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

5 Yan Zhang[1,2,3¶], Shantao Li[1¶], Alexej Abyzov[4*], Mark B. Gerstein[1,2,5*]

7 [1]Program in Computational Biology and Bioinformatics, Yale University, New
8 Haven, CT 06520
9 [2]Department of Molecular Biophysics and Biochemistry, School of Medicine, Yale
10 University, New Haven, CT 06520
11 [3]Department of Biomedical Informatics, College of Medicine, The Ohio State
12 University, Columbus, OH 43210
13 [4]Department of Health Sciences Research, Center for Individualized Medicine, Mayo
14 Clinic, Rochester, Minnesota, MN 55905
15 [5]Department of Computer Science, Yale University, New Haven, CT 06520

17 [*]Corresponding authors
18 E-mails: abyzov.alexej@mayo.edu (AA); mark@gersteinlab.org (MBG)

20 [¶]These authors contributed equally to this work.

## Abstract

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

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

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

## Author Summary

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

- 2 -

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

## 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 retrotransponsons, 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

## 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    (FDR) using decoy exon junction libraries. As a result, we have called a total of
121    15,642 retroduplications from 2,533 individuals (with two outlier samples excluded)
122    for 503 unique parent genes (**S1-S2 Fig.**, **S1 Table**, and **S2 File**). On average, an
123    individual has 6 novel retroduplications identified based on exon-exon junctions.
124    Next, we identified retroduplication insertion sites for 152 of the parent genes based
125    on 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

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

We have identified novel retroduplications, which are insertions relative to the reference genome. There are also retroduplications that are deletions relative to the reference genome (i.e. absent in the individuals but present in the reference genome). These events can be detected by overlapping known processed pseudogenes in the reference genome with 1000 Genomes Phase 3 deletion set. We carried out this in the supplement, finding 68 such deletion events (**S4 File**). These 68 regions are present in hg19 as processed pseudogenes, but reported as deletions in 1000 Genomes Phase 3. This type of events is far less common than retroduplication insertions, thus we suggest focusing on retroduplication insertions in the study.

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

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

- 5 -

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

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

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

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

we tested the association of genomic features with insertion sites. Inheritable retroduplication events occurred in germline so we focused on gametes. We found that retroduplication insertions sites are enriched within hypomethylated regions in sperm (2.0-fold, empirical p-value < 0.0012). It is likely that retroduplication events exhibit certain preference in insertion sites associated with open chromatin. Furthermore, we characterized nucleosome positioning [37,38] around insertion sites. Overall, insertion sites show high regularity of nucleosome location (empirical p-value from permutation test $2\times10^{-4}$) (**Fig. 3A**). Highly nucleosome regularity often indicates the presence of chromatin remodeling and DNA binding proteins [39], which creates favorable loosely packed microenviroment for insertion.

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

Insertions points could be supported by discordant reads from both sides or just one side around the insertion point. There is no fundamental preference for retroduplicated DNA segments to insert into other retroelements such as L1 elements. However, L1 involved in retroduplication is sometime co-duplicated and co-inserted next to the retroduplicated segment. This type of co-insertion event can be detected by looking at the insertion sites that only have discordant-read support on one side. In these cases, we found co-inserted L1 tend to belong to young L1 subfamilies, represented by L1HS (4.7-fold, p-value < 0.001) and L1PA (1.9-fold, p-value < 0.001) (**Fig. 3B**). Contrastingly, for insertion sites without evidence for co-insertion (i.e. insertion sites that are supported by both sides) we did not observe such young L1 preference (p-value > 0.05). Enrichment of young and active L1 subfamilies involving in speculated L1 transductions suggests these novel retroduplication variants happened very recently.

