

1 Landscape and Variation of Novel Retroduplications in 26 Human
2 Populations

3

4

5 Yan Zhang^{1,2,3¶}, Shantao Li^{1¶}, Alexej Abyzov^{4*}, Mark B. Gerstein^{1,2,5*}

6

7 ¹Program in Computational Biology and Bioinformatics, Yale University, New
8 Haven, CT 06520

9 ²Department of Molecular Biophysics and Biochemistry, School of Medicine, Yale
10 University, New Haven, CT 06520

11 ³Department of Biomedical Informatics, College of Medicine, The Ohio State
12 University, Columbus, OH 43210

13 ⁴Department of Health Sciences Research, Center for Individualized Medicine, Mayo
14 Clinic, Rochester, Minnesota, MN 55905

15 ⁵Department of Computer Science, Yale University, New Haven, CT 06520

16

17 *Corresponding authors

18 E-mails: abyzov.alexej@mayo.edu (AA); mark@gersteinlab.org (MBG)

19

20 ¶These authors contributed equally to this work.

21

22

23

24

25

26

27

28

29

30

31

32

33 **Abstract**

34 Retroduplications come from reverse transcription of mRNAs and their
35 insertion back into the genome. In this study, we performed comprehensive discovery
36 and analysis of retroduplications in unprecedented 2,535 individuals from 26 human
37 populations. We developed an integrated approach to discover novel retroduplications
38 from both high-coverage exome sequencing and low-coverage whole genome
39 sequencing data, utilizing information of both exon-exon junctions and discordant
40 locations of paired-end reads. We detected 503 parent genes having novel
41 retroduplications absent in the human reference genome. The set reveals the high-
42 resolution landscape of human germline retroduplication polymorphism. It gives us
43 the power to perform extensive analysis of retroduplication variation.

44 We successfully constructed phylogenetic trees of human populations solely
45 based on retroduplication variations, which confidently represents the superpopulation
46 structure, and indicates that variable retroduplications are effective markers of human
47 populations. We further identified 43 retroduplication parent genes that can
48 differentiate superpopulations. We have also detected several interesting intragenic
49 insertion events, including SLMO2 retroduplication and insertion into CAV3, which
50 worth further investigation for disease propensity. By investigating local genomic
51 features at retroduplication insertion sites, we observed that novel retroduplications
52 insertion sites are associated with nucleosome positioning and co-inserted L1
53 elements belong to young L1 families, indicating recent retroduplication activity
54 occurred in human migration.

55 Our investigation provides valuable insight into retroduplication functional
56 impact and their association with genomic elements. We anticipate our
57 retroduplication discovery approach and analytical methodology to have broader
58 applications in biomedical researches, where exome sequencing data is abundant.

59

60 **Author Summary**

61 We developed an approach and performed comprehensive discovery of
62 retroduplications from 26 human populations, utilizing whole exome/genome
63 sequencing data. Our high-resolution landscape of retroduplications reveals that
64 variable retroduplications are effective markers of human populations and can track
65 population divergence. We observed that novel retroduplications come from genes

Shantao 3/20/2017 5:29 PM

Deleted: Retroduplications are common in mammal genomes.

Shantao 3/20/2017 5:37 PM

Deleted: have

Shantao 3/20/2017 5:37 PM

Deleted: have

Shantao 3/20/2017 5:42 PM

Deleted: come from genes with relatively high expression level

Shantao 3/20/2017 5:44 PM

Comment [1]: Need to rewrite...this needs to be like pop. Science. Not snippets from abstract.

72 with relatively high expression level and co-inserted L1 elements belong to young L1
73 families, indicating recent retroduplication activity occurred in human migration. We
74 have also detected several interesting intragenic insertion events, including SLMO2
75 retroduplication and insertion into CAV3, which worth further investigation for
76 disease propensity.

77

78 **Introduction**

79 Retrotransposons are class I transposable elements. In retrotransposition
80 events, they are first transcribed into RNA and then reverse transcribed back into
81 DNA, which are eventually inserted into a new position in the genome. It has been
82 found that L1 retrotransposons, the only autonomous mobile elements in human
83 genome, also occasionally pick up cellular mRNAs as templates for reverse
84 transcription and insertion [1–3]. Although RNA-mediated retroduplication is less
85 common and widespread than DNA-mediated duplication [4], recent studies have
86 revealed extensive retroduplication polymorphism in human genomes [5–7].

87 Retroduplication of genes contributes to new gene generation and genome
88 evolution [4,8,9]. While most of the retroduplications suffer from lack of promoters,
89 5' truncation, mutations, inactive local chromatin environment and other unfavorable
90 factors that hinder the expression of functional protein products, they do exhibit
91 functional impact at times. In some cases, cellular environment change, such as
92 cancer initiation, can “activate” retroduplications, and both transcription and
93 translation evidence of such cases have been observed [10–12]. In other cases,
94 transcription products play a role in the expression regulation of their parent genes
95 [13,14]. Two known transcriptional level regulatory mechanisms are RNA
96 interference [15–17], and transcription products serving as competitive miRNAs
97 binding targets [18,19]. Sometimes retroduplications can have high impact on
98 genomic functions when inserting into functional regions. Studies have confirmed
99 cases in which germline intragenic retroduplications result in liver cancer
100 susceptibility [20] and primary immunodeficiency [21]. Besides germline events, a
101 number of studies have reported massive somatic retroduplication events and their
102 critical roles in tumor development [20,22–25] and neuron development [26,27].

103 Retroduplications carry several distinctive features: exon-exon junctions,
104 genome locations distant to parent genes, poly-A tails, and L1 transposition markers

105 such as target-site duplications (TSDs) and human L1 endonuclease preferential
106 cleavage sites. In this study, we developed an integrative approach to exploit these
107 features for novel variable retroduplication identification, and successfully applied it
108 to 2,535 individuals from 26 populations sequenced by the 1000 Genomes Project
109 Phase 3 [28–30]. Our study adds an additional category of genetic variation to the
110 released Phase 3 categories [29,30]. We further performed extensive population
111 genetic analysis, association analysis, event mechanism inference, and functional
112 analysis of retroduplications. Our study is indicative of human migration and
113 evolution history, and provides valuable insight into retroduplication functional
114 impact and their association with genomic elements.
115

