Additional file 2: Figure S1) to obtain all RIPs from our samples. The retrotransposon
18
19
20 310 insertions with a frequency equal to 1 were removed from our non-reference RIPs. The
21
22
23 311 “reference” RIPs were defined as the reference genome-specific retrotransposon insertions
24
25
26 312 compared with each individual of the YH90 group. These reference RIPs were selected from
27
28
29 313 the dataset of YH90 deletions, and only the RIPs absent in chimpanzee were retained.
30

31 314 AFS was not only influenced by the natural selection but also by demographic history. For
32
33
34 315 example, a low-frequency bias for the majority of mutations can also be obtained if the
35
36
37 316 population recently experienced a bottleneck [36].
38

39 317 To perform the neutral test more accurately, we took the demographic history into
40
41
42 318 consideration (Additional file 2: Text S10). We simulated the following two different
43
44
45 319 demographic scenarios: a two-epoch population with a recent contraction and a three-epoch
46
47
48 320 bottleneck-shaped history containing a reduction of effective population size in the past
49
50
51 321 followed by a recent phase of size recovery (Fig. 4a). We tested the different assumptions with
52
53
54 322 the SNP dataset (Fig. 4b and Additional file 2: Table S12), which supported that the
55
56
57 323 three-epoch model was the best model.
58

59 324 Next, we explored the possibility of using RIP information to perform population evolution
60

1 325 analysis. Based on the genotyping result of the merged RIP dataset, we described the RIP
2
3 326 allele frequency spectrum (AFS) (Fig. 4c and Additional file 2: Text S11). The neutral model
4
5
6 327 expectation can be calculated using the formula θ/i , where θ is the insertion diversity
7
8
9 328 parameter and i (180) is the allele count in a fixed number of samples n (90) [36]. The
10
11 329 spectrum was skewed toward low-allele frequency compared with the distribution of the
12
13
14 330 expected neutral model, indicating possible negative selection pressure on retrotransposon
15
16
17 331 insertions.

18
19
20 332 To investigate the influence of the demographic history on RIP AFS, we performed
21
22 333 demographic correction and re-analyzed the RIP AFS under different selection models (Fig. 4d
23
24
25 334 and Additional file 2: Figure S8-9). The classification of neutral with negative and positive
26
27
28 335 selection indicates that a proportion of RIPs was neutral, and a proportion of RIPs was under
29
30
31 336 negative selection. In addition, other RIPs were under positive selection (m1), neutral with
32
33
34 337 negative selection (m2), neutral with positive selection (m3), negative selection (m4), positive
35
36 338 selection (m5), and neutral selection (m6). We further calculated the selection coefficient (S')
37
38
39 339 under each best-fit model with the determination of an approximately neutral selection effect
40
41
42 340 threshold ($S' < 0.01\%$) [37]. Models m1 and m2 were the most fitted models with the observed
43
44
45 341 RIP AFS (Additional file 2: Table S13). The best-fit result of model m1 demonstrated that
46
47
48 342 approximately 75% RIPs were under negative selection with $s = 0.0290\%$, which indicates that
49
50
51 343 these RIPs are weakly deleterious. In addition, 10% were under positive selection, whereas 15%
52
53 344 were neutral. Under model m2, the best-fit result demonstrated that 70% of RIPs were under
54
55
56 345 negative selection with $s = 0.0396\%$. In addition, 30% of RIPs were neutral. The selection
57
58
59 346 coefficient was 0.0079% under the all negative selection model, indicating an approximately

1 347 neutral selection effect.

2
3 348 The distribution of fitness effects of retrotransposon subfamilies (L1, SVA and Alu) was
4
5
6 349 also estimated under the same demographic model. Assuming that all RIPs of different
7
8
9 350 subfamilies were under negative selection (model m1), the selection coefficient models were
10
11
12 351 various among three subfamilies of RIPs ($S' = -0.0143\%$, $S' = -0.0172\%$, $S' = -0.0068\%$ for L1,
13
14 352 SVA and Alu, respectively), suggesting that there is more natural selection pressure on L1 and
15
16
17 353 SVA (weakly negative selection) compared with Alu (nearly neutral selection).

20 354 **Phylogenetic analysis**

22 355 To investigate whether RIP information can be used to separate the Northern and Southern
23
24
25 356 Chinese groups, we performed principal component analysis (PCA) using the RIPs detected
26
27
28 357 from the YH90 dataset, which provided well-resolved Northern and Southern Chinese groups
29
30
31 358 (Fig. 5a and Additional file 2: Text S12). Compared with the PCA result derived from the SNPs
32
33
34 359 detected from the same dataset (Fig. 5b), there seemed to be more overlapping observations,
35
36 360 indicating SNPs might be more informative in resolving the two distinctive populations. Next,
37
38
39 361 we determined whether it is possible to perform phylogenetic analysis using RIP information
40
41
42 362 detected from the YH90 dataset. Two phylogenetic trees were constructed using RIPs and
43
44
45 363 SNPs, separately (Fig. 5c and 5d; for details, see Additional file 2: Text S13). Similar to the
46
47
48 364 PCA result, increased mixing between Northern and Southern Chinese individuals was
49
50
51 365 observed for the phylogenetic tree derived from the RIP information. Interestingly, HG00534,
52
53 366 an isolated Southern Chinese individual located in a northern cluster in the phylogenetic tree
54
55
56 367 established using the SNP information, clustered largely with Southern Chinese individuals in
57
58
59 368 the phylogenetic tree derived from the RIP information. Future studies are warranted to

