# Population-wide Sampling of Retrotransposon Insertion Polymorphisms Using Deep Sequencing and Efficient Detection

Qichao Yu[1,2,†], Wei Zhang[1,2,†], Xiaolong Zhang[2], Yongli Zeng[2], Yeming Wang[2], Yanhui Wang[2],

Liqin Xu[2], Xiaoyun Huang[2], Nannan Li[2], Xinlan Zhou[2], Jie Lu[3], Xiaosen Guo[2], Guibo Li[2,4], Yong

Hou[2,4], Shiping Liu[2,5,*] and Bo Li[2,6,*]


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

China

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

[†] Equal contributors

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

**Emails of all authors:**

Qichao Yu: yuqichao@genomics.cn (ORCID: 0000-0003-2158-8424); Wei Zhang:

zhangwei7@genomics.cn (ORCID: 0000-0002-5792-7662); Xiaolong Zhang: 13528497060@163.com;

Yongli Zeng: zeoly100@163.com; Yeming Wang: 1398738509@qq.com (ORCID:

0000-0002-1521-2140); Yanhui Wang: 839584901@qq.com; Liqin Xu: xuliqin@genomics.cn; Nannan Li:

linannan@genomics.cn (ORCID: 0000-0003-3632-7964); Xinlan Zhou: zhouxinlan@genomics.cn

(ORCID: 0000-0001-9293-0894); Xiaoyun Huang: huangxiaoyun@genomics.cn (ORCID:

23 0000-0002-3389-9759); Jie Lu: lujie1@genomics.cn (ORCID: 0000-0001-7304-2023); Xiaosen Guo:

24 guoxs@genomics.cn (ORCID: 0000-0003-1317-2760); Guibo Li: liguibo@genomics.cn (ORCID:

25 0000-0002-6141-4931); Yong Hou: houyong@genomics.cn (ORCID: 0000-0002-0420-0726); Bo Li:

26 libo@genomics.cn; Shiping Liu: liushiping@genomics.cn (ORCID: 0000-0003-0019-619X).

27

28 **Abstract**

29 **Background:** Active retrotransposons play important roles during evolution and continue to

30 shape our genomes today, especially in genetic polymorphisms underlying a diverse set of

31 diseases. However, studies of human retrotransposon insertion polymorphisms (RIPs) based

32 on whole-genome deep sequencing at the population level have not been sufficiently

33 undertaken, despite the obvious need for a thorough characterization of RIPs in the general

34 population.

35 **Findings:** Herein, we present a novel and efficient computational tool named Specific

36 Insertions Detector (SID) for the detection of non-reference RIPs. We demonstrate that SID is

37 suitable for high depth whole-genome sequencing (WGS) data using paired-end reads

38 obtained from simulated and real datasets. We construct a comprehensive RIP database

39 using a large population of 90 Han Chinese individuals with a mean 68× depth per individual.

40 In total, we identify 9342 recent RIPs, and 8433 of these RIPs are novel compared with dbRIP,

41 including 5826 Alu, 2169 long interspersed nuclear element 1 (L1), 383 SVA, and 55 long

42 terminal repeats (LTR). Among the 9342 RIPs, 4828 were located in gene regions and five

43 were located in protein-coding regions. We demonstrate that RIPs can, in principle, be an

44 informative resource to perform population evolution and phylogenetic analyses. Taking the

45 demographic effects into account, we identify a weak negative selection on SVA and L1 but

46 approximately neutral selection for Alu elements based on the frequency spectrum of RIPs.

47 **Conclusions:** SID is a powerful open-source program for the detection of non-reference RIPs.

48 We built a non-reference RIP dataset that greatly enhanced the diversity of RIPs detected in

49 the general population and should be invaluable to researchers interested in many aspects of

50    human evolution, genetics, and disease. As a proof-of-concept, we demonstrate that the RIPs

51    can be used as biomarkers in a similar way as single nucleotide polymorphisms (SNPs).

52    **Keywords:** Transposable element, retrotransposon insertion polymorphism, next-generation

53    sequencing, whole-genome sequencing

54

55

56    **Findings**

57    **Introduction**

58    Transposable elements (TEs) are genomic sequences that can replicate within the genome

59    either autonomously or in conjunction with other TEs, resulting in insertion polymorphisms.

60    Over the evolutionary timescale, this process leads to drastic changes in genomic structure.

61    Current estimates suggest that approximately half of the human genome is derived from TEs

62    [1]. Retrotransposons, which constitute ~93% of TEs [2], can be subdivided into those

63    sequences containing LTRs and those that do not (non-LTR). The majority of human TEs

64    result from the activity of non-LTR retrotransposons, including the L1, Alu and SVA elements,

65    which collectively account for approximately one-third of the human genome [1]. Although

66    most retrotransposons are inactive remnants prevalent among the human population, younger

67    retrotransposons account for much of the structural variation among individual genomes [3].

68    Only a small proportion of total L1s are highly active [4]. The current rate of retrotransposition

69    in humans has been approximately estimated as 1 for every 20 births for Alu, 1 for every 200

70    births for L1 and 1 for every 900 births for SVA [5, 6].

71        Retrotransposon insertion is a disease-causing mechanism [7], and next-generation

72    sequencing (NGS) technology has been widely used to explore the association between

73    retrotransposon insertions and disease, such as cancer [8-10]. In this respect, a

74    comprehensive RIP dataset of a healthy population is necessary to serve as a reference for

75    the identification of disease-related RIPs. Based on the database of the 1000 Genomes

76    Project (1000GP), researchers performed RIP detection on an unprecedented scale and

77    detected thousands of novel RIPs [11-14]. This finding implies that an insertion allele present

78    in multiple individuals would effectively receive high coverage across the pooled dataset,

3

79 leading to a detection bias toward common insertions. It was previously estimated that at least

80 30× coverage of sequencing is needed to detect heterozygous RIPs with high sensitivity using

81 WGS [15].

82    Here, we developed the software SID to detect RIPs, which fulfilled our needs regarding

83 detection efficiency, accuracy and sensitivity. We also generated a non-reference TE insertion

84 polymorphism database by employing SID to analyze the whole-genome sequences of 90 Han

85 Chinese individuals (YH90) acquired at a mean depth of 68×.

**Materials and methods**

86

**Samples and whole genome sequencing**

87

88 We obtained B-lymphocyte cell lines from 90 Han Chinese individuals at the Coriell Institute

89 (Camden, New Jersey, USA). These individuals were selected from Beijing, Hunan province

90 and Fujian province, respectively. We broadly separated the samples into "Northern group" (45

91 samples) and "Southern group" (45 samples). DNA was extracted from the B-lymphocyte cells

92 of each individual, and libraries were then constructed following the manufacturer's

93 instructions. High-coverage paired-end 100 bp WGS libraries were sequenced on the Illumina

94 HiSeq 2000 Platform. For more on this dataset see the Data Note describing its production

95 published alongside this paper [16]. In addition, we also used a Chinese sample [17] for which

96 the data were previously released in the European Nucleotide Archive (ENA) repository

97 (Additional file 1: Table S1). The Institutional Review Board on Bioethics and Biosafety of BGI

98 (BGI-IRB) approved the study.