116 **Results and Discussion**

117 First, we performed retroduplication discovery for each individual, using the
118 exon-exon junction strategy on high-coverage whole exome sequencing (WES) data
119 (see **Supplementary Methods**, and **Fig. 1**). We controlled the false discovery rate
120 using decoy exon junction libraries. As a result, we have called a total of 15,642
121 retroduplications from 2,533 individuals (with two outlier samples excluded) for 503
122 unique parent genes (**S1-S2 Fig.**, **S1 Table**, and **S2 File**). On average, an individual
123 has 6 novel retroduplications identified based on exon-exon junctions. Next, we
124 identified retroduplication insertion sites for 152 of the parent genes based on
125 discordant paired-end reads, using shallow-sequenced whole genome sequencing
126 (WGS) data pooled by population (**Fig. 1**, and **S3 File**). Multiple genomic features are
127 exploited in this pipeline, in order to achieve high sensitivity in calling, while
128 conservatively controlling the false discovery rate. The retroduplications identified in
129 our study adds an additional category of genetic variation to the released Phase 3
130 categories [29,30].
131

132 **Fig. 1. Overview of the retroduplication calling pipeline.** A – A simplified flow chart of our calling
133 pipeline. B – A schematic diagram of our strategies. We first align unmapped reads to exon junction
134 libraries and use decoy libraries to control false discovery rate (FDR). Then, we collect discordant
135 paired-end reads, and cluster the reads that are mapped distal to the parent genes. The location of
136 clustered distal reads indicates retroduplication insertion site.
137

138 Compared to previous studies of human germline retroduplications, which
139 relied on about 1,000 shallow-sequenced individuals [5–7] from 1000 Genomes
140 Project Phase 1[31], the population set and sequencing coverage in Phase 3 has scaled
141 up about 10-fold combined (S3 Fig.). Besides the retroduplication calls shared among
142 callsets, there are also large number of calls unique to our callset, which is likely due
143 to newly enrolled populations in Phase 3 data, and the higher sensitivity of our
144 methods. Similarly, we resolved a higher fraction than before[5,7] of insertion sites of
145 retroduplications, 152/503 (30.2%). Functional enrichment analysis for the 503
146 unique parent genes shows that functions related to ribosome/structural molecule
147 activity, intracellular organelle lumen/nucleoplasm, and protein complex assembly are
148 among the most enriched. This observation is in accordance with previous study [5],
149 indicating retrotransposition is coupled with cell division.

150 We have identified novel retroduplications, which are insertions relative to the
151 reference genome. There are also retroduplications that are deletions relative to the
152 reference genome (i.e. absent in the individuals but present in the reference genome).
153 These events can be detected by overlapping known processed pseudogenes in the
154 reference genome with Phase 3 deletion set. We carried out this in the supplement,
155 finding 68 such deletion events (S4 File). These 68 regions are present in Hg19 as
156 processed pseudogenes, but reported as deletions in 1000 Genomes Phase 3. This type
157 of events is far less common than retroduplication insertions, thus we suggest
158 focusing on retroduplication insertions in the study.

159 The high-resolution landscape of germline retroduplication polymorphism
160 presented by our callset gives us the power to perform extensive analyses of
161 retroduplication variation. Among all 503 parent genes with novel retroduplications,
162 361 (71.8%) are exclusively identified in a single population, while only 29 (5.8%)
163 are commonly identified in more than 10 populations (S2 Fig. B). Retroduplications
164 are larger events than SNPs. It is known that individual structural variations are more
165 likely to lead to phenotypic differences than individual SNPs [32]; thus,
166 retroduplications might be more influential population-specific events than SNPs.
167 From all identified parent genes, we identified 43 that can differentiate
168 superpopulations (with significantly large fixation index F_{ST} , adjusted empirical p-
169 value < 0.001, see S2 Table).

170 From the frequency spectrum of retroduplication parent genes (Fig. 2AB, S5
171 Fig.), we observed expected and confident cluster cohesion of superpopulations

172 (African, Ad Mixed American, East Asian, European, and South Asian). We
173 hypothesize that **many of the exclusive retroduplications emerged after population**
174 **divergence**. We further constructed phylogenetic trees of human populations solely
175 based on novel retroduplication variations (**Fig. 2C**). The phylogenetic trees can
176 confidently represent the superpopulation structure and also show mixed populations
177 (Ad Mixed American) mingling with other superpopulations. These observed
178 population relationships are consistent with human migration history, which
179 reconfirms the effectiveness of retroduplications as population markers as well as
180 validates our approach to retroduplication discovery.

181

182 **Fig. 2. Common retroduplication frequency spectrum and phylogenetic tree.** A – Frequency
183 spectrum of 29 retroduplication events that are detected in more than 10 populations. Hierarchical
184 clustering was used. B – PCA biplot of the populations based on frequencies of these 29
185 retroduplication events. Different colors indicate five superpopulations, i.e. AFR (African), AMR (Ad
186 Mixed American), EAS (East Asian), EUR (European), and SAS (South Asian). Arrows represent
187 loadings of parents genes. Ad Mixed Americans are marked with '*'. C – Consensus phylogenetic tree
188 built based on novel retroduplications from all 26 populations enrolled in the 1000 Genome Project
189 Phase 3. Bootstrap probability (BP) value is computed from ordinary bootstrap resampling. It is the
190 frequency of the cluster appearing in bootstrap replicates. Approximately unbiased (AU) probability
191 value is calculated from multiscale bootstrap resampling [33,34]. AU is less biased than BP. Bootstrap
192 resampling was performed 1,000 times for generating the trees that are summarized in the consensus
193 tree. Manhattan distance and average linkage is used in hierarchical clustering.

194

195 For each population enrolled in the Geuvadis RNA-sequencing project (i.e.
196 CEU, FIN, GBR, TSI, and YRI) [35], **we tested whether having novel**
197 **retroduplication(s) is associated with the parent gene's expression level. We did not**
198 **see any significant association from this analysis (S6 File), i.e. no retroduplication**
199 **event was identified as an eQTL**. However, while comparing expression level of
200 retroduplication parent genes to all genes, we see a weak but ubiquitous trend that
201 novel retroduplications came from the genes with relatively high expression level (p-
202 value $< 1.4 \times 10^{-5}$ for each population, calculated from omnibus tests, see **S7 File**). It
203 is consistent with our expectation that the more mRNAs a gene has made, the higher
204 probability that it will be converted into complementary DNA and inserted back into
205 the genome.