In order to investigate the functional impact of retroduplication insertions on genomic functions, we tested the significance of overlap between retroduplication

Shantao 3/23/2017 2:37 AM
**Deleted:** enrichment

Shantao 3/23/2017 2:37 AM
**Deleted:** overlap between the

Shantao 3/23/2017 2:37 AM
**Deleted:** and

Shantao 3/23/2017 2:39 AM
**Deleted:** In order to further

Shantao 3/23/2017 2:39 AM
**Deleted:**

Shantao 3/23/2017 2:40 AM
**Deleted:** , we aggregated nucleosome positions around insertion sites

Shantao 3/23/2017 2:42 AM
**Deleted:** either

Shantao 3/23/2017 2:43 AM
**Deleted:** that is

Shantao 3/23/2017 2:44 AM
**Deleted:** see a preference

Shantao 3/23/2017 2:44 AM
**Deleted:** for the

Shantao 3/23/2017 2:45 AM
**Deleted:** a

Shantao 3/23/2017 2:45 AM
**Deleted:** L1

Shantao 3/23/2017 2:46 AM
**Deleted:** where we do not detect a

Shantao 3/23/2017 2:45 AM
**Deleted:** o

Shantao 3/23/2017 2:46 AM
**Deleted:** have

insertion sites and genomic elements compared to random genomic background (**Fig. 3C**). As expected, ultraconserved regions are significantly depleted (p-value < 0.001). This observation is consistent with our knowledge that in general population, variable retroduplications should not interrupt with evolutionary or functionally constrained regions. Unexpectedly, we observed that intron regions are also depleted (p-value < 0.01), which might be due to negative selection that maintains conserved alternative splicing by avoiding interruption from insertion into introns.

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

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

Shantao 3/23/2017 2:49 AM
**Comment [5]:** Is this really surprising? In general, all intragenic insertions are bad

Shantao 3/23/2017 2:50 AM
**Deleted:** interesting

Shantao 3/23/2017 2:52 AM
**Deleted:** retroduplication

Shantao 3/23/2017 2:52 AM
**Deleted:**

Shantao 3/23/2017 2:55 AM
**Deleted:** using whole exome/genome sequencing data

## Materials and Methods

**Data resources**

Whole exome sequencing and whole genome sequencing data of 2,535 individuals from 26 populations were generated by the 1000 Genomes Project Phase 3 (whole-genome sequencing with mean depth 7.4x and read length of 100bp; targeted exome sequencing with mean depth 65.7x and read length of 76bp) [28–30]. Population description can be found at http://www.1000genomes.org/category/frequently-asked-questions/population. Protein-coding gene expression data (Peer-factor normalized RPKM) is obtained from the Geuvadis RNA-sequencing project [35], which generated RNA sequencing data from lymphoblastoid cell lines of 462 individuals from 5 populations (CEU, FIN, GBR, TSI and YRI) enrolled in the 1000 Genomes Project. We use human reference genome build 37 [44] and GENCODE v19 human genome annotation [45] in the study.


**Calling pipeline**

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

Shantao 3/20/2017 5:47 PM
**Deleted:** and

Shantao 3/20/2017 5:47 PM
**Deleted:** S8 File

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

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

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

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

369 **Report novel retroduplication calls.** Multiple "previously unmapped" reads
370 (unmapped to the reference genome) might be mapped to the same exon-exon
371 junction, supporting the existence of the novel exon-exon junction. Furthermore,
372 multiple exon-exon junctions with mapped reads might support the existence of a
373 gene retroduplication event. We report a gene having novel retroduplications, when it
374 has at least two non-overlapping supporting exon-exon junctions, and at least one
375 junction is supported by at least two mapped reads. The genes (also called parent
376 genes) with novel retroduplications are called for each person individually. We
377 noticed that the 1000 Genomes Project Phase 3 provides paired-end sequencing data
378 for all individuals but NA19318. We include this individual into our analysis, as
379 single-end sequencing does not seem to affect the performance of this pipeline.