**Processing of the WGS data**

99

100 Reads were aligned to the human genome reference (HG19, Build37) using *BWA* (BWA ,

101 RRID:SCR_010910)[18]. Duplications were removed using Picard tools, and the quality values

102 of each reads were recalibrated using the Genome Analysis Toolkit (GATK)( GATK ,

103 RRID:SCR_001876)[19]. The resulting Binary Alignment/Map (BAM) files were used as input

104 for SID (Additional file 2: Text S1).

**The specific insertion detector pipeline**

105

4

106    SID is compiled in Perl and includes the following two steps: discordant reads detection and

107    reads clustering. Generally, the first step collects informative reads and generates other

108    necessary files, whereas the second step discovers the specific insertion sites and exports the

109    final results into plain text.

110    *Detection of discordant reads.* The "discordant reads" were extracted for the subsequent

111    clustering step. Paired-end reads were determined as "discordant reads" if they met one of the

112    following criteria: a. one read mapped to HG19 uniquely and the other read mapped to the

113    retrotransposon library (multi-mapped or unmapped to HG19); b. one read mapped to HG19

114    uniquely and the other soft-clipped read mapped to HG19, and the clipped sequence could be

115    mapped to the retrotransposon library; c. one soft-clipped read mapped to HG19, and the

116    clipped sequence could be mapped to the retrotransposon library. The other read mapped to

117    the retrotransposon library (multi-mapped or unmapped to HG19). The retrotransposon library

118    includes objective TE classes, such as L1, Alu, and SVA. In this study, the TE reference

119    database contains known TE sequences collected from RepBase version 17.07 [20], dbRIP

120    [21] and Hot L1s [4]. To reduce the long processing time due to large volumes of WGS data,

121    we implemented a parallel approach to process each bam files of samples simultaneously in

122    the discordant reads detection step.

123    *Reads clustering and detection of breakpoints.* First, the "discordant reads" were scanned and

124    clustered into blocks that supported potential RIPs based on the Maximal Valid Clusters

125    algorithm [22]. Second, we extracted all reads located within the cluster regions and

126    determined the breakpoints. Although high-depth, data-enabled RIP detection with high

127    sensitivity was possible given that more soft-clipped reads neighboring target site duplication

5

128 (TSD) could be detected, alignments neighboring the TSDs had apparently lower depth

129 compared with the mean sequencing depth of the whole genome due to occasional

130 sequencing and system errors. This feature made breakpoint detection difficult and increased

131 the false discovery rate (FDR). Thus, we added the recalibration process of clipped points to

132 determine breakpoints. Each read located within the cluster regions flanking potential

133 breakpoints was used to confirm the precise location of the breakpoints. Small deletions were

134 extracted to perform breakpoint recalibration, and the mismatched bases were removed from

135 the deletion sequences.

136     The clipped sequences were realigned to local regions on HG19 to determine the actual

137 breakpoints. Breakpoints were assigned as "clips" if greater than half of the new clipped

138 sequences were discordant with the reference sequence and the length of gap within the new

139 clipped sequence was less than 30%. The point would not be a candidate unless it was a "clip"

140 and the mismatch was less than 5 bp or contained poly-A/T.

141     Some terminals of reads containing mismatched bases may be the clipped parts because

142 these bases were treated as mismatches rather than clips. The breakpoints candidates were

143 re-estimated by SID if mismatches accounted for greater than half of the read terminals.

144     Notably, we implemented the "Asynchronous Scanning" algorithm (Additional file 2: Text

145 S2). Using this algorithm, once the program clustered one possible insertion region by

146 scanning unique reads, the process of breakpoint detection in this region was immediately

147 performed, rendering it possible to detect TE insertions in one chromosome in only a few

148 minutes. The detailed algorithm for RIP candidate determination is provided in Additional file 2:

149 Text S2.

6

150 **Annotation of TE insertions**

151 *Orientation annotation for the TE insertions*. We annotated the orientation of TE insertions

152 based on the BLAST results [23]. First, we extracted the discordant repeat anchored mate

153 (RAM) reads and clipped reads that supported the TE insertion and made the reads'

154 orientations the same as HG19. Then, we realigned the supporting reads against the

155 consensus sequences of known active retrotransposons to identify the mapped orientation in

156 known active retrotransposons. The orientations of TE insertions were judged by the reads'

157 orientation (for details see Additional file 2: Text S3). The accuracy of orientation annotation

158 was assessed by comparing 396 matched insertions from dbRIP and 21 fully sequenced

159 insertions from PCR validation experiments (Additional file 1: Table S2). In total, 326 insertions

160 were verified, and the FDR of orientation annotation was 21.82%.

161 *Subfamily annotation for RIPs.* The subfamily annotation of RIPs was performed according to

162 known active retrotransposons. We first constructed a comprehensive retrotransposons

163 sequence library. Alu subfamily consensus sequences were acquired from RepBase 17.07

164 [20]. L1 subfamily consensus sequences were acquired from Eunjung Lee [10]. SVA and LTR

165 consensus sequences were acquired from Baillie [24]. Next, we performed multiple subfamily

166 sequence alignment for each type of retrotransposon and discovered the diagnostic nucleotide

167 for each subfamily (for details see Additional file 1: Table S3-5). Specially, we discovered the

168 diagnostic nucleotide of L1 from previous studies [25-28]. We then assembled the "discordant

169 reads" of each RIP into contigs using CAP3 [29] and realigned them against all of the

170 subfamily sequences using BLAST (NCBI BLAST , RRID:SCR_004870)[30] (Additional file 2:

171 Text S3-4).

7

172 *Length annotation for RIPs.* During mapping the contigs to subfamily sequences, we identified

173 the first mapped site of the 5' and 3' ends of the subfamily sequence and accordingly counted

174 the lengths from the initial site ($L_{min}$ and $L_{max}$). The length of inserted retrotransposon ($L_{\text{retro}}$)

175 was calculated as the difference between the maximum and the minimum length of the aligned

176 sequence, as follows:

177
$$L_{retro} = L_{max} - L_{min} + 1.$$

178 **Simulation of RIP data**

179 In total, 761 TEs were randomly selected from our reference TE database (see Materials and

180 methods: Annotation of TE insertions) and inserted into HG19 autosomes randomly to

181 generate a new human genome (for details see Additional file 1: Table S6). The pIRS [31]

182 software was used to generate approximately 60× paired-end 100 bp reads; then, we mapped

183 these reads to the HG19 genome by BWA. Then, we used SID to detect these RIPs in the

184 simulated genome. By repeating this process, we obtained results from simulated data with

185 different depths to assess the sensitivity and specificity of RIP detection in sequence data with

186 distinct depth using SID.