206 To investigate local genomic features around insertion sites which might
207 explain insertion localization preference and imply retroduplication mechanism [36],

208 we tested the enrichment of overlap between the genomic features and insertion sites.
209 Inheritable retroduplication events occurred in germline. We found that
210 retroduplication insertion sites are enriched within hypomethylated regions in sperm
211 (2.0-fold, empirical p-value < 0.0012). It is likely that retroduplication events exhibit
212 certain preference in insertion sites associated with open chromatin. In order to further
213 characterize nucleosome positioning [37,38], we aggregated nucleosome positions
214 around insertion sites. Overall, insertion sites show high regularity of nucleosome
215 location (empirical p-value from permutation test 2×10^{-4}) (Fig. 3A). Highly
216 nucleosome regularity often indicates the presence of chromatin remodeling and DNA
217 binding proteins [39], which creates favorable loosely packed microenvironment for
218 insertion.

219

220 **Fig. 3. Overlap between retroduplication insertion sites and genomic features/functional elements.**
221 A – Aggregation plot around insertion sites with strongly positioned nucleosomes. B – Association
222 between discordant reads clusters that only have support on one side and L1 element subfamilies. Fold
223 change and empirical p-values were obtained from permutations tests. *** indicates adjusted p-value <
224 0.001. C – Overlap between genomic elements and retroduplication insertion sites. The enrichment of
225 overlap is expressed as log₂ fold change of the observed overlap statistic versus the mean of its null
226 distribution. Positive (negative) log₂ fold change indicates enriched (depleted) genomic element-
227 insertion overlap, compared to random background. * indicates empirical p-value ≤ 0.002.

228

229 Insertions points could be supported by discordant reads from both sides or
230 either one side. There is no fundamental preference for retroduplicated DNA
231 segments to insert into other retroelements such as L1 elements. However, L1 that is
232 involved in retroduplication is sometime co-duplicated and inserted next to the
233 retroduplicated segment. This type of co-insertion event can be detected by looking at
234 the insertion sites that only have discordant-read support on one side. In these cases,
235 we see a preference for the co-inserted L1 to be a young L1, represented by LIHS
236 (4.7-fold, p-value < 0.001) and L1PA (1.9-fold, p-value < 0.001) (Fig. 3B).
237 Contrastingly, for insertion sites where we do not detect a co-insertion (i.e. insertion
238 sites that are supported by both sides) we do not observe such young L1 preference
239 (p-value > 0.05). Enrichment of young and active L1 subfamilies involving in
240 speculated L1 transductions suggests novel retroduplication variants have happened
241 very recently.

242 In order to investigate the functional impact of retroduplication insertions on
243 genomic functions, we tested the significance of overlap between retroduplication
244 insertion sites and genomic elements compared to random genomic background (**Fig.**
245 **3C**). As expected, ultraconserved regions are significantly depleted (p-value < 0.001).
246 This observation is consistent with our knowledge that in **general population**, variable
247 retroduplications should not interrupt with evolutionary or functionally constrained
248 regions. Unexpectedly, we observed that intron regions are also depleted (p-value <
249 0.01), which might be due to negative selection that maintains conserved alternative
250 splicing by avoiding interruption from insertion into introns.

251 Among the 43 parent genes that differentiate superpopulations (top 43 genes
252 in **S2 Table**), we have detected several interesting intragenic insertion events. For
253 example, we observed that SLMO2 (slowmo homolog 2, ENSG00000101166) is
254 retroduplicated and inserts into the last intron of CAV3 (caveolin 3,
255 ENSG00000182533). SLMO2 retroduplication insertion sweeps through all seven
256 African populations almost exclusively. Based on exon-exon junction evidence, we
257 found 30 cases in African populations and only one case in MXL (Ad Mixed
258 American, **S5 File**). CAV3 variants are strongly associated with cardiac dysrhythmia,
259 such as long QT syndrome [40] and sudden infant death syndrome [41].
260 Epidemiological studies have shown that African descendant is a risk factor for
261 prolongation of QT interval [42] and sudden infant death syndrome [43]. Such
262 insertion events might worth further investigation for susceptibility of diseases. **We**
263 **have identified a total of 12 retroduplication insertion events related to disease, and**
264 **report the full list and affected populations in S3 Table.**

265 In summary, we developed an integrative approach for variable
266 retroduplication discovery and successfully applied it to sequencing data of 2,535
267 individuals from 26 populations. We have shown the power of using whole
268 exome/genome sequencing data in retroduplication identification. Furthermore, we
269 performed comprehensive analysis of our large retroduplication dataset, which reveals
270 variational landscape of novel retroduplications, and shed a light on population
271 differentiation, and functional impact of retroduplications on the genome.

272

273 **Materials and Methods**

274 **Data resources**

275 Whole exome sequencing and whole genome sequencing data of 2,535
276 individuals from 26 populations were generated by the 1000 Genomes Project Phase 3
277 (whole-genome sequencing with mean depth 7.4x and read length of 100bp; targeted
278 exome sequencing with mean depth 65.7x and read length of 76bp) [28–30].

279 Population description can be found at
280 <http://www.1000genomes.org/category/frequently-asked-questions/population>.

281 Protein-coding gene expression data (Peer-factor normalized RPKM) is obtained from
282 the Geuvadis RNA-sequencing project [35], which generated RNA sequencing data
283 from lymphoblastoid cell lines of 462 individuals from 5 populations (CEU, FIN,
284 GBR, TSI and YRI) enrolled in the 1000 Genomes Project. We use human reference
285 genome build 37 [44] and GENCODE v19 human genome annotation [45] in the
286 study.

287

288 **Calling pipeline**

289 The calling pipeline is developed and customized for generating
290 retroduplication calls from high-coverage exome sequencing data. A simplified
291 flowchart of the current pipeline is shown in **Fig. 1**. We also provide the code for
292 download ([URL](#)).

Shantao 3/20/2017 5:47 PM

Deleted: and

Shantao 3/20/2017 5:47 PM

Deleted: S8 File

295 **Build true and decoy exon junction libraries.** For calling retroduplications from
296 whole exome sequencing data, we need to build exon junction libraries from
297 annotated protein coding exons. The true exon junction library is built by joining pairs
298 of protein coding exon segments within the same genes, while maintaining exons'
299 order on the strand. Exon segments of length 100 bases adjacent to the joining splice
300 sites are combined (S4 Fig). We also build five decoy exon junction libraries for the
301 purpose of controlling false call rate. The decoy exon junction libraries contain fake
302 exon junctions, in which exon annotations are shifted by e base(s) on both sides (i.e.
303 start location + e, end location - e). e is taken as 1, 2, 3, 6, and 12 for each decoy exon
304 library, respectively.