380 **Detect retroduplication insertion sites.** In the insertion site detection step, we pool
381 low-coverage whole genome sequencing data by population, and call insertion sites
382 for each population. We search for discordant paired-end reads (with a minimum
383 quality score of 15) with one read correctly mapped to the parent gene, and the other
384 read mapped to a different chromosome or at least 1 kb away from the gene. In order
385 to avoid false discovery, we limit our searching scope to the parent genes identified
386 from the exon-exon junctions.
387    Read pairs with proper orientations are clustered using average linkage
388 clustering. It can be shown that this linkage criterion is not likely affected by the local
389 coverage. Assuming uniform distribution of reads, it can be shown mathematically
390 that the expected distance between reads supporting the same insertion point is

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

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

- 11 -

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

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

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

411

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

413 **Generate retroduplication frequency matrix.** Some retroduplication parent genes
414 are called commonly among multiple populations, while some others are called
415 exclusively in a single population. Besides, parent genes are called at different
416 frequencies within a population. This information can be used for measuring distance
417 between populations, while taking into account different retroduplication frequencies.
418 We define a retroduplication frequency matrix, from which distance measures can be
419 calculated.

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

425 **Bootstrap phylogenetic trees.** We use Manhattan distance as the distance measure
426 between each pair of populations (i.e. Manhattan distance between two columns in *A*).
427 Average linkage is used in hierarchical clustering for generating each tree. 1000
428 bootstrap replications are performed, and the uncertainty is assessed using Pvclust
429 [33]. The reported AU (Approximately Unbiased) probability values [33,34] are used
430 to indicate the certainty of sub-tree structures generated from multi-scale bootstrap
431 resampling [47–49]. The higher the AU probability value, the more confident the
432 substructure is.
433

434 **Detect population differentiation due to retroduplication polymorphism**
435 We check population differentiation due to retroduplication polymorphism,
436 based on retroduplication frequencies in different superpopulations. Herein we pool
437 the 26 populations into 5 superpopulations (African, Ad Mixed American, East Asian,
438 European, and South Asian) as defined by the 1000 Genomes Project. For each given
439 retroduplication parent gene, we calculate the population differentiation measure
440 equivalent to the fixation index [50]. We define the test statistic

441
$$F_{ST} = \frac{p(1-p) - \sum_{i=1}^{5} c_i p_i (1-p_i)}{p(1-p)},$$

442 in which $i = 1, ..., 5$ corresponds to the $i$th superpopulation, $p$ is the retroduplication
443 frequency of a given parent gene in the total population, $p_i$ is the retroduplication
444 frequency of the same parent gene in the $i$th superpopulation, and $c_i$ is the relative
445 population size of the $i$th superpopulation. $c_i$ is calculated as the number of
446 individuals in the $i$th superpopulation divided by the number of individuals in the total
447 population. The larger the $F_{ST}$, the more different the retroduplication frequencies
448 among superpopulations. One-tailed empirical p-value is calculated comparing the
449 observed $F_{ST}$ versus the null distribution of $F_{ST}$. The null distribution is calculated
450 from 1000 fake population sets generated by shuffling individual labels, while
451 maintaining the size unchanged for each population. By the significance of $F_{ST}$, i.e.
452 the p-value adjusted by Benjamini-Hochberg procedure [51], we can detect the
453 retroduplications that can differentiate populations.
454

**Analyze association between retroduplication and gene expression**

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

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

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

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

480 **Expression level of retroduplication parent genes compared to all genes.** For a
481 certain population, we test whether the retroduplication parent genes are highly
482 expressed among all the genes measured in the Geuvadis data set. We take medium
483 expression value over all individuals for each gene as the representative expression
484 value. One-tailed empirical p-value is calculated comparing the expression value of
485 each parent gene versus the null distribution of expression values of all genes. It
486 indicates the significance of each retroduplication parent gene having high expression
487 value among all genes. Fisher's combined probability test is performed on the
488 empirical p-values. If the combined p-value is less than 0.05, that means in general
489 the parent genes are significantly highly expressed among all genes. The results are
490 available in S7 File.
491