187 **Reference RIP detection**

188 The reference RIPs were detected as a subset of deletions of the samples relative to the HG19

189 reference (Additional file 2: Figure S1). These deletions were selected from the results of

190 structural variation (SV) detection of YH90 samples, and the RIPs were annotated based on

191 matched deletion coordinates to HG19 annotation of RepeatMasker (greater than 90% of them

192 overlap with each other) [32].

193     The reference RIPs should be absent in the chimpanzee genome. The alignments of

8

194 chimpanzee mapped to the human genome were downloaded from UCSC

195 (http://hgdownload.cse.ucsc.edu). One reference RIP candidate should correspond to a gap

196 with an overlap of greater than 90% to each other, and no gaps were present in the

197 chimpanzee genome at this locus. The RIP candidates were filtered if no polymorphisms were

198 present in the YH90 samples (i.e., the allele frequency was equal to 180).

199 **Results**

200 **Establishment of SID**

201 To detect non-reference RIPs from WGS data accurately and in a time-efficient manner, we

202 developed SID, which can detect non-reference RIPs easily and quickly through discordant

203 reads detection and reads clustering. In the first step, three types of informative discordant

204 reads were selected for further analysis (Fig. 1a). Then, the reads that had mismatched bases

205 at the terminals (Fig. 1b, 1c) were used for judging heterozygosity. The clipped reads were

206 used to confirm the sequence of TSD and the precise insertion site of certain TEs.

207 **Non-reference retrotransposon insertion calling**

208 To investigate the influence of sequencing depth on RIP detection sensitivity and accuracy, we

209 simulated sequence data at different depths. Detection sensitivity dramatically increased with

210 increasing sequencing depth and achieved 95% (730/761) when the sequencing depth was

211 greater than 30×. By contrast, detection accuracy slightly changed with increasing sequencing

212 depth (Fig. 2a).

213 We next estimated the RIP detection sensitivity using two real sequencing datasets. One

214 dataset was the CEU trio data, which was deep-sequenced (> 75×) Illumina HiSeq data

215 generated by the Broad Institute (father NA12891, mother NA12892 and the female offspring

9

216 NA12878) from the 1000GP. We first used SID to detect the RIPs of each individual in the CEU

217 dataset and evaluated the sensitivity by comparing the detection results with the

218 PCR-validated datasets from Stewart et al. [12]. For Alu, the mean sensitivity reached 96.3%

219 among individuals. We also obtained a mean sensitivity of 80.3% and 83.3% for L1 and SVA,

220 respectively (Additional file 1: Table S7).

221     The other dataset, including NA18571, NA18572 and NA18537, was also recruited in

222 1000GP. The RIP datasets of these three individuals detected by SID were larger and covered

223 70.08% of the same sample's results in 1000GP on average (Additional file 2: Figure S2). We

224 estimated RIP detection accuracy using the sequencing data from a lymphocytic cell line

225 (YH_CL, ~52×) obtained from an Asian individual. These data represent the first Asian diploid

226 genome dataset, and we performed PCR validation. We randomly selected 103 detected RIPs,

227 and 93/96 (7 loci were removed because of the poor primer specificity) loci were successfully

228 validated, indicating that SID had an accuracy of 90.29% - 96.88% (Additional file 1: Table S8

229 and Additional file 2: Figure S3 and Text S5). We also used the PCR validation result to access

230 the accuracy of genotyping, which was approximately 93.55% (87/93, Fig. 2b, Additional file 2:

231 Text S6).

232     We next compared the RIP detection efficiency of different methods (SID, RetroSeq [11]

233 and TEA [10]) using YH_CL and three samples (NA18571, NA18572 and NA18537) from

234 YH90 (Additional file 2: Text S7). The run time of SID was approximately 3-fold reduced

235 compared with the other two methods, suggesting that SID was the most time-saving method

236 among the three methods (Additional file 2: Table S9). SID and TEA had comparable

237 sensitivities that were increased compared with RetroSeq (Additional file 2: Figure S4). We

10

238 also validated the uniquely detected RIPs by PCR (Additional file 1: Table S10) with an

239 accuracy of 75.86% (22/29) and 77.78% (7/9) for Alu and L1, respectively, revealing a higher

240 RIP detection accuracy (Alu: 42.10% (8/19) and 82.61% (19/23) and L1: 66.67% (2/3) and

241 66.67% (2/3) for RetroSeq and TEA, respectively).

**A comprehensive RIP landscape of the Han Chinese population**

243 We then performed RIP detection on a much larger scale. We sequenced 90 Han Chinese

244 individuals and generated Illumina paired-end sequence data at an average depth of 68× for

245 each sample (Additional file 1: Table S1). Using SID, the high depth of the dataset (much more

246 than 30×) allowed us to build a comprehensive non-reference RIP landscape with high

247 confidence[16].

248 In total, we identified 9342 non-reference RIPs in autosome regions, including 6483 Alu

249 elements, 2398 L1s, 61 LTRs and 400 SVAs (Fig. 3a; for details, see Additional file 1: Table

250 S11 and Additional file 2: Text S8). Of this dataset, 8433 RIPs, including 5826 Alu elements,

251 2169 L1s, 383 SVAs, and 55 LTRs, were novel compared with dbRIP (Fig. 3b). The average

252 number of non-reference RIPs per individual was 1394 (ranging from 1304 to 1493, Fig. 3c),

253 including 1110.80 Alu elements, 231.34 L1s, 43.14 SVAs and 9.01 LTRs, and each type of RIP

254 had a similar proportion ($P = 0.6364$, $P = 0.2711$, $P = 0.2128$, $P = 0.5582$, respectively,

255 Wilcoxon signed-rank test). We compared pair-wise individuals of all 90 samples, and the

256 average specific loci number was 672.79, which is approximately half (48.25%) of

257 non-reference RIPs of one individual.