305 **Generate unmapped read alignments.** We generate reduplication calls for each
306 individual. Unmapped reads can be utilized for calling novel retroduplications that are
307 absent in the reference genome. We use SAMtools [46] to extract unmapped reads
308 from exome bam files, then use BWA-0.7.7 to align the unmapped reads to all of true
309 and decoy exon junction libraries (S4 Fig). d1 and d2 are the number of bases that the
310 read maps to either exon segment. $\min(d1, d2) \geq d$ is required for a newly mapped
311 read to be reported from our pipeline. We also calculate the mismatch rate r for each
312 mapped read. d and r are parameters automatically tuned in the range [1, 15] and
313 [0.00, 0.05], respectively, ensuring the most number of calls from the true exon
314 junction library while satisfying no false calls from any decoy library.

315 **Estimate FDR of the exon-exon junction callset.** We optimize the calling
316 parameters so that no calls are detected in any decoy library, still this does not
317 guarantee that the generated retroduplication calls are free of false positives. Let us
318 assume that per sample FDR is λ . For simplicity, but without losing generality, we
319 assume that λ is uniform across all samples. Then, the count of false calls per sample
320 follows a Poisson distribution. The chance of having zero false calls per sample is
321 $\exp(-\lambda)$. Since we never detect false calls in the 2,533 samples, $\exp(-\lambda)^{2533}$ is the
322 chance of observing no false calls. For 95% confidence level, this probability is equal
323 to 0.05. This yields per sample FDR λ of 1.2×10^{-3} . Similarly, for 99% confidence
324 level, λ is 2.7×10^{-3} . This projects to 3 (at 95% confidence) and 7 (at 99% confidence)
325 false calls over the entire callset. Thus, for the 503 unique parent genes with variable
326 retroduplications, we estimate <2% FDR with 99% confidence.

327 Moreover, as we always try to move further to more restricted calling criteria
328 after no call is detected in decoy libraries, our FDR estimation above is conservative.
329 Using additional simulated decoy libraries with different shifting coordinates as test
330 libraries, we do not detect any false positive call under our final calling parameters.
331 This further supports our very low FDR estimation.

332 **Report novel retroduplication calls.** Multiple “previously unmapped” reads
333 (unmapped to the reference genome) might be mapped to the same exon-exon
334 junction, supporting the existence of the novel exon-exon junction. Furthermore,
335 multiple exon-exon junctions with mapped reads might support the existence of a
336 gene retroduplication event. We report a gene having novel retroduplications, when it
337 has at least two non-overlapping supporting exon-exon junctions, and at least one
338 junction is supported by at least two mapped reads. The genes (also called parent
339 genes) with novel retroduplications are called for each person individually. We
340 noticed that the 1000 Genomes Project Phase 3 provides paired-end sequencing data
341 for all individuals but NA19318. We include this individual into our analysis, as
342 single-end sequencing does not seem to affect the performance of this pipeline.

343 **Detect retroduplication insertion sites.** In the insertion site detection step, we pool
344 low-coverage whole genome sequencing data by population, and call insertion sites
345 for each population. We search for discordant paired-end reads (with a minimum
346 quality score of 15) with one read correctly mapped to the parent gene, and the other
347 read mapped to a different chromosome or at least 1 kb away from the gene. In order
348 to avoid false discovery, we limit our searching scope to the parent genes identified
349 from the exon-exon junctions.

350 Read pairs with proper orientations are clustered using average linkage
351 clustering. It can be shown that this linkage criterion is not likely affected by the local
352 coverage. Assuming uniform distribution of reads, it can be shown mathematically
353 that the expected distance between reads supporting the same insertion point is

$$\frac{2(IS - RL) + 1}{3},$$

354 where IS is the insertion size and RL is the read length. As the insertion size in most
355 cases is around 200-400 bp and the read length is about 70-100 bp, we choose 500 bp
356 as the cut-off for average linkage distance to stop clustering. This cut-off not only
357 takes the deviations of insertion size into consideration, but also allows sufficient

358 space for target site duplications (TSDs). A valid insertion site must have at least two
359 reads on both sides (i.e. stands). Overlapped insertion sites with identical parent gene
360 and orientation are further merged across populations, as these sites should represent
361 one single event.

362 In our insertion site detection step, we have discovered single-side clusters that
363 have sufficient number of supporting reads. We require at least four reads on one side
364 and no reads on the other side to call those incomplete single-side events. Single-side
365 events across populations are merged by requiring identical parent gene, same
366 orientation, and within 500 bp distance using locations defined by the cluster of one
367 end. Also we only use insertion sites on chromosomes (i.e. exclude alternative locus).

368 **Detect retroduplication deletions.** Retroduplication deletions (relative to the
369 reference genome) are the variable retroduplications that are absent in the individuals
370 but present in the reference genome. We detect the retroduplication deletions by
371 overlapping known processed pseudogenes in the GENCODE v19 with 1000
372 Genomes Phase 3 deletion set, requiring the processed pseudogene region overlaps at
373 least 50% of the deletion region. The results are available in **S4 File**.

374

375 **Build population phylogenetic trees based on novel retroduplication calls**

376 **Generate retroduplication frequency matrix.** Some retroduplication parent genes
377 are called commonly among multiple populations, while some others are called
378 exclusively in a single population. Besides, parent genes are called at different
379 frequencies within a population. This information can be used for measuring distance
380 between populations, while taking into account different retroduplication frequencies.
381 We define a retroduplication frequency matrix, from which distance measures can be
382 calculated.

383 Suppose there are N populations, and M unique parent genes are identified in
384 these populations. The retroduplication frequency matrix A is defined as an $M \times N$
385 matrix, with each element $A_{m,n}$ ($m=1,2,\dots,M; n=1,2,\dots,N$) being a value in $[0, 1]$,
386 representing the percentage of individuals in population n having this unique parent
387 gene m called.

388 **Bootstrap phylogenetic trees.** We use Manhattan distance as the distance measure
389 between each pair of populations (i.e. Manhattan distance between two columns in A).
390 Average linkage is used in hierarchical clustering for generating each tree. 1000
391 bootstrap replications are performed, and the uncertainty is assessed using Pvcust
392 [33]. The reported AU (Approximately Unbiased) probability values [33,34] are used
393 to indicate the certainty of sub-tree structures generated from multi-scale bootstrap
394 resampling [47–49]. The higher the AU probability value, the more confident the
395 substructure is.