1 369 explore whether combining SNPs with RIP results in the construction of a more accurate
2
3 370 phylogenetic tree.
4
5
6

7 371 **Conclusions**

8
9
10 372 In this paper, we developed the computer program SID to detect the non-reference RIPs of 90
11
12 373 healthy Han Chinese individuals using high-depth WGS. We described the landscape of RIP
13
14 374 distribution on population genomes and annotated the subfamily, orientation, and length of
15
16 375 RIPs. We demonstrated that the RIPs could be used as a normal baseline for
17
18 376 retrotransposon-related disease research.
19
20
21

22
23
24 377 To our knowledge, this dataset is the largest Han Chinese dataset to date. Compared with
25
26 378 1000GP results from the same samples, approximately half (mean 48.05%; Additional file 2:
27
28 379 Figure S2) of RIPs in our dataset were previously observed, suggesting that our
29
30 380 deep-sequenced data exhibited increased detection sensitivity compared with low coverage
31
32 381 data. For example, serum ACE levels were determined by the Alu insertion/deletion (I/D)
33
34 382 polymorphism in the following order: DD > ID > II [38]. The D allele of the ACE gene was
35
36 383 associated with essential hypertension in different populations [39-42]. We found that the ACE
37
38 384 gene harbored an Alu insertion in the 15th intron with a frequency of 81/90 in our 90 Chinese
39
40 385 genomes compared with a considerably reduced frequency (7/63) in CEPH individuals [12],
41
42 386 which was supported by a previous study [43]. To our surprise, no RIP ACEs were present in
43
44 387 Han Chinese samples from the 1000GP dataset, which is a high-frequency inserted gene in
45
46 388 our RIP data. ACE-specific PCR validation (Additional file 2: Figure S10) and a previous ACE
47
48 389 study [44] indicated that our results were consistent with the real values. This finding suggests
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1 390 that adequate sequencing depth is important to investigate RIP frequency and that our data
2
3 391 present a result that is consistent with the actual situation. The highly sensitive and accurate
4
5
6 392 RIP dataset provided a perfect opportunity to perform RIP fitness analysis. This study is the
7
8
9 393 first to evaluate the natural selection effect on retrotransposon insertions at the population
10
11
12 394 level. As a type of long fragment insertion, RIPs are under approximately neutral selection.
13
14
15 395 This finding is consistent with our result that retrotransposon insertions are mostly relatively
16
17
18 396 inconsequential because the harbored genes are always relatively unimportant. Regarding
19
20
21 397 different types of RIPs in addition to Alu, the longer insertion elements L1 and SVA exhibit
22
23
24 398 weakly positive selection pressure.

25 399 This dataset can be compared with others to provide guidance in research of the
26
27
28 400 disease-causing mechanisms in certain populations and to successfully determine the
29
30
31 401 insertion time of a specific locus. This dataset can also be used as a standard for other RIP
32
33
34 402 research and can serve as a baseline to filter irrelevant RIPs in disease-causing
35
36
37 403 retrotransposon research. Genome-wide association studies (GWAS) have proven their utility
38
39
40 404 in identifying genomic variants associated with the risk for numerous diseases. Unlike SNPs
41
42
43 405 and copy number variations (CNVs) that are widely used in GWAS, RIPs have generally been
44
45
46 406 overlooked as a major contributor to human variation. Significantly, this dataset provides a
47
48
49 407 valuable resource to perform GWAS and identify more markers related to complex diseases.

50 408 The high cost of WGS at high depth is still a major limitation, preventing it from being
51
52
53 409 widely used in TE research. Furthermore, the large amount of data yielded by high-depth WGS
54
55
56 410 makes it difficult to undertake bioinformatic analysis. With the development of biotechnology
57
58
59 411 and IT, this situation should improve soon.

1 412 The next step is to research RIPs at the transcriptome level. The impact of RIPs on gene
2
3 413 expression remains unclear. Combining the genome and transcriptome would provide a
4
5
6 414 comprehensive picture about the regulation of RIPs. Thus, we can further expound the
7
8
9 415 position of the retrotransposon in the course of human evolution.

10 11 416 **Availability and requirements**

- 12
13
14 417 ● Project name: Specific Insertions Detector (SID)
- 15
16
17 418 ● Project home page: <https://github.com/Jonathanyu2014/SID>
- 18
19
20 419 ● Operating system(s): Linux
- 21
22
23 420 ● Programming language: Perl
- 24
25
26 421 ● Other requirements: Perl 5.14 or later, BLAST v2.2.25 or later, Samtools v1.0 or later
- 27
28
29 422 ● License: Apache License 2.0
- 30
31 423 ● Any restrictions to use by non-academics: None

32 33 34 35 424 **Additional files**

- 36
37
38 425 Additional file 1: Supplementary tables. Data description and the results of RIPs calling. (XLSX
39
40
41 426 1991 kb)
- 42
43
44 427 Additional file 2: Supplementary texts, figures and tables. (PDF 1010 kb)

45 46 47 428 **Abbreviations**

- 48
49
50
51 429 RIP, retrotransposon insertion polymorphism; TE, transposable element; LTR, long terminal
52
53
54 430 repeat; L1, long interspersed nuclear element 1; WGS, whole-genome sequencing; NGS,
55
56
57 431 next-generation sequencing; SID, specific insertions detector; TSD, target site duplication;
58
59
60 432 CNV, copy number variation; SNP, single nucleotide polymorphism; ENA, European