492 **Explore association between local genomic features and retroduplications**
493 **insertion sites**
494       To test the association between sperm methylation patterns and
495 retroduplication insertion sites, we intersect out insertion sites with hypomethylated
496 regions in sperms [53]. L1 annotation (RepeatMask), ENCODE HESC DNase I
497 hypersensitive data and genomic GC contents are downloaded from the USCS
498 Genome Browser [54]. Well-positioned nucleosome data is obtained from a recent
499 study on multiple individuals [55].
500       We randomly shuffle the locations of insertion sites for 10,000 times on the
501 same chromosome, excluding the gap regions, to obtain an empirical distribution of
502 the null hypothesis. For fold changes, we use the mean of this distribution as the best
503 estimate of the expected value. Calculation of p-value is empirical in order to be
504 conservative. We use Bonferroni correction in case of multiple hypothesis testing.
505 Unless specified otherwise, we only report corrected p-value. In order to avoid any
506 effect of the difference of location precision across different insertion sites, we
507 enlarge the insertion site region to 500 bp while keeping the middle point of insertions
508 unchanged. We also exclude insertion points on alternative locus in the genome.
509       For aggregation plot on well-positioned nucleosome and GC content, we use
510 200 bp bins to calculate the base overlap, and the final plot was further window-
511 smoothed with window size of 10. Normalization is performed by taking mean value

512  of the first and last 20 bins as background. We use the GC contents from UCSC

513  browser track, which is binned in 5 bp.

514

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

516  We test the significance of overlap between retroduplication insertion sites and

517  genomic elements, including gene, CDS, exon, UTR, intron, pseudogene and

518  lincRNA annotated in GENCODE v19, and ultraconserved regions (evolutionary

519  constraint regions across species), ultrasensitive non-coding regions (regions

520  particularly sensitive to disruptive mutations) and TF (transcription factor) peak

521  regions obtained from ENCODE RNA-seq data [10] and literature [30,56–59]. The

522  overlap between a genomic element type and the insertion sites is measured by the

523  partial overlap statistic, which is the count of genomic elements that have at least 1 bp

524  overlap with the detected insertion sites.

525  We randomly shuffle the locations of insertion sites for 1,000 times on the

526  same chromosome, excluding the Hg19 gap regions, to obtain an empirical

527  distribution of the null hypothesis. In the permutation tests, the null distribution of the

528  overlap measures is calculated from true genomic elements intersecting randomly

529  shuffled insertion locations. The enrichment of overlap is represented by log2 fold

530  change of the observed overlap statistic versus the mean of its null distribution.

531  Empirical p-value is calculated.

532  In order to avoid any effect from different location precisions, we enlarge the

533  insertion intervals uniformly to 1000 bp, while keeping the middle point of insertions.

534  We only use insertion sites on chromosomes (i.e. exclude alternative locus) in the

535  analysis.

536

537  **Functional enrichment analysis**

538  We use DAVID [60] to annotate functional terms for retroduplication parent

539  genes, and survey functional term enrichment.

540

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

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

## Acknowledgements

## References

1. Esnault C, Maestre J, Heidmann T. Human LINE retrotransposons generate processed pseudogenes. Nat Genet. 2000;24: 363–7. doi:10.1038/74184

2. Wei W, Gilbert N, Ooi SL, Lawler JF, Ostertag EM, Kazazian HH, et al. Human L1 retrotransposition: cis preference versus trans complementation. Mol Cell Biol. 2001;21: 1429–39. doi:10.1128/MCB.21.4.1429-1439.2001

3. Mandal PK, Ewing AD, Hancks DC, Kazazian HH. Enrichment of processed pseudogene transcripts in L1-ribonucleoprotein particles. Hum Mol Genet. 2013;22: 3730–48. doi:10.1093/hmg/ddt225

4. Kaessmann H. Origins, evolution, and phenotypic impact of new genes. Genome Res. Cold Spring Harbor Lab; 2010;20: 1313–1326. doi:10.1101/gr.101386.109