396

397 **Detect population differentiation due to retroduplication polymorphism**

398 We check population differentiation due to retroduplication polymorphism,
399 based on retroduplication frequencies in different superpopulations. Herein we pool
400 the 26 populations into 5 superpopulations (African, Ad Mixed American, East Asian,
401 European, and South Asian) as defined by the 1000 Genomes Project. For each given
402 retroduplication parent gene, we calculate the population differentiation measure
403 equivalent to the fixation index [50]. We define the test statistic

$$404 \quad F_{ST} = \frac{p(1-p) - \sum_{i=1}^5 c_i p_i (1-p_i)}{p(1-p)},$$

405 in which $i = 1, \dots, 5$ corresponds to the i th superpopulation, p is the retroduplication
406 frequency of a given parent gene in the total population, p_i is the retroduplication
407 frequency of the same parent gene in the i th superpopulation, and c_i is the relative
408 population size of the i th superpopulation. c_i is calculated as the number of
409 individuals in the i th superpopulation divided by the number of individuals in the total
410 population. The larger the F_{ST} , the more different the retroduplication frequencies
411 among superpopulations. One-tailed empirical p-value is calculated comparing the
412 observed F_{ST} versus the null distribution of F_{ST} . The null distribution is calculated
413 from 1000 fake population sets generated by shuffling individual labels, while
414 maintaining the size unchanged for each population. By the significance of F_{ST} , i.e.
415 the p-value adjusted by Benjamini-Hochberg procedure [51], we can detect the
416 retroduplications that can differentiate populations.

417

418 **Analyze association between retroduplication and gene expression**

419 We utilize our retroduplication callset and the Geuvadis gene expression data
420 (Peer-factor normalized RPKM) [35] to analyze the association between
421 retroduplication occurrence and gene expression. Matching data of the individuals
422 enrolled in both the 1000 Genomes Project and the Geuvadis project are used. The
423 association tests are performed for each population, respectively, in order to rule out
424 the confounding by population stratification.

425 **Retroduplication eQTL analysis.** For a certain population, we perform the
426 association test within the set of retroduplication parent genes: test whether having
427 novel retroduplication(s) or not is associated with the parent gene's expression level.

428 First, differential expression of each parent gene is tested between the group of
429 individuals that have novel retroduplications of this gene and the group of individuals
430 that do not. Two-sided Wilcoxon rank sum test is used. P-values are adjusted by
431 Benjamini-Hochberg procedure [51]. A gene is reported to be differentially expressed
432 in the parent gene set if its adjusted p-value is less than 0.05. Furthermore, the global
433 differential expression of all the parent gene set is tested using Fisher's combined
434 probability test [52] on unadjusted p-values. This omnibus test can test the combined
435 effect of multiple parent genes, whose individual effects are not necessarily strong. If
436 the combined p-value is less than 0.05, we can conclude that the association between
437 retroduplication variance and parent gene expression is significant. The results are
438 available in **S6 File**.

439 To re-confirm the result, we also perform two-sided Wilcoxon signed rank
440 test. For each gene, medium expressions of both groups (having the novel
441 retroduplication or not) are paired. The test result is consistent with that of the
442 Fisher's method.

443 **Expression level of retroduplication parent genes compared to all genes.** For a
444 certain population, we test whether the retroduplication parent genes are highly
445 expressed among all the genes measured in the Geuvadis data set. We take medium
446 expression value over all individuals for each gene as the representative expression
447 value. One-tailed empirical p-value is calculated comparing the expression value of
448 each parent gene versus the null distribution of expression values of all genes. It
449 indicates the significance of each retroduplication parent gene having high expression
450 value among all genes. Fisher's combined probability test is performed on the
451 empirical p-values. If the combined p-value is less than 0.05, that means in general
452 the parent genes are significantly highly expressed among all genes. The results are
453 available in **S7 File**.
454

455 **Explore association between local genomic features and retroduplications**
456 **insertion sites**

457 To test the association between sperm methylation patterns and
458 retroduplication insertion sites, we intersect out insertion sites with hypomethylated
459 regions in sperms [53]. L1 annotation (RepeatMask), ENCODE HESC DNase I
460 hypersensitive data and genomic GC contents are downloaded from the USCS
461 Genome Browser [54]. Well-positioned nucleosome data is obtained from a recent
462 study on multiple individuals [55].

463 We randomly shuffle the locations of insertion sites for 10,000 times on the
464 same chromosome, excluding the gap regions, to obtain an empirical distribution of
465 the null hypothesis. For fold changes, we use the mean of this distribution as the best
466 estimate of the expected value. Calculation of p-value is empirical in order to be
467 conservative. We use Bonferroni correction in case of multiple hypothesis testing.
468 Unless specified otherwise, we only report corrected p-value. In order to avoid any
469 effect of the difference of location precision across different insertion sites, we
470 enlarge the insertion site region to 500 bp while keeping the middle point of insertions
471 unchanged. We also exclude insertion points on alternative locus in the genome.

472 For aggregation plot on well-positioned nucleosome and GC content, we use
473 200 bp bins to calculate the base overlap, and the final plot was further window-
474 smoothed with window size of 10. Normalization is performed by taking mean value

475 of the first and last 20 bins as background. We use the GC contents from UCSC
476 browser track, which is binned in 5 bp.

477

478 **Investigate impact of retroduplication insertions on genomic functions**

479 We test the significance of overlap between retroduplication insertion sites and
480 genomic elements, including gene, CDS, exon, UTR, intron, pseudogene and
481 lincRNA annotated in GENCODE v19, and ultraconserved regions (evolutionary
482 constraint regions across species), ultrasensitive non-coding regions (regions
483 particularly sensitive to disruptive mutations) and TF (transcription factor) peak
484 regions obtained from ENCODE RNA-seq data [10] and literature [30,56–59]. The
485 overlap between a genomic element type and the insertion sites is measured by the
486 partial overlap statistic, which is the count of genomic elements that have at least 1 bp
487 overlap with the detected insertion sites.

488 We randomly shuffle the locations of insertion sites for 1,000 times on the
489 same chromosome, excluding the Hg19 gap regions, to obtain an empirical
490 distribution of the null hypothesis. In the permutation tests, the null distribution of the
491 overlap measures is calculated from true genomic elements intersecting randomly
492 shuffled insertion locations. The enrichment of overlap is represented by log2 fold
493 change of the observed overlap statistic versus the mean of its null distribution.
494 Empirical p-value is calculated.