1 433 Nucleotide Archive; GWAS, genome-wide association study.
2
3

4 434 **Acknowledgments**
5
6
7

8 435 We are grateful for Zengli Yan, Nan Li, Na Li and Runze Jiang for optimizing and testing the
9
10 436 SID program. We thank Haoxiang Lin and Wenjuan Zhu for providing technical assistance to
11
12 437 us. We thank Liang Wu and Xulian Shi for polishing the manuscript. We acknowledge the
13
14 438 support by the 1000 Genomes Project Consortium. This work was supported by the Shenzhen
15
16 439 Municipal Government of China [JSGG20140702161347218] and
17
18 440 [KQCX20150330171652450].
19
20
21
22
23
24

25 441 **Availability of data and materials**
26
27

28 442 The source code of SID is available from the GitHub repository[45]. The human (Homo
29
30 443 sapiens) reference genome sequence (HG19) and its annotation files were downloaded from
31
32 444 UCSC Genome Bioinformatics (<http://genome.ucsc.edu/>). The raw sequence data of YH_CL is
33
34 445 available from the ENA repository (accession number ERA000005) [46]. All the YH90 raw
35
36 446 sequences have been released to the ENA repository (accession number ERA496654).
37
38
39
40
41
42

43 447 **Authors' contributions**
44
45
46

47 448 BL, SL and YH initiated this project and reviewed the manuscript. QY, XZ, YZ and XH drafted
48
49 449 the manuscript. XH and JL edited the manuscript. QY, WZ, XZ and YW performed the data
50
51 450 analysis and drew the pictures. YZ and YW designed and developed the SID program. NL, XZ
52
53 451 and GL conducted the experiment for sequencing. LX designed the primers and performed
54
55 452 PCR validation. YH, BL, SL, XZ, XG and XH provided fruitful discussions.
56
57
58
59
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

453 **Competing interests**

454 The authors declare that they have no competing interests.

455 **Author details**

456 ¹ BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,

457 China. ² BGI-Shenzhen, Shenzhen 518083, China. ³ BGI College, Shenzhen 518083, China.

458 ⁴ Department of Biology, University of Copenhagen, Copenhagen 1599, Denmark. ⁵ School of

459 Life Sciences, Sun Yat-sen University, Guangzhou 510006, China. ⁶ BGI-Forensics,

460 Shenzhen 518083, China.

461 **Ethics, consent and permissions**

462 This study was approved by BGI-IRB (NO. 16101).

463 **Consent to publish**

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

465

466 **References**

467 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle
468 M, FitzHugh W *et al*: **Initial sequencing and analysis of the human genome**. *Nature* 2001,
469 **409**(6822):860-921.

470 2. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution**.
471 *Nature reviews Genetics* 2009, **10**(10):691-703.

1 472 3. Kidd JM, Graves T, Newman TL, Fulton R, Hayden HS, Malig M, Kallicki J, Kaul R, Wilson
2
3 473 RK, Eichler EE: **A human genome structural variation sequencing resource reveals**
4
5
6 474 **insights into mutational mechanisms.** *Cell* 2010, **143**(5):837-847.
7
8
9 475 4. Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran JV, Kazazian HH, Jr.: **Hot**
10
11 476 **L1s account for the bulk of retrotransposition in the human population.** *Proceedings of the*
12
13
14 477 *National Academy of Sciences of the United States of America* 2003, **100**(9):5280-5285.
15
16
17 478 5. Xing J, Zhang Y, Han K, Salem AH, Sen SK, Huff CD, Zhou Q, Kirkness EF, Levy S, Batzer
18
19
20 479 MA *et al*: **Mobile elements create structural variation: analysis of a complete human**
21
22 480 **genome.** *Genome research* 2009, **19**(9):1516-1526.
23
24
25 481 6. Cordaux R, Hedges DJ, Herke SW, Batzer MA: **Estimating the retrotransposition rate of**
26
27
28 482 **human Alu elements.** *Gene* 2006, **373**:134-137.
29
30
31 483 7. Hancks DC, Kazazian HH, Jr.: **Active human retrotransposons: variation and disease.** *Curr*
32
33
34 484 *Opin Genet Dev* 2012, **22**(3):191-203.
35
36
37 485 8. Shukla R, Upton KR, Munoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, Brennan PM,
38
39 486 Baillie JK, Collino A, Ghisletti S *et al*: **Endogenous retrotransposition activates oncogenic**
40
41
42 487 **pathways in hepatocellular carcinoma.** *Cell* 2013, **153**(1):101-111.
43
44
45 488 9. Solyom S, Ewing AD, Rahrmann EP, Doucet T, Nelson HH, Burns MB, Harris RS, Sigmon DF,
46
47
48 489 Casella A, Erlanger B *et al*: **Extensive somatic L1 retrotransposition in colorectal tumors.**
49
50 490 *Genome research* 2012, **22**(12):2328-2338.
51
52
53 491 10. Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ 3rd, Lohr JG, Harris CC, Ding L,
54
55
56 492 Wilson RK *et al*: **Landscape of somatic retrotransposition in human cancers.** *Science* 2012,
57
58
59 493 **337**(6097):967-971.