5. Abyzov A, Iskow R, Gokcumen O, Radke DW, Balasubramanian S, Pei B, et al. Analysis of variable retroduplications in human populations suggests coupling of retrotransposition to cell division. Genome Res. 2013;23: 2042–2052. doi:10.1101/gr.154625.113

6. Ewing AD, Ballinger TJ, Earl D, Harris CC, Ding L, Wilson RK, et al. Retrotransposition of gene transcripts leads to structural variation in mammalian genomes. Genome Biol. 2013;14: R22. doi:10.1186/gb-2013-14-3-r22

7. Schrider DR, Navarro FCP, Galante PAF, Parmigiani RB, Camargo AA, Hahn MW, et al. Gene copy-number polymorphism caused by retrotransposition in humans. PLoS Genet. 2013;9: e1003242. doi:10.1371/journal.pgen.1003242

8. Ciomborowska J, Rosikiewicz W, Szklarczyk D, Makałowski W, Makałowska I. "Orphan" retrogenes in the human genome. Mol Biol Evol. 2013;30: 384–96. doi:10.1093/molbev/mss235

9. Long M, VanKuren NW, Chen S, Vibranovski MD. New gene evolution: little did we know. Annu Rev Genet. 2013;47: 307–33. doi:10.1146/annurev-genet-111212-133301

580 10. Bernstein BE, Birney E, Dunham I, Green ED, Gunter C, Snyder M. An
581 integrated encyclopedia of DNA elements in the human genome. Nature.
582 Nature Publishing Group, a division of Macmillan Publishers Limited. All
583 Rights Reserved.; 2012;489: 57–74. doi:10.1038/nature11247

584 11. Pei B, Sisu C, Frankish A, Howald C, Habegger L, Mu XJ, et al. The
585 GENCODE pseudogene resource. Genome Biol. 2012;13: R51.
586 doi:10.1186/gb-2012-13-9-r51

587 12. Sisu C, Pei B, Leng J, Frankish A, Zhang Y, Balasubramanian S, et al.
588 Comparative analysis of pseudogenes across three phyla. Proc Natl Acad Sci U
589 S A. 2014;111: 13361–6. doi:10.1073/pnas.1407293111

590 13. Sasidharan R, Gerstein M. Genomics: protein fossils live on as RNA. Nature.
591 Nature Publishing Group; 2008;453: 729–31. doi:10.1038/453729a

592 14. Salmena L, Poliseno L, Tay Y, Kats L, Pandolfi PP. A ceRNA hypothesis: the
593 Rosetta Stone of a hidden RNA language? Cell. 2011;146: 353–8.
594 doi:10.1016/j.cell.2011.07.014

595 15. Tam OH, Aravin AA, Stein P, Girard A, Murchison EP, Cheloufi S, et al.
596 Pseudogene-derived small interfering RNAs regulate gene expression in mouse
597 oocytes. Nature. Nature Publishing Group; 2008;453: 534–8.
598 doi:10.1038/nature06904

599 16. Watanabe T, Totoki Y, Toyoda A, Kaneda M, Kuramochi-Miyagawa S, Obata
600 Y, et al. Endogenous siRNAs from naturally formed dsRNAs regulate
601 transcripts in mouse oocytes. Nature. 2008;453: 539–43.
602 doi:10.1038/nature06908

603 17. Wen Y-Z, Zheng L-L, Liao J-Y, Wang M-H, Wei Y, Guo X-M, et al.
604 Pseudogene-derived small interference RNAs regulate gene expression in
605 African Trypanosoma brucei. Proc Natl Acad Sci U S A. 2011;108: 8345–50.
606 doi:10.1073/pnas.1103894108

607 18. Betrán E, Emerson JJ, Kaessmann H, Long M. Sex chromosomes and male
608 functions: where do new genes go? Cell Cycle. 2004;3: 873–5.