495 In order to avoid any effect from different location precisions, we enlarge the
496 insertion intervals uniformly to 1000 bp, while keeping the middle point of insertions.
497 We only use insertion sites on chromosomes (i.e. exclude alternative locus) in the
498 analysis.

499

500 **Functional enrichment analysis**

501 We use DAVID [60] to annotate functional terms for retroduplication parent
502 genes, and survey functional term enrichment.

503

504 **Search for literature supported disease-associated insertion events**

505 We generate a list of genes where the novel retroduplication insert into. We
506 then search these genes in the DISEASES database [61] to find disease-gene
507 associations reported in literature.

508

509 **Acknowledgements**

510 The authors would like to thank Arif O. Harmanci, Jieming Chen and Yao Fu
511 for discussion on useful datasets, and Baikang Pei for discussion on processed
512 pseudogenes.

513

514 **References**

- 515 1. Esnault C, Maestre J, Heidmann T. Human LINE retrotransposons generate
516 processed pseudogenes. *Nat Genet.* 2000;24: 363–7. doi:10.1038/74184
- 517 2. Wei W, Gilbert N, Ooi SL, Lawler JF, Ostertag EM, Kazazian HH, et al.
518 Human L1 retrotransposition: cis preference versus trans complementation.
519 *Mol Cell Biol.* 2001;21: 1429–39. doi:10.1128/MCB.21.4.1429-1439.2001
- 520 3. Mandal PK, Ewing AD, Hancks DC, Kazazian HH. Enrichment of processed
521 pseudogene transcripts in L1-ribonucleoprotein particles. *Hum Mol Genet.*
522 2013;22: 3730–48. doi:10.1093/hmg/ddt225
- 523 4. Kaessmann H. Origins, evolution, and phenotypic impact of new genes.
524 *Genome Res. Cold Spring Harbor Lab;* 2010;20: 1313–1326.
525 doi:10.1101/gr.101386.109
- 526 5. Abyzov A, Iskow R, Gokcumen O, Radke DW, Balasubramanian S, Pei B, et
527 al. Analysis of variable retroduplications in human populations suggests
528 coupling of retrotransposition to cell division. *Genome Res.* 2013;23: 2042–
529 2052. doi:10.1101/gr.154625.113
- 530 6. Ewing AD, Ballinger TJ, Earl D, Harris CC, Ding L, Wilson RK, et al.
531 Retrotransposition of gene transcripts leads to structural variation in
532 mammalian genomes. *Genome Biol.* 2013;14: R22. doi:10.1186/gb-2013-14-3-
533 r22
- 534 7. Schrider DR, Navarro FCP, Galante PAF, Parmigiani RB, Camargo AA, Hahn
535 MW, et al. Gene copy-number polymorphism caused by retrotransposition in
536 humans. *PLoS Genet.* 2013;9: e1003242. doi:10.1371/journal.pgen.1003242
- 537 8. Ciomborowska J, Rosikiewicz W, Szklarczyk D, Makałowski W, Makałowska
538 I. “Orphan” retrogenes in the human genome. *Mol Biol Evol.* 2013;30: 384–96.
539 doi:10.1093/molbev/mss235
- 540 9. Long M, VanKuren NW, Chen S, Vibranovski MD. New gene evolution: little
541 did we know. *Annu Rev Genet.* 2013;47: 307–33. doi:10.1146/annurev-genet-
542 111212-133301

- 543 10. Bernstein BE, Birney E, Dunham I, Green ED, Gunter C, Snyder M. An
544 integrated encyclopedia of DNA elements in the human genome. *Nature*.
545 Nature Publishing Group, a division of Macmillan Publishers Limited. All
546 Rights Reserved.; 2012;489: 57–74. doi:10.1038/nature11247
- 547 11. Pei B, Sisu C, Frankish A, Howald C, Habegger L, Mu XJ, et al. The
548 GENCODE pseudogene resource. *Genome Biol.* 2012;13: R51.
549 doi:10.1186/gb-2012-13-9-r51
- 550 12. Sisu C, Pei B, Leng J, Frankish A, Zhang Y, Balasubramanian S, et al.
551 Comparative analysis of pseudogenes across three phyla. *Proc Natl Acad Sci U*
552 *S A.* 2014;111: 13361–6. doi:10.1073/pnas.1407293111
- 553 13. Sasidharan R, Gerstein M. Genomics: protein fossils live on as RNA. *Nature*.
554 Nature Publishing Group; 2008;453: 729–31. doi:10.1038/453729a
- 555 14. Salmena L, Poliseno L, Tay Y, Kats L, Pandolfi PP. A ceRNA hypothesis: the
556 Rosetta Stone of a hidden RNA language? *Cell.* 2011;146: 353–8.
557 doi:10.1016/j.cell.2011.07.014
- 558 15. Tam OH, Aravin AA, Stein P, Girard A, Murchison EP, Cheloufi S, et al.
559 Pseudogene-derived small interfering RNAs regulate gene expression in mouse
560 oocytes. *Nature.* Nature Publishing Group; 2008;453: 534–8.
561 doi:10.1038/nature06904
- 562 16. Watanabe T, Totoki Y, Toyoda A, Kaneda M, Kuramochi-Miyagawa S, Obata
563 Y, et al. Endogenous siRNAs from naturally formed dsRNAs regulate
564 transcripts in mouse oocytes. *Nature.* 2008;453: 539–43.
565 doi:10.1038/nature06908
- 566 17. Wen Y-Z, Zheng L-L, Liao J-Y, Wang M-H, Wei Y, Guo X-M, et al.
567 Pseudogene-derived small interference RNAs regulate gene expression in
568 African *Trypanosoma brucei*. *Proc Natl Acad Sci U S A.* 2011;108: 8345–50.
569 doi:10.1073/pnas.1103894108
- 570 18. Betrán E, Emerson JJ, Kaessmann H, Long M. Sex chromosomes and male
571 functions: where do new genes go? *Cell Cycle.* 2004;3: 873–5.
- 572 19. Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. A
573 coding-independent function of gene and pseudogene mRNAs regulates tumour
574 biology. *Nature.* 2010;465: 1033–8. doi:10.1038/nature09144
- 575 20. Shukla R, Upton KR, Muñoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, et
576 al. Endogenous retrotransposition activates oncogenic pathways in
577 hepatocellular carcinoma. *Cell.* 2013;153: 101–11.
578 doi:10.1016/j.cell.2013.02.032
- 579 21. de Boer M, van Leeuwen K, Geissler J, Weemaes CM, van den Berg TK,
580 Kuijpers TW, et al. Primary immunodeficiency caused by an exonized
581 retroposed gene copy inserted in the CYBB gene. *Hum Mutat.* 2014;35: 486–
582 96. doi:10.1002/humu.22519
- 583 22. Solyom S, Ewing AD, Rahrmann EP, Doucet T, Nelson HH, Burns MB, et al.
584 Extensive somatic L1 retrotransposition in colorectal tumors. *Genome Res.*
585 2012;22: 2328–38. doi:10.1101/gr.145235.112