1 494 11. Keane TM, Wong K, Adams DJ: **RetroSeq: transposable element discovery from**
2
3 495 **next-generation sequencing data.** *Bioinformatics* 2013, **29**(3):389-390.
4
5
6 496 12. Stewart C, Kural D, Stromberg MP, Walker JA, Konkel MK, Stutz AM, Urban AE, Grubert F,
7
8
9 497 Lam HY, Lee WP *et al*: **A comprehensive map of mobile element insertion polymorphisms**
10
11 498 **in humans.** *PLoS Genet* 2011, **7**(8):e1002236.
12
13
14 499 13. Ewing AD, Kazazian HH, Jr.: **Whole-genome resequencing allows detection of many rare**
15
16 500 **LINE-1 insertion alleles in humans.** *Genome research* 2011, **21**(6):985-990.
17
18
19 501 14. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye K,
20
21
22 502 Jun G, Hsi-Yang Fritz M *et al*: **An integrated map of structural variation in 2,504 human**
23
24 503 **genomes.** *Nature* 2015, **526**(7571):75-81.
25
26
27 504 15. Xing J, Witherspoon DJ, Jorde LB: **Mobile element biology: new possibilities with**
28
29 505 **high-throughput sequencing.** *Trends in genetics : TIG* 2013, **29**(5):280-289.
30
31
32 506 16. Wang J, Wang W, Li R, Li Y, Tian G, Goodman L, Fan W, Zhang J, Li J, Zhang J *et al*: **The**
33
34 507 **diploid genome sequence of an Asian individual.** *Nature* 2008, **456**(7218):60-65.
35
36
37 508 17. Li H, Durbin R: **Fast and accurate short read alignment with Burrows-Wheeler transform.**
38
39 509 *Bioinformatics* 2009, **25**(14):1754-1760.
40
41
42 510 18. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
43
44 511 Altshuler D, Gabriel S, Daly M *et al*: **The Genome Analysis Toolkit: a MapReduce**
45
46 512 **framework for analyzing next-generation DNA sequencing data.** *Genome research* 2010,
47
48 513 **20**(9):1297-1303.
49
50
51 514 19. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J: **Repbase Update,**
52
53 515 **a database of eukaryotic repetitive elements.** *Cytogenet Genome Res* 2005,

1 516 110(1-4):462-467.

2

3 517 20. Wang J, Song L, Grover D, Azrak S, Batzer MA, PL: **dbRIP: a highly integrated database of**

4

5

6 518 **retrotransposon insertion polymorphisms in humans.** *Hum Mutat* 2006, **27**(4):323-329.

7

8

9 519 21. Mount DW: **Using the Basic Local Alignment Search Tool (BLAST).** *CSH Protoc* 2007,

10

11 520 **2007:pdb top17.**

12

13

14 521 22. Baillie JK, Barnett MW, Upton KR, Gerhardt DJ, Richmond TA, De Sapio F, Brennan PM,

15

16

17 522 Rizzu P, Smith S, Fell M *et al*: **Somatic retrotransposition alters the genetic landscape of the**

18

19

20 523 **human brain.** *Nature* 2011, **479**(7374):534-537.

21

22

23 524 23. Boissinot S, Chevret P, AV F: **L1 (LINE-1) retrotransposon evolution and amplification in**

24

25 525 **recent human history.** *Mol Biol Evol* 2000, **17**(6):915-928.

26

27

28 526 24. Dombroski BA, Mathias SL, Nanthakumar E, Scott AF, Jr KH: **Isolation of an active human**

29

30 527 **transposable element.** *Science* 1991, **254**(5039):1805-1808.

31

32

33

34 528 25. Ovchinnikov I, Rubin A, GD S: **Tracing the LINEs of human evolution.** *Proceedings of the*

35

36 529 *National Academy of Sciences of the United States of America* 2002, **99**(16):10522-10527.

37

38

39 530 26. Ovchinnikov I, Troxel AB, GD S: **Genomic characterization of recent human LINE-1**

40

41 531 **insertions: evidence supporting random insertion.** *Genome research* 2001,

42

43 532 **11**(12):2050-2058.

44

45

46

47 533 27. Huang X, Madan A: **CAP3: A DNA sequence assembly program.** *Genome research* 1999,

48

49 534 **9**(9):868-877.

50

51

52

53 535 28. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.**

54

55 536 *Journal of molecular biology* 1990, **215**(3):403-410.

56

57

58 537 29. Hu X, Yuan J, Shi Y, Lu J, Liu B, Li Z, Chen Y, Mu D, Zhang H, Li N: **pIRS: Profile-based**

59

60

61

62

63

64

65

1 538 **Illumina pair-end reads simulator.** *Bioinformatics* 2012, **28**(11):1533-1535.

2

3 539 30. Tarailo-Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in**

4

5

6 540 **genomic sequences.** *Curr Protoc Bioinformatics* 2009, **Chapter 4**:Unit 4 10.

7

8

9 541 31. Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ, 3rd, Lohr JG, Harris CC, Ding L,

10

11 542 Wilson RK *et al*: **Landscape of somatic retrotransposition in human cancers.** *Science* 2012,

12

13

14 543 **337**(6097):967-971.

15

16

17 544 32. Hormozdiari F, Alkan C, Ventura M, Hajirasouliha I, Malig M, Hach F, Yorukoglu D, Dao P,

18

19

20 545 Bakhshi M, Sahinalp SC *et al*: **Alu repeat discovery and characterization within human**

21

22 546 **genomes.** *Genome research* 2011, **21**(6):840-849.

23

24

25 547 33. Batzer MA, Deininger PL: **Alu repeats and human genomic diversity.** *Nature reviews*

26

27

28 548 *Genetics* 2002, **3**(5):370-379.

29

30

31 549 34. Burns KH, Boeke JD: **Human transposon tectonics.** *Cell* 2012, **149**(4):740-752.

32

33

34 550 35. Wang J, Song L, Grover D, Azrak S, Batzer MA, Liang P: **dbRIP: a highly integrated**

35

36 551 **database of retrotransposon insertion polymorphisms in humans.** *Hum Mutat* 2006,

37

38

39 552 **27**(4):323-329.