258 We next compared our results with the 1000GP SV dataset. In total, 34.94% (3264/9342)

259 of RIPs in YH90 were also found in the 1000GP dataset. The Pearson correlation coefficient

11

260    was 0.7998 ($P < 2.2 \times 10^{-16}$) between YH90 and all the 26 populations in 1000GP SV dataset.

261    The Pearson correlation coefficient was 0.8856 between YH90 and the East Asian (EAS)

262    population in 1000GP, which was higher than other populations ($r = 0.7662$, $r = 0.5741$, $r =$

263    0.7025 and $r = 0.7627$ for American (AMR), African (AFR), European (EUR) and South Asian

264    (SAS) populations, respectively. Additional file 2: Text S9)[14].

265        Specific insert location information enabled us to investigate genome-wide sequence

266    patterns of these non-reference RIPs. We observed that the non-reference RIPs varied among

267    chromosomes (Fig. 3d, e). Notably, we found that the two different subpopulations (from

268    southern and northern China) had similar patterns of RIP distribution ($r = 0.782$, Fig. 3e and for

269    details see Additional file 2: Figure S5). However, the distribution of non-reference RIPs was

270    not obviously correlated with GC content, fixed RIPs, or SNPs of the same sample within 10M

271    non-N bins (Additional file 2: Figure S6).

272        To further investigate the distribution of non-reference RIPs in the functional region, we

273    annotated all the inserted loci (Fig. 3f). Greater than half of RIPs (4828/9342) were located in

274    gene regions, and the majority of these were located in introns. Only 5/9342 RIPs were located

275    in protein-coding regions, including three genes, C1orf66 (Alu-inserted), SNX31 (Alu-inserted)

276    and APH1B (SVA-inserted), with low frequency (1/90) and two genes, ADORA3 (Alu-inserted)

277    and Slco1b3 (L1-inserted), with higher frequency (44/90 and 12/90, respectively). In addition to

278    gene regions, we also found that on average 9.78% and 4.93% RIPs were located in enhancer

279    regions and promoter regions per sample, respectively (Fig. 3f).

280        Furthermore, we annotated the subfamily, orientation and sequence length of all detected

281    inserted retrotransposons based on regional sequence assembly and remapping to the

12

282 retrotransposon library. The AluY sub-family constituted essentially all non-reference Alu

283 insertions, in which AluYa5 and AluYb8 were mostly active (Additional file 1: Table S11),

284 supporting conclusions from previous studies [26, 33, 34].

285     The orientation of one RIP is determined from the mapping orientation of contigs to a

286 retrotransposon reference and the existence of poly-A or poly-T tails of the inserted sequence

287 (Additional file 1: Table S11). Previous studies have reported that the gene-inserted RIP had a

288 greater influence on gene expression if it was inserted on the same orientation as the target

289 gene [2, 35]. However, we detected a comparable number of direct and reverse events (0.475

290 and 0.525, respectively), arguing against an obvious natural selection on the RIPs with

291 consistent orientation with the inserted gene.

292     Along with subfamily and orientation annotation, we also calculated the length of each

293 insertion sequence. We found that different types of TE insertions had different length

294 distributions (Additional file 2: Figure S7). Greater than half of Alu elements (~70%) were

295 full-length, whereas the length of the L1 was distributed more discretely. Most L1s (> 80%)

296 were fractured during the process of retrotransposon, which is consistent with a previous study

297 [13].

298 **RIPs of a healthy population**

299 The pure and comprehensive RIP dataset can be used as a baseline of healthy people for

300 other disease-related research, especially single-gene diseases. The candidate

301 disease-related retrotransposon insertions found in this dataset were filtered. We explicitly

302 measured the overlap between our dataset and the disease-related retrotransposon insertion

303 data in dbRIP (http://dbrip.org) [36]. None of the insertion sites existed in our dataset,

304    indicating the accuracy of the database. We also tested some cancer research data. We

305    tested the dataset of candidate cancer-related somatic retrotransposon insertions that was

306    strictly generated from data of The Cancer Genome Atlas (TCGA) Pan-Cancer Project for 11

307    tumor types. No overlapping RIPs were detected, whereas 43.36% germline retrotransposons

308    were detected. According to the comparison of colon cancer-specific data [9], we identified two

309    L1 insertions consistent with our dataset with frequency of 51/90 and 50/90. These two L1

310    insertions were germline retrotransposon insertions that were further validated by PCR

311    validation in Solyom's research. We also tested the candidate hepatocellular

312    carcinoma-specific insertions [8] and identified one L1 insertion that was also present in our

313    dataset with a frequency of 9/90. This site was finally validated as a germline insertion by PCR

314    in that research. In conclusion, our data provide a reference panel to exclude false positive

315    insertions related to cancer.

316    **Population evolution analysis**

317    To perform the population evolution analysis of RIPs, we first merged the non-reference RIP

318    dataset with the "reference" retrotransposon insertions that were polymorphic in YH90

319    samples (Additional file 2: Figure S1) to obtain all RIPs from our samples. The retrotransposon

320    insertions with a frequency equal to 1 were removed from our non-reference RIPs. The

321    "reference" RIPs were defined as the reference genome-specific retrotransposon insertions

322    compared with each individual of the YH90 group. These reference RIPs were selected from

323    the dataset of YH90 deletions, and only the RIPs absent in chimpanzee were retained.

324        Allele frequency spectrum (AFS) was not only influenced by the natural selection but also

325    by demographic history. For example, a low-frequency bias for the majority of mutations can

14

326 also be obtained if the population recently experienced a bottleneck [37].

327    To perform the neutral test more accurately, we took the demographic history into

328 consideration (Additional file 2: Text S10). We simulated the following two different

329 demographic scenarios: a two-epoch population with a recent contraction and a three-epoch

330 bottleneck-shaped history containing a reduction of effective population size in the past

331 followed by a recent phase of size recovery (Fig. 4a). We tested the different assumptions with

332 the SNP dataset (Fig. 4b and Additional file 2: Table S12), which supported that the

333 three-epoch model was the best model.