- 586 23. Cooke SL, Shlien A, Marshall J, Pipinikas CP, Martincorena I, Tubio JMC, et
587 al. Processed pseudogenes acquired somatically during cancer development.
588 Nat Commun. 2014;5: 3644. doi:10.1038/ncomms4644
- 589 24. Tubio JMC, Li Y, Ju YS, Martincorena I, Cooke SL, Tojo M, et al. Extensive
590 transduction of nonrepetitive DNA mediated by L1 retrotransposition in cancer
591 genomes. Science (80-). 2014;345: 1251343–1251343.
592 doi:10.1126/science.1251343
- 593 25. Helman E, Lawrence MS, Stewart C, Sougnez C, Getz G, Meyerson M.
594 Somatic retrotransposition in human cancer revealed by whole-genome and
595 exome sequencing. Genome Res. 2014;24: 1053–63.
596 doi:10.1101/gr.163659.113
- 597 26. Richardson SR, Salvador-Palomeque C, Faulkner GJ. Diversity through
598 duplication: whole-genome sequencing reveals novel gene retrocopies in the
599 human population. Bioessays. 2014;36: 475–81. doi:10.1002/bies.201300181
- 600 27. Evrony GD, Lee E, Mehta BK, Benjamini Y, Johnson RM, Cai X, et al. Cell
601 Lineage Analysis in Human Brain Using Endogenous Retroelements. Neuron.
602 2015;85: 49–59. doi:10.1016/j.neuron.2014.12.028
- 603 28. The 1000 Genomes Project [Internet]. [cited 29 Oct 2015]. Available:
604 <http://www.1000genomes.org/>
- 605 29. Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti
606 A, et al. A global reference for human genetic variation. Nature. Nature
607 Publishing Group, a division of Macmillan Publishers Limited. All Rights
608 Reserved.; 2015;526: 68–74. doi:10.1038/nature15393
- 609 30. Sudmant PHPH, Rausch T, Gardner EJEJ, Handsaker RERE, Abyzov A,
610 Huddleston J, et al. An integrated map of structural variation in 2,504 human
611 genomes. Nature. Nature Publishing Group, a division of Macmillan Publishers
612 Limited. All Rights Reserved.; 2015;526: 75–81. doi:10.1038/nature15394
- 613 31. Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE,
614 et al. An integrated map of genetic variation from 1,092 human genomes.
615 Nature. Nature Publishing Group, a division of Macmillan Publishers Limited.
616 All Rights Reserved.; 2012;491: 56–65. doi:10.1038/nature11632
- 617 32. Stranger BE, Forrest MS, Dunning M, Ingle CE, Beazley C, Thorne N, et al.
618 Relative impact of nucleotide and copy number variation on gene expression
619 phenotypes. Science. American Association for the Advancement of Science;
620 2007;315: 848–53. doi:10.1126/science.1136678
- 621 33. Suzuki R, Shimodaira H. Pvcust: an R package for assessing the uncertainty in
622 hierarchical clustering. Bioinformatics. 2006;22: 1540–2.
623 doi:10.1093/bioinformatics/btl117
- 624 34. Shimodaira H, Hasegawa M. CONSEL: for assessing the confidence of
625 phylogenetic tree selection. Bioinformatics. 2001;17: 1246–1247.
626 doi:10.1093/bioinformatics/17.12.1246
- 627 35. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas
628 MA, et al. Transcriptome and genome sequencing uncovers functional
629 variation in humans. Nature. Nature Publishing Group, a division of Macmillan

- 630 Publishers Limited. All Rights Reserved.; 2013;501: 506–11.
631 doi:10.1038/nature12531
- 632 36. Abyzov A, Li S, Kim DR, Mohiyuddin M, Stütz AM, Parrish NF, et al.
633 Analysis of deletion breakpoints from 1,092 humans reveals details of mutation
634 mechanisms. *Nat Commun.* 2015;6: 7256. doi:10.1038/ncomms8256
- 635 37. Baller JA, Gao J, Stamenova R, Curcio MJ, Voytas DF. A nucleosomal surface
636 defines an integration hotspot for the *Saccharomyces cerevisiae* Ty1
637 retrotransposon. *Genome Res.* 2012;22: 704–13. doi:10.1101/gr.129585.111
- 638 38. Mularoni L, Zhou Y, Bowen T, Gangadharan S, Wheelan SJ, Boeke JD.
639 Retrotransposon Ty1 integration targets specifically positioned asymmetric
640 nucleosomal DNA segments in tRNA hotspots. *Genome Res.* 2012;22: 693–
641 703. doi:10.1101/gr.129460.111
- 642 39. Segal E, Widom J. What controls nucleosome positions? *Trends Genet.*
643 2009;25: 335–43. doi:10.1016/j.tig.2009.06.002
- 644 40. Vatta M, Ackerman MJ, Ye B, Makielski JC, Ughanze EE, Taylor EW, et al.
645 Mutant caveolin-3 induces persistent late sodium current and is associated with
646 long-QT syndrome. *Circulation.* 2006;114: 2104–12.
647 doi:10.1161/CIRCULATIONAHA.106.635268
- 648 41. Cronk LB, Ye B, Kaku T, Tester DJ, Vatta M, Makielski JC, et al. Novel
649 mechanism for sudden infant death syndrome: persistent late sodium current
650 secondary to mutations in caveolin-3. *Heart Rhythm.* 2007;4: 161–6.
651 doi:10.1016/j.hrthm.2006.11.030
- 652 42. Williams ES, Thomas KL, Broderick S, Shaw LK, Velazquez EJ, Al-Khatib
653 SM, et al. Race and gender variation in the QT interval and its association with
654 mortality in patients with coronary artery disease: results from the Duke
655 Databank for Cardiovascular Disease (DDCD). *Am Heart J.* 2012;164: 434–41.
656 doi:10.1016/j.ahj.2012.05.024
- 657 43. Hakeem GF, Oddy L, Holcroft CA, Abenhaim HA. Incidence and determinants
658 of sudden infant death syndrome: a population-based study on 37 million
659 births. *World J Pediatr.* 2014; doi:10.1007/s12519-014-0530-9
- 660 44. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al.
661 Initial sequencing and analysis of the human genome. *Nature.* Macmillian
662 Magazines Ltd.; 2001;409: 860–921. doi:10.1038/35057062
- 663 45. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et
664 al. GENCODE: the reference human genome annotation for The ENCODE
665 Project. *Genome Res.* 2012;22: 1760–74. doi:10.1101/gr.135350.111
- 666 46. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The
667 Sequence Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25:
668 2078–9. doi:10.1093/bioinformatics/btp352
- 669 47. Efron B, Halloran E, Holmes S. Bootstrap confidence levels for phylogenetic
670 trees. *Proc Natl Acad Sci U S A.* 1996;93: 13429–13434.
671 doi:10.1073/pnas.93.23.13429
- 672 48. Shimodaira H. An approximately unbiased test of phylogenetic tree selection.