40

41

42 553 36. Tajima F: **Statistical method for testing the neutral mutation hypothesis by DNA**

43

44 554 **polymorphism.** *Genetics* 1989, **123**(3):585-595.

45

46

47 555 37. Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, Adams

48

49

50 556 MD, Schmidt S, Sninsky JJ, Sunyaev SR *et al*: **Assessing the evolutionary impact of amino**

51

52 557 **acid mutations in the human genome.** *PLoS Genet* 2008, **4**(5):e1000083.

53

54

55 558 38. Rigat B, Hubert C, Alhenc-Gelas F, Cambien F, Corvol P, F S: **An insertion/deletion**

56

57 559 **polymorphism in the angiotensin I-converting enzyme gene accounting for half the**

1 560 **variance of serum enzyme levels.** *J Clin Invest* 1990, **86**(4):1343-1346.

2

3 561 39. Jeng JR, Harn HJ, Jeng CY, Yueh KC, SM S: **Angiotensin I converting enzyme gene**

4

5

6 562 **polymorphism in Chinese patients with hypertension.** *Am J Hypertens* 1997,

7

8

9 563 **10**(5Pt1):558-561.

10

11 564 40. Zee RY, Lou YK, Griffiths LR, BJ M: **Association of a polymorphism of the angiotensin**

12

13

14 565 **I-converting enzyme gene with essential hypertension.** *Biochem Biophys Res Commun* 1992,

15

16

17 566 **184**(1):9-15.

18

19

20 567 41. Asamoah A, Yanamandra K, Thurmon TF, Richter R, Green R, Lakin T, C M: **A deletion in the**

21

22 568 **angiotensin converting enzyme (ACE) gene is common among African Americans with**

23

24

25 569 **essential hypertension.** *Clin Chim Acta* 1996, **254**(1):41-46.

26

27

28 570 42. Duru K, Farrow S, Wang JM, Lockette W, T K: **Frequency of a deletion polymorphism in the**

29

30

31 571 **gene for angiotensin converting enzyme is increased in African-Americans with**

32

33

34 572 **hypertension.** *Am J Hypertens* 1994, **7**(8):759-762.

35

36 573 43. Anand SS, Yusuf S, Vuksan V, Devanesen S, Teo KK, Montague PA, Kelemen L, Yi C, Lonn E,

37

38

39 574 Gerstein H *et al*: **Differences in risk factors, atherosclerosis, and cardiovascular disease**

40

41

42 575 **between ethnic groups in Canada: the Study of Health Assessment and Risk in Ethnic**

43

44

45 576 **groups (SHARE).** *Lancet* 2000, **356**(9226):279-284.

46

47 577 44. Batzer MA, Stoneking M, Alegria-Hartman M, Bazan H, Kass DH, Shaikh TH, Novick GE,

48

49

50 578 Ioannou PA, Scheer WD, RJ H: **African origin of human-specific polymorphic Alu**

51

52

53 579 **insertions.** *Proceedings of the National Academy of Sciences of the United States of America*

54

55

56 580 1994, **91**(25):12288-12292.

57

58 581 45. Yu Q: **Specific Insertions Detector.** In: *Zenodo.* 2016.

59

60

61

62

63

64

65

1 582 46. Zong C, Lu S, Chapman AR, Xie XS: **Genome-wide detection of single-nucleotide and**
2
3 583 **copy-number variations of a single human cell.** *Science* 2012, **338**(6114):1622-1626.
4

5 584
6
7 585
8

9
10 586 **Figure legends**

11
12
13
14 587 **Fig. 1** The principle of retrotransposon insertion detection. **(a)** Schematic diagram of using SID
15
16 588 for RIP detection in the genome. TSD: target site duplication. SID: Specific Insertions Detector.
17
18
19 589 **(b)** An example of reads mapping for predicted homozygous insertions. **(c)** An example of
20
21
22 590 reads mapping for predicted heterozygous insertions. In **(b)** and **(c)**, the red bases indicate the
23
24
25 591 mismatches, and the sequences with an orange background represent the clipped part of the
26
27
28 592 reads. The clipped reads are derived from one allele with inserted retrotransposons, and the
29
30
31 593 normal reads are derived from the other allele with the same reference. The three reads with
32
33
34 594 asterisks indicate no clipped part but the presence of terminal mismatches, which can also
35
36 595 support the breakpoint and exhibit consistency with the clipped reads.

37
38
39
40 596 **Fig. 2** Assessing the SID results. **(a)** Detecting accuracy and sensitivity estimation along
41
42
43 597 cumulating sequencing depth of simulated data. **(b)** RIP genotyping of YH_CL. PCR validation
44
45 598 results are marked. HEE: estimated heterozygous site. HOE: estimated homozygous site.
46
47
48 599 HEV: validated heterozygous site. HOV: validated homozygous site. The dash line indicates
49
50
51 600 the estimated boundary between heterozygous and heterozygous sites. Note that some of the
52
53
54 601 validated RIPs are present in the same locus in the plot figure.

55
56
57 602 **Fig. 3** Comprehensive landscape of non-reference RIPs of YH90. **(a)** Proportions of novel
58
59

1 603 insertions identified for each type of retrotransposon. **(b)** Comparison of YH90 non-reference
2
3 604 RIP results with dbRIP. Adjacent 100-bp regions of RIPs were taken into consideration. **(c)** TE
4
5 605 distribution of each YH90 sample. **(d)** Box plots of non-reference RIP distribution among
6
7 606 autosomes. **(e)** TE frequency distribution among YH90 samples. Rings from outer to inner
8
9 607 indicate Alu insertions frequency, L1 insertion frequency, SVA insertion frequency, LTR
10
11 608 insertion frequency and cytoband structure. The inside frequency of the rings indicates the
12
13 609 insertion frequency for the Northern Chinese group, and the outside frequency represents that
14
15 610 of the Southern Chinese group. **(f)** RIP distribution in different functional regions of the
16
17 611 genome.