609 19. Poliseno L, Salmena L, Zhang J, Carver B, Haveman WJ, Pandolfi PP. A
610 coding-independent function of gene and pseudogene mRNAs regulates tumour
611 biology. Nature. 2010;465: 1033–8. doi:10.1038/nature09144

612 20. Shukla R, Upton KR, Muñoz-Lopez M, Gerhardt DJ, Fisher ME, Nguyen T, et
613 al. Endogenous retrotransposition activates oncogenic pathways in
614 hepatocellular carcinoma. Cell. 2013;153: 101–11.
615 doi:10.1016/j.cell.2013.02.032

616 21. de Boer M, van Leeuwen K, Geissler J, Weemaes CM, van den Berg TK,
617 Kuijpers TW, et al. Primary immunodeficiency caused by an exonized
618 retroposed gene copy inserted in the CYBB gene. Hum Mutat. 2014;35: 486–
619 96. doi:10.1002/humu.22519

620 22. Solyom S, Ewing AD, Rahrmann EP, Doucet T, Nelson HH, Burns MB, et al.
621 Extensive somatic L1 retrotransposition in colorectal tumors. Genome Res.
622 2012;22: 2328–38. doi:10.1101/gr.145235.112

623  23.  Cooke SL, Shlien A, Marshall J, Pipinikas CP, Martincorena I, Tubio JMC, et
624       al. Processed pseudogenes acquired somatically during cancer development.
625       Nat Commun. 2014;5: 3644. doi:10.1038/ncomms4644

626  24.  Tubio JMC, Li Y, Ju YS, Martincorena I, Cooke SL, Tojo M, et al. Extensive
627       transduction of nonrepetitive DNA mediated by L1 retrotransposition in cancer
628       genomes. Science (80- ). 2014;345: 1251343–1251343.
629       doi:10.1126/science.1251343

630  25.  Helman E, Lawrence MS, Stewart C, Sougnez C, Getz G, Meyerson M.
631       Somatic retrotransposition in human cancer revealed by whole-genome and
632       exome sequencing. Genome Res. 2014;24: 1053–63.
633       doi:10.1101/gr.163659.113

634  26.  Richardson SR, Salvador-Palomeque C, Faulkner GJ. Diversity through
635       duplication: whole-genome sequencing reveals novel gene retrocopies in the
636       human population. Bioessays. 2014;36: 475–81. doi:10.1002/bies.201300181

637  27.  Evrony GD, Lee E, Mehta BK, Benjamini Y, Johnson RM, Cai X, et al. Cell
638       Lineage Analysis in Human Brain Using Endogenous Retroelements. Neuron.
639       2015;85: 49–59. doi:10.1016/j.neuron.2014.12.028

640  28.  The 1000 Genomes Project [Internet]. [cited 29 Oct 2015]. Available:
641       http://www.1000genomes.org/

642  29.  Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti
643       A, et al. A global reference for human genetic variation. Nature. Nature
644       Publishing Group, a division of Macmillan Publishers Limited. All Rights
645       Reserved.; 2015;526: 68–74. doi:10.1038/nature15393

646  30.  Sudmant PHPH, Rausch T, Gardner EJEJ, Handsaker RERE, Abyzov A,
647       Huddleston J, et al. An integrated map of structural variation in 2,504 human
648       genomes. Nature. Nature Publishing Group, a division of Macmillan Publishers
649       Limited. All Rights Reserved.; 2015;526: 75–81. doi:10.1038/nature15394

650  31.  Abecasis GR, Auton A, Brooks LD, DePristo MA, Durbin RM, Handsaker RE,
651       et al. An integrated map of genetic variation from 1,092 human genomes.
652       Nature. Nature Publishing Group, a division of Macmillan Publishers Limited.
653       All Rights Reserved.; 2012;491: 56–65. doi:10.1038/nature11632