334    Next, we explored the possibility of using RIP information to perform population evolution

335 analysis. Based on the genotyping result of the merged RIP dataset, we described the RIP

336 AFS (Fig. 4c and Additional file 2: Text S11). The neutral model expectation can be calculated

337 using the formula $\theta/i$, where $\theta$ is the insertion diversity parameter and $i$ (180) is the allele

338 count in a fixed number of samples $n$ (90) [37]. The spectrum was skewed toward low-allele

339 frequency compared with the distribution of the expected neutral model, indicating possible

340 negative selection pressure on retrotransposon insertions.

341    To investigate the influence of the demographic history on RIP AFS, we performed

342 demographic correction and re-analyzed the RIP AFS under different selection models (Fig. 4d

343 and Additional file 2: Figure S8-9). The classification of neutral with negative and positive

344 selection indicates that a proportion of RIPs was neutral, and a proportion of RIPs was under

345 negative selection. In addition, other RIPs were under positive selection (m1), neutral with

346 negative selection (m2), neutral with positive selection (m3), negative selection (m4), positive

347 selection (m5), and neutral selection (m6). We further calculated the selection coefficient ($S'$)

15

348 under each best-fit model with the determination of an approximately neutral selection effect

349 threshold ($S'$ < 0.01%) [38]. Models m1 and m2 were the most fitted models with the observed

350 RIP AFS (Additional file 2: Table S13). The best-fit result of model m1 demonstrated that

351 approximately 75% RIPs were under negative selection with s = 0.0290%, which indicates that

352 these RIPs are weakly deleterious. In addition, 10% were under positive selection, whereas 15%

353 were neutral. Under model m2, the best-fit result demonstrated that 70% of RIPs were under

354 negative selection with s = 0.0396%. In addition, 30% of RIPs were neutral. The selection

355 coefficient was 0.0079% under the all negative selection model, indicating an approximately

356 neutral selection effect.

357     The distribution of fitness effects of retrotransposon subfamilies (L1, SVA and Alu) was

358 also estimated under the same demographic model. Assuming that all RIPs of different

359 subfamilies were under negative selection (model m1), the selection coefficient models were

360 various among three subfamilies of RIPs ($S'$ = -0.0143%, $S'$ = -0.0172%, $S'$ = -0.0068% for L1,

361 SVA and Alu, respectively), suggesting that there is more natural selection pressure on L1 and

362 SVA (weakly negative selection) compared with Alu (nearly neutral selection).

363 **Phylogenetic analysis**

364 To investigate whether RIP information can be used to separate the Northern and Southern

365 Chinese groups, we performed principal component analysis (PCA) using the RIPs detected

366 from the YH90 dataset, which provided well-resolved Northern and Southern Chinese groups

367 (Fig. 5a and Additional file 2: Text S12). Compared with the PCA result derived from the SNPs

368 detected from the same dataset (Fig. 5b), there seemed to be more overlapping observations,

369 indicating SNPs might be more informative in resolving the two distinctive populations. Next,

16

370    we determined whether it is possible to perform phylogenetic analysis using RIP information

371    detected from the YH90 dataset. Two phylogenetic trees were constructed using RIPs and

372    SNPs, separately (Fig. 5c and 5d; for details, see Additional file 2: Text S13). Similar to the

373    PCA result, increased mixing between Northern and Southern Chinese individuals was

374    observed for the phylogenetic tree derived from the RIP information. Interestingly, HG00534,

375    an isolated Southern Chinese individual located in a northern cluster in the phylogenetic tree

376    established using the SNP information, clustered largely with Southern Chinese individuals in

377    the phylogenetic tree derived from the RIP information. Future studies are warranted to

378    explore whether combining SNPs with RIP results in the construction of a more accurate

379    phylogenetic tree.

380    **Conclusions**

381    In this paper, we developed the computer program SID to detect the non-reference RIPs of 90

382    healthy Han Chinese individuals using high-depth WGS. We described the landscape of RIP

383    distribution on population genomes and annotated the subfamily, orientation, and length of

384    RIPs. We demonstrated that the RIPs could be used as a normal baseline for

385    retrotransposon-related disease research.

386        To our knowledge, this is the largest Han Chinese genomics dataset to date. Compared

387    with 1000GP results from the same samples, approximately half (mean 48.05%; Additional file

388    2: Figure S2) of RIPs in our dataset were previously observed, suggesting that our

389    deep-sequenced data exhibited increased detection sensitivity compared with low coverage

390    data. For example, serum ACE levels were determined by the Alu insertion/deletion (I/D)

17

391　polymorphism in the following order: DD > ID > II [39]. The D allele of the ACE gene was

392　associated with essential hypertension in different populations [40-43]. We found that the ACE

393　gene harbored an Alu insertion in the 15th intron with a frequency of 81/90 in our 90 Chinese

394　genomes compared with a considerably reduced frequency (7/63) in CEPH individuals [12],

395　which was supported by a previous study [44]. To our surprise, no RIP ACEs were present in

396　Han Chinese samples from the 1000GP dataset, which is a high-frequency inserted gene in

397　our RIP data. ACE-specific PCR validation (Additional file 2: Figure S10) and a previous ACE

398　study [45] indicated that our results were consistent with the real values. This finding suggests

399　that adequate sequencing depth is important to investigate RIP frequency and that our data

400　present a result that is consistent with the actual situation. The highly sensitive and accurate

401　RIP dataset provided a perfect opportunity to perform RIP fitness analysis. This study

402　evaluates the natural selection effect on retrotransposon insertions at the population level. As

403　a type of long fragment insertion, RIPs are under approximately neutral selection. This finding

404　is consistent with our result that retrotransposon insertions are mostly relatively

405　inconsequential because the harbored genes are always relatively unimportant. Regarding

406　different types of RIPs in addition to Alu, the longer insertion elements L1 and SVA exhibit

407　weakly positive selection pressure.

408　　　This dataset can be compared with others to provide guidance in research of the

409　disease-causing mechanisms in certain populations and to successfully determine the

410　insertion time of a specific locus. This dataset can also be used as a standard for other RIP

411　research and can serve as a baseline to filter irrelevant RIPs in disease-causing

412　retrotransposon research. Genome-wide association studies (GWAS) have proven their utility

18

413 in identifying genomic variants associated with the risk for numerous diseases. Unlike SNPs

414 and copy number variations (CNVs) that are widely used in GWAS, RIPs have generally been

415 overlooked as a major contributor to human variation. Significantly, this dataset provides a

416 valuable resource to perform GWAS and identify more markers related to complex diseases.