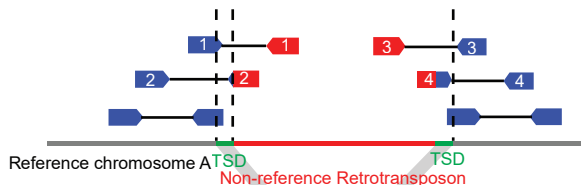
18
19
20
21
22
23
24
25
26 612 **Fig. 4** Population genetics analysis based on YH90. **(a)** A two-epoch population with a recent
27
28 613 contraction; a three-epoch bottleneck-shaped history, which contained a reduction of the
29
30 614 effective population size in the past followed by a recent phase of size recovery. Details of the
31
32 615 parameters for all models are provided in Additional file 2: Table S12. **(b)** The observed SNP
33
34 616 frequency spectra and expected neutral SNP frequency spectra under different demographic
35
36 617 models. **(c)** Observed and expected RIP site frequency spectra before demographic correction
37
38 618 of each subfamily. **(d)** Assessing the evolutionary impact of RIPs in the human genome. The
39
40 619 allele frequency distribution of RIPs was compared among observed, neutral models and
41
42 620 negative models after demographic correction.

43
44
45
46
47
48
49
50
51
52 621 **Fig. 5** Phylogenetic analysis using RIPs and SNPs. **(a)** The detected RIPs were used for PCA.
53
54 622 Each dot represents a sample from YH90 and is plotted as scatterplot using PC1 and PC2.
55
56 623 Red indicates samples from individuals from northern China, and blue indicates individuals
57
58
59
60
61
62
63
64
65

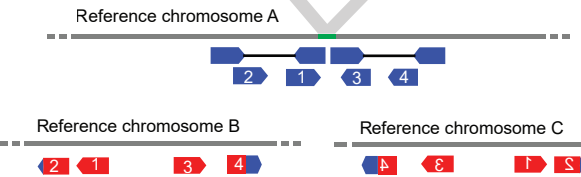
1 624 from southern China. **(b)** The detected SNPs were used for PCA. The plot layout and legend
2
3 625 are the same as those presented in **(a)**. **(c)** Phylogenetic tree constructed using the detected
4
5
6 626 RIPs. HG19 (green) is used as a control. Red indicates samples from individuals from northern
7
8
9 627 China, and blue indicates samples from individuals from southern China. **(d)** Phylogenetic tree
10
11 628 constructed using the detected SNPs. HG19 (green) is used as a control. Plot layout and
12
13
14 629 legend are same as that presented in **(c)**.
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

Figure 1

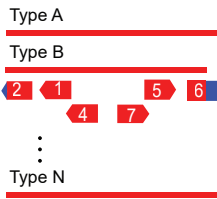
Sample genome



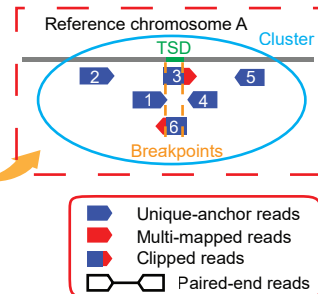
Hg19 genome



Retrotransposon library



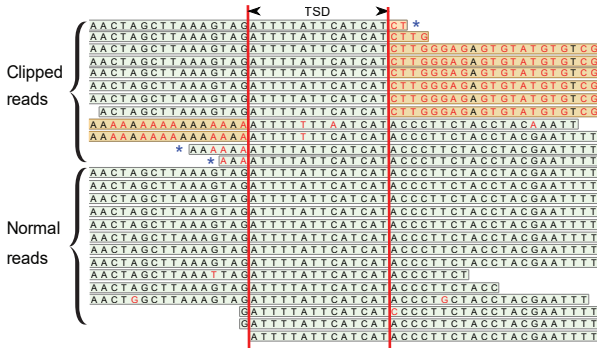
Detection result



B

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

Reference AACTAGCTTAAAGTAGATTTTATT CAT CATA CCCTTCTACCTACGAAATTT



C

Reference GGTTTACCAATTAGTCTCCCTTAAAGGGCACTGTTTATGATCATCACCACAAATGGATGCA

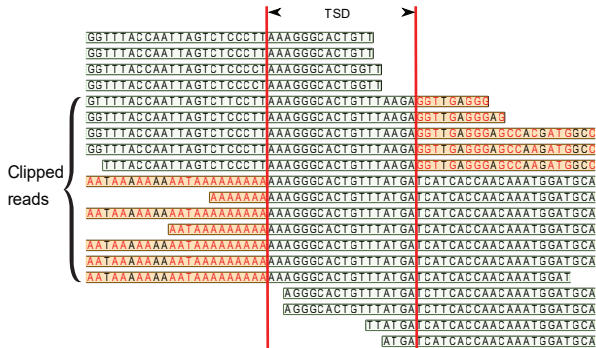
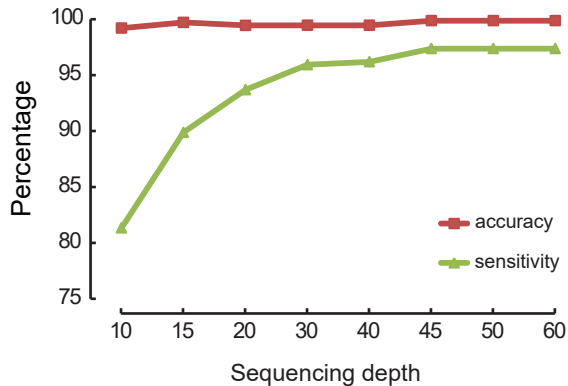