654  32.  Stranger BE, Forrest MS, Dunning M, Ingle CE, Beazley C, Thorne N, et al.
655       Relative impact of nucleotide and copy number variation on gene expression
656       phenotypes. Science. American Association for the Advancement of Science;
657       2007;315: 848–53. doi:10.1126/science.1136678

658  33.  Suzuki R, Shimodaira H. Pvclust: an R package for assessing the uncertainty in
659       hierarchical clustering. Bioinformatics. 2006;22: 1540–2.
660       doi:10.1093/bioinformatics/btl117

661  34.  Shimodaira H, Hasegawa M. CONSEL: for assessing the confidence of
662       phylogenetic tree selection. Bioinformatics. 2001;17: 1246–1247.
663       doi:10.1093/bioinformatics/17.12.1246

664  35.  Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PAC, Monlong J, Rivas
665       MA, et al. Transcriptome and genome sequencing uncovers functional
666       variation in humans. Nature. Nature Publishing Group, a division of Macmillan

667　　　　　Publishers Limited. All Rights Reserved.; 2013;501: 506–11.
668　　　　　doi:10.1038/nature12531

669　36.　Abyzov A, Li S, Kim DR, Mohiyuddin M, Stütz AM, Parrish NF, et al.
670　　　　Analysis of deletion breakpoints from 1,092 humans reveals details of mutation
671　　　　mechanisms. Nat Commun. 2015;6: 7256. doi:10.1038/ncomms8256

672　37.　Baller JA, Gao J, Stamenova R, Curcio MJ, Voytas DF. A nucleosomal surface
673　　　　defines an integration hotspot for the Saccharomyces cerevisiae Ty1
674　　　　retrotransposon. Genome Res. 2012;22: 704–13. doi:10.1101/gr.129585.111

675　38.　Mularoni L, Zhou Y, Bowen T, Gangadharan S, Wheelan SJ, Boeke JD.
676　　　　Retrotransposon Ty1 integration targets specifically positioned asymmetric
677　　　　nucleosomal DNA segments in tRNA hotspots. Genome Res. 2012;22: 693–
678　　　　703. doi:10.1101/gr.129460.111

679　39.　Segal E, Widom J. What controls nucleosome positions? Trends Genet.
680　　　　2009;25: 335–43. doi:10.1016/j.tig.2009.06.002

681　40.　Vatta M, Ackerman MJ, Ye B, Makielski JC, Ughanze EE, Taylor EW, et al.
682　　　　Mutant caveolin-3 induces persistent late sodium current and is associated with
683　　　　long-QT syndrome. Circulation. 2006;114: 2104–12.
684　　　　doi:10.1161/CIRCULATIONAHA.106.635268

685　41.　Cronk LB, Ye B, Kaku T, Tester DJ, Vatta M, Makielski JC, et al. Novel
686　　　　mechanism for sudden infant death syndrome: persistent late sodium current
687　　　　secondary to mutations in caveolin-3. Heart Rhythm. 2007;4: 161–6.
688　　　　doi:10.1016/j.hrthm.2006.11.030

689　42.　Williams ES, Thomas KL, Broderick S, Shaw LK, Velazquez EJ, Al-Khatib
690　　　　SM, et al. Race and gender variation in the QT interval and its association with
691　　　　mortality in patients with coronary artery disease: results from the Duke
692　　　　Databank for Cardiovascular Disease (DDCD). Am Heart J. 2012;164: 434–41.
693　　　　doi:10.1016/j.ahj.2012.05.024

694　43.　Hakeem GF, Oddy L, Holcroft CA, Abenhaim HA. Incidence and determinants
695　　　　of sudden infant death syndrome: a population-based study on 37 million
696　　　　births. World J Pediatr. 2014; doi:10.1007/s12519-014-0530-9

697　44.　Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al.
698　　　　Initial sequencing and analysis of the human genome. Nature. Macmillian
699　　　　Magazines Ltd.; 2001;409: 860–921. doi:10.1038/35057062