417    The high cost of WGS at high depth is still a major limitation, preventing it from being

418 widely used in TE research. Furthermore, the large amount of data yielded by high-depth WGS

419 makes it difficult to undertake bioinformatic analysis. With the development of biotechnology

420 and IT, this situation should improve soon.

421    The next step is to research RIPs at the transcriptome level. The impact of RIPs on gene

422 expression remains unclear. Combining the genome and transcriptome would provide a

423 comprehensive picture about the regulation of RIPs. Thus, we can further expound the

424 position of the retrotransposon in the course of human evolution.

425

426 **Availability and requirements**

427  ● Project name: Specific Insertions Detector (SID)

428  ● Project home page: https://github.com/Jonathanyu2014/SID

429  ● Operating system(s): Linux

430  ● Programming language: Perl

431  ● Other requirements: Perl 5.14 or later, BLAST v2.2.25 or later, Samtools v1.0 or later

432  ● License: Apache License 2.0

433  ● Any restrictions to use by non-academics: None

434 **Additional files**

435  Additional file 1: Supplementary tables. Data description and the results of RIPs calling. (XLSX

436  1991 kb)

437  Additional file 2: Supplementary texts, figures and tables. (PDF 956 kb)

438  **Abbreviations**

439  CNV, copy number variation; ENA, European Nucleotide Archive; GWAS, genome-wide

440  association study; LTR, long terminal repeat; L1, long interspersed nuclear element 1; NGS,

441  next-generation sequencing; PCA, principal component analysis; RIP, retrotransposon

442  insertion polymorphism; SID, specific insertions detector; SNP, single nucleotide

443  polymorphism; TCGA, The Cancer Genome Atlas; TE, transposable element; TSD, target site

444  duplication；WGS, whole-genome sequencing.

445  **Acknowledgments**

452  **Availability of data and materials**

453  The source code of SID is available from the GitHub and Zenodo repositories[46]. The human

454  (Homo sapiens) reference genome sequence (HG19) and its annotation files were

20

455　downloaded from UCSC Genome Bioinformatics (http://genome.ucsc.edu/). The raw

456　sequence data of YH_CL is available from the ENA repository (accession number ERA000005)

457　[47]. All the YH90 raw sequences have been released to the ENA repository (bioproject

458　number: PRJEB11005) and the processed data is also available from the *GigaScience*

459　GigaDB repository [48]. Snapshots of the code, alignments, and results files are also hosted in

460　GigaDB[49]. Protocols used for simulating reads for SNP Indel calling and detection of

461　transportable element insertions are also hosted in the protocols.io repository[50, 51].

462　**Authors' contributions**

463　BL, SL and YH initiated this project and reviewed the manuscript. QY, XZ, YZ and XH drafted

464　the manuscript. XH and JL edited the manuscript. QY, WZ, XZ and YW performed the data

465　analysis and drew the pictures. YZ and YW designed and developed the SID program. NL, XZ

466　and GL conducted the experiment for sequencing. LX designed the primers and performed

467　PCR validation. YH, BL, SL, XZ, XG and XH provided fruitful discussions.

468　**Competing interests**

469　The authors declare that they have no competing interests.

470　**Author details**

471　[1] BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083,

472　China. [2] BGI-Shenzhen, Shenzhen 518083, China. [3] BGI College, Shenzhen 518083, China.

21

473    [4] Department of Biology, University of Copenhagen, Copenhagen 1599, Denmark.. [5] School of

474    Biology and Biological Engineering, South China University of Technology, Guangzhou

475    510641, China. [6] BGI-Forensics, Shenzhen 518083, China.

**Ethics, consent and permissions**

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

**Consent to publish**

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

480

**References**

482    1.    Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing

483          and analysis of the human genome. Nature. 2001;409 6822:860-921.

484    2.    Cordaux R and Batzer MA. The impact of retrotransposons on human genome evolution.

485          Nature reviews Genetics. 2009;10 10:691-703. doi:10.1038/nrg2640.

486    3.    Kidd JM, Graves T, Newman TL, Fulton R, Hayden HS, Malig M, et al. A human genome

487          structural variation sequencing resource reveals insights into mutational mechanisms. Cell.

488          2010;143 5:837-47. doi:10.1016/j.cell.2010.10.027.

489    4.    Brouha B, Schustak J, Badge RM, Lutz-Prigge S, Farley AH, Moran JV, et al. Hot L1s account

490          for the bulk of retrotransposition in the human population. Proceedings of the National

491          Academy of Sciences of the United States of America. 2003;100 9:5280-5.

492          doi:10.1073/pnas.0831042100.

493   5.   Xing J, Zhang Y, Han K, Salem AH, Sen SK, Huff CD, et al. Mobile elements create structural

494       variation: analysis of a complete human genome. Genome research. 2009;19 9:1516-26.

495       doi:10.1101/gr.091827.109.

496   6.   Cordaux R, Hedges DJ, Herke SW and Batzer MA. Estimating the retrotransposition rate of

497       human Alu elements. Gene. 2006;373:134-7. doi:10.1016/j.gene.2006.01.019.

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

499       Opin Genet Dev. 2012;22 3:191-203. doi:10.1016/j.gde.2012.02.006.

500   8.   Shukla R, Upton KR, Munoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, et al. Endogenous

501       retrotransposition activates oncogenic pathways in hepatocellular carcinoma. Cell. 2013;153

502       1:101-11. doi:10.1016/j.cell.2013.02.032.

503   9.   Solyom S, Ewing AD, Rahrmann EP, Doucet T, Nelson HH, Burns MB, et al. Extensive somatic

504       L1 retrotransposition in colorectal tumors. Genome research. 2012;22 12:2328-38.

505       doi:10.1101/gr.145235.112.

506   10.   Lee E, Iskow R, Yang L, Gokcumen O, Haseley P, Luquette LJ, 3rd, et al. Landscape of somatic

507       retrotransposition in human cancers. Science. 2012;337 6097:967-71.

508       doi:10.1126/science.1222077.

509   11.   Keane TM, Wong K and Adams DJ. RetroSeq: transposable element discovery from

510       next-generation sequencing data. Bioinformatics. 2013;29 3:389-90.

511       doi:10.1093/bioinformatics/bts697.

512   12.   Stewart C, Kural D, Stromberg MP, Walker JA, Konkel MK, Stutz AM, et al. A comprehensive

513       map of mobile element insertion polymorphisms in humans. PLoS Genet. 2011;7 8:e1002236.

514       doi:10.1371/journal.pgen.1002236.