Figure 2

[Click here to download Figure Fig. 2.pdf](#)

A



B

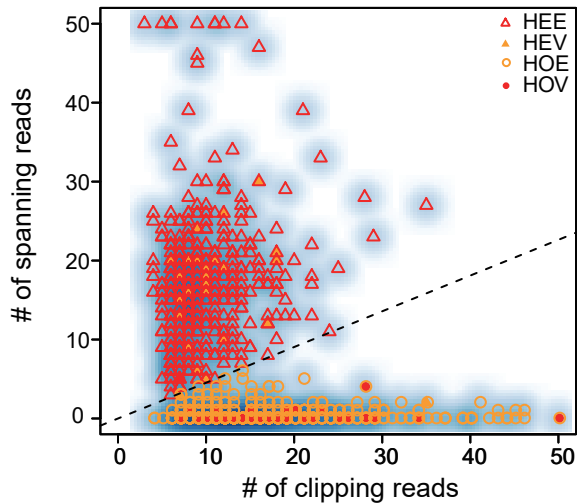
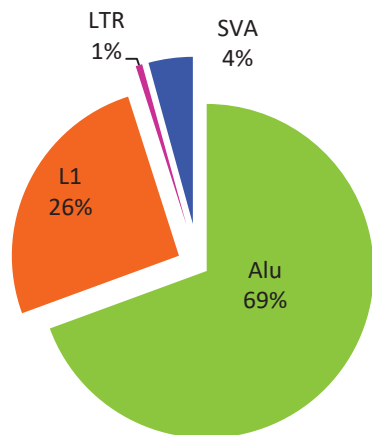