- 673 Syst Biol. 2002;51: 492–508. doi:10.1080/10635150290069913
- 674 49. Shimodaira H. Approximately unbiased tests of regions using multistep-
675 multiscale bootstrap resampling. *Ann Stat. Institute of Mathematical Statistics;*
676 2004;32: 2616–2641.
- 677 50. Holsinger KE, Weir BS. Genetics in geographically structured populations:
678 defining, estimating and interpreting F(ST). *Nat Rev Genet.* 2009;10: 639–50.
679 doi:10.1038/nrg2611
- 680 51. Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical
681 and Powerful Approach to Multiple Testing. *J R Stat Soc Ser B. Blackwell*
682 *Publishers;* 1995;57: 289–300.
- 683 52. Fisher RA. Statistical methods for research workers. Boyd OA, editor.
684 *Biological monographs and manuals. Oliver and Boyd;* 1925.
- 685 53. Molaro A, Hodges E, Fang F, Song Q, McCombie WR, Hannon GJ, et al.
686 Sperm Methylation Profiles Reveal Features of Epigenetic Inheritance and
687 Evolution in Primates. *Cell.* 2011;146: 1029–1041.
688 doi:10.1016/j.cell.2011.08.016
- 689 54. Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, et al.
690 The UCSC Genome Browser database: 2014 update. *Nucleic Acids Res.*
691 2014;42: D764-70. doi:10.1093/nar/gkt1168
- 692 55. Gaffney DJ, McVicker G, Pai AA, Fondufe-Mittendorf YN, Lewellen N,
693 Michelini K, et al. Controls of nucleosome positioning in the human genome.
694 *PLoS Genet.* 2012;8: e1003036. doi:10.1371/journal.pgen.1003036
- 695 56. Bejerano G, Pheasant M, Makunin I, Stephen S, Kent WJ, Mattick JS, et al.
696 Ultraconserved elements in the human genome. *Science (80-).* 2004;304:
697 1321–1325.
- 698 57. Fu Y, Liu Z, Lou S, Bedford J, Mu XJ, Yip KY, et al. FunSeq2: A framework
699 for prioritizing noncoding regulatory variants in cancer. *Genome Biol.*
700 2014;15: 480. doi:10.1186/s13059-014-0480-5
- 701 58. Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, et al.
702 Integrative annotation of variants from 1092 humans: application to cancer
703 genomics. *Science (80-).* 2013;342: 1235587.
- 704 59. Ha H, Song J, Wang S, Kapusta A, Feschotte C, Chen KC, et al. A
705 comprehensive analysis of piRNAs from adult human testis and their
706 relationship with genes and mobile elements. *BMC Genomics.* 2014;15: 545.
- 707 60. Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of
708 large gene lists using DAVID bioinformatics resources. *Nat Protoc.* 2009;4:
709 44–57. doi:10.1038/nprot.2008.211
- 710 61. Pletscher-Frankild S, Pallegà A, Tsafou K, Binder JX, Jensen LJ. DISEASES:
711 Text mining and data integration of disease-gene associations. *Methods.* 2014;
712 doi:10.1016/j.ymeth.2014.11.020
- 713

714 **Supporting Information**

715 **S1 File. Supplementary file.** This file contains supplementary figures and supplementary tables.

716

717 **S2 File. Retroduplication callset derived from indicative exon-exon junctions.** Retroduplication
718 calls from each person are listed. Each row contains the following information: the junction location
719 represented by the interval between a pair of exons being joined (Chrom: chromosome, Start: end site
720 of the upstream exon, End: start site of the downstream exon), Parent Gene ID, the person's ID in the
721 1000 Genomes Project, and the population abbreviation.

722

723 **S3 File. Detected retroduplication insertion sites.** The file contains the confidence regions of
724 detected insertion sites.

725

726 **S4 File. Detected retroduplication deletions.** The file reports overlaps between deletions (DEL) and
727 processed pseudogenes where the processed pseudogene region overlaps at least 50% of the deletion
728 regions. The first six columns are the information for each DEL region (chromosome, start site, end
729 site, structural variation type, allele frequency, ID in Phase 3). The last three columns are the
730 information for overlapping pseudogenes (chromosome, start site, end site).

731

732 **S5 File. Retroduplication counts and frequencies in five superpopulations.** The file contains the
733 retroduplication counts (in terms of the number of individuals having the retroduplication in a
734 superpopulation), and the retroduplication frequencies, for all the 503 unique parent genes detected in
735 the whole callset.

736

737 **S6 File. Retroduplication eQTL results.** The file contains retroduplication eQTL results for five
738 populations (CEU, FIN, GBR, TSI, YRI). Each sheet contains the result of one population. Each row
739 (except the last) contains the following information: Parent Gene ID, the statistic from two-sided
740 Wilcoxon rank sum test, the original p-value from the test, and the p-value adjusted by Benjamini-
741 Hochberg procedure. The last row contains the combined p-value from the omnibus test.

742

743 **S7 File. Expression level of retroduplication parent genes compared to all genes.** The file contains
744 gene expression level comparison results for five populations (CEU, FIN, GBR, TSI, YRI). Each sheet
745 contains the result of one population. Each row (except the last) contains the following information:
746 Parent Gene ID, the observed statistic (medium of the expression level of the parent gene), quantile of
747 the observed statistic compared to null distribution, the empirical p-value, and the p-value adjusted by
748 Benjamini-Hochberg procedure. The last row contains the combined p-value from the omnibus test.

749

750 **S8 File. The code of retroduplication calling pipeline.** The file contains the zipped code.