515    13.    Ewing AD and Kazazian HH, Jr. Whole-genome resequencing allows detection of many rare

516          LINE-1 insertion alleles in humans. Genome research. 2011;21 6:985-90.

517          doi:10.1101/gr.114777.110.

518    14.    Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An

519          integrated map of structural variation in 2,504 human genomes. Nature. 2015;526 7571:75-81.

520          doi:10.1038/nature15394.

521    15.    Xing J, Witherspoon DJ and Jorde LB. Mobile element biology: new possibilities with

522          high-throughput sequencing. Trends in genetics : TIG. 2013;29 5:280-9.

523          doi:10.1016/j.tig.2012.12.002.

524    16.    Lan, T; Lin, H; Asker Melchior Tellier, L, C; Zhu, W; Yang, M; Liu, X; Wang, J; Wang, J; Yang,

525          H; Xu, X; Guo, X (2017): Deep whole-genome sequencing of 90 Han Chinese genomes.

526          GigaScience. In Press.

527    17.    Wang J, Wang W, Li R, Li Y, Tian G, Goodman L, et al. The diploid genome sequence of an

528          Asian individual. Nature. 2008;456 7218:60-5. doi:10.1038/nature07484.

529    18.    Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.

530          Bioinformatics. 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.

531    19.    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome

532          Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.

533          Genome research. 2010;20 9:1297-303. doi:10.1101/gr.107524.110.

534    20.    Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O and Walichiewicz J. Repbase

535          Update, a database of eukaryotic repetitive elements. Cytogenet Genome Res. 2005;110

536          1-4:462-7. doi:10.1159/000084979.

537   21.   Wang J, Song L, Grover D, Azrak S, Batzer MA and P L. dbRIP: a highly integrated database of

538         retrotransposon insertion polymorphisms in humans. Hum Mutat. 2006;27 4:323-9.

539   22.   Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, et al. Next-generation

540         VariationHunter: combinatorial algorithms for transposon insertion discovery. Bioinformatics.

541         2010;26 12:i350-7. doi:10.1093/bioinformatics/btq216.

542   23.   Mount DW. Using the Basic Local Alignment Search Tool (BLAST). CSH Protoc.

543         2007;2007:pdb top17. doi:10.1101/pdb.top17.

544   24.   Baillie JK, Barnett MW, Upton KR, Gerhardt DJ, Richmond TA, De Sapio F, et al. Somatic

545         retrotransposition alters the genetic landscape of the human brain. Nature. 2011;479 7374:534-7.

546         doi:10.1038/nature10531.

547   25.   Boissinot S, Chevret P and AV F. L1 (LINE-1) retrotransposon evolution and amplification in

548         recent human history. Mol Biol Evol. 2000;17 6:915-28.

549   26.   Dombroski BA, Mathias SL, Nanthakumar E, Scott AF and Jr KH. Isolation of an active human

550         transposable element. Science. 1991;254 5039:1805-8.

551   27.   Ovchinnikov I, Rubin A and GD S. Tracing the LINEs of human evolution. Proceedings of the

552         National Academy of Sciences of the United States of America. 2002;99 16:10522-7.

553   28.   Ovchinnikov I, Troxel AB and GD S. Genomic characterization of recent human LINE-1

554         insertions: evidence supporting random insertion. Genome research. 2001;11 12:2050-8.

555   29.   Huang X and Madan A. CAP3: A DNA sequence assembly program. Genome research. 1999;9

556         9:868-77.

557   30.   Altschul SF, Gish W, Miller W, Myers EW and Lipman DJ. Basic local alignment search tool.

558         Journal of molecular biology. 1990;215 3:403-10.

559    31.    Hu X, Yuan J, Shi Y, Lu J, Liu B, Li Z, et al. pIRS: Profile-based Illumina pair-end reads

560           simulator. Bioinformatics. 2012;28 11:1533-5.

561    32.    Tarailo-Graovac M and Chen N. Using RepeatMasker to identify repetitive elements in

562           genomic sequences. Curr Protoc Bioinformatics. 2009;Chapter 4:Unit 4 10.

563           doi:10.1002/0471250953.bi0410s25.

564    33.    Hormozdiari F, Alkan C, Ventura M, Hajirasouliha I, Malig M, Hach F, et al. Alu repeat

565           discovery and characterization within human genomes. Genome research. 2011;21 6:840-9.

566           doi:10.1101/gr.115956.110.

567    34.    Batzer MA and Deininger PL. Alu repeats and human genomic diversity. Nature reviews

568           Genetics. 2002;3 5:370-9. doi:10.1038/nrg798.

569    35.    Burns KH and Boeke JD. Human transposon tectonics. Cell. 2012;149 4:740-52.

570           doi:10.1016/j.cell.2012.04.019.

571    36.    Wang J, Song L, Grover D, Azrak S, Batzer MA and Liang P. dbRIP: a highly integrated

572           database of retrotransposon insertion polymorphisms in humans. Hum Mutat. 2006;27 4:323-9.

573           doi:10.1002/humu.20307.

574    37.    Tajima F. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism.

575           Genetics. 1989;123 3:585-95.

576    38.    Boyko AR, Williamson SH, Indap AR, Degenhardt JD, Hernandez RD, Lohmueller KE, et al.

577           Assessing the evolutionary impact of amino acid mutations in the human genome. PLoS Genet.

578           2008;4 5:e1000083. doi:10.1371/journal.pgen.1000083.

579    39.    Rigat B, Hubert C, Alhenc-Gelas F, Cambien F, Corvol P and F S. An insertion/deletion

580           polymorphism in the angiotensin I-converting enzyme gene accounting for half the variance of

26

581    serum enzyme levels. J Clin Invest. 1990;86 4:1343-6.

582    40.    Jeng JR, Harn HJ, Jeng CY, Yueh KC and SM S. Angiotensin I converting enzyme gene

583    polymorphism in Chinese patients with hypertension. Am J Hypertens. 1997;10 5Pt1:558-61.

584    41.    Zee RY, Lou YK, Griffiths LR and BJ M. Association of a polymorphism of the angiotensin

585    I-converting enzyme gene with essential hypertension. Biochem Biophys Res Commun.

586    1992;184 1:9-15.

587    42.    Asamoah A, Yanamandra K, Thurmon TF, Richter R, Green R, Lakin T, et al. A deletion in the

588    angiotensin converting enzyme (ACE) gene is common among African Americans with

589    essential hypertension. Clin Chim Acta. 1996;254 1:41-6.