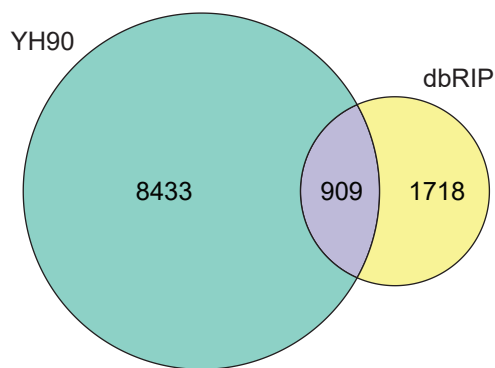
Figure 3

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

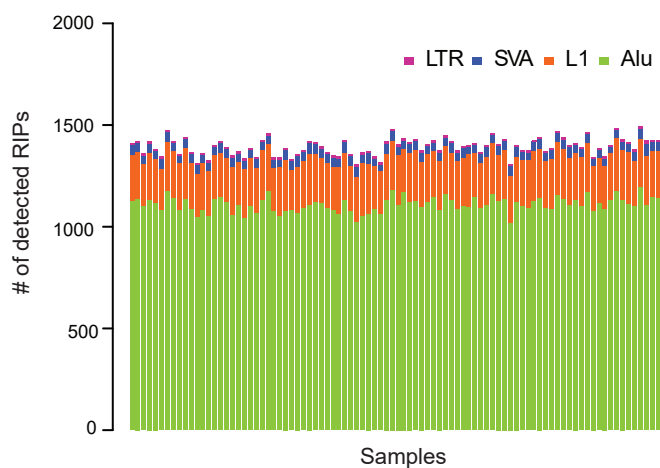
A



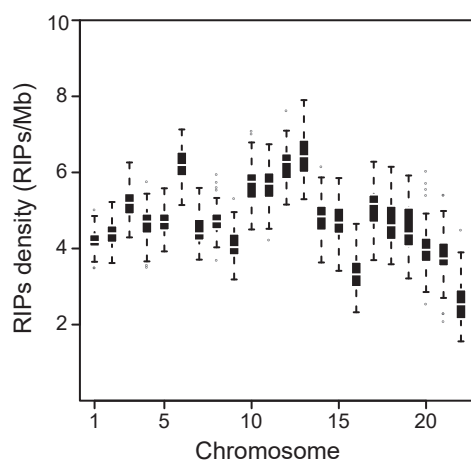
B



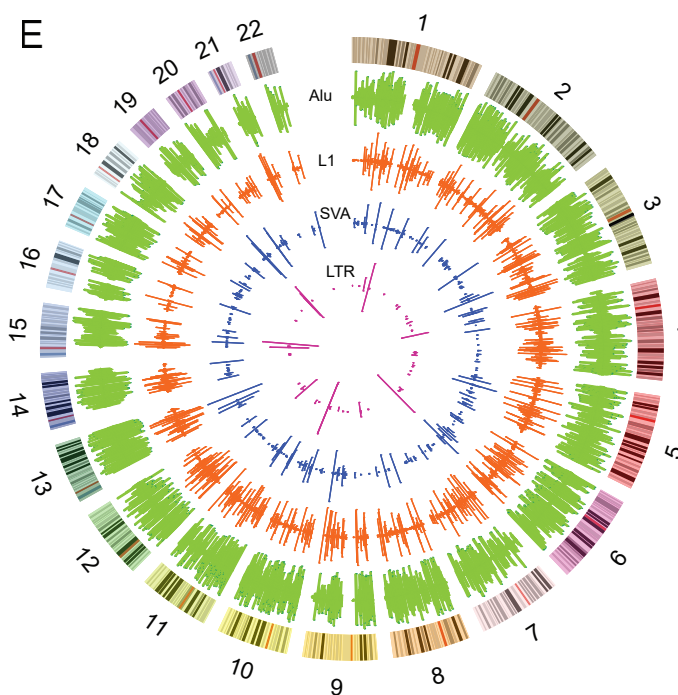
C



D



E



F

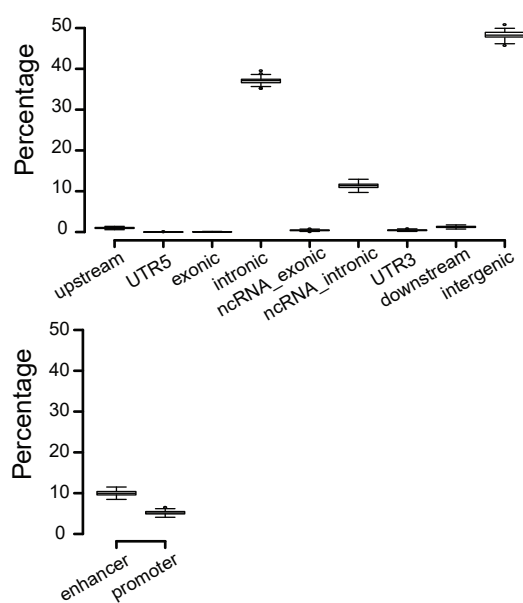
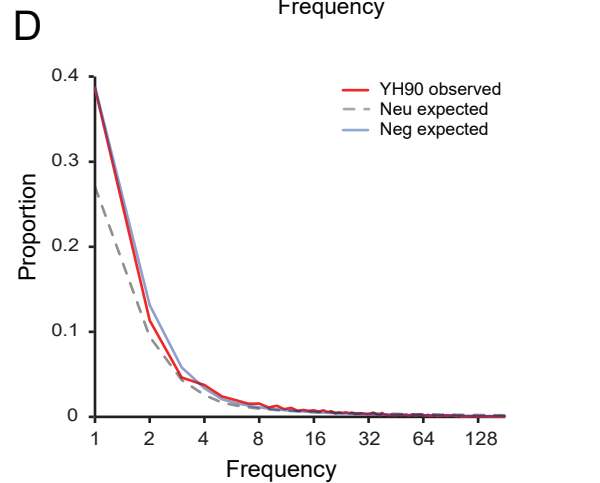
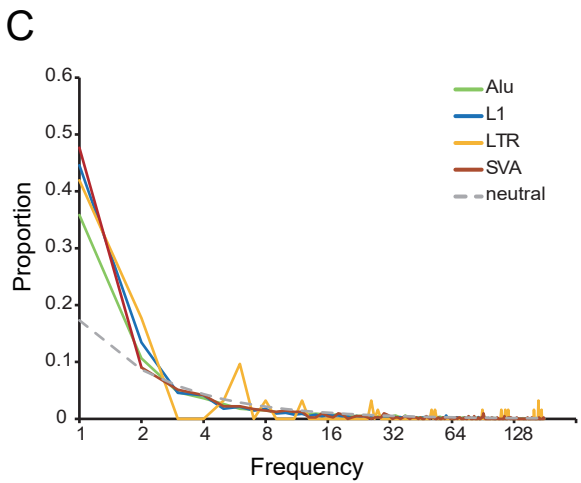
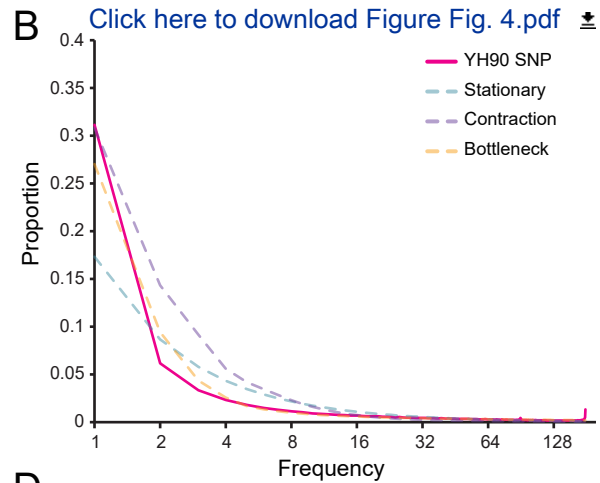
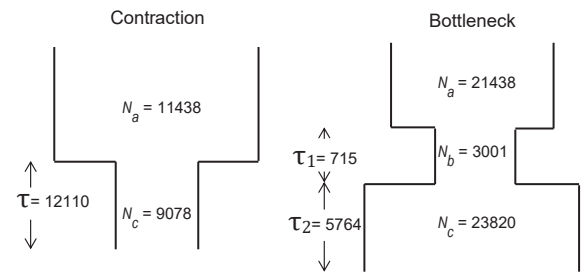
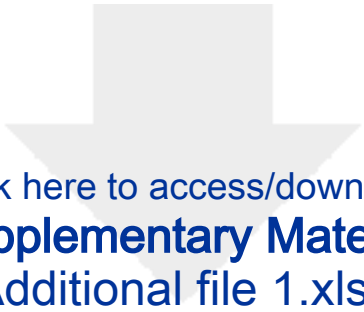




Figure 4





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





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