700　45.　Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et
701　　　　al. GENCODE: the reference human genome annotation for The ENCODE
702　　　　Project. Genome Res. 2012;22: 1760–74. doi:10.1101/gr.135350.111

703　46.　Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The
704　　　　Sequence Alignment/Map format and SAMtools. Bioinformatics. 2009;25:
705　　　　2078–9. doi:10.1093/bioinformatics/btp352

706　47.　Efron B, Halloran E, Holmes S. Bootstrap confidence levels for phylogenetic
707　　　　trees. Proc Natl Acad Sci U S A. 1996;93: 13429–13434.
708　　　　doi:10.1073/pnas.93.23.13429

709　48.　Shimodaira H. An approximately unbiased test of phylogenetic tree selection.

710        Syst Biol. 2002;51: 492–508. doi:10.1080/10635150290069913

711  49.  Shimodaira H. Approximately unbiased tests of regions using multistep-
712        multiscale bootstrap resampling. Ann Stat. Institute of Mathematical Statistics;
713        2004;32: 2616–2641.

714  50.  Holsinger KE, Weir BS. Genetics in geographically structured populations:
715        defining, estimating and interpreting F(ST). Nat Rev Genet. 2009;10: 639–50.
716        doi:10.1038/nrg2611

717  51.  Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical
718        and Powerful Approach to Multiple Testing. J R Stat Soc Ser B. Blackwell
719        Publishers; 1995;57: 289–300.

720  52.  Fisher RA. Statistical methods for research workers. Boyd OA, editor.
721        Biological monographs and manuals. Oliver and Boyd; 1925.

722  53.  Molaro A, Hodges E, Fang F, Song Q, McCombie WR, Hannon GJ, et al.
723        Sperm Methylation Profiles Reveal Features of Epigenetic Inheritance and
724        Evolution in Primates. Cell. 2011;146: 1029–1041.
725        doi:10.1016/j.cell.2011.08.016

726  54.  Karolchik D, Barber GP, Casper J, Clawson H, Cline MS, Diekhans M, et al.
727        The UCSC Genome Browser database: 2014 update. Nucleic Acids Res.
728        2014;42: D764-70. doi:10.1093/nar/gkt1168

729  55.  Gaffney DJ, McVicker G, Pai AA, Fondufe-Mittendorf YN, Lewellen N,
730        Michelini K, et al. Controls of nucleosome positioning in the human genome.
731        PLoS Genet. 2012;8: e1003036. doi:10.1371/journal.pgen.1003036

732  56.  Bejerano G, Pheasant M, Makunin I, Stephen S, Kent WJ, Mattick JS, et al.
733        Ultraconserved elements in the human genome. Science (80- ). 2004;304:
734        1321–1325.

735  57.  Fu Y, Liu Z, Lou S, Bedford J, Mu XJ, Yip KY, et al. FunSeq2: A framework
736        for prioritizing noncoding regulatory variants in cancer. Genome Biol.
737        2014;15: 480. doi:10.1186/s13059-014-0480-5

738  58.  Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, et al.
739        Integrative annotation of variants from 1092 humans: application to cancer
740        genomics. Science (80- ). 2013;342: 1235587.

741  59.  Ha H, Song J, Wang S, Kapusta A, Feschotte C, Chen KC, et al. A
742        comprehensive analysis of piRNAs from adult human testis and their
743        relationship with genes and mobile elements. BMC Genomics. 2014;15: 545.

744  60.  Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of
745        large gene lists using DAVID bioinformatics resources. Nat Protoc. 2009;4:
746        44–57. doi:10.1038/nprot.2008.211

747  61.  Pletscher-Frankild S, Pallejà A, Tsafou K, Binder JX, Jensen LJ. DISEASES:
748        Text mining and data integration of disease-gene associations. Methods. 2014;
749        doi:10.1016/j.ymeth.2014.11.020

750

## Supporting Information

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

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

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

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

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

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

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

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