590    43.    Duru K, Farrow S, Wang JM, Lockette W and T K. Frequency of a deletion polymorphism in

591    the gene for angiotensin converting enzyme is increased in African-Americans with

592    hypertension. Am J Hypertens. 1994;7 8:759-62.

593    44.    Anand SS, Yusuf S, Vuksan V, Devanesen S, Teo KK, Montague PA, et al. Differences in risk

594    factors, atherosclerosis, and cardiovascular disease between ethnic groups in Canada: the Study

595    of Health Assessment and Risk in Ethnic groups (SHARE). Lancet. 2000;356 9226:279-84.

596    45.    Batzer MA, Stoneking M, Alegria-Hartman M, Bazan H, Kass DH, Shaikh TH, et al. African

597    origin of human-specific polymorphic Alu insertions. Proceedings of the National Academy of

598    Sciences of the United States of America. 1994;91 25:12288-92.

599    46.    Qichao Yu. (2016, September 1). Specific Insertions Detector. Zenodo.

600    http://doi.org/10.5281/zenodo.212115

601    47.    Zong C, Lu S, Chapman AR and Xie XS. Genome-wide detection of single-nucleotide and

602    copy-number variations of a single human cell. Science. 2012;338 6114:1622-6.

603     doi:10.1126/science.1229164.

604     48.     Lan, T; Lin, H; Asker Melchior Tellier, L, C; Zhu, W; Yang, M; Liu, X; Wang, J; Wang, J; Yang,

605             H; Xu, X; Guo, X (2017): Supporting data for "Deep whole-genome sequencing of 90 Han

606             Chinese genomes" GigaScience Database. http://dx.doi.org/10.5524/100302

607     49.     Yu, Q; Zhang, W; Zeng, Y; Zhang, X; Wang, Y; Wang, Y; Xu, L; Huang, X; Li, N; Zhou, X; Lu,

608             J; Guo, X; Li, G; Hou, Y; Liu, S; Li, B (2017): Supporting data for "Population-wide Sampling

609             of Retrotransposon Insertion Polymorphisms Using Deep Sequencing and Efficient Detection"

610             GigaScience Database. http://dx.doi.org/10.5524/100318

611     50      Haoxiang Lin: SNP INDEL calling. protocols.io dx.doi.org/10.17504/protocols.io.grkbv4w

612     51.     GigaScience Database: Simulating reads for detection of transportable element insertions.

613             protocols.io. 2017. dx.doi.org/10.17504/protocols.io.imrcc56

614

615

616     **Figure legends**

617     **Fig. 1** The principle of retrotransposon insertion detection. (**a**) Schematic diagram of using SID

618     for RIP detection in the genome. TSD: target site duplication. SID: Specific Insertions Detector.

619     (**b**) An example of reads mapping for predicted homozygous insertions. (**c**) An example of

620     reads mapping for predicted heterozygous insertions. In (**b**) and (**c**), the red bases indicate the

621     mismatches, and the sequences with an orange background represent the clipped part of the

622     reads. The clipped reads are derived from one allele with inserted retrotransposons, and the

623     normal reads are derived from the other allele with the same reference. The three reads with

624     asterisks indicate no clipped part but the presence of terminal mismatches, which can also

625 support the breakpoint and exhibit consistency with the clipped reads.

626 **Fig. 2** Assessing the SID results. (**a**) Detecting accuracy and sensitivity estimation along

627 cumulating sequencing depth of simulated data. (**b**) RIP genotyping of YH_CL. PCR validation

628 results are marked. HEE: estimated heterozygous site. HOE: estimated homozygous site.

629 HEV: validated heterozygous site. HOV: validated homozygous site. The dash line indicates

630 the estimated boundary between heterozygous and heterozygous sites. Note that some of the

631 validated RIPs are present in the same locus in the plot figure.

632 **Fig. 3** Comprehensive landscape of non-reference RIPs of YH90. (**a**) Proportions of novel

633 insertions identified for each type of retrotransposon. (**b**) Comparison of YH90 non-reference

634 RIP results with dbRIP. Adjacent 100-bp regions of RIPs were taken into consideration. (**c**) TE

635 distribution of each YH90 sample. (**d**) Box plots of non-reference RIP distribution among

636 autosomes. (**e**) TE frequency distribution among YH90 samples. Rings from outer to inner

637 indicate Alu insertions frequency, L1 insertion frequency, SVA insertion frequency, LTR

638 insertion frequency and cytoband structure. The inside frequency of the rings indicates the

639 insertion frequency for the Northern Chinese group, and the outside frequency represents that

640 of the Southern Chinese group. (**f**) RIP distribution in different functional regions of the

641 genome.

642 **Fig. 4** Population genetics analysis based on YH90. (**a**) A two-epoch population with a recent

643 contraction; a three-epoch bottleneck-shaped history, which contained a reduction of the

644 effective population size in the past followed by a recent phase of size recovery. Details of the

645 parameters for all models are provided in Additional file 2: Table S12. (**b**) The observed SNP

646    frequency spectra and expected neutral SNP frequency spectra under different demographic

647    models. (**c**) Observed and expected RIP site frequency spectra before demographic correction

648    of each subfamily. (**d**) Assessing the evolutionary impact of RIPs in the human genome. The

649    allele frequency distribution of RIPs was compared among observed, neutral models and

650    negative models after demographic correction.

651    **Fig. 5** Phylogenetic analysis using RIPs and SNPs. (**a**) The detected RIPs were used for PCA.

652    Each dot represents a sample from YH90 and is plotted as scatterplot using PC1 and PC2.

653    Red indicates samples from individuals from northern China, and blue indicates individuals

654    from southern China. (**b**) The detected SNPs were used for PCA. The plot layout and legend

655    are the same as those presented in (**a**). (**c**) Phylogenetic tree constructed using the detected

656    RIPs. HG19 (green) is used as a control. Red indicates samples from individuals from northern

657    China, and blue indicates samples from individuals from southern China. (**d**) Phylogenetic tree

658    constructed using the detected SNPs. HG19 (green) is used as a control. Plot layout and

659    legend are same as that presented in (**c**).

Figure 1

Figure 2

Figure 3

Figure 3

Figure 4

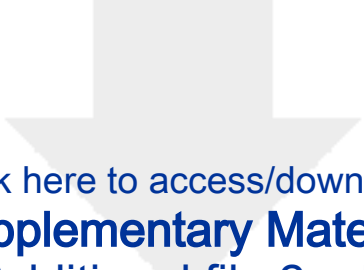Figure 5

Additional file 1

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

Additional file 2